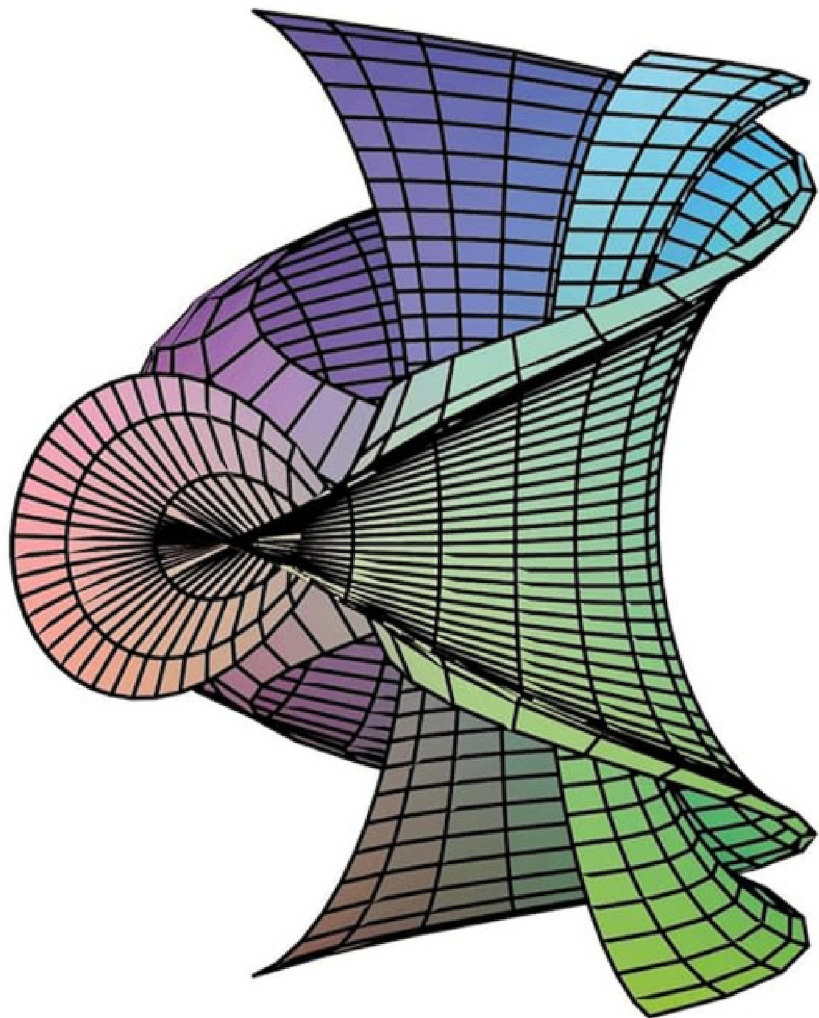


DE GRUYTER

GRADUATE

*Gerard Walschap*

# MULTIVARIABLE CALCULUS AND DIFFERENTIAL GEOMETRY





Gerard Walschap

**Multivariable Calculus and Differential Geometry**





Gerard Walschap

# **Multivariable Calculus and Differential Geometry**

---

**DE GRUYTER**

ISBN 978-3-11-036949-6

**Library of Congress Cataloging-in-Publication Data**

A CIP catalog record for this book has been applied for at the Library of Congress.

**Bibliographic information published by the Deutsche Nationalbibliothek**

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available on the Internet at <http://dnb.dnb.de>.

© 2015 Walter de Gruyter GmbH, Berlin/Boston

Typesetting: PTP-Berlin, Protago-TeX-Production GmbH, Berlin

Printing and binding: CPI books GmbH, Leck

⊗ Printed on acid-free paper

Printed in Germany

[www.degruyter.com](http://www.degruyter.com)

# Preface

The purpose of this text is to introduce the reader to the basic concepts of differential geometry with a minimum of prerequisites. The only absolute requirement is a solid background in single variable calculus, although some mathematical maturity would also be helpful. All the material from multivariable calculus, linear algebra, and basic analysis that is needed in the study of elementary differential geometry has been included. Many students who have completed a standard course in multivariable calculus still struggle to grasp manifold theory. There are several factors contributing to this; chief among them, at least in the author's view, is the lack of integration of linear algebra with calculus. Although most texts on calculus of several variables introduce the reader to matrices, few emphasize concepts such as linear transformations. The latter are admittedly of little use to someone who is only interested in computing, say, partial derivatives of a map, but they are a cornerstone of both calculus and differential geometry. These areas are concerned with retrieving information about a smooth map between Euclidean spaces from its derivative, and the derivative of such a map at a point is essentially synonymous with the notion of linear transformation. We hope that even the reader who is familiar with linear algebra will find it useful to have access to this material should it prove necessary. It constitutes the core of the first chapter. We emphasize this is not a course in linear algebra: only those topics that will be needed in studying geometry are tackled. Furthermore, the presentation and proofs of the material are carried out through the prism of calculus rather than in full generality. For example, several properties that hold in general vector spaces are only established for inner product spaces if their proof is shorter or easier in that context. The remaining part of the chapter is devoted to metric properties of Euclidean spaces, and to the notions of limits and continuity of maps between those spaces.

The second chapter introduces the reader to differentiation of maps between Euclidean spaces, as well as to further concepts from linear algebra that are relevant to inner product spaces. Special emphasis is given to vector fields and their Lie brackets.

The third chapter discusses the spaces that are studied in differential geometry, namely differentiable manifolds. Instead of defining them abstractly, we only consider submanifolds of Euclidean space, which, as geometers know, does not constitute – at least in principle – a restriction. The reason for this is twofold: on the one hand, it is much more intuitive to consider a surface in 3-space and make the leap of imagination to higher dimensions than to study an abstract topological space; on the other, the abstract approach requires concepts from general topology, such as paracompactness, which lie outside the scope of this text. We next proceed to calculus on manifolds, and introduce the basic concepts of covariant derivatives, geodesics, and curvature.

Chapter 4 discusses integration of functions of several variables. We first look at functions with domain in some Euclidean space, determine conditions for integrability, and highlight some of the main tools available, such as Fubini's theorem and the

change of variable formula. The special cases of cylindrical and spherical change of coordinates is discussed in detail. Several applications of multiple integrals to concepts from physics are also examined.

Chapter 5 briefly introduces tensors and then discusses differential forms on manifolds and their integration. This leads us to the modern formulation of Stokes' theorem, and how it unifies the classical versions of the fundamental theorem of Calculus, Green's theorem, and Stokes' theorem. As a practical example, we discuss how Green's theorem explains the principle behind the polar planimeter, a device that calculates the area of a plane region by merely tracing out its boundary.

The next chapter examines manifolds as spaces where a distance between any two points can be defined; when the space is complete, this distance is given by the length of the shortest curve joining the points. Further properties of these metric spaces are discussed, and the chapter ends with an illustration of how curvature affects the shape of space.

The last chapter examines all these concepts in the special case of hypersurfaces; i.e., manifolds whose dimension is one less than that of the ambient Euclidean space. In this setting, there are additional tools such as the Gauss map that provide further insight into the structure of the space. Many features, in particular geodesics, convexity, and curvature become more transparent in this context. We also briefly discuss the geometry of some classical surfaces.

# Contents

Preface — v

## 1 Euclidean Space — 1

- 1.1 Vector spaces — 1
- 1.2 Linear transformations — 6
- 1.3 Determinants — 12
- 1.4 Euclidean spaces — 19
- 1.5 Subspaces of Euclidean space — 25
- 1.6 Determinants as volume — 27
- 1.7 Elementary topology of Euclidean spaces — 30
- 1.8 Sequences — 36
- 1.9 Limits and continuity — 41
- 1.10 Exercises — 48

## 2 Differentiation — 57

- 2.1 The derivative — 57
- 2.2 Basic properties of the derivative — 62
- 2.3 Differentiation of integrals — 67
- 2.4 Curves — 69
- 2.5 The inverse and implicit function theorems — 75
- 2.6 The spectral theorem and scalar products — 81
- 2.7 Taylor polynomials and extreme values — 89
- 2.8 Vector fields — 94
- 2.9 Lie brackets — 103
- 2.10 Partitions of unity — 108
- 2.11 Exercises — 110

## 3 Manifolds — 117

- 3.1 Submanifolds of Euclidean space — 117
- 3.2 Differentiable maps on manifolds — 124
- 3.3 Vector fields on manifolds — 129
- 3.4 Lie groups — 137
- 3.5 The tangent bundle — 141
- 3.6 Covariant differentiation — 143
- 3.7 Geodesics — 148
- 3.8 The second fundamental tensor — 153
- 3.9 Curvature — 156
- 3.10 Sectional curvature — 160
- 3.11 Isometries — 163
- 3.12 Exercises — 168

<b>4</b>	<b>Integration on Euclidean space — 177</b>
4.1	The integral of a function over a box — 177
4.2	Integrability and discontinuities — 181
4.3	Fubini's theorem — 187
4.4	Sard's theorem — 195
4.5	The change of variables theorem — 198
4.6	Cylindrical and spherical coordinates — 202
4.6.1	Cylindrical coordinates — 202
4.6.2	Spherical coordinates — 206
4.7	Some applications — 210
4.7.1	Mass — 211
4.7.2	Center of mass — 211
4.7.3	Moment of inertia — 213
4.8	Exercises — 214
<b>5</b>	<b>Differential Forms — 221</b>
5.1	Tensors and tensor fields — 221
5.2	Alternating tensors and forms — 224
5.3	Differential forms — 232
5.4	Integration on manifolds — 236
5.5	Manifolds with boundary — 240
5.6	Stokes' theorem — 243
5.7	Classical versions of Stokes' theorem — 246
5.7.1	An application: the polar planimeter — 249
5.8	Closed forms and exact forms — 252
5.9	Exercises — 257
<b>6</b>	<b>Manifolds as metric spaces — 267</b>
6.1	Extremal properties of geodesics — 267
6.2	Jacobi fields — 271
6.3	The length function of a variation — 275
6.4	The index form of a geodesic — 278
6.5	The distance function — 283
6.6	The Hopf-Rinow theorem — 285
6.7	Curvature comparison — 289
6.8	Exercises — 292
<b>7</b>	<b>Hypersurfaces — 301</b>
7.1	Hypersurfaces and orientation — 301
7.2	The Gauss map — 304
7.3	Curvature of hypersurfaces — 308
7.4	The fundamental theorem for hypersurfaces — 313

7.5	Curvature in local coordinates —	<b>316</b>
7.6	Convexity and curvature —	<b>318</b>
7.7	Ruled surfaces —	<b>320</b>
7.8	Surfaces of revolution —	<b>323</b>
7.9	Exercises —	<b>328</b>

**Appendix A — 339**

**Appendix B — 345**

**Index — 351**





# 1 Euclidean Space

We begin by exploring some of the properties of Euclidean space that will be needed to study Calculus. They roughly fall into two categories, algebraic and topological. The former are closely related to the concept of vector spaces.

## 1.1 Vector spaces

As a set,  $k$ -dimensional Euclidean space consists of all  $k$ -tuples  $(a_1, \dots, a_k)$  of real numbers  $a_i$ ,  $i = 1, \dots, k$  (the notation Euclidean space usually requires some additional structure as we shall see in section 4; we will, for now, ignore this out of laziness and for lack of a better name). Terms such as “dimension” and “Euclidean” will be examined in more detail later in the chapter. The number  $a_i$  is called the  $i$ -th *coordinate* of the point, and the map  $u^i : \mathbb{R}^k \rightarrow \mathbb{R}$  which assigns to a point its  $i$ -th coordinate is called the  $i$ -th *projection*. When  $k$  equals 2 or 3,  $\mathbb{R}^k$  can be visualized geometrically as the plane or 3-space, with  $(a_1, a_2, a_3)$  representing the point in 3-space whose orthogonal projection onto the  $x$ ,  $y$ , and  $z$ -axes equals  $a_1$ ,  $a_2$ , and  $a_3$  respectively (the plane may be identified with those points in  $\mathbb{R}^3$  with zero third coordinate). The motivation behind the notation  $\mathbb{R}^k$  is that it is a kind of product of  $\mathbb{R}$  with itself  $k$  times: the *cartesian product*  $A \times B$  of two sets  $A$  and  $B$  is defined to be the set of all pairs  $(a, b)$  where  $a$  and  $b$  range over  $A$  and  $B$  respectively, and  $\mathbb{R}^k$  is then the product  $\mathbb{R} \times \dots \times \mathbb{R}$  of  $\mathbb{R}$  with itself  $k$  times.

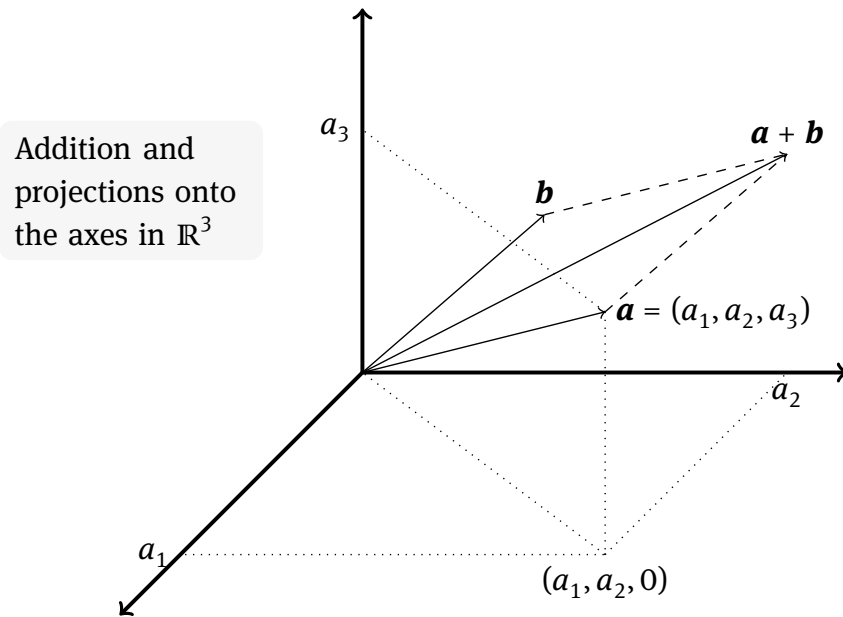
$\mathbb{R}^k$  comes with two operations, *addition*  $+$ :  $\mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}^k$ , given by

$$(a_1, \dots, a_k) + (b_1, \dots, b_k) = (a_1 + b_1, \dots, a_k + b_k), \quad a_i \in \mathbb{R},$$

and *scalar multiplication*  $\cdot$ :  $\mathbb{R} \times \mathbb{R}^k \rightarrow \mathbb{R}^k$ ,

$$c \cdot (a_1, \dots, a_k) = (ca_1, \dots, ca_k), \quad c, a_i \in \mathbb{R}.$$

We will mostly dispense with the dot in the scalar multiplication, and often abbreviate  $(a_1, \dots, a_k)$  by  $\mathbf{a}$ . If one visualizes a point  $\mathbf{a}$  in 3-space as a directed line segment from the origin  $\mathbf{0}$  to  $\mathbf{a}$ , then the sum of  $\mathbf{a}$  and  $\mathbf{b}$  is the diagonal of the parallelogram determined by  $\mathbf{a}$  and  $\mathbf{b}$  (unless  $\mathbf{a}$  and  $\mathbf{b}$  are parallel, in which case  $\mathbf{a} + \mathbf{b}$  is also parallel).



Because of the properties of scalar addition and product in  $\mathbb{R}$ , one easily verifies that  $\mathbb{R}^k$  with these operations has the following structure:

**Definition 1.1.1.** A (real) *vector space* is a set  $V$  together with two operations  $+$ :  $V \times V \rightarrow V$  and  $\cdot$ :  $\mathbb{R} \times V \rightarrow V$  that satisfy for all  $\mathbf{u}, \mathbf{v}, \mathbf{w} \in V$ ,  $a, b \in \mathbb{R}$ :

- (1)  $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$ ;
- (2)  $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$ ;
- (3) There exists an element  $\mathbf{0} \in V$ , called the *zero vector*, satisfying  $\mathbf{u} + \mathbf{0} = \mathbf{u}$  for all  $\mathbf{u}$ ;
- (4) Any  $\mathbf{u}$  in  $V$  has an additive inverse; i.e., an element  $-\mathbf{u}$  such that  $\mathbf{u} + (-\mathbf{u}) = \mathbf{0}$ .
- (5)  $(a + b) \cdot \mathbf{u} = a \cdot \mathbf{u} + b \cdot \mathbf{u}$ ;
- (6)  $a \cdot (\mathbf{u} + \mathbf{v}) = a \cdot \mathbf{u} + a \cdot \mathbf{v}$ ;
- (7)  $1 \cdot \mathbf{u} = \mathbf{u}$ ;
- (8)  $a \cdot (b \cdot \mathbf{u}) = (ab) \cdot \mathbf{u}$

Elements of  $V$  are called *vectors*. As in  $\mathbb{R}^k$ , we will for the most part omit the dot in scalar multiplication.

**Examples 1.1.1.** (i) An  $m \times n$  *matrix*  $A$  is a rectangular array of  $mn$  real numbers organized in  $m$  rows and  $n$  columns:

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

The element  $a_{ij}$  which lies on the  $i$ -th row and  $j$ -th column is referred to as the  $(i, j)$ -th entry of  $A$ . The sum of two such matrices  $A$  and  $B$  is defined to be the matrix of the same size whose  $(i, j)$ -th entry is  $a_{ij} + b_{ij}$ . Given  $c \in \mathbb{R}$ ,  $cA$  is the matrix whose

$(i, j)$ -th entry equals  $ca_{ij}$ . The collection  $M_{m,n}$  of all  $m \times n$  matrices is a vector space under these operations. In fact, when  $m = 1$ , we recover  $\mathbb{R}^n$ . Much of the time, though, we will identify  $\mathbb{R}^n$  with column matrices (i.e., with  $M_{n,1}$ ) rather than with row matrices.

- (ii) The collection  $\mathcal{F}(\mathbb{R})$  of all functions  $f : \mathbb{R} \rightarrow \mathbb{R}$  with addition and scalar multiplication defined in the usual way,

$$(f + g)(x) = f(x) + g(x), \quad (cf)(x) = cf(x), \quad c, x \in \mathbb{R}, \quad f, g \in \mathcal{F}(\mathbb{R}),$$

is a vector space.

- (iii) The collection  $\mathcal{C}(\mathbb{R})$  of continuous functions  $f : \mathbb{R} \rightarrow \mathbb{R}$  is a vector space with the operations from (ii).

Notice that a zero vector in a vector space  $V$  is necessarily unique: for if  $\mathbf{0}$  and  $\mathbf{0}'$  are two zero vectors, then  $\mathbf{0} + \mathbf{0}' = \mathbf{0}$  because  $\mathbf{0}'$  is a zero vector, and  $\mathbf{0} + \mathbf{0}' = \mathbf{0}'$  because  $\mathbf{0}$  is also one. A similar argument shows that a vector has a unique inverse. Furthermore, for any  $\mathbf{v} \in V$ ,  $0\mathbf{v} = \mathbf{0}$ , since  $0\mathbf{v} = (0 + 0)\mathbf{v} = 0\mathbf{v} + 0\mathbf{v}$ , and the conclusion follows by adding the inverse of  $0\mathbf{v}$  to both sides.

Example (iii) above is a subset of (ii). This is a quite common occurrence: A nonempty subset of a vector space  $V$  is called a *subspace* of  $V$  if it is a vector space with the operations inherited from  $V$ . Clearly, a necessary condition for  $W \subset V$  to be a subspace is that it be closed under addition and scalar multiplication; i.e., that the sum of vectors in  $W$  is again in  $W$ , and similarly for scalar multiplication. It turns out this is also a sufficient condition:

**Proposition 1.1.1.** *A nonempty subset  $W$  of a vector space  $V$  is a subspace of  $V$  if*

- (1)  $\mathbf{v} + \mathbf{w} \in W$  for any  $\mathbf{v}, \mathbf{w} \in W$ , and
- (2)  $c\mathbf{v} \in W$  for any  $c \in \mathbb{R}, \mathbf{v} \in W$ .

*Proof.* All the axioms for vector space, except possibly the existence of zero vector and inverses, are satisfied because they hold for elements of  $V$  and therefore also for elements of  $W$ . The remaining two axioms follow from closure under scalar multiplication: pick any  $\mathbf{v} \in W$ ; then  $\mathbf{0} = 0\mathbf{v} \in W$ . For the other axiom, observe that  $\mathbf{0} = 0\mathbf{v} = (1 - 1)\mathbf{v} = \mathbf{v} + (-1)\mathbf{v}$ , and by uniqueness of inverses,  $-\mathbf{v} = (-1)\mathbf{v}$ . This shows that  $-\mathbf{v} \in W$ .  $\square$

**Examples 1.1.2.** (i) The transpose  $A^T$  of an  $m \times n$  matrix  $A$  is the  $n \times m$  matrix obtained by interchanging the rows and columns of  $A$ ; i.e., the  $(i, j)$ -th element of  $A^T$  is  $a_{ji}$ . An  $n \times n$  matrix  $A$  is said to be symmetric (respectively skew-symmetric) if  $A = A^T$  (resp.  $A = -A^T$ ). The easily checked identities

$$(A + B)^T = A^T + B^T, \quad (cA)^T = cA^T,$$

imply that the set of symmetric matrices and that of skew-symmetric ones are both subspaces of  $M_{n,n}$ .

- (ii) The *trace* of an  $n \times n$  matrix  $A$  is the sum  $\text{tr } A = \sum_{i=1}^n a_{ii}$  of the elements on the diagonal. The collection of matrices with trace equal to 0 is a subspace of  $M_{n,n}$ , but those with trace any other number is not.
- (iii) A *linear combination* of vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k \in V$  is a vector of the form  $c_1\mathbf{v}_1 + \dots + c_k\mathbf{v}_k$ , where  $c_i \in \mathbb{R}$ . The collection  $\text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  of all such linear combinations is called the *span* of  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ . It follows immediately from Proposition 1.1.1 that the span of a subset  $S \subset V$  is a subspace of  $V$ . For instance, if  $V = \mathbb{R}^3$ , then the span of  $\{(1, 2, 0), (-1, 1, 2)\}$  is the space

$$\{s(1, 2, 0) + t(-1, 1, 2) \mid s, t \in \mathbb{R}\} = \{(s - t, 2s + t, 2t) \in \mathbb{R}^3 \mid s, t \in \mathbb{R}\}.$$

**Definition 1.1.2.** Vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k \in V$  are said to be *linearly dependent* if there exist scalars  $a_1, \dots, a_k$ , *not all zero*, such that  $a_1\mathbf{v}_1 + \dots + a_k\mathbf{v}_k = \mathbf{0}$ . Otherwise, they are said to be *linearly independent*.

For vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$ , the equation  $x_1\mathbf{v}_1 + \dots + x_k\mathbf{v}_k = \mathbf{0}$  in the variables  $x_i$  always has the trivial solution  $x_1 = \dots = x_k = 0$ . In fact, the collection of solutions  $(x_1, \dots, x_k)$  is a subspace of  $\mathbb{R}^k$ . The  $\mathbf{v}_i$ 's are linearly independent iff the trivial solution is the only one.

Another way to characterize linear dependence is by means of linear combinations (see Example 1.1.2 (iii)): A set  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is linearly dependent if and only if one of the vectors is a linear combination of the others. Indeed, if  $\mathbf{v}_i = \sum_{j \neq i} a_j\mathbf{v}_j$ , then  $\sum_l x_l\mathbf{v}_l = \mathbf{0}$  for  $x_l = a_l$  when  $l \neq i$ , and  $x_i = -1 \neq 0$ . Conversely, if the set is linearly dependent, then  $\sum a_i\mathbf{v}_i = \mathbf{0}$  for some scalars  $a_i$ , with at least one of them, say  $a_j$ , different from zero. Then  $\mathbf{v}_j = -\sum_{i \neq j} (a_i/a_j)\mathbf{v}_i$  is a linear combination of the others.

Notice that if a spanning set for a vector space is linearly dependent, then any vector that is a linear combination of the others can be discarded, and the remaining ones still span the space. The result is a minimal spanning set:

**Definition 1.1.3.** A *basis* of a vector space  $V$  is a linearly independent spanning set.

**Examples 1.1.3.** (i) Let  $\mathbf{e}_i$ ,  $1 \leq i \leq n$ , denote the point in  $\mathbb{R}^n$  with  $u^j(\mathbf{e}_i) = 1$  if  $i = j$  and 0 otherwise (recall that  $u^j$  is the  $j$ -th projection). The identity  $(a_1, \dots, a_n) = \sum_i a_i\mathbf{e}_i$  implies that the  $\mathbf{e}_i$  are both independent and span  $\mathbb{R}^n$ . This is the so-called *standard basis* of  $\mathbb{R}^n$ .

- (ii) The collection of points  $(x, y, z) \in \mathbb{R}^3$  such that  $x + y + z = 0$  is a subspace  $V$  of  $\mathbb{R}^3$ . The vectors  $\mathbf{u} = (1, 0, -1)$  and  $\mathbf{v} = (0, 1, -1)$  belong to  $V$ , and are linearly independent: the equation  $a\mathbf{u} + b\mathbf{v} = \mathbf{0}$  yields

$$a(1, 0, -1) + b(0, 1, -1) = (a, b, -a - b) = (0, 0, 0),$$

so that  $a = b = 0$ . They also span  $V$ , since  $(x, y, z) \in V$  iff  $z = -x - y$ ; i.e., if and only if

$$(x, y, z) = (x, y, x - y) = (x, 0, x) + (0, y, -y) = x\mathbf{u} + y\mathbf{v}.$$

Thus, they form a basis of  $V$ .

**Theorem 1.1.1.** *If  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a basis of  $V$ , then any subset of  $V$  that contains more than  $n$  vectors is linearly dependent.*

*Proof.* Let  $S = \{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  be a subset with  $k > n$  elements. If one of the  $\mathbf{u}_i$  is the zero vector, then  $S$  is certainly linearly dependent (see Exercise 1.4), so we may assume this is not the case. Since the  $\mathbf{v}_i$  span  $V$ , we may express  $\mathbf{u}_1$  as a linear combination  $\mathbf{u}_1 = \sum_i a_i \mathbf{v}_i$ , where at least one of the scalar coefficients is nonzero. By renumbering the basis elements if necessary, we may assume  $a_1 \neq 0$ , so that the above equation can be solved for  $\mathbf{v}_1$ :

$$\mathbf{v}_1 = (1/a_1)(\mathbf{u}_1 - a_2 \mathbf{v}_2 - \dots - a_n \mathbf{v}_n).$$

This implies that  $V$  is spanned by  $\{\mathbf{u}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ . Continuing in this way, we can replace the  $\mathbf{v}$ 's one by one with  $\mathbf{u}$ 's. More precisely, suppose  $V$  has been shown to be spanned by  $\{\mathbf{u}_1, \dots, \mathbf{u}_l, \mathbf{v}_{l+1}, \dots, \mathbf{v}_n\}$ . Then  $\mathbf{u}_{l+1}$  can be expressed as a linear combination

$$\mathbf{u}_{l+1} = \sum_{i \leq l} b_i \mathbf{u}_i + \sum_{i > l} c_i \mathbf{v}_i.$$

where at least one of the scalars is nonzero. We may assume one of the  $c_i$  is nonzero, for if they are all zero, then  $\mathbf{u}_{l+1}$  is a linear combination of the  $\mathbf{u}$ 's and  $S$  is then linearly dependent. Renumbering the vectors if necessary, we may suppose that  $c_{l+1} \neq 0$ , and we can write  $\mathbf{v}_{l+1}$  as a linear combination of  $\{\mathbf{u}_1, \dots, \mathbf{u}_{l+1}, \mathbf{v}_{l+2}, \dots, \mathbf{v}_n\}$ . This means the latter set spans  $V$ . By induction, we conclude that  $V$  is spanned by  $\mathbf{u}_1, \dots, \mathbf{u}_n$ . But then any  $\mathbf{u}_{n+i}$  is a linear combination of  $\mathbf{u}_1, \dots, \mathbf{u}_n$ , and the result follows.  $\square$

**Corollary 1.1.1.** *If  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  are two bases of  $V$ , then  $k = n$ .*

*Proof.* Since the  $\mathbf{v}_i$  are linearly independent, Theorem 1.1.1 implies that  $n \leq k$ . By symmetry,  $k \leq n$ .  $\square$

**Definition 1.1.4.** If  $V$  has a basis consisting of  $n$  elements, the *dimension*  $\dim V$  of  $V$  is defined to be  $n$ .

A trivial vector space (one that consists of a single element, which is then necessarily the zero vector) has no nonempty linearly independent spanning set. We therefore say it has the empty set as basis, and its dimension is zero. A vector space with dimension  $n$  for some  $n$  is called a *finite-dimensional* space. Unless specified otherwise, we will concern ourselves here only with finite-dimensional vector spaces, even though it should be noted that many spaces have infinite bases: for example, one such is the space of all polynomials with the usual addition of polynomials and scalar multiplication. The reader is invited to check that the collection  $\{p_k(x) \mid k = 0, 1, \dots\}$ , where  $p_k(x) = x^k$ , is a basis of this space.

One particularly useful feature of bases is that they can be used to uniquely identify vectors:

**Proposition 1.1.2.** *If  $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a basis of  $V$ , then any  $\mathbf{v} \in V$  can be written in one and only one way as a linear combination of the basis elements.*

*Proof.* That  $\mathbf{v}$  can be written as a linear combination follows from the fact that the  $\mathbf{v}_i$  span  $V$ . That the expression is unique follows from linear independence: if  $\sum a_i \mathbf{v}_i = \sum b_i \mathbf{v}_i$ , then  $\sum (a_i - b_i) \mathbf{v}_i = \mathbf{0}$ , so that each  $a_i - b_i$  must be zero.  $\square$

The uniqueness part of the above proposition enables us to make the following:

**Definition 1.1.5.** Let  $\mathcal{B} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  be an ordered basis of  $V$ . The *coordinate vector* of  $\mathbf{v} = a_1 \mathbf{v}_1 + \dots + a_n \mathbf{v}_n \in V$  with respect to  $\mathcal{B}$  is the element  $[\mathbf{v}]_{\mathcal{B}} = [a_1, \dots, a_n]^T \in M_{n,1}$ .

For example, if  $V = \mathbb{R}^n$ , then the coordinate vector of  $\mathbf{v}$  with respect to the standard basis is just  $\mathbf{v}$  itself, written as a column rather than a row. A more interesting example is the plane  $V$  with the ordered basis  $(\mathbf{u}, \mathbf{v})$  from Examples 1.1.3 (ii): The vector  $(1, -1, 0)$  belongs to  $V$ , and since it equals  $\mathbf{u} - \mathbf{v}$ , its coordinate vector in that basis is  $[1, -1]^T$ .

Notice that the order in which the basis elements are listed is crucial. We will, when referring to an ordered basis, dispense with braces when listing the elements, and either use parentheses as in the above definition or no symbols at all.

## 1.2 Linear transformations

**Definition 1.2.1.** A map  $L : V \rightarrow W$  between vector spaces  $V$  and  $W$  is said to be *linear* if it preserves the vector space operations; i.e., if

- (1)  $L(\mathbf{u} + \mathbf{v}) = L(\mathbf{u}) + L(\mathbf{v})$  for all  $\mathbf{u}, \mathbf{v} \in V$ ;
- (2)  $L(c\mathbf{u}) = cL(\mathbf{u})$  for all  $c \in \mathbb{R}, \mathbf{u} \in V$ .

When there is a single vector  $\mathbf{u}$  in the argument, it is customary to write  $L\mathbf{u}$  instead of  $L(\mathbf{u})$ . Notice that the  $+$  on the right of the equality sign in (1) denotes the addition in  $W$ , whereas the one on the left refers to addition in  $V$ . A similar observation involving scalar multiplication holds for (2). If  $L$  is invertible, i.e., if there exists a map  $L^{-1} : W \rightarrow V$  such that  $L \circ L^{-1} = 1_W$  and  $L^{-1} \circ L = 1_V$  ( $1_V$  is the identity map on  $V$ ), then  $L$  is said to an *isomorphism*; in this case  $V$  and  $W$  are said to be *isomorphic*. Observe that the inverse of an isomorphism is again linear: given  $\mathbf{w}_i \in W, i = 1, 2$ , there exist unique  $\mathbf{v}_i \in V$  such that  $L\mathbf{v}_i = \mathbf{w}_i$ . Then

$$\begin{aligned} L^{-1}(\mathbf{w}_1 + \mathbf{w}_2) &= L^{-1}(L\mathbf{v}_1 + L\mathbf{v}_2) = L^{-1}(L(\mathbf{v}_1 + \mathbf{v}_2)) = \mathbf{v}_1 + \mathbf{v}_2 \\ &= L^{-1}\mathbf{w}_1 + L^{-1}\mathbf{w}_2, \end{aligned}$$

and a similar argument shows that  $L^{-1}$  preserves scalar multiplication. Isomorphic vector spaces are identical as far as their algebraic structure is concerned, and can thus be identified.

**Examples 1.2.1.** (i) The map  $[\ ]_{\mathcal{S}} : \mathbb{R}^n \rightarrow M_{n,1}$  that assigns to a point its coordinate vector with respect to the standard basis (see Definition 1.1.5) is an isomorphism. This formally justifies why we can use either row or column matrices to denote elements of  $\mathbb{R}^n$ .

- (ii) More generally, the map  $M_{m,n} \rightarrow M_{n,m}$  that sends a matrix to its transpose is an isomorphism.
- (iii) Let  $V$  be any  $n$ -dimensional vector space with basis  $\mathcal{B}$ . As in (i), we see that the map  $[\ ]_{\mathcal{B}} : V \rightarrow M_{n,1}$  that assigns to a vector in  $V$  its coordinate vector with respect to  $\mathcal{B}$  is an isomorphism. Together with (i), this shows that any  $n$ -dimensional vector space is isomorphic to  $\mathbb{R}^n$ . This isomorphism is not canonical, however, in the sense that it depends on the chosen basis.

More generally, any two vector spaces of the same dimension are isomorphic: if  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a basis of  $V$ , and  $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$  is a basis of  $W$ , define  $L : V \rightarrow W$  as follows: given  $\mathbf{v} \in V$ , it can be uniquely written as  $\mathbf{v} = \sum a_i \mathbf{v}_i$  for some choices of  $a_i \in \mathbb{R}$ . Set  $L\mathbf{v} = \sum a_i \mathbf{w}_i$ . It is straightforward to check that  $L$  is a well-defined linear transformation, and that it is entirely determined by the fact that  $L\mathbf{v}_i = \mathbf{w}_i$ ,  $i = 1, \dots, n$ . In fact, if  $W$  is a vector space of any dimension, then for any vectors  $\mathbf{w}_i \in W$ ,  $i = 1, \dots, n$ , there exists a unique linear map  $L : V \rightarrow W$  such that  $L(\mathbf{v}_i) = \mathbf{w}_i$  for all  $i = 1, \dots, n$ , see Exercise 1.10.  $L$  will not, in general, be an isomorphism, though.

- (iv) Let  $C^\infty(\mathbb{R})$  denote the space of functions from  $\mathbb{R}$  to itself that have derivatives of any order. The map  $D : C^\infty(\mathbb{R}) \rightarrow C^\infty(\mathbb{R})$ ,  $Df = f'$ , is linear. Why is it not an isomorphism?

Linear maps between finite-dimensional Euclidean spaces are most conveniently expressed in terms of matrices. If  $A = [a_{ij}]$  is an  $m \times n$  matrix, and  $B = [b_{ij}]$  is an  $n \times k$  matrix, the *product* matrix is the  $m \times k$  matrix  $AB$  whose  $(i, j)$ -th entry is  $\sum_{l=1}^n a_{il}b_{lj}$ . For example,

$$\begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & -1 \\ -1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 0 & 1 \\ -1 & 0 \end{bmatrix} = \begin{bmatrix} -2 & 4 \\ 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad \begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}.$$

One easily verifies that matrix multiplication has the following properties:

**Proposition 1.2.1.** For appropriately sized matrices  $A, B, C$ , and  $c \in \mathbb{R}$ ,

- (1)  $A(BC) = (AB)C$ ;
- (2)  $(A + B)C = AC + BC$ , and  $A(B + C) = AB + AC$ ;
- (3)  $c(AB) = (cA)B = A(cB)$ .

The expression “appropriately sized” means that the matrices are assumed to have the correct size for addition and multiplication; for instance, in the first identity, the number of columns of  $A$  equals the number of rows of  $B$ , and the number of columns of  $B$  equals the number of rows of  $C$ . Unlike scalar multiplication, though,  $AB \neq BA$  in general (even when both products make sense). There is nevertheless a matrix that

plays the same role that 1 does in scalar multiplication: the  $n \times n$  *identity matrix*

$$I_n = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

is the matrix whose  $(i, j)$ -th entry is 1 if  $i = j$  and 0 otherwise. It follows from the definition of multiplication that if  $A$  is  $m \times n$ , then  $AI_n = A$  and if  $B$  is  $n \times m$ , then  $I_n B = B$ . A (necessarily square)  $n \times n$  matrix  $A$  is said to be *invertible* if there exists a matrix  $B$  such that  $AB = BA = I_n$ . If such a  $B$  exists (and in general it need not), then it is unique: indeed, if  $B$  and  $B'$  are two such, then  $B' = B'I_n = B'(AB) = (B'A)B = I_n B = B$ . We then call  $B$  the *inverse* of  $A$ , and denote it  $A^{-1}$ .

**Examples 1.2.2.** (i) The reader is invited to verify that the matrices

$$A = \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix} \text{ and } B = \begin{bmatrix} -1 & 1 \\ 2 & -1 \end{bmatrix}$$

satisfy  $AB = BA = I_2$ , so that  $B = A^{-1}$ .

(ii) The matrix

$$A = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}$$

is not invertible. One way of seeing this is by noting that  $AA$  equals the zero matrix  $\mathbf{0}$  whose entries are all 0. If  $A$  did have an inverse, then we would have that  $A = I_2 A = (A^{-1}A)A = A^{-1}(AA) = A^{-1}\mathbf{0} = \mathbf{0}$ , which is not true.

Matrix multiplication provides many examples of linear transformations:

**Example 1.2.3.** Let  $A$  be an  $m \times n$  matrix. For any  $\mathbf{u} \in \mathbb{R}^n$ ,  $A\mathbf{u}^T \in M_{1,m}$ , so that  $(A\mathbf{u}^T)^T \in \mathbb{R}^m$ . By Proposition 1.2.1 (2), (3), and the fact that transposing is linear, *left multiplication*  $L_A$  by  $A$ ,

$$\begin{aligned} L_A : \mathbb{R}^n &\rightarrow \mathbb{R}^m, \\ \mathbf{u} &\mapsto (A\mathbf{u}^T)^T \end{aligned}$$

is a linear transformation. We will shortly see that *every* linear transformation between Euclidean spaces is of this form.

A word about terminology: left multiplication by  $A$  may sound like a poor choice of words, since  $L_A \mathbf{u} = \mathbf{u}A^T$ , so that  $L_A$  is really right multiplication by the transpose of  $A$ . Nevertheless, as observed earlier, it is customary to identify elements of  $\mathbb{R}^n$  with column rather than row matrices, and when doing so, we may write  $L_A \mathbf{u} = A\mathbf{u}$ . We will be deliberately vague as to which representation we use, but a general rule of thumb is that an element of  $\mathbb{R}^n$  considered as a point is represented by a row matrix, whereas an element considered as a vector is represented as a column.



**Definition 1.2.2.** Let  $L : V \rightarrow W$  be a linear transformation, and suppose  $\mathcal{B} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$  and  $\mathcal{C} = (\mathbf{v}_1, \dots, \mathbf{v}_m)$  are ordered bases of  $V$  and  $W$  respectively. The *matrix of  $L$  with respect to  $\mathcal{B}$  and  $\mathcal{C}$*  is the  $m \times n$  matrix

$$[L]_{\mathcal{B},\mathcal{C}} = \begin{bmatrix} [L\mathbf{u}_1]_{\mathcal{C}} & \cdots & [L\mathbf{u}_n]_{\mathcal{C}} \end{bmatrix}$$

whose  $i$ -th column is the coordinate vector of  $L\mathbf{u}_i$  with respect to  $\mathcal{C}$ , cf. Examples 1.2.1 (iii).

**Example 1.2.4.** Let  $P_n$  denote the  $(n + 1)$ -dimensional space of polynomials of degree  $\leq n$  with its standard basis  $\mathcal{S} = \{1, x, \dots, x^n\}$ , and consider the derivative operator  $D : P_2 \rightarrow P_1$ . If  $\mathcal{B}$  is the basis  $\{1, 1 + x\}$  of  $P_1$ , then  $D1 = 0$ ,  $Dx = 1$ , and  $Dx^2 = 2x = -2 \cdot 1 + 2(1 + x)$ . Thus,

$$[D]_{\mathcal{S},\mathcal{B}} = \begin{bmatrix} 0 & 1 & -2 \\ 0 & 0 & 2 \end{bmatrix}.$$

The following result implies that any linear transformation is entirely determined by its matrix with respect to any given bases; in fact, it says that once bases are fixed, the vector space of all linear transformations from  $V$  to  $W$  is isomorphic to the space of  $m \times n$  matrices, where  $\dim V = n$  and  $\dim W = m$ , see Exercise 1.16:

**Theorem 1.2.1.** Let  $L : V \rightarrow W$  be a linear transformation, and suppose  $\mathcal{B} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  and  $\mathcal{C} = \{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  are bases of  $V$  and  $W$  respectively. Then for any  $\mathbf{u} \in V$ ,

$$[L\mathbf{u}]_{\mathcal{C}} = [L]_{\mathcal{B},\mathcal{C}}[\mathbf{u}]_{\mathcal{B}}.$$

*Proof.* We begin by observing that if  $A$  is an  $m \times n$  matrix with columns  $A_i$ , and  $\mathbf{u} = [x_1 \dots x_n]^T \in \mathbb{R}^n$ , then by definition of matrix multiplication,  $A\mathbf{u} = x_1A_1 + \dots + x_nA_n$ . Now, if  $\mathbf{v} = \sum a_i\mathbf{v}_i \in V$ , then  $[\mathbf{v}]_{\mathcal{B}} = [a_1 \dots a_n]^T$ , and together with Examples 1.2.1 (iii),

$$[L\mathbf{v}]_{\mathcal{C}} = [L(\sum a_i\mathbf{v}_i)]_{\mathcal{C}} = [\sum a_iL\mathbf{v}_i]_{\mathcal{C}} = \sum a_i[L\mathbf{v}_i]_{\mathcal{C}} = [L]_{\mathcal{B},\mathcal{C}}[\mathbf{v}]_{\mathcal{B}},$$

as claimed.  $\square$

**Corollary 1.2.1.** (1) If  $A$  is the matrix of a linear transformation  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  with respect to the standard bases, then  $L = L_A$  (see Example 1.2.3).

(2) Let  $L : V_1 \rightarrow V_2$ ,  $T : V_2 \rightarrow V_3$  be linear, and  $\mathcal{B}_i$  be a basis of  $V_i$ ,  $1 \leq i \leq 3$ . Then  $[T \circ L]_{\mathcal{B}_1,\mathcal{B}_3} = [T]_{\mathcal{B}_2,\mathcal{B}_3}[L]_{\mathcal{B}_1,\mathcal{B}_2}$ . In particular, if  $L : V_1 \rightarrow V_2$  is an isomorphism, then  $[L^{-1}]_{\mathcal{B}_2,\mathcal{B}_1} = [L]_{\mathcal{B}_1,\mathcal{B}_2}^{-1}$ .

(3) (Change of basis) If  $\mathcal{B}$  and  $\mathcal{C}$  are bases of  $V$ , then  $[\mathbf{v}]_{\mathcal{C}} = [1_V]_{\mathcal{B},\mathcal{C}}[\mathbf{v}]_{\mathcal{B}}$  for any  $\mathbf{v} \in V$ .

*Proof.* (1) We denote by the same letter  $\mathcal{S}$  the standard bases in both Euclidean spaces. Recalling that a vector in Euclidean space equals its coordinate vector in the standard basis, Theorem 1.2.1 implies that for any  $\mathbf{v} \in \mathbb{R}^n$ ,

$$L\mathbf{v} = [L\mathbf{v}]_{\mathcal{S}} = [L]_{\mathcal{S},\mathcal{S}}[\mathbf{v}]_{\mathcal{S}} = A\mathbf{v} = L_A\mathbf{v}.$$

(2) Given  $\mathbf{u} \in V_1$ ,

$$[(T \circ L)\mathbf{u}]_{\mathcal{B}_3} = [T(L\mathbf{u})]_{\mathcal{B}_3} = [T]_{\mathcal{B}_2, \mathcal{B}_3} [L\mathbf{u}]_{\mathcal{B}_2} = [T]_{\mathcal{B}_2, \mathcal{B}_3} [L]_{\mathcal{B}_1, \mathcal{B}_2} [\mathbf{u}]_{\mathcal{B}_1}.$$

The second statement follows by taking  $T = L^{-1}$ , observing that  $L^{-1} \circ L$  is the identity map, and the matrix  $[1_{V_1}]_{\mathcal{B}, \mathcal{B}}$  of the identity map with respect to any basis  $\mathcal{B}$  is the identity matrix.

(3) is an immediate consequence of Theorem 1.2.1.  $\square$

In light of Theorem 1.2.1 (3), we refer to  $[1_V]_{\mathcal{B}, \mathcal{C}}$  as the *change of basis matrix* from the basis  $\mathcal{B}$  to  $\mathcal{C}$ .

**Examples 1.2.5.** (i) Theorem 1.2.1 has many other applications that can be proved by arguments similar to the ones used above. Suppose, for example, that  $L : V \rightarrow V$  is a linear transformation, and  $\mathcal{B}, \mathcal{C}$  are two bases of  $V$ . How are the matrices of  $L$  with respect to the two bases related? To answer this, let  $\mathbf{v} \in V$ , and observe that

$$[L\mathbf{v}]_{\mathcal{C}} = [(1_V \circ L \circ 1_V)\mathbf{v}]_{\mathcal{C}} = [1_V]_{\mathcal{C}, \mathcal{B}} [L]_{\mathcal{C}, \mathcal{C}} [1_V]_{\mathcal{B}, \mathcal{C}} [\mathbf{v}]_{\mathcal{B}},$$

so that if  $P$  denotes the change of basis matrix  $[1_V]_{\mathcal{B}, \mathcal{C}}$  from  $\mathcal{B}$  to  $\mathcal{C}$ , then

$$[L]_{\mathcal{C}, \mathcal{C}} = P^{-1} [L]_{\mathcal{B}, \mathcal{B}} P. \quad (1.2.1)$$

For the sake of brevity, the matrix  $[L]_{\mathcal{B}, \mathcal{B}}$  of  $L$  in a given basis  $\mathcal{B}$  will often be denoted by  $[L]_{\mathcal{B}}$ , so that (1.2.1) reads

$$[L]_{\mathcal{C}} = P^{-1} [L]_{\mathcal{B}} P.$$

Two  $n \times n$  matrices  $A$  and  $B$  are said to be *similar* if there exists an invertible matrix  $P$  such that  $A = P^{-1}BP$ . This is easily seen to be an equivalence relation for the class of all  $n \times n$  matrices (meaning (1) any  $A$  is similar to itself, (2) if  $A$  is similar to  $B$ , then  $B$  is similar to  $A$ , and (3) if  $A$  is similar to  $B$  and  $B$  is similar to  $C$ , then  $A$  is similar to  $C$ ). (1.2.1) says that the matrices of a linear transformation  $L : V \rightarrow V$  with respect to two different bases are similar.

(ii) The set  $\mathcal{B} = (\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3)$  where

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 2 \\ 1 \\ -1 \end{bmatrix}, \quad \mathbf{v}_3 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix},$$

is an ordered basis of  $\mathbb{R}^3$ . Furthermore,  $\mathbf{e}_1 = (1/2)\mathbf{v}_1 + (1/4)\mathbf{v}_2 - (1/4)\mathbf{v}_3$ ,  $\mathbf{e}_2 = (-1/2)\mathbf{v}_1 + (1/4)\mathbf{v}_2 - (1/4)\mathbf{v}_3$ , and  $\mathbf{e}_3 = (1/2)\mathbf{v}_1 - (1/4)\mathbf{v}_2 + (1/4)\mathbf{v}_3$ , so that the change of basis matrix from the standard basis  $\mathcal{S}$  to  $\mathcal{B}$  is

$$[1_{\mathbb{R}^3}]_{\mathcal{S}, \mathcal{B}} = [\mathbf{e}_1]_{\mathcal{B}} \quad [\mathbf{e}_2]_{\mathcal{B}} \quad [\mathbf{e}_3]_{\mathcal{B}} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{4} & -\frac{1}{4} \\ -\frac{1}{4} & \frac{3}{4} & \frac{1}{4} \end{bmatrix}.$$

This means for example that the coordinate vector of  $\begin{bmatrix} 1 & 2 & 1 \end{bmatrix}^T$  in the basis  $\mathcal{B}$  equals

$$\begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}_{\mathcal{B}} = [\mathbf{1}_{\mathbb{R}^3}]_{\mathcal{B}, \mathcal{S}} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}_{\mathcal{S}} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{4} & -\frac{1}{4} \\ -\frac{1}{4} & \frac{3}{4} & \frac{1}{4} \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ \frac{1}{2} \\ \frac{3}{2} \end{bmatrix},$$

which reflects the fact that  $\begin{bmatrix} 1 & 2 & 1 \end{bmatrix}^T = (1/2)\mathbf{v}_2 + (3/2)\mathbf{v}_3$ .

**Definition 1.2.3.** The *kernel* or *nullspace*  $\ker L$  of a linear map  $L : V \rightarrow W$  is the set of all  $\mathbf{v} \in V$  that are mapped to  $\mathbf{0} \in W$  by  $L$ . The *image*  $\text{Im } L$  of  $L$  is the collection of all elements of  $W$  that can be expressed as  $L\mathbf{v}$  for some  $\mathbf{v} \in V$ . It is often denoted  $L(V)$ .

It is easily seen that kernel and image are subspaces of  $V$  and  $W$  respectively.

**Theorem 1.2.2.** If  $L : V \rightarrow W$  is linear, then  $\dim V = \dim \ker L + \dim \text{Im } L$ .

*Proof.* Let  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  denote a basis of  $\ker L$ , and  $\{\mathbf{w}_1, \dots, \mathbf{w}_r\}$  one of  $\text{Im } L$ . For each  $i = 1, \dots, r$ , choose some  $\mathbf{v}_i \in V$  such that  $L\mathbf{v}_i = \mathbf{w}_i$ . The claim follows once we establish that the set  $S = \{\mathbf{u}_1, \dots, \mathbf{u}_k, \mathbf{v}_1, \dots, \mathbf{v}_r\}$  is a basis of  $V$ . To see that it spans  $V$ , take an arbitrary element  $\mathbf{v} \in V$ . By assumption,  $L\mathbf{v} = \sum_i a_i \mathbf{w}_i$  for some scalars  $a_i$ . Since  $L(\mathbf{v} - \sum_i a_i \mathbf{v}_i) = L\mathbf{v} - \sum_i a_i \mathbf{w}_i = \mathbf{0}$ , the vector  $\mathbf{v} - \sum_i a_i \mathbf{v}_i$  lies in the kernel of  $L$  and can therefore be expressed as a linear combination of the  $\mathbf{u}_i$ . In other words,  $\mathbf{v}$  lies in the span of  $S$ . For linear independence, consider the equation

$$\sum_{i=1}^k a_i \mathbf{u}_i + \sum_{j=1}^r b_j \mathbf{v}_j = \mathbf{0}.$$

Applying  $L$  to both sides, and recalling that the first summand lies in the kernel of  $L$ , we obtain  $\sum b_j \mathbf{w}_j = \mathbf{0}$ . By linear independence of the  $\mathbf{w}_j$ , each  $b_j$  must vanish. But then only the first summand remains in the above equation, and by linear independence of the  $\mathbf{u}_i$ , each  $a_i = 0$ .  $\square$

The dimension of the kernel of  $L$  is called the *nullity* of  $L$ , that of its image the *rank* of  $L$ .

**Example 1.2.6.** Let  $L : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  be given by

$$L \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x + z \\ y + 2z \\ 2x + y + 4z \end{bmatrix}.$$

$L$  is linear, since it is left multiplication by the matrix

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 2 \\ 2 & 1 & 4 \end{bmatrix}.$$

The columns of  $A$  are the image via  $L$  of the standard coordinate vectors, and they are linearly dependent: in fact,  $L\mathbf{e}_3 = L\mathbf{e}_1 + 2L\mathbf{e}_2$ . Thus, the image of  $L$  is spanned by the

first 2 columns of  $A$ . Since they are linearly independent,  $L$  has rank 2. On the other hand,  $[xyz]^T$  belongs to the kernel of  $L$  iff  $x + z = 0$ ,  $y + 2z = 0$ , and  $2x + y + 4z = 0$ . The last equation is just the sum of the second one with two times the first one, and can be discarded. The first two equations enable us to express  $x$  and  $y$  in terms of  $z$ , so that the kernel consists of all vectors of the form

$$\begin{bmatrix} -z \\ -2z \\ z \end{bmatrix} = z \begin{bmatrix} -1 \\ -2 \\ 1 \end{bmatrix}, \quad z \in \mathbb{R},$$

and is therefore one-dimensional, as predicted by Theorem 1.2.2.

### 1.3 Determinants

Determinants play a crucial role in Linear Algebra, and by extension, in Calculus and Differential Geometry. Their properties are closely related to those of permutations: let  $J_n = \{1, \dots, n\}$ . A *permutation* of  $J_n$  is a bijection  $\sigma : J_n \rightarrow J_n$ . It is often represented by

$$[\sigma(1) \ \sigma(2) \ \dots \ \sigma(n)].$$

For example,

$$\sigma = [2 \ 4 \ 1 \ 3]$$

is the permutation of  $J_4$  such that  $\sigma(1) = 2$ ,  $\sigma(2) = 4$ ,  $\sigma(3) = 1$ , and  $\sigma(4) = 3$ . The composition  $\sigma \circ \tau$  of two permutations  $\sigma$  and  $\tau$  is again a permutation, denoted  $\sigma\tau$  for brevity, and the collection of all permutations of  $J_n$  is denoted  $S_n$ . It is easy to see that  $S_n$  consists of  $n!$  elements, where

$$n! = n(n-1)(n-2) \cdots 2,$$

since there are  $n$  possible choices in assigning the value of  $\sigma$  at, say, 1, then only  $n-1$  remaining choices for  $\sigma(2)$ , and so on.

A *transposition* is a permutation that interchanges two elements and leaves all the others fixed. We denote by  $(i, j)$  the transposition that interchanges  $i$  and  $j$ .

**Lemma 1.3.1.** *Every permutation is a product of transpositions.*

*Proof.* The term product in the statement actually means composition; it is suggested by the notation in use. The argument will be by induction on  $n$ . For  $n = 1$ , there is nothing to prove. So assume that every permutation of  $J_{n-1}$  is a product of transpositions, and consider  $\sigma \in S_n$ . Let  $k = \sigma(n)$ . Define  $\tau$  to be the transposition  $(k, n)$  if  $k \neq n$ , and the identity if  $k = n$ . Then  $\tau\sigma(n) = \tau(k) = n$ , so  $\tau\sigma \in S_{n-1}$  and by the induction hypothesis, there exist transpositions  $\tau_1, \dots, \tau_j$  such that  $\tau\sigma = \tau_1 \cdots \tau_j$ . Since  $\tau$  equals its own inverse,

$$\sigma = \tau\tau_1 \cdots \tau_j$$

as claimed. □

The argument explicitly yields an algorithm for expressing a permutation as a product of transpositions. For example, the permutation

$$\sigma = [2 \quad 4 \quad 1 \quad 3]$$

considered earlier can be written as  $(3, 4)(3, 1)(1, 2)$ . There is no unique representation of a permutation as a product of transpositions, but the number of transpositions that appear is either always even or always odd. To see this, consider the polynomial  $p$  in  $n$  variables

$$p(x_1, \dots, x_n) = \prod_{1 \leq i < j \leq n} (x_i - x_j).$$

For example, when  $n = 3$ ,

$$p(x_1, x_2, x_3) = (x_1 - x_2)(x_1 - x_3)(x_2 - x_3).$$

Define a map

$$\begin{aligned} F : S_n \times \{\pm p\} &\rightarrow \{\pm p\}, \\ (\sigma, \pm p) &\mapsto \pm \prod_{i < j} (x_{\sigma(i)} - x_{\sigma(j)}). \end{aligned}$$

It is easily checked that  $F(\text{id}, \pm p) = \pm p$ , and  $F(\sigma\tau, \pm p) = F(\sigma, F(\tau, \pm p))$  for any  $\tau, \sigma \in S_n$  (such a map is called an *action* of  $S_n$  on the set  $\{\pm p\}$ ). Observe that for a transposition  $\tau$ ,  $F(\tau, p) = -p$ . Thus, if  $\sigma \in S_n$  can be represented as the product of an even number of transpositions, then  $F(\sigma, p) = p$ , and consequently any other representation of  $\sigma$  also involves an even number of transpositions.

**Definition 1.3.1.** A permutation is said to be *even* if it can be written as a product of an even number of transpositions. A permutation that is not even is said to be *odd*. The *sign*  $\varepsilon(\sigma)$  of  $\sigma$  is defined to be  $+1$  if  $\sigma$  is even, and  $-1$  if it is odd.

A map  $M : V_1 \times \dots \times V_k \rightarrow V$  from a Cartesian product of vector spaces  $V_i$  to a vector space  $V$  is said to be *multilinear* if it is linear in each component; i.e., if for any  $1 \leq i \leq k$  and  $\mathbf{v}_j \in V_j, j \neq i$ , the map

$$\begin{aligned} V_i &\rightarrow V \\ \mathbf{v} &\mapsto M(\mathbf{v}_1, \dots, \mathbf{v}_{i-1}, \mathbf{v}, \mathbf{v}_{i+1}, \dots, \mathbf{v}_k) \end{aligned}$$

is linear. A multilinear map  $M$  is said to be *alternating* if

$$M(\mathbf{v}_1, \dots, \mathbf{v}_i, \dots, \mathbf{v}_j, \dots, \mathbf{v}_k) = -M(\mathbf{v}_1, \dots, \mathbf{v}_j, \dots, \mathbf{v}_i, \dots, \mathbf{v}_k),$$

for all  $1 \leq i < j \leq k$ ,  $\mathbf{v}_l \in V_l, l = 1, \dots, k$ ; equivalently,  $M$  is alternating if and only if  $M(\mathbf{v}_1, \dots, \mathbf{v}_i, \dots, \mathbf{v}_j, \dots, \mathbf{v}_k) = 0$  whenever  $\mathbf{v}_i = \mathbf{v}_j$  for some  $1 \leq i < j \leq k$ : clearly, if  $M$  is alternating, then  $M$  evaluates to zero when applied to a list of vectors at least one of which is repeated. Conversely, if  $M$  evaluates to zero under these conditions, then

$$M(\mathbf{v}_1, \dots, \mathbf{v}_i + \mathbf{v}_j, \dots, \mathbf{v}_i + \mathbf{v}_j, \dots, \mathbf{v}_k) = 0.$$

By multilinearity, the expression on the left is a sum of four terms:  $M(\mathbf{v}_1, \dots, \mathbf{v}_i, \dots, \mathbf{v}_i, \dots, \mathbf{v}_k)$ , which by assumption vanishes; a similar term with  $\mathbf{v}_j$  replacing  $\mathbf{v}_i$  which also vanishes; and

$$M(\mathbf{v}_1, \dots, \mathbf{v}_i, \dots, \mathbf{v}_j, \dots, \mathbf{v}_k) + M(\mathbf{v}_1, \dots, \mathbf{v}_j, \dots, \mathbf{v}_i, \dots, \mathbf{v}_k).$$

Since this must equal zero, the claim follows.

With these preliminaries out of the way, we are able to prove the following:

**Theorem 1.3.1.** *There exists a unique map  $\det : (\mathbb{R}^n)^n = \mathbb{R}^n \times \dots \times \mathbb{R}^n \rightarrow \mathbb{R}$  satisfying the following:*

- (1) *det is multilinear;*
- (2) *det is alternating;*
- (3)  $\det(\mathbf{e}_1, \dots, \mathbf{e}_n) = 1$ .

*Proof.* We begin with uniqueness. Let  $\mathbf{v}_i = \sum_k a_{ki} \mathbf{e}_k$ ,  $i = 1, \dots, n$ . Since det is multilinear,

$$\begin{aligned} \det(\mathbf{v}_1, \dots, \mathbf{v}_n) &= \det(a_{11}\mathbf{e}_1 + \dots + a_{n1}\mathbf{e}_n, \dots, a_{1n}\mathbf{e}_1 + \dots + a_{nn}\mathbf{e}_n) \\ &= \sum_{\sigma} a_{\sigma(1)1} \cdots a_{\sigma(n)n} \det(\mathbf{e}_{\sigma(1)}, \dots, \mathbf{e}_{\sigma(n)}), \end{aligned}$$

where the sum runs over all maps  $\sigma : J_n \rightarrow J_n$ . Now, if  $\sigma$  is not one-to-one, i.e., if  $\sigma(i) = \sigma(j)$  for some  $i \neq j$ , then  $\det(\mathbf{e}_{\sigma(1)}, \dots, \mathbf{e}_{\sigma(n)}) = 0$  because det is alternating. Thus, the sum actually runs over all permutations of  $J_n$ . Finally, by properties (2), (3), and the definition of the sign of a permutation,  $\det(\mathbf{e}_{\sigma(1)}, \dots, \mathbf{e}_{\sigma(n)}) = \varepsilon(\sigma)$ . Summarizing, if det is to satisfy the three properties listed, then

$$\det \left( \begin{bmatrix} a_{11} \\ \vdots \\ a_{n1} \end{bmatrix}, \dots, \begin{bmatrix} a_{1n} \\ \vdots \\ a_{nn} \end{bmatrix} \right) = \sum_{\sigma \in S_n} \varepsilon(\sigma) a_{\sigma(1)1} \cdots a_{\sigma(n)n}, \quad (1.3.1)$$

which establishes uniqueness. To determine existence, it is enough to show that the above equation defines a map satisfying the three stated properties. For the first one, we compute

$$\begin{aligned} &\det \left( \begin{bmatrix} a_{11} \\ \vdots \\ a_{n1} \end{bmatrix}, \dots, \begin{bmatrix} a_{1i} \\ \vdots \\ a_{ni} \end{bmatrix} + c \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}, \dots, \begin{bmatrix} a_{1n} \\ \vdots \\ a_{nn} \end{bmatrix} \right) \\ &= \det \left( \begin{bmatrix} a_{11} \\ \vdots \\ a_{n1} \end{bmatrix}, \dots, \begin{bmatrix} a_{1i} + cb_1 \\ \vdots \\ a_{ni} + cb_n \end{bmatrix} + \cdots \begin{bmatrix} a_{1n} \\ \vdots \\ a_{nn} \end{bmatrix} \right) \\ &= \sum_{\sigma \in S_n} \varepsilon(\sigma) a_{\sigma(1)1} \cdots (a_{\sigma(i)i} + cb_{\sigma(i)}) \cdots a_{\sigma(n)n} \end{aligned}$$

$$\begin{aligned}
&= \sum_{\sigma \in S_n} \varepsilon(\sigma) a_{\sigma(1)1} \cdots a_{\sigma(i)i} \cdots a_{\sigma(n)n} + c \sum_{\sigma \in S_n} \varepsilon(\sigma) a_{\sigma(1)1} \cdots b_{\sigma(i)} \cdots a_{\sigma(n)n} \\
&= \det \left( \begin{bmatrix} a_{11} \\ \vdots \\ a_{n1} \end{bmatrix}, \dots, \begin{bmatrix} a_{1i} \\ \vdots \\ a_{ni} \end{bmatrix}, \dots, \begin{bmatrix} a_{1n} \\ \vdots \\ a_{nn} \end{bmatrix} \right) \\
&\quad + c \det \left( \begin{bmatrix} a_{11} \\ \vdots \\ a_{n1} \end{bmatrix}, \dots, \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}, \dots, \begin{bmatrix} a_{1n} \\ \vdots \\ a_{nn} \end{bmatrix} \right),
\end{aligned}$$

which shows that  $\det$  is multilinear. For the second property, let  $\tau \in S_n$ , and  $\mathbf{v}_i = [a_{1i} \ \dots \ a_{ni}]^T$ . Then, with notation as above, and noting that  $\varepsilon^2(\tau) = 1$ , we have

$$\begin{aligned}
\det(\mathbf{v}_{\tau(1)}, \dots, \mathbf{v}_{\tau(n)}) &= \sum_{\sigma \in S_n} \varepsilon(\sigma) a_{\sigma \circ \tau(1)1} \cdots a_{\sigma \circ \tau(n)n} \\
&= \varepsilon(\tau) \sum_{\sigma \in S_n} \varepsilon(\sigma) \varepsilon(\tau) a_{\sigma \circ \tau(1)1} \cdots a_{\sigma \circ \tau(n)n} \\
&= \varepsilon(\tau) \sum_{\sigma \circ \tau \in S_n} \varepsilon(\sigma \circ \tau) a_{\sigma \circ \tau(1)1} \cdots a_{\sigma \circ \tau(n)n} \\
&= \varepsilon(\tau) \det(\mathbf{v}_1, \dots, \mathbf{v}_n).
\end{aligned}$$

Notice that for the third equality, we used the fact that  $S_n = \{\sigma \circ \tau \mid \sigma \in S_n\}$ . Taking  $\tau$  to be a transposition now shows that  $\det$  is alternating. The last property is immediate.  $\square$

If  $A$  is an  $n \times n$  matrix, we define its *determinant* to be the number  $\det A = \det(\mathbf{a}_1, \dots, \mathbf{a}_n)$ , where  $\mathbf{a}_i$  is the  $i$ -th column of  $A$ . Thus, by (1.3.1),

$$\det \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \cdots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} = \sum_{\sigma \in S_n} \varepsilon(\sigma) a_{\sigma(1)1} \cdots a_{\sigma(n)n}. \quad (1.3.2)$$

For example, in the case of a  $2 \times 2$  matrix,  $S_2$  consists only of the identity and the transposition  $(1, 2)$ , so that

$$\det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad - bc.$$

The definition given here is not very convenient for computing determinants of larger matrices. In order to describe a different approach for  $n > 1$ , let us denote by  $A_{ij}$  the  $(n-1) \times (n-1)$  matrix obtained by deleting row  $i$  and column  $j$  from  $A$ .

**Theorem 1.3.2** (Expansion along the  $i$ -th row).

$$\det A = (-1)^{i+1} a_{i1} \det A_{i1} + \cdots + (-1)^{i+n} a_{in} \det A_{in}.$$

*Proof.* According to Theorem 1.3.1, it suffices to show that the function defined by the right side of the above identity satisfies the three properties stated in the theorem. To check linearity in the  $k$ -th column of  $A$ , let  $\tilde{A}$  denote the matrix obtained by adding  $\mathbf{b} \in \mathbb{R}^n$  to the  $k$ -th column of  $A$ , and  $\bar{A}$  the one obtained by replacing the  $k$ -th column of  $A$  by  $\mathbf{b}$ . We must show that

$$\sum_j (-1)^{i+j} \tilde{a}_{ij} \det \tilde{A}_{ij} = \sum_j (-1)^{i+j} a_{ij} \det A_{ij} + \sum_j (-1)^{i+j} \bar{a}_{ij} \det \bar{A}_{ij}. \quad (1.3.3)$$

Now, if  $j \neq k$ , then  $\tilde{A}_{ij}$  is obtained by adding  $(b_1, \dots, b_{i-1}, b_{i+1}, \dots, b_n)^T$  to a column of  $A_{ij}$ , whereas  $\bar{A}_{ij}$  is obtained by replacing that column by the same vector. Thus,  $\det \tilde{A}_{ij} = \det A_{ij} + \det \bar{A}_{ij}$ ; furthermore,  $a_{ij} = \tilde{a}_{ij} = \bar{a}_{ij}$ , so that

$$\tilde{a}_{ij} \det \tilde{A}_{ij} = a_{ij} \det A_{ij} + \bar{a}_{ij} \det \bar{A}_{ij}, \quad j \neq k. \quad (1.3.4)$$

When  $j = k$  on the other hand,  $A_{ik} = \tilde{A}_{ik} = \bar{A}_{ik}$ , and  $\tilde{a}_{ik} = a_{ik} + b_k = a_{ik} + \bar{a}_{ik}$ , so that

$$\tilde{a}_{ik} \det \tilde{A}_{ik} = a_{ik} \det A_{ik} + \bar{a}_{ik} \det \bar{A}_{ik}. \quad (1.3.5)$$

Identities (1.3.4) and (1.3.5) then imply (1.3.3). A similar argument shows that if  $B$  denotes the matrix obtained by multiplying the  $k$ -th column of  $A$  by  $c$ , then

$$\sum_j (-1)^{i+j} b_{ij} \det B_{ij} = c \sum_j (-1)^{i+j} a_{ij} \det A_{ij}.$$

This shows that the expansion along the  $i$ -th row is multilinear.

To see that the expansion is alternating, consider first the matrix  $B$  obtained by interchanging two adjacent columns from  $A$ , say columns  $k$  and  $k + 1$ . Then for  $j \neq k, k + 1$ ,  $B_{ij}$  is obtained from  $A_{ij}$  by interchanging two columns, so that  $\det B_{ij} = -\det A_{ij}$ , and

$$b_{ij} \det B_{ij} = a_{ij} \det B_{ij} = -a_{ij} \det A_{ij}.$$

For the remaining two cases, observe that  $A_{i(k+1)} = B_{ik}$ ,  $B_{i(k+1)} = A_{ik}$ , and similar identities hold for lower case letters. Thus,

$$b_{ik} \det B_{ik} = a_{i(k+1)} \det A_{i(k+1)}, \quad a_{ik} \det A_{ik} = b_{i(k+1)} \det B_{i(k+1)}.$$

However, the terms have opposite sign in the corresponding sums, since one is in column  $k$ , and the other in column  $k + 1$ . Thus,  $\det B = -\det A$ . Finally, if  $B$  is obtained by interchanging two non-adjacent columns  $i$  and  $j$  from  $A$ , it is also obtained by interchanging adjacent columns an odd number of times: assuming, without loss of generality, that  $i < j$ , interchange column  $i$  with the column to its right  $j - i$  times in succession, then interchange column  $j$  with the column to its left  $j - i - 1$  times. This shows that the sum is alternating.

It remains to check that the right side of the identity equals 1 when  $A$  is the  $n \times n$  identity matrix  $I_n$ . Since  $a_{ij} = 0$  unless  $i = j$ , the only remaining term in the sum is  $(-1)^{2i} a_{ii} \det A_{ii} = \det I_{n-1}$ . Applying this repeatedly yields  $\det I_n = \det I_{n-1} = \dots = \det I_1 = 1$ .  $\square$



**Theorem 1.3.3.** *If  $A$  is a square matrix, then  $\det A = \det(A^T)$ .*

*Proof.* Recall that for an  $n \times n$  matrix  $A$ ,

$$\det A = \sum_{\sigma \in S_n} \varepsilon(\sigma) a_{\sigma(1)1} \cdots a_{\sigma(n)n}.$$

Now, if  $\sigma(i) = j$ , then  $i = \sigma^{-1}(j)$ , and  $a_{\sigma(i)i} = a_{j\sigma^{-1}(j)}$ . Furthermore,  $a_{j\sigma^{-1}(j)}$  occurs exactly once for each integer  $j$ , so that

$$a_{\sigma(1)1} \cdots a_{\sigma(n)n} = a_{1\sigma^{-1}(1)} \cdots a_{n\sigma^{-1}(n)}.$$

Finally,  $\varepsilon(\sigma^{-1}) = \varepsilon(\sigma)$  (since  $\varepsilon(\sigma)\varepsilon(\sigma^{-1}) = \varepsilon(\sigma \circ \sigma^{-1}) = \varepsilon(\text{id}) = 1$ ), and the map that assigns to each  $\sigma \in S_n$  its inverse is a bijection. Thus,

$$\begin{aligned} \det A &= \sum_{\sigma^{-1} \in S_n} \varepsilon(\sigma^{-1}) a_{1\sigma^{-1}(1)} \cdots a_{n\sigma^{-1}(n)} = \sum_{\sigma \in S_n} \varepsilon(\sigma) a_{1\sigma(1)} \cdots a_{n\sigma(n)} \\ &= \det(A^T). \end{aligned} \quad \square$$

**Corollary 1.3.1** (Expansion along  $j$ -th column).

$$\det A = (-1)^{1+j} a_{1j} \det(A_{1j}) + \cdots + (-1)^{n+j} a_{nj} \det(A_{nj}).$$

*Proof.* Let  $B$  denote the transpose of  $A$ . Thus,  $B_{ji}$  is the matrix obtained by deleting row  $j$  and column  $i$  from  $A^T$ ; i.e., it is the transpose of the matrix obtained by deleting row  $i$  and column  $j$  from  $A$ . By Theorem 1.3.3,  $\det B_{ji} = \det A_{ij}$ . Using Theorem 1.3.3 once again, and expanding along the  $j$ -th row, we have

$$\begin{aligned} \det A &= \det B = (-1)^{1+j} b_{j1} \det(B_{j1}) + \cdots + (-1)^{n+j} b_{jn} \det(B_{jn}) \\ &= (-1)^{1+j} a_{1j} \det(A_{1j}) + \cdots + (-1)^{n+j} a_{nj} \det(A_{nj}). \end{aligned} \quad \square$$

For example, to compute the determinant of the three by three matrix

$$A = \begin{bmatrix} 1 & 2 & -1 \\ 0 & 1 & 2 \\ 1 & 4 & 1 \end{bmatrix},$$

we can expand along the second row (which has the only zero) to obtain

$$1 \det \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} - 2 \det \begin{bmatrix} 1 & 2 \\ 1 & 4 \end{bmatrix} = 2 - 2 \cdot 2 = -2.$$

Better yet, the multilinear and alternating properties imply that adding a multiple of a row to another row (or a multiple of a column to another column) does not change the determinant. In the previous example, we may therefore subtract row 1 from row 3, and expand along the first column to get

$$\det A = \det \begin{bmatrix} 1 & 2 & -1 \\ 0 & 1 & 2 \\ 0 & 2 & 2 \end{bmatrix} = 1 \det \begin{bmatrix} 1 & 2 \\ 2 & 2 \end{bmatrix} = -2.$$

Even though the determinant of a sum is not the sum of the determinants, products are another matter:

**Theorem 1.3.4.**  $\det(AB) = \det A \det B$ .

*Proof.* Suppose both matrices are  $n \times n$ . As observed in the proof of Theorem 1.2.1, the  $k$ -th column of  $AB$  equals

$$A\mathbf{b}_k = b_{1k}\mathbf{a}_1 + \cdots + b_{nk}\mathbf{a}_n,$$

where  $\mathbf{b}_k$  denotes the  $k$ -th column of  $B$ , and similarly,  $\mathbf{a}_i$  is the  $i$ -th column of  $A$ . Since the determinant is multilinear and alternating, we obtain, just as in the proof of Theorem 1.3.1,

$$\begin{aligned} \det(AB) &= \det(b_{11}\mathbf{a}_1 + \cdots + b_{n1}\mathbf{a}_n, \dots, b_{1n}\mathbf{a}_1 + \cdots + b_{nn}\mathbf{a}_n) \\ &= \sum_{\sigma} b_{\sigma(1)1} \cdots b_{\sigma(n)n} \det(\mathbf{a}_{\sigma(1)}, \dots, \mathbf{a}_{\sigma(n)}) \\ &= \sum_{\sigma} \varepsilon(\sigma) b_{\sigma(1)1} \cdots b_{\sigma(n)n} \det(\mathbf{a}_1, \dots, \mathbf{a}_n) \\ &= \det B \det A. \end{aligned} \quad \square$$

As a special case, if  $A$  is an invertible matrix, then  $\det(A^{-1}) = 1/\det A$ , since  $AA^{-1} = I_n$  and the identity matrix has determinant 1. Another important consequence is that similar matrices (see Example 1.2.5) have the same determinant. In particular, the determinant of the matrix  $[L]_{\mathcal{B},\mathcal{B}}$  of a linear transformation  $L : V \rightarrow V$  in a basis  $\mathcal{B}$  equals that of the matrix of  $L$  in any other basis. This justifies the following:

**Definition 1.3.2.** The *determinant* of a linear transformation  $L : V \rightarrow V$  is the determinant of the matrix  $[L]_{\mathcal{B},\mathcal{B}}$  of  $L$  with respect to any given basis  $\mathcal{B}$  of  $V$ .

Determinants are often useful for, well, determining whether a linear map is an isomorphism:

**Theorem 1.3.5.** *Let  $V$  be a finite-dimensional vector space. A linear transformation  $L : V \rightarrow V$  is an isomorphism if and only if  $\det L \neq 0$ .*

*Proof.* If  $L$  is an isomorphism, then it has an inverse  $L^{-1}$ , and the product  $\det L \det(L^{-1}) = \det(L \circ L^{-1}) = \det(1_V) = 1$ , so  $L$  cannot have vanishing determinant. If, on the other hand,  $L$  is not an isomorphism, then by Theorem 1.2.2, there exists a nonzero  $\mathbf{v} \in V$  such that  $L\mathbf{v} = \mathbf{0}$ . Let  $\mathcal{B} = \mathbf{v}_1, \dots, \mathbf{v}_n$  denote an ordered basis of  $V$ , and write  $\mathbf{v}$  as a linear combination  $\sum_i c_i \mathbf{v}_i$  of the basis elements. Notice that not all coefficients vanish. Then

$$\mathbf{0} = [L\mathbf{v}]_{\mathcal{B}} = [L(\sum_i c_i \mathbf{v}_i)]_{\mathcal{B}} = \sum_i c_i [L\mathbf{v}_i]_{\mathcal{B}}.$$

This says that the columns  $[L\mathbf{v}_i]_{\mathcal{B}}$  of the matrix  $[L]_{\mathcal{B},\mathcal{B}}$  of  $L$  in the basis  $\mathcal{B}$  are linearly dependent. But any square matrix with linearly dependent columns (or rows) has vanishing determinant: Indeed, suppose  $A$  has linearly dependent columns  $\mathbf{a}_i$  (for rows,

consider the transpose of  $A$  instead). Then, one of the columns, say,  $\mathbf{a}_j$ , can be written as a linear combination

$$\mathbf{a}_j = \sum_{i \neq j} \alpha_i \mathbf{a}_i$$

of the others. By linearity of  $\det$  in each column,

$$\det A = \det(\mathbf{a}_1, \dots, \mathbf{a}_j, \dots, \mathbf{a}_n) = \sum_{i \neq j} \alpha_i \det(\mathbf{a}_1, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n).$$

Each determinant in this sum is the determinant of a matrix with two equal columns (the  $i$ -th and the  $j$ -th). Such a determinant must vanish by the alternating property (interchanging columns  $i$  and  $j$  changes the sign of the determinant, but the matrix remains the same). This establishes the result.  $\square$

The proof simplifies considerably if one uses the fact that given any nonzero  $\mathbf{v} \in V$ , there is a basis of  $V$  containing  $\mathbf{v}$  (this will be proved in the next section). In fact, if  $L$  is not an isomorphism, consider any nonzero vector in the kernel of  $L$ , and extend it to a basis. The matrix of  $L$  in that basis then has a zero column and therefore also zero determinant.

**Theorem 1.3.6.** *Let  $V$  be a finite-dimensional vector space. For any linear transformation  $L : V \rightarrow V$ , there exists a linear map  $\tilde{L} : V \rightarrow V$  such that*

$$L \circ \tilde{L} = (\det L)1_V.$$

*In particular, if  $L$  is invertible, then  $L^{-1} = (1/\det L)\tilde{L}$ , and if not, then  $L \circ \tilde{L} = 0$ .*

*Proof.* Consider any ordered basis  $\mathcal{B}$  of  $V$  and denote by  $A = (a_{ij})$  the matrix of  $L$  with respect to  $\mathcal{B}$ . For each  $i$  and  $j$  between 1 and  $n = \dim V$ , let  $A_{ij}$  be the  $(n-1) \times (n-1)$  matrix obtained by deleting row  $i$  and column  $j$  from  $A$ , and define an  $n \times n$  matrix  $\tilde{A} = (\tilde{a}_{ij})$  by setting

$$\tilde{a}_{ij} = (-1)^{i+j} \det A_{ji}.$$

A straightforward computation yields

$$A\tilde{A} = (\det A)I_n. \tag{1.3.6}$$

The result now follows if we let  $\tilde{L}$  be the linear map whose matrix in the basis  $\mathcal{B}$  equals  $\tilde{A}$ .  $\square$

The map  $\tilde{L}$  constructed in the above proof will be called the *linear map adjugate to  $L$* . It can be constructed without using a basis, but this requires extra work.

## 1.4 Euclidean spaces

An elementary concept that is used throughout Calculus is the distance  $|a - b|$  between two points  $a$  and  $b$  in the real line. We now generalize this concept to  $\mathbb{R}^n$ . The

(standard) *inner product* of  $\mathbf{a}$  and  $\mathbf{b} \in \mathbb{R}^n$  is the number

$$\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^n u^i(\mathbf{a})u^i(\mathbf{b}).$$

This is nothing but the entry of the  $1 \times 1$  matrix  $\mathbf{a}^T \cdot \mathbf{b}$ . The vector space  $\mathbb{R}^n$  together with this inner product is called *n-dimensional Euclidean space*. Some authors use this terminology for a larger class of spaces, which we call inner product spaces, to be introduced shortly. The *norm* of  $\mathbf{a}$  is

$$|\mathbf{a}| = \langle \mathbf{a}, \mathbf{a} \rangle^{1/2}.$$

When  $n = 1$ , the norm of  $\mathbf{a}$  is just its absolute value.

**Theorem 1.4.1.** *If  $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^n$ , and  $a \in \mathbb{R}$ , then*

- (1)  $\langle \mathbf{a}, \mathbf{b} \rangle = \langle \mathbf{b}, \mathbf{a} \rangle$ ;
- (2)  $\langle a\mathbf{a} + \mathbf{b}, \mathbf{c} \rangle = a\langle \mathbf{a}, \mathbf{b} \rangle + \langle \mathbf{b}, \mathbf{c} \rangle$ ;
- (3)  $|\mathbf{a}| \geq 0$ , and  $|\mathbf{a}| = 0$  if and only if  $\mathbf{a} = \mathbf{0}$ ;
- (4)  $|a\mathbf{a}| = |a||\mathbf{a}|$ ;
- (5) (Cauchy-Schwarz inequality)  $|\langle \mathbf{a}, \mathbf{b} \rangle| \leq |\mathbf{a}||\mathbf{b}|$ ;
- (6) (Triangle inequality)  $|\mathbf{a} + \mathbf{b}| \leq |\mathbf{a}| + |\mathbf{b}|$ .

*Proof.* The first four statements are obvious. For the Cauchy-Schwarz inequality, let  $a_i = u^i(\mathbf{a})$ ,  $b_i = u^i(\mathbf{b})$ . We may assume that  $\mathbf{b} \neq \mathbf{0}$ , since otherwise the conclusion is trivial. Now,

$$\begin{aligned} 0 &\leq \sum_i (|\mathbf{b}|^2 a_i - \langle \mathbf{a}, \mathbf{b} \rangle b_i)^2 \\ &= |\mathbf{b}|^4 \sum a_i^2 + \langle \mathbf{a}, \mathbf{b} \rangle^2 \sum b_i^2 - 2|\mathbf{b}|^2 \langle \mathbf{a}, \mathbf{b} \rangle \sum a_i b_i \\ &= |\mathbf{b}|^4 |\mathbf{a}|^2 + |\mathbf{b}|^2 \langle \mathbf{a}, \mathbf{b} \rangle^2 - 2|\mathbf{b}|^2 \langle \mathbf{a}, \mathbf{b} \rangle^2 \\ &= |\mathbf{b}|^2 (|\mathbf{a}|^2 |\mathbf{b}|^2 - \langle \mathbf{a}, \mathbf{b} \rangle^2). \end{aligned}$$

Since  $|\mathbf{b}|^2 > 0$ ,  $|\mathbf{a}|^2 |\mathbf{b}|^2 - \langle \mathbf{a}, \mathbf{b} \rangle^2 \geq 0$ , which is the desired inequality.

The triangle inequality is a direct consequence of the Cauchy-Schwarz inequality:

$$\begin{aligned} |\mathbf{a} + \mathbf{b}|^2 &= \langle \mathbf{a} + \mathbf{b}, \mathbf{a} + \mathbf{b} \rangle = |\mathbf{a}|^2 + |\mathbf{b}|^2 + 2\langle \mathbf{a}, \mathbf{b} \rangle \leq |\mathbf{a}|^2 + |\mathbf{b}|^2 + 2|\mathbf{a}||\mathbf{b}| \\ &= (|\mathbf{a}| + |\mathbf{b}|)^2. \end{aligned} \quad \square$$

**Corollary 1.4.1.**  $||\mathbf{a}| - |\mathbf{b}|| \leq |\mathbf{a} - \mathbf{b}|$  for  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ .

*Proof.* It is enough to show that  $|\mathbf{a}| - |\mathbf{b}| \leq |\mathbf{a} - \mathbf{b}|$  (why?). But this follows from the triangle inequality, since

$$|\mathbf{a}| = |\mathbf{a} - \mathbf{b} + \mathbf{b}| \leq |\mathbf{a} - \mathbf{b}| + |\mathbf{b}|. \quad \square$$

**Definition 1.4.1.** The *operator norm* of a linear transformation  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is

$$|L| = \sup\{|L\mathbf{u}| \mid |\mathbf{u}| = 1\},$$

where sup denotes the supremum or least upper bound of a set as introduced in Section 1.7.

Notice that the norm of  $L$  is always a finite number: If  $\mathbf{u} = \sum_i a_i \mathbf{e}_i$  has norm 1, then each  $|a_i| \leq 1$ , so that

$$|L\mathbf{u}| = \left| \sum_i a_i L\mathbf{e}_i \right| \leq \sum_i |a_i| |L\mathbf{e}_i| \leq \sum_i |L\mathbf{e}_i|,$$

and consequently  $|L| \leq \sum_i |L\mathbf{e}_i|$ .

It also follows from the definition that for all  $\mathbf{u} \in \mathbb{R}^n$ ,

$$|L\mathbf{u}| \leq |L||\mathbf{u}|. \quad (1.4.1)$$

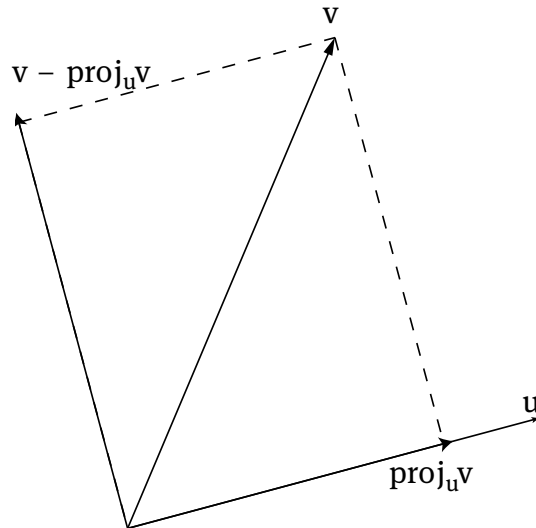
More generally, an *inner product* on a vector space  $V$  is a map  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$  that satisfies (1), (2), and (3) in Theorem 1.4.1. In this case,  $(V, \langle \cdot, \cdot \rangle)$  is called an *inner product space*.  $\mathbf{v}, \mathbf{w} \in V$  are said to be *orthogonal* if  $\langle \mathbf{v}, \mathbf{w} \rangle = 0$ . Given a nonzero vector  $\mathbf{u} \in V$ , the *projection of  $\mathbf{v} \in V$  along  $\mathbf{u}$*  is the vector

$$\text{proj}_{\mathbf{u}} \mathbf{v} = \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\langle \mathbf{u}, \mathbf{u} \rangle} \mathbf{u}.$$

Any  $\mathbf{v} \in V$  decomposes uniquely as the sum

$$\mathbf{v} = \text{proj}_{\mathbf{u}} \mathbf{v} + (\mathbf{v} - \text{proj}_{\mathbf{u}} \mathbf{v})$$

of a vector parallel to  $\mathbf{u}$  (namely its projection along  $\mathbf{u}$ ) and a vector orthogonal to  $\mathbf{u}$ .



A basis  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  of an inner product space is said to be *orthonormal* if  $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = \delta_{ij}$  (the *Kronecker delta*  $\delta_{ij}$  is the symbol that equals 1 if  $i = j$  and 0 otherwise); in

other words, a basis is orthonormal if it consists of unit vectors any two of which are orthogonal.

The following theorem provides an algorithm for obtaining an orthogonal basis from an arbitrary one. To get an orthonormal one, it only remains to divide each basis element by its length.

**Theorem 1.4.2** (Gram-Schmidt orthogonalization). *Suppose  $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  denotes a basis of an inner product space  $V$ . Then  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ , where*

$$\begin{aligned}\mathbf{v}_1 &= \mathbf{u}_1, \\ \mathbf{v}_2 &= \mathbf{u}_2 - \text{proj}_{\mathbf{v}_1} \mathbf{u}_2, \\ \mathbf{v}_3 &= \mathbf{u}_3 - \text{proj}_{\mathbf{v}_1} \mathbf{u}_3 - \text{proj}_{\mathbf{v}_2} \mathbf{u}_3, \\ &\vdots \\ \mathbf{v}_n &= \mathbf{u}_n - \sum_{i=1}^{n-1} \text{proj}_{\mathbf{v}_i} \mathbf{u}_n,\end{aligned}$$

is an orthogonal basis of  $V$ .

*Proof.* By construction, each  $\mathbf{v}_i$  is orthogonal to  $\mathbf{v}_1, \dots, \mathbf{v}_{i-1}$ , so that any two of them are orthogonal. Moreover,

$$\mathbf{u}_i = \mathbf{v}_i + \sum_{l=1}^{i-1} \text{proj}_{\mathbf{v}_l} \mathbf{u}_i \in \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_i\},$$

so that  $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  spans  $V$ . But then  $\mathcal{B}$  is linearly independent, for otherwise a strict subset of  $\mathcal{B}$  would span  $V$ , contradicting the fact that  $\dim V = n$ . Thus,  $\mathcal{B}$  is an orthogonal basis of  $V$ .  $\square$

**Example 1.4.1.** Let  $\mathbf{u}_1 = [1 \ 0 \ 1]^T$ ,  $\mathbf{u}_2 = [2 \ 1 \ 0]^T$ , and  $V$  the subspace of  $\mathbb{R}^3$  spanned by  $\mathbf{u}_1$  and  $\mathbf{u}_2$ . Then  $\text{proj}_{\mathbf{u}_1} \mathbf{u}_2 = (\langle \mathbf{u}_1, \mathbf{u}_2 \rangle / |\mathbf{u}_1|^2) \mathbf{u}_1 = \mathbf{u}_1$ , and  $\mathbf{v}_2 = \mathbf{u}_2 - \mathbf{u}_1 = [1 \ 1 \ -1]^T$  is a vector in  $V$  orthogonal to  $\mathbf{u}_1$ . Thus, the vectors

$$\mathbf{w}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \quad \mathbf{w}_2 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}$$

form an orthonormal basis of  $V$ . If we wish to extend this basis to an orthonormal basis of  $\mathbb{R}^3$ , we need only find a vector that does not belong to  $V$ , say,  $\mathbf{e}_3$ , and apply the Gram-Schmidt process to it:

$$\mathbf{v}_3 = \mathbf{e}_3 - \text{proj}_{\mathbf{w}_1} \mathbf{e}_3 - \text{proj}_{\mathbf{w}_2} \mathbf{e}_3 = \mathbf{e}_3 - \mathbf{w}_1 + \mathbf{w}_2 = \begin{bmatrix} -1/6 \\ 1/3 \\ 1/6 \end{bmatrix}.$$

$\{\mathbf{w}_1, \mathbf{w}_2, \mathbf{v}_3/|\mathbf{v}_3|\}$  is then an orthonormal basis of  $\mathbb{R}^3$ .

One reason orthonormal bases are useful is that it is quite easy to compute coordinate vectors in them:

**Theorem 1.4.3.** *If  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is an orthonormal basis of  $V$ , then for any  $\mathbf{v} \in V$ ,*

$$\mathbf{v} = \langle \mathbf{v}, \mathbf{v}_1 \rangle \mathbf{v}_1 + \dots + \langle \mathbf{v}, \mathbf{v}_n \rangle \mathbf{v}_n.$$

*Thus, the coordinate vector of  $\mathbf{v}$  in the ordered basis  $(\mathbf{v}_1, \dots, \mathbf{v}_n)$  has  $\langle \mathbf{v}, \mathbf{v}_i \rangle$  as  $i$ -th entry.*

*Proof.* By assumption  $\mathbf{v} = \sum a_i \mathbf{v}_i$  for unique scalars  $a_1, \dots, a_n$ . Taking the inner product on both sides with  $\mathbf{v}_j$  yields

$$\langle \mathbf{v}, \mathbf{v}_j \rangle = \sum_{i=1}^n a_i \langle \mathbf{v}_i, \mathbf{v}_j \rangle = \sum_{i=1}^n a_i \delta_{ij} = a_j. \quad \square$$

It is now easy to see that the Cauchy-Schwarz inequality is valid in any inner product space: just replace  $a_i$  in the proof by  $\langle \mathbf{a}, \mathbf{v}_i \rangle$ , where  $\mathbf{v}_1, \dots, \mathbf{v}_n$  is an orthonormal basis, and similarly for  $b_i$ .

**Definition 1.4.2.** Let  $A$  be a nonempty subset of an inner product space  $V$ . The *orthogonal complement* of  $A$  is the set

$$A^\perp = \{\mathbf{v} \in V \mid \langle \mathbf{v}, \mathbf{a} \rangle = 0 \text{ for every } \mathbf{a} \in A\}.$$

Linearity of the inner product in each entry implies that  $A^\perp$  is a subspace of  $V$ . For the same reason,  $(\text{span } A)^\perp = A^\perp$ .

In general, if  $W_1, W_2$  are subspaces of a vector space  $V$ , then the set  $W_1 + W_2 = \{\mathbf{v}_1 + \mathbf{v}_2 \mid \mathbf{v}_i \in W_i\}$  is easily seen to be a subspace of  $V$ , called the *sum of  $W_1$  and  $W_2$* . If  $W_1 \cap W_2 = \{0\}$ , this sum is called a *direct sum*, and is denoted  $W_1 \oplus W_2$ . In this case, any  $\mathbf{v} \in W_1 \oplus W_2$  can be written in one and only one way as a sum of a vector in  $W_1$  and a vector in  $W_2$ : indeed if  $\mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}'_1 + \mathbf{v}'_2$  with  $\mathbf{v}_i, \mathbf{v}'_i \in W_i$ , then  $\mathbf{v}_1 - \mathbf{v}'_1 = \mathbf{v}_2 - \mathbf{v}'_2$ . Since the left side is in  $W_1$  and the right side in  $W_2$ , both terms vanish, and  $\mathbf{v}_i = \mathbf{v}'_i$  as claimed.

**Proposition 1.4.1.** *If  $A$  is a nonempty subset of an inner product space  $V$ , then  $V = (\text{span } A) \oplus A^\perp$ .*

*Proof.* If  $\mathbf{v} \in (\text{span } A) \cap A^\perp$ , then  $\mathbf{v}$  is orthogonal to itself, and must then equal  $\mathbf{0}$ . It therefore remains to show that  $V \subset (\text{span } A) + A^\perp$ . Given  $\mathbf{v} \in V$ , and an orthonormal basis  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  of  $\text{span } A$ , write

$$\mathbf{v} = \mathbf{u}_1 + \mathbf{u}_2, \quad \mathbf{u}_1 = \sum_{i=1}^k \text{proj}_{\mathbf{v}_i} \mathbf{v}, \quad \mathbf{u}_2 = \mathbf{v} - \mathbf{u}_1. \quad (1.4.2)$$

By definition  $\mathbf{u}_1 \in \text{span } A$ . For any  $j = 1, \dots, k$ ,

$$\langle \mathbf{u}_2, \mathbf{v}_j \rangle = \langle \mathbf{v} - \sum_{i=1}^k \text{proj}_{\mathbf{v}_i} \mathbf{v}, \mathbf{v}_j \rangle = \langle \mathbf{v}, \mathbf{v}_j \rangle - \langle \text{proj}_{\mathbf{v}_j} \mathbf{v}, \mathbf{v}_j \rangle = 0,$$

so that  $\mathbf{u}_2$  is orthogonal to every  $\mathbf{v}_j$  and therefore to every element in  $\text{span } A$ ; i.e.,  $\mathbf{u}_2 \in A^\perp$ .  $\square$

The (necessarily unique) vector  $\mathbf{u}_1$  in (1.4.2) is called the *orthogonal projection of  $\mathbf{v}$  onto  $\text{span } A$* . When  $A$  consists of a single vector  $\mathbf{u}$ , it coincides with what we defined earlier as the projection of  $\mathbf{v}$  along  $\mathbf{u}$ .

Incidentally, we have in essence established the following:

**Proposition 1.4.2.** *If  $A$  is a linearly independent subset of  $\mathbb{R}^n$ , then there exists a basis of  $\mathbb{R}^n$  that contains  $A$ .*

*Proof.* Let  $\mathcal{B}$  denote a basis of  $A^\perp$ . Then  $A \cup \mathcal{B}$  is a basis as in the statement.  $\square$

It should be noted that this is true in any vector space, since such a space is isomorphic to some Euclidean space. In fact, any two spaces with the same dimension are isomorphic, as observed in Examples 1.2.1 (iii). Such an isomorphism depends of course on the choice of bases. More to the point, the above proposition can be proved without introducing an inner product. The argument given here is, however, sufficient for our purposes.

The following concept too does not require the existence of an inner product:

**Definition 1.4.3.** The *dual space* of a vector space  $V$  is the collection  $V^*$  of all linear maps  $V \rightarrow \mathbb{R}$ .

An element of  $V^*$  is also called a *one-form on  $V$* . These maps can be added in the usual way, as well as multiplied by scalars, thereby inducing a vector space structure on  $V^*$ . If  $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  denotes a basis of  $V$ , then the map

$$\begin{aligned} V^* &\rightarrow M_{n,1}, \\ \alpha &\mapsto [\alpha(\mathbf{v}_1) \quad \cdots \quad \alpha(\mathbf{v}_n)] \end{aligned}$$

is by definition linear, surjective, and has zero kernel. In particular,  $V^*$  has the same dimension as  $V$  and is therefore isomorphic to  $V$ . Although, as noted earlier, isomorphisms between spaces of the same dimension depend in general on the choice of bases, in this case there is a canonical one, provided  $V$  is an inner product space:

**Theorem 1.4.4.** *Let  $V$  be a finite-dimensional inner product space. The map  $\flat : V \rightarrow V^*$ , given by  $\mathbf{u}^\flat(\mathbf{v}) = \langle \mathbf{u}, \mathbf{v} \rangle$ , is an isomorphism.*

*Proof.* The map is clearly linear, and since the spaces have the same dimension, it suffices to check its kernel is zero. But if  $\mathbf{u} \in \ker \flat$ , then  $\langle \mathbf{u}, \mathbf{v} \rangle = 0$  for all  $\mathbf{v} \in V$ , and in particular,  $\|\mathbf{u}\| = \langle \mathbf{u}, \mathbf{u} \rangle^{1/2} = 0$ , so that  $\mathbf{u} = \mathbf{0}$ .  $\square$

The expression  $\mathbf{u}^\flat$  reads “ $\mathbf{u}$  flat”.  $\flat$  and its inverse  $\sharp$  (the existence of which is guaranteed by Theorem 1.4.4) are called the *musical isomorphisms* associated to the inner product. An important consequence of Theorem 1.4.4 is the following:



**Corollary 1.4.2.** *If  $V$  is a finite-dimensional inner product space, then for any  $\alpha \in V^*$ , there exists a unique  $\mathbf{u} \in V$  such that*

$$\langle \mathbf{u}, \mathbf{v} \rangle = \alpha(\mathbf{v}), \quad \mathbf{v} \in V.$$

*Proof.*  $\mathbf{u} = \alpha^\sharp$ . □

**Examples 1.4.2.** (i) If  $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a basis of  $V$ , then there exists, for each  $j = 1, \dots, n$ , a unique element  $\alpha^j \in V^*$  satisfying  $\alpha^j(\sum_{i=1}^n a_i \mathbf{v}_i) = a_j$ . In fact, the matrix of  $\alpha^j$  in the basis  $\mathcal{B}$  is  $\mathbf{e}_j^T$ , which also shows that the  $\alpha^j$  are linearly independent. By dimension considerations, they form a basis of  $V^*$ , called the *basis dual* to  $\mathcal{B}$ . (It is also easy to see directly that they span  $V^*$ , since any  $\alpha \in V^*$  may be written as  $\alpha = \sum_{i=1}^n \alpha(\mathbf{v}_i) \alpha^i$ .) For example, if  $V = \mathbb{R}^n$ , then the basis of  $V^*$  dual to the standard basis is  $u^1, \dots, u^n$ , because  $u^j(\sum_i a_i \mathbf{e}_i) = a_j$ . Superscripts are traditionally used with dual elements.

(ii) A simple yet illustrative example of a musical isomorphism is the one corresponding to the standard inner product of  $\mathbb{R}^n$ . If  $\alpha$  is a one-form, and  $[\alpha]$  denotes its  $n \times 1$  matrix in the standard basis, then, recalling that any vector is its own coordinate vector in that basis, we have

$$\alpha(\mathbf{v}) = [\alpha]\mathbf{v} = \langle [\alpha]^T, \mathbf{v} \rangle, \quad \mathbf{v} \in \mathbb{R}^n.$$

Thus,  $\alpha^\sharp$  is the transpose of the matrix of  $\alpha$  in the standard basis. For example, if  $\alpha : \mathbb{R}^3 \rightarrow \mathbb{R}$  is given by  $\alpha(x, y, z) = 2x - 3y + z$ , then

$$\alpha^\sharp = \begin{bmatrix} 2 \\ -3 \\ 1 \end{bmatrix}.$$

It should be noted that the musical isomorphism  $\flat$  maps the standard basis of  $\mathbb{R}^n$  to its dual basis in the sense of (i). This is because the standard basis is orthonormal, and the above claim holds in any inner product space: If  $\mathbf{v}_1, \dots, \mathbf{v}_n$  is an orthonormal basis of an inner product space  $V$ , then

$$\mathbf{v}_j^\flat \left( \sum_{i=1}^n a_i \mathbf{v}_i \right) = \langle \mathbf{v}_j, \sum_{i=1}^n a_i \mathbf{v}_i \rangle = a_j.$$

## 1.5 Subspaces of Euclidean space

We wish to take a closer look at subspaces of dimension 1 or  $n - 1$  of  $\mathbb{R}^n$ , since these represent all proper subspaces in the important special case when  $n = 3$ .

If  $V$  is a subspace of dimension  $n - 1$ , then its orthogonal complement  $V^\perp$  is one-dimensional, and an element of  $V^\perp$  is called a *normal vector* of  $V$ . If  $\mathbf{n} \neq \mathbf{0}$  is one such, then any other is a multiple of it, and  $V$  consists of the set of all  $\mathbf{v} \in \mathbb{R}^n$  such

that  $\langle \mathbf{v}, \mathbf{n} \rangle = 0$ . Equivalently,  $V$  is the kernel of left multiplication  $L_{\mathbf{n}^T}$  by  $\mathbf{n}^T$ . If  $\mathbf{n}^T = [a_1 \ \dots \ a_n]$ , then  $V$  consists of all  $\mathbf{v} = [x_1 \ \dots \ x_n]^T$  such that  $a_1x_1 + \dots + a_nx_n = 0$ .  $V$  is called a *hyperplane* (or plane when  $n = 3$ ) through the origin.

An *affine hyperplane* of  $\mathbb{R}^n$  is a set of the form  $\mathbf{u} + V$ , where  $V$  is a hyperplane and  $\mathbf{u} \in \mathbb{R}^n$ . In other words, it is the subspace  $V$  “parallel translated” by  $\mathbf{u}$ . If  $\mathbf{u} = [u_1 \ \dots \ u_n]^T$  and  $V$  has equation  $a_1x_1 + \dots + a_nx_n = 0$  as above, then the equation of the affine hyperplane is

$$\{(x_1, \dots, x_n) \mid a_1x_1 + \dots + a_nx_n = b\}, \quad \text{where } b = a_1u_1 + \dots + a_nu_n.$$

Indeed, the affine hyperplane consists of all points  $\mathbf{u} + \mathbf{v}$ , where  $\langle \mathbf{v}, \mathbf{n} \rangle = 0$ . Letting  $\mathbf{x} = \mathbf{u} + \mathbf{v}$ , this is equivalent to the set of all  $\mathbf{x}$  such that  $\langle \mathbf{x} - \mathbf{u}, \mathbf{n} \rangle = 0$ ; setting  $x_i = u^i(\mathbf{x})$ , this corresponds to all  $[x_1 \ \dots \ x_n]^T$  such that  $\sum a_i(x_i - u_i) = 0$ , as claimed.

Even though we are primarily interested in the cases when  $\dim V = 1$  or  $n - 1$ , a similar approach can be adopted to describe a subspace  $V$  of dimension  $n - k$  for any  $k$ : If  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  is a basis of  $V^\perp$ , then  $\mathbf{a} \in \mathbb{R}^n$  belongs to  $V$  if and only if  $\langle \mathbf{a}, \mathbf{u}_i \rangle = 0$  for  $i = 1, \dots, k$ . Each of these equations represents a hyperplane, and  $V$  is then the intersection of these  $k$  hyperplanes. As above, we obtain an affine subspace by translating  $V$  away from the origin.

A one-dimensional subspace  $V$  of  $\mathbb{R}^n$  is called a *line* through the origin. We could describe it as an intersection of  $n - 1$  hyperplanes, but when  $n$  is large, it is usually more convenient to resort to a vector  $\mathbf{a}$  that spans  $V$ .  $\mathbf{v} \in \mathbb{R}^n$  belongs to  $V$  if and only if  $\mathbf{v} = t\mathbf{a}$  for some  $t \in \mathbb{R}$ . Letting  $x_i = u^i(\mathbf{v})$ ,  $a_i = u^i(\mathbf{a})$ , we obtain so-called *parametric equations* of the line

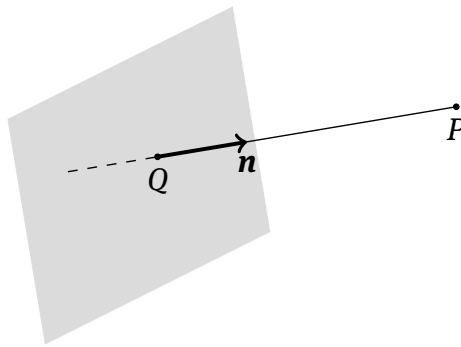
$$\{(x_1, \dots, x_n) \mid x_1 = ta_1, \dots, x_n = ta_n, t \in \mathbb{R}\}. \quad (1.5.1)$$

When all  $a_i \neq 0$ , one can eliminate the parameter  $t$  in (1.5.1) and describe the line by the equation  $x_1/a_1 = \dots = x_n/a_n$ .

As in the case of hyperplanes, one can parallel translate one-dimensional subspaces to describe all possible lines. Thus, the line parallel to the one in (1.5.1) that passes through  $\mathbf{p} = (p_1, \dots, p_n)$  has parametric equations

$$x_1 = p_1 + ta_1, \dots, x_n = p_n + ta_n, \quad t \in \mathbb{R},$$

which can be written  $(x_1 - p_1)/a_1 = \dots = (x_n - p_n)/a_n$  if every  $a_i \neq 0$ .



**Example 1.5.1.** Suppose we are asked to find the distance from the point  $P = (1, 2, -1)$  to the plane  $x + y - z = 1$  in  $\mathbb{R}^3$ . This distance equals the length  $|PQ|$  of the line segment  $PQ$ , where  $Q$  is the point of intersection of the plane with the line passing through  $P$  perpendicular to the plane. Since a vector normal to the plane is  $\mathbf{n} = [1 \ 1 \ -1]^T$ , the line in question has vector equation

$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} + t \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}.$$

Thus, a point on the line has coordinates  $(1 + t, 2 + t, -1 - t)$ , and the intersection with the plane occurs for the value of  $t$  for which

$$x + y - z = 1 + t + 2 + t + 1 + t = 1,$$

or  $t = -1$ . This means that  $Q = (0, 1, 0)$ . Elementary geometry implies that the distance between vectors  $\mathbf{p}$  and  $\mathbf{q}$  in  $\mathbb{R}^3$  is  $|\mathbf{q} - \mathbf{p}|$ , so that

$$|PQ| = \left( (0 - 1)^2 + (1 - 2)^2 + (0 - (-1))^2 \right)^{1/2} = \sqrt{3}.$$

We will later find another solution to this problem, one that uses Calculus.

## 1.6 Determinants as volume

Any two vectors  $\mathbf{a}$  and  $\mathbf{b}$  in  $\mathbb{R}^3$  determine a linear map

$$\begin{aligned} \mathbb{R}^3 &\rightarrow \mathbb{R}, \\ \mathbf{u} &\mapsto \det [\mathbf{a} \ \mathbf{b} \ \mathbf{u}], \end{aligned}$$

and thus an element of the dual space of  $\mathbb{R}^3$ . By Corollary 1.4.2, there is a unique vector  $\mathbf{a} \times \mathbf{b} \in \mathbb{R}^3$ , called the *cross product* of  $\mathbf{a}$  and  $\mathbf{b}$ , such that

$$\langle \mathbf{a} \times \mathbf{b}, \mathbf{u} \rangle = \det [\mathbf{a} \ \mathbf{b} \ \mathbf{u}], \quad \mathbf{u} \in \mathbb{R}^3.$$

The following theorem is an immediate consequence of properties of the determinant:

**Theorem 1.6.1.** Given  $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^3, t \in \mathbb{R}$ ,

- $\mathbf{b} \times \mathbf{a} = -\mathbf{a} \times \mathbf{b}$ ;
- $\mathbf{a} \times \mathbf{b}$  is orthogonal to  $\mathbf{a}$  and  $\mathbf{b}$ ;
- $(t\mathbf{a}) \times \mathbf{b} = \mathbf{a} \times (t\mathbf{b}) = t(\mathbf{a} \times \mathbf{b})$ ;
- $\mathbf{a} \times (\mathbf{b} + \mathbf{c}) = \mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}$ .

It is straightforward to compute the components of the vector  $\mathbf{a} \times \mathbf{b}$ : for example, if  $a_i = u^i(\mathbf{a})$  and  $b_i = u^i(\mathbf{b})$ , then the first component equals

$$\begin{aligned} \langle \mathbf{a} \times \mathbf{b}, \mathbf{e}_1 \rangle &= \det [\mathbf{a} \ \mathbf{b} \ \mathbf{e}_1] = \det \begin{bmatrix} \mathbf{a}^T \\ \mathbf{b}^T \\ \mathbf{e}_1^T \end{bmatrix} = \det \begin{bmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \\ 1 & 0 & 0 \end{bmatrix} \\ &= \det \begin{bmatrix} a_2 & a_3 \\ b_2 & b_3 \end{bmatrix}, \end{aligned}$$

and similar formulas hold for the other components, replacing  $\mathbf{e}_1$  by  $\mathbf{e}_2$  and  $\mathbf{e}_3$ . A common way of writing these formulas is

$$\mathbf{a} \times \mathbf{b} = \det \begin{bmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{bmatrix},$$

where the meaningless right side is meant to be “expanded” along the first row; i.e.,

$$\mathbf{a} \times \mathbf{b} = \det \begin{bmatrix} a_2 & a_3 \\ b_2 & b_3 \end{bmatrix} \mathbf{e}_1 - \det \begin{bmatrix} a_1 & a_3 \\ b_1 & b_3 \end{bmatrix} \mathbf{e}_2 + \det \begin{bmatrix} a_1 & a_2 \\ b_1 & b_2 \end{bmatrix} \mathbf{e}_3. \quad (1.6.1)$$

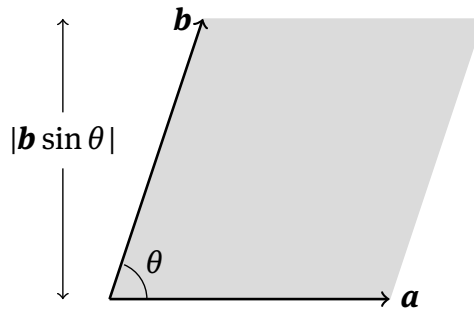
Using the latter formula, it is straightforward to compute that the norm squared of a cross product equals

$$\begin{aligned} |\mathbf{a} \times \mathbf{b}|^2 &= (a_1^2 + a_2^2 + a_3^2)(b_1^2 + b_2^2 + b_3^2) - (a_1b_1 + a_2b_2 + a_3b_3)^2 \\ &= |\mathbf{a}|^2 |\mathbf{b}|^2 - \langle \mathbf{a}, \mathbf{b} \rangle^2 = |\mathbf{a}|^2 |\mathbf{b}|^2 - |\mathbf{a}|^2 |\mathbf{b}|^2 \cos^2 \theta \\ &= |\mathbf{a}|^2 |\mathbf{b}|^2 \sin^2 \theta, \end{aligned}$$

with  $\theta$  denoting the angle between  $\mathbf{a}$  and  $\mathbf{b}$ . Thus,

$$|\mathbf{a} \times \mathbf{b}| = |\mathbf{a}| |\mathbf{b}| |\sin \theta|. \quad (1.6.2)$$

This means that  $|\mathbf{a} \times \mathbf{b}|$  represents the area of the parallelogram spanned by  $\mathbf{a}$  and  $\mathbf{b}$ : notice that the area equals that of a rectangle with same base length  $|\mathbf{a}|$  and height  $|\mathbf{b}| \sin \theta$ .



If, on the other hand, we take the two vectors to lie in the  $xy$ -plane, so that their third component is zero, then by (1.6.1),

$$\mathbf{a} \times \mathbf{b} = \det \begin{bmatrix} a_1 & a_2 \\ b_1 & b_2 \end{bmatrix} \mathbf{e}_3.$$

Combining this with the previous remark, we conclude: *The absolute value of a  $2 \times 2$  determinant equals the area of the parallelogram spanned by the columns (or the rows).*

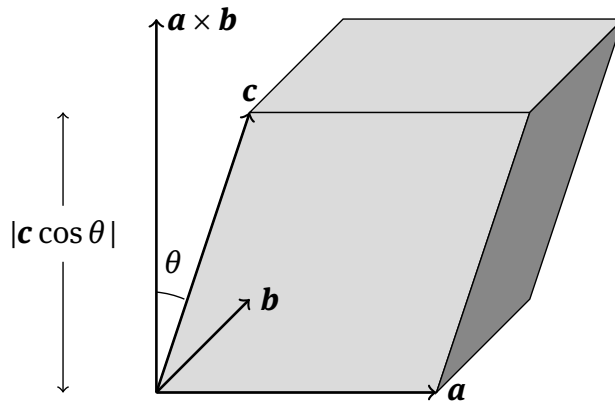
A similar interpretation exists for  $3 \times 3$  determinants. If  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c} \in \mathbb{R}^3$ , then the volume of the parallelepiped

$$\{\mathbf{a}\mathbf{a} + \mathbf{b}\mathbf{b} + \mathbf{c}\mathbf{c} \mid 0 \leq a, b, c \leq 1\}$$

spanned by the three vectors equals the area of the base times the height. The former is  $|\mathbf{a} \times \mathbf{b}|$ , and the latter is  $|\mathbf{c}| \cos \theta$ , where  $\theta$  is the angle between  $\mathbf{a} \times \mathbf{b}$  and  $\mathbf{c}$ . Thus, the volume equals

$$|\langle \mathbf{a} \times \mathbf{b}, \mathbf{c} \rangle| = |\det [\mathbf{a} \ \mathbf{b} \ \mathbf{c}]|,$$

so that *the absolute value of a  $3 \times 3$  determinant equals the volume of the parallelepiped spanned by the three columns (or rows).* Incidentally, we define this to be also the volume of the “open” (see the following section) parallelepiped  $\{\mathbf{a}\mathbf{a} + \mathbf{b}\mathbf{b} + \mathbf{c}\mathbf{c} \mid 0 < a, b, c < 1\}$ . This is because the open solid may be expressed as an increasing union of “closed” ones together with continuity of volume: for example,  $(0, 1) = \bigcup_{i=1}^{\infty} [1/n, 1 - 1/n]$ , so that the volume or length of  $(0, 1)$  equals  $\lim_{n \rightarrow \infty} 1 - 1/n - 1/n = 1$ .



Later, we will discuss volume in  $n$ -dimensional space. For now, the previous observations provide some justification for the following:

**Definition 1.6.1.** The *volume* of the (open or closed) parallelepiped in  $\mathbb{R}^n$  spanned by the vectors  $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^n$  is defined to be the absolute value of

$$\det [\mathbf{a}_1 \cdots \mathbf{a}_n].$$

## 1.7 Elementary topology of Euclidean spaces

A fundamental application of the inner product introduced in Section 1.3 is the concept of distance: The *distance* between  $\mathbf{a}$  and  $\mathbf{b} \in \mathbb{R}^n$  is  $d(\mathbf{a}, \mathbf{b}) = |\mathbf{a} - \mathbf{b}|$ . The distance function  $d$  satisfies the following properties for any  $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^n$ :

- $d(\mathbf{a}, \mathbf{b}) = d(\mathbf{b}, \mathbf{a})$ ;
- $d(\mathbf{a}, \mathbf{b}) \geq 0$ , and  $d(\mathbf{a}, \mathbf{b}) = 0$  if and only if  $\mathbf{a} = \mathbf{b}$ ;
- $d(\mathbf{a}, \mathbf{b}) \leq d(\mathbf{a}, \mathbf{c}) + d(\mathbf{c}, \mathbf{b})$ .

All these properties follow from corresponding properties of the norm; for example, the third one is a consequence of the triangle inequality (and is actually responsible for that name, since it says that in a triangle, the length of one side cannot exceed the sum of the other two):

$$d(\mathbf{a}, \mathbf{b}) = |\mathbf{a} - \mathbf{b}| = |(\mathbf{a} - \mathbf{c}) + (\mathbf{c} - \mathbf{b})| \leq |(\mathbf{a} - \mathbf{c})| + |(\mathbf{c} - \mathbf{b})| = d(\mathbf{a}, \mathbf{c}) + d(\mathbf{c}, \mathbf{b}).$$

More generally, a set  $X$  together with a function  $d : X \times X \rightarrow \mathbb{R}$  that satisfies the above three properties is called a *metric space*. Thus, every inner product space is a metric space. Notice that the inner product itself is not essential, but rather the norm associated to it.

**Definition 1.7.1.** A *norm* on a vector space  $E$  is a map  $|| \cdot || : E \rightarrow \mathbb{R}$  such that for all  $\mathbf{a}, \mathbf{b} \in E$  and  $\alpha \in \mathbb{R}$ ,

- (1)  $|\mathbf{a}| \geq 0$ , and  $|\mathbf{a}| = 0$  if and only if  $\mathbf{a} = \mathbf{0}$ ;
- (2)  $|\mathbf{a} + \mathbf{b}| \leq |\mathbf{a}| + |\mathbf{b}|$ ;
- (3)  $|\alpha \mathbf{a}| = |\alpha| |\mathbf{a}|$ .

A *normed vector space* is a vector space together with a norm. It follows that any normed space is a metric space if one defines the distance by  $d(\mathbf{a}, \mathbf{b}) = |\mathbf{a} - \mathbf{b}|$ . One important example is the space of all linear transformations  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  with the operator norm from Definition 1.4.1.

**Definition 1.7.2.** Let  $(X, d)$  be a metric space.

- (1) The *open ball of radius  $r > 0$  around  $\mathbf{a}$*  is the set  $B_r(\mathbf{a})$  of all points at distance less than  $r$  from  $\mathbf{a}$ ;
- (2) A set  $U \subset X$  is said to be a *neighborhood* of  $\mathbf{a} \in X$  if  $U$  contains some open ball around  $\mathbf{a}$ ; in this case, we say  $\mathbf{a}$  is an *interior point* of  $U$ ;
- (3) A set  $V \subset X$  is said to be *open* if it is a neighborhood of each and every one of its elements.
- (4) A set  $C \subset X$  is said to be *closed* if its complement  $X \setminus C$  is open.

It follows from the definition that both  $X$  and the empty set are open. Being complements of each other, they are also closed. Although we will mostly deal with the case  $X = \mathbb{R}^n$ , it is useful, whenever possible, to state properties in terms of abstract

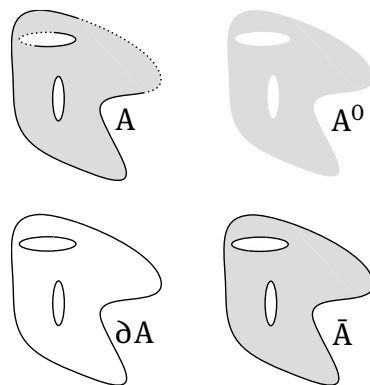
metric spaces. Observe for example that any subset of  $\mathbb{R}^n$  becomes a metric space when endowed with the restriction of the metric on  $\mathbb{R}^n$ .

The triangle inequality guarantees that open balls are indeed open: for if  $\mathbf{b} \in B_r(\mathbf{a})$ , then  $\varepsilon := r - d(\mathbf{a}, \mathbf{b}) > 0$ , and the ball of radius  $\varepsilon$  around  $\mathbf{b}$  is contained inside  $B_r(\mathbf{a})$ . In fact, given  $\mathbf{c} \in B_\varepsilon(\mathbf{b})$ ,  $d(\mathbf{a}, \mathbf{c}) \leq d(\mathbf{a}, \mathbf{b}) + d(\mathbf{b}, \mathbf{c}) < d(\mathbf{a}, \mathbf{b}) + \varepsilon = r$ .

Notice that a union of open sets is again open. Since the complement of a union is the intersection of the complements, an arbitrary intersection of closed sets is closed. The intersection of a finite number of open sets is open (and therefore a finite union of closed sets is closed): if  $\mathbf{a} \in U_i$  for  $i = 1, \dots, k$  and each  $U_i$  is open, then there exists  $r_i > 0$  such that the ball of radius  $r_i$  around  $\mathbf{a}$  is contained in  $U_i$ . Consequently the ball of radius  $r$  around  $\mathbf{a}$ , where  $r = \min\{r_1, \dots, r_k\}$ , is contained in every  $U_i$ , and therefore in their intersection. On the other hand, arbitrary intersections of open sets need not be open: for example, the intersection of all open balls of radius  $1/k$ ,  $k \in \mathbb{N}$ , about a point  $\mathbf{p} \in \mathbb{R}^n$ , consists of the single point  $\mathbf{p}$ , and one-point sets are closed.

Given real numbers  $a < b$ , the interval  $(a, b)$  is open in  $\mathbb{R}$ . Viewing the real line as a subset of  $\mathbb{R}^2$  (i.e., identifying  $\mathbb{R}$  with  $\mathbb{R} \times \{0\} \subset \mathbb{R}^2$ )  $(a, b) \times \{0\}$  is no longer open. In order to be able to say that it is open as a subset of  $\mathbb{R} \times \{0\}$ , we introduce the following:

**Definition 1.7.3.** Let  $A$  be a subset of a metric space  $X$ . A subset  $B$  of  $A$  is said to be *open* (resp. *closed*) in  $A$  or *relative to  $A$*  if  $B = U \cap A$ , where  $U$  is open (resp. closed) in  $X$ .



**Definition 1.7.4.** Let  $A \subset X$ .

- (1) The *interior*  $A^0$ , also denoted  $\text{int } A$ , of  $A$  is the set of all interior points of  $A$ .
- (2) A point is said to be a *boundary point* of  $A$  if every neighborhood of that point intersects both  $A$  and the complement of  $A$ . The *boundary*  $\partial A$  of  $A$  is the collection of all boundary points of  $A$ .
- (3) The *closure* of  $A$  is the set  $\bar{A} = A \cup \partial A$ .
- (4)  $A$  is said to be *bounded* if it is contained in some metric ball (of finite radius).

It is clear that  $A^0$  is open,  $\partial A$  and  $\bar{A}$  are closed, and  $A^0 \subset A \subset \bar{A}$ . In fact,  $A^0$  is the largest open set that is contained in  $A$ , and  $\bar{A}$  is the smallest closed set that contains  $A$ , see Exercise 1.20.

**Definition 1.7.5.** An *open cover*  $\mathcal{O}$  of  $A \subset X$  is a collection of open sets whose union contains  $A$ . A subcollection of  $\mathcal{O}$  is called a *subcover* if it is also a cover of  $A$ .  $A$  is said to be *compact* if any open cover of  $A$  contains a finite subcover.

Any finite set is of course compact, but finiteness is not a requisite: If  $\{a_n\}$  is any sequence of real numbers that converges to some  $a$ , then the set  $K = \{a\} \cup \{a_n \mid n \in \mathbb{N}\}$  is compact, since any open interval containing  $a$  will contain all but finitely many elements of  $K$ .

**Theorem 1.7.1.** (1) *Any compact set is closed;*  
 (2) *Any closed subset of a compact set is compact.*

*Proof.* For (1), suppose  $K$  is compact. We will show that the complement of  $K$  is open. Given  $\mathbf{a} \notin K$ , choose for every  $\mathbf{b} \in K$  open balls  $U_{\mathbf{b}}$  and  $V_{\mathbf{b}}$  centered at  $\mathbf{b}$  and  $\mathbf{a}$  respectively of small enough radius that they don't intersect. By compactness of  $K$ , finitely many of these, say,  $U_{\mathbf{b}_1}, \dots, U_{\mathbf{b}_k}$  cover  $K$ . Then  $V_{\mathbf{b}_1} \cap \dots \cap V_{\mathbf{b}_k}$  is an open ball around  $\mathbf{a}$  that does not intersect any  $U_{\mathbf{b}_j}$ ,  $j = 1, \dots, k$ , and therefore does not intersect  $K$ .

For (2), suppose  $C \subset K$ , where  $C$  is closed and  $K$  is compact, and consider an open cover  $\mathcal{O}$  of  $C$ . Then  $\mathcal{O}$  together with the complement of  $C$  is an open cover of  $K$  and we may extract a finite subcover. If  $X \setminus C$  is a member of this subcover, remove it. What remains is a finite subcollection of  $\mathcal{O}$  that still covers  $C$ .  $\square$

Compact sets are important enough to warrant a more concrete characterization. This alternative description does not, however, hold in arbitrary metric spaces. Recall that  $\alpha \in \mathbb{R}$  is said to be an *upper bound* of a set  $A$  if  $\alpha \geq x$  for any  $x \in A$ , and the *least upper bound* of  $A$  if it no larger than any other upper bound; i.e., if  $\alpha \leq \beta$  for any upper bound  $\beta$  of  $A$ . In this case, we write  $\alpha = \sup A$ . Notice that even though  $\alpha$  itself need not belong to  $A$ , any neighborhood of  $\alpha$  must contain an element of  $A$ : for if, say,  $(\alpha - \varepsilon, \alpha + \varepsilon)$  did not intersect  $A$  for some  $\varepsilon > 0$ , then  $\alpha - \varepsilon$  would be an upper bound of  $A$  smaller than  $\alpha$ . A key property of the real numbers is that any nonempty set of reals that is bounded above has a least upper bound. Appendix A explores these concepts in further detail.

A *box* or *rectangle* in  $\mathbb{R}^n$  is a cartesian product of  $n$  intervals. In the definition of an open set, we could have replaced open metric balls by open boxes: the ball  $B_r(\mathbf{p})$  contains the open box

$$\left(p_1 - \frac{r}{\sqrt{n}}, p_1 + \frac{r}{\sqrt{n}}\right) \times \dots \times \left(p_n - \frac{r}{\sqrt{n}}, p_n + \frac{r}{\sqrt{n}}\right)$$

(with  $p_i = u^i(\mathbf{p})$ ) centered at  $\mathbf{p}$ , and this box in turn contains the ball  $B_{r/\sqrt{n}}(\mathbf{p})$ . In particular, if  $U$  is any open set containing a point  $\mathbf{p}$ , then there exists an open box  $R$  centered at  $\mathbf{p}$  which lies inside  $U$ . Conversely, if  $U$  is a set with the property that any point of  $U$  admits an open box centered at that point which is contained in  $U$ , then any  $\mathbf{p}$  in  $U$  also admits an open ball centered at  $\mathbf{p}$  which is contained in  $U$ , so that  $U$  is



open. Our first goal is to show that a closed and bounded box is compact. We begin in dimension one:

**Theorem 1.7.2 (Heine-Borel).** *For any  $a < b$ , the interval  $[a, b]$  is compact.*

*Proof.* Given an open cover  $\mathcal{O}$  of  $[a, b]$ , denote by  $A$  the set of all  $x \in [a, b]$  such that  $[a, x]$  is covered by finitely many sets in  $\mathcal{O}$ . Then  $a \in A$ ,  $A$  is bounded above by  $b$ , and so  $A$  has a least upper bound  $\alpha \in [a, b]$ . We first observe that  $\alpha \in A$ : indeed,  $\alpha$  belongs to some element  $U$  of  $\mathcal{O}$ . Since  $\alpha$  is the least upper bound of  $A$ ,  $U$  must contain some element  $x \in A$ . By assumption,  $[a, x]$  is then covered by finitely many sets in  $\mathcal{O}$ , and  $[x, \alpha]$  by one, namely  $U$ . Thus  $[a, \alpha] = [a, x] \cup [x, \alpha]$  is also covered by finitely many sets, and  $\alpha \in A$ . We now conclude the proof by showing that  $\alpha = b$ . If  $\alpha < b$ , then the set  $U$  above must contain some  $c \in (\alpha, b)$ , and the above argument shows that  $c \in A$ , contradicting the fact that  $\alpha$  is an upper bound.  $\square$

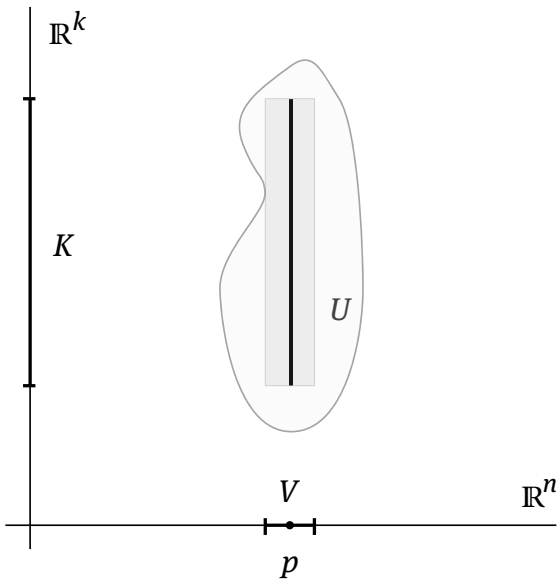


Fig. 1.1: A tube about  $\{p\} \times K$  contained in  $U$

In order to show that a Cartesian product of compact sets is compact, we will need the following:

**Lemma 1.7.1 (the tube lemma).** *Suppose  $p \in \mathbb{R}^n$ , and  $K$  is a compact subset of  $\mathbb{R}^k$ . If  $U$  is an open set in  $\mathbb{R}^{n+k}$  that contains  $\{p\} \times K$ , then  $U$  contains  $V \times K$  for some open neighborhood  $V$  of  $p$  in  $\mathbb{R}^n$ .*

*Proof.* For any  $a \in K$ , there exists an open box inside  $U$  that contains the point  $(p, a)$ . Such a box is a product  $V_a \times W_a$  of open boxes in  $\mathbb{R}^n$  and  $\mathbb{R}^k$  containing  $p$  and  $a$  respectively. Now,  $K$  is compact, and is therefore covered by finitely many  $W_{a_1}, \dots, W_{a_l}$ .

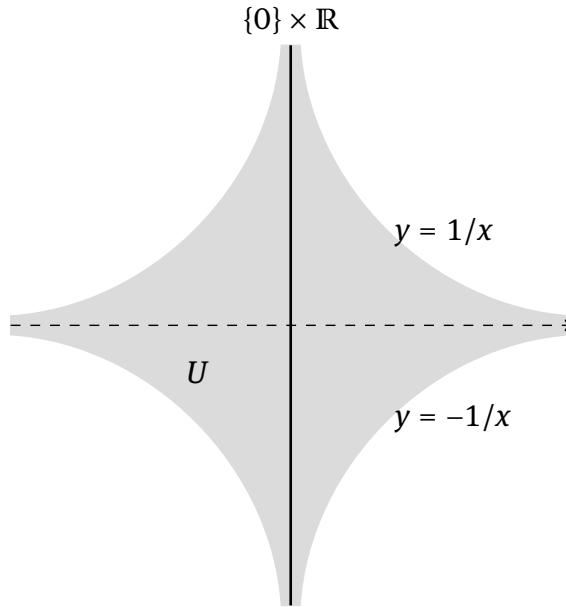
If  $V = V_{a_1} \cap \cdots \cap V_{a_l}$ , then  $V$  is an open neighborhood of  $p$ , and

$$V \times K \subset V \times \left( \bigcup_{i=1}^l W_{a_i} \right) = \bigcup_{i=1}^l (V \times W_{a_i}) \subset \bigcup_{i=1}^l (V_{a_i} \times W_{a_i}) \subset U,$$

as claimed.  $V \times K$  is called a *tube* about  $\{p\} \times K$ . □

Notice that the set  $\{p\} \times K$  in the above lemma is compact if  $K$  is: indeed, if  $\pi : \mathbb{R}^n \times \mathbb{R}^k \rightarrow \mathbb{R}^k$  denotes the projection onto the second factor, then  $\pi$  is an *open map*; i.e.,  $\pi$  maps open sets to open sets – in fact,  $\pi$  maps balls of a given radius onto balls of the same radius. Thus, if  $\mathcal{O}$  is an open cover of  $\{p\} \times K$ , then the sets  $\pi(U)$ ,  $U \in \mathcal{O}$ , form an open cover of  $K$ , and a finite subcover  $\pi(U_1), \dots, \pi(U_l)$  of  $K$  may be chosen. This means that  $U_1 \cup \cdots \cup U_l \supset \{p\} \times K$ .

It is easily seen that compactness of  $K$  is essential in the tube lemma: if, for example  $n = k = 1$ ,  $p = 0$ , and  $K = \mathbb{R}$ , then  $U = \{(x, y) \mid |xy| < 1\}$  is an open subset of  $\mathbb{R}^2$  containing  $\{0\} \times \mathbb{R}$ , but it cannot contain any tube around it.



**Theorem 1.7.3.** *A Cartesian product of compact sets is compact. In particular, if  $a_i < b_i$ ,  $i = 1, \dots, n$ , then  $[a_1, b_1] \times \cdots \times [a_n, b_n]$  is compact.*

*Proof.* The second statement follows from the first one together with Theorem 1.7.2. For the first one, it is enough to show that  $A \times B$  is compact whenever  $A$  and  $B$  are. Let  $\mathcal{O}$  denote an open cover of  $A \times B$ . For each  $p \in A$ , the set  $\{p\} \times B$ , being compact, can be covered by a finite subcollection  $\mathcal{O}_p$  of  $\mathcal{O}$ . By the tube lemma, the union of the sets in  $\mathcal{O}_p$  contains  $U_p \times B$  for some neighborhood  $U_p$  of  $p$ , and  $A$  can be covered by finitely many of these, say,  $U_{p_1}, \dots, U_{p_l}$ . But then  $\mathcal{O}_{p_1} \cup \cdots \cup \mathcal{O}_{p_l}$  is a finite cover of  $A \times B$ . □

**Theorem 1.7.4.** *A subset  $K$  of  $\mathbb{R}^n$  is compact if and only if it is closed and bounded.*

*Proof.* If  $K$  is compact, then it is closed by Theorem 1.7.1. It is also bounded, for if  $\mathbf{p} \in K$ , then the collection of balls  $B_k(\mathbf{p})$ ,  $k \in \mathbb{N}$ , is an open cover of  $K$ , and admits a finite subcover. The ball of largest radius in that subcover contains  $K$ . Conversely, suppose  $K$  is closed and bounded. Being bounded, it is contained inside some closed rectangle. The latter is compact by Theorem 1.7.3. Since  $K$  is closed, it is compact by Theorem 1.7.1.  $\square$

By definition, if  $E$  is not compact, then there exists some open cover of  $E$  with no finite subcover. However, if  $E$  lies in Euclidean space, there will always be a countable subcover:

**Theorem 1.7.5.** *If  $E \subset \mathbb{R}^n$ , then any open cover  $\{U_\alpha\}_{\alpha \in J}$  of  $E$  has a countable subcover; i.e., there exists a subset  $A$  of natural numbers, and a map  $f : A \rightarrow J$  such that*

$$E \subset \bigcup_{k \in A} U_{f(k)}.$$

*Proof.* The argument uses the notion of countable set and dense set, both of which are examined in Appendix A. For each  $\mathbf{a} \in E$ , choose some  $U_\alpha$  that contains it. Since this set is open, there exists some  $r > 0$  such that the ball  $B_r(\mathbf{a}) \subset U_\alpha$ . The fact that rationals are dense in  $\mathbb{R}$  (see Appendix A) is easily extended to  $\mathbb{Q}^n$  being dense in  $\mathbb{R}^n$ . Thus, there exists a ball  $B_{\mathbf{a}}$  with center  $\mathbf{a} \in \mathbb{Q}^n$  and rational radius such that  $\mathbf{a} \in B_{\mathbf{a}} \subset U_\alpha$ . But  $\mathbb{Q}^n \times \mathbb{Q}$  is countable, and therefore so is the collection

$$\{B_{\mathbf{a}} \mid \mathbf{a} \in E\} = \{V_1, V_2, \dots\}.$$

By assumption, each  $V_k$  is contained in some  $U_\alpha$ . Choosing some such  $\alpha$  and setting  $\alpha = f(k)$ , we conclude that  $E \subset \bigcup_k V_k \subset \bigcup_k U_{f(k)}$ .  $\square$

The least upper bound of a set  $A$  of real numbers is also called the *supremum* of  $A$ , and is denoted  $\sup A$ . A similar notion can be introduced for lower bounds: the *infimum* or *greatest lower bound*  $\alpha = \inf A$  of  $A$  is a lower bound of  $A$  (i.e.,  $\alpha \leq x$  for any  $x \in A$ ) that is greater than or equal to any other lower bound of  $A$ . Notice that  $\alpha$  is a lower bound of  $A$  if and only if  $-\alpha$  is an upper bound of  $-A = \{-a \mid a \in A\}$ . It easily follows that  $\inf A = -\sup(-A)$ , in the sense that if one of the two exists, then so does the other, and they are equal. In particular, any nonempty set of real numbers that is bounded below has a greatest lower bound.

A useful observation is that if  $A$  is closed, nonempty, and bounded above (respectively below), then  $A$  contains its least upper bound (resp. greatest lower bound): indeed, by definition, given any  $\varepsilon > 0$ , there must be a point  $x \in A$  such that

$$\sup A - \varepsilon < x \leq \sup A,$$

since otherwise  $\sup A - \varepsilon$  would be an upper bound of  $A$ . Thus, any open interval of radius  $\varepsilon$  around  $\sup A$  contains points of  $A$  (as was just pointed out) and also points

outside  $A$  (any  $y \in (\sup A, \sup A + \varepsilon)$ ), so that  $\sup A$  is a boundary point of  $A$ . But closed sets contain their boundary points. The argument for infimum is similar.

The concepts of infimum and supremum also enable us to prove a property of intervals that will be quite useful in the sequel:

**Proposition 1.7.1.** *If  $I$  is an interval of real numbers, then any subset that is both open and closed in  $I$  is either empty or equals all of  $I$ .*

*Proof.* We argue by contradiction: suppose  $A$  is a nonempty subset of  $I$  that is both open and closed, but there exists some  $c \in I$  that does not belong to  $A$ . Consider any  $t_0 \in A$ . Then  $t_0 < c$  or  $t_0 > c$ . The argument is similar in both cases, so we only consider the former. Let

$$B = \{t \in I \mid [t_0, t] \subset A\}.$$

Then  $B$  is nonempty (since it contains  $t_0$ ), bounded above (by  $c$ ), and so admits a supremum  $\alpha$ . By definition,  $\alpha$  is a boundary point of  $A$ . Since  $A$  is closed in  $I$ ,  $A$  equals the intersection of  $I$  with some closed set  $C$ , and  $\alpha \in C$ : indeed, if some neighborhood of  $\alpha$  is contained in  $C$ , then certainly  $\alpha \in C$ . Otherwise, every neighborhood of  $\alpha$  contains points outside  $C$ . It must also contain points in  $C$  since it contains points of  $A \subset C$ . This means  $\alpha$  is a boundary point of  $C$ , and since  $C$  is closed,  $\alpha \in C$  as claimed. But  $\alpha$  also belongs to  $I$ , because the latter is an interval that contains a point to the left of  $\alpha$  (e.g.,  $t_0$ ) and one to the right of it (e.g.,  $c$ ). Thus,  $\alpha \in A = C \cap I$ , and so  $\alpha < c$ .  $A$  is open in  $I$ , however, so there exists some  $\varepsilon > 0$  such that  $(\alpha, \alpha + \varepsilon) \cap I \subset A$ . Thus, if  $\beta$  is any point in the nonempty set  $(\alpha, \alpha + \varepsilon) \cap I$ , then  $[t_0, \beta] \subset A$ , contradicting the assumption that  $\alpha$  is an upper bound of  $B$ .  $\square$

## 1.8 Sequences

The reader is already familiar with sequences of real numbers from Calculus. The generalization to metric spaces is straightforward: A *sequence* in a metric space  $(X, d)$  is a map from the set  $\mathbb{N}$  of natural numbers into  $X$ . The value of the map at  $k \in \mathbb{N}$  is usually denoted  $\mathbf{a}_k$  (or some other letter such as  $\mathbf{p}_k$ ) when  $X = \mathbb{R}^n$ ,  $n > 1$ , and the sequence itself by  $\{\mathbf{a}_k\}$  – although we will often use  $\mathbf{a}_k$  to denote either one when there is no risk of confusion. When  $n = 1$  or  $X$  is not Euclidean space, we use the regular font  $a_k$  instead of the bold one.

**Definition 1.8.1.** A sequence  $\{\mathbf{a}_k\}$  is said to *converge* to  $\mathbf{a}$  if for any  $\varepsilon > 0$ , there is a positive integer  $N$  such that  $d(\mathbf{a}_k, \mathbf{a}) < \varepsilon$  whenever  $k \geq N$ . In this case,  $\mathbf{a}$  is called the *limit* of the sequence, and we write  $\mathbf{a}_k \rightarrow \mathbf{a}$  or  $\lim_{k \rightarrow \infty} \mathbf{a}_k = \mathbf{a}$ . If no such  $\mathbf{a}$  exists, the sequence is said to *diverge*. The sequence is said to be *bounded* if the set  $\{\mathbf{a}_k \mid k \in \mathbb{N}\}$  of values is bounded; i.e., if there exists some  $R > 0$  such that  $d(\mathbf{a}_1, \mathbf{a}_k) \leq R$  for all  $k$ .

Another way to describe convergence  $\mathbf{a}_k \rightarrow \mathbf{a}$  is to say that any neighborhood of  $\mathbf{a}$  contains  $\mathbf{a}_k$  except perhaps for finitely many values of  $k$ . The one in the above definition is of course just the open ball of radius  $\varepsilon$  about  $\mathbf{a}$ . This also justifies using the terminology “the limit”, for a sequence can have at most one limit: if  $\mathbf{a} \neq \mathbf{b}$ , then these two points admit disjoint neighborhoods (such as the open metric balls of radius half the distance between them), and then both neighborhoods cannot contain  $\mathbf{a}_k$  for all sufficiently large  $k$ . It also shows that any convergent sequence is bounded: there are at most finitely many points of the sequence outside the ball of radius, say, 1 about the limit, and a finite set is always bounded: specifically, if  $\mathbf{a}_k \in B_1(\mathbf{a})$  for all  $k \geq N$ , let  $R = \max\{1, d(\mathbf{a}_1, \mathbf{a}), \dots, d(\mathbf{a}_{N-1}, \mathbf{a})\}$ . Then  $d(\mathbf{a}_k, \mathbf{a}) \leq R$  for all  $k$ . It is also clear from the above discussion that changing finitely many terms in a sequence does not affect convergence; i.e., if  $\mathbf{a}_k$  and  $\mathbf{b}_k$  are two sequences such that  $\mathbf{a}_k = \mathbf{b}_k$  for all sufficiently large  $k$ , then they either both converge to the same limit, or they both diverge.

We first consider convergence of sequences of real numbers, since it is very much related to that of sequences in higher-dimensional Euclidean spaces. Even though the reader is probably familiar with the contents of the following theorem from a previous Calculus course, it has been included because its proof is often not covered there.

**Theorem 1.8.1.** *Let  $\{x_k\}$  and  $\{y_k\}$  be sequences of real numbers that converge to  $x$  and  $y$  respectively. Then*

- (1)  $x_k + y_k \rightarrow x + y$ ;
- (2)  $cx_k \rightarrow cx$  for any  $c \in \mathbb{R}$ ;
- (3)  $x_k y_k \rightarrow xy$ ;
- (4) If  $x_k, x \neq 0$ , then  $1/x_k \rightarrow 1/x$ ;
- (5) If  $y_k, y \neq 0$ , then  $(x_k/y_k) \rightarrow (x/y)$ .

*Proof.* (1) Let  $\varepsilon > 0$  be given. By assumption, there exist positive integers  $N_1$  and  $N_2$  such that  $|x_k - x| < \varepsilon/2$  for all  $k \geq N_1$ , and  $|y_k - y| < \varepsilon/2$  for all  $k \geq N_2$ . So, if  $N$  is the largest of the two numbers  $N_1$  and  $N_2$ , then for any  $k \geq N$ ,

$$|(x_k + y_k) - (x + y)| = |(x_k - x) + (y_k - y)| \leq |x_k - x| + |y_k - y| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.$$

(2) This is an immediate consequence of (3), taking  $y_k = c$  for all  $k$ .

(3) Since  $\{x_k\}$  converges, there exists some  $M > 0$  such that  $|x_k| \leq M$  for all  $k$ . Choose  $N_1$  so that  $|y_k - y| < \varepsilon/(2M)$  for  $k \geq N_1$ , and  $N_2$  so that  $|x_k - x| < \varepsilon/(2|y|)$  if  $k \geq N_2$  (unless  $y = 0$ , in which case take  $N_2 = 1$ ). If  $N$  is the larger of  $N_1$  and  $N_2$ , then for  $k \geq N$ ,

$$|x_k y_k - xy| = |x_k(y_k - y) + y(x_k - x)| \leq |x_k||y_k - y| + |y||x_k - x|.$$

The first term on the right side of the above inequality is smaller than  $\varepsilon/2$ . The second term is either zero (if  $y = 0$ ) or less than  $\varepsilon/2$  (if  $y \neq 0$ ). In either case, the left side is less than  $\varepsilon$ , which establishes the claim.

(4) Take  $\varepsilon = |x|/2$  in the definition of convergence of  $\{x_k\}$  to conclude that there exists  $N_1$  such that  $|x_k| > |x|/2$  whenever  $n \geq N_1$ . Next, for a given  $\varepsilon > 0$ , there exists  $N_2$  such

that  $|x_k - x| < \varepsilon|x|^2/2$  whenever  $k \geq N_2$ . Let  $N = \max\{N_1, N_2\}$ . If  $k \geq N$ , then

$$\left| \frac{1}{x_k} - \frac{1}{x} \right| = \left| \frac{x_k - x}{x_k x} \right| < \frac{2|x_k - x|}{|x|^2} < \varepsilon.$$

(5) is an easy consequence of (3) and (4).  $\square$

**Examples and Remarks 1.8.1.** (i) A sequence  $\{a_k\}$  of real numbers is said to be *increasing* (resp. *decreasing*) if  $a_k \leq a_{k+1}$  (resp.  $a_k \geq a_{k+1}$ ) for all  $k$ . An increasing sequence that is bounded above converges: to see this, let  $\alpha$  denote the supremum of  $\{a_k \mid k \in \mathbb{N}\}$ . Given  $\varepsilon > 0$ , there exists some  $N \in \mathbb{N}$  such that  $0 \leq \alpha - a_N < \varepsilon$ . Since the sequence is increasing,  $0 \leq \alpha - a_k \leq \alpha - a_N < \varepsilon$  for all  $k \geq N$ , and  $a_k \rightarrow \alpha$ . This also implies that a decreasing sequence that is bounded below converges to the infimum of its set of values, because if  $\{a_k\}$  is decreasing bounded below, then  $\{-a_k\}$  is increasing and bounded above.

(ii) Another useful tool for proving convergence of a real sequence is the so-called squeeze theorem: if  $a_k \leq b_k \leq c_k$ , and both  $\{a_k\}$  and  $\{c_k\}$  converge to the same limit  $L$ , then so does the middle sequence  $\{b_k\}$ . To see this, let  $\varepsilon > 0$ , and choose  $N_1$  such that  $|a_k - L| < \varepsilon$  whenever  $k \geq N_1$ . Similarly, let  $N_2$  be such that  $|c_k - L| < \varepsilon$  for  $k \geq N_2$ . If  $k$  is larger than  $N = \max\{N_1, N_2\}$ , then  $|b_k - L| < \varepsilon$ , because

$$-\varepsilon < a_k - L \leq b_k - L \leq c_k - L < \varepsilon.$$

(iii) One of the simplest applications of (i) and (ii) is the so-called *geometric sequence*  $r^k$ , where  $0 < r < 1$ . It is a decreasing sequence, bounded below by zero, and therefore converges. This does not quite tell us what the limit is, but if  $r^k \rightarrow L$ , then  $a_k \rightarrow L$ , where  $a_k = r^{k+1}$ : indeed, given  $\varepsilon > 0$ , choose  $N$  such that  $|r^k - L| < \varepsilon$  for  $k > N$ ; then  $|a_k - L| = |r^{k+1} - L|$  is also less than  $\varepsilon$  for  $k > N$  because  $k + 1 > k > N$ . On the other hand, by Theorem 1.8.1 (2),  $a_k \rightarrow rL$ . Thus,  $L = rL$ , and since  $r \neq 1$ ,  $L = 0$ . In conclusion,  $r^k \rightarrow 0$ .

We have, in fact, that  $r^k \rightarrow 0$  for any  $r \in (-1, 1)$ ; indeed, we know that  $|r|^k \rightarrow 0$ . By Theorem 1.8.1 (2),  $-|r|^k \rightarrow 0$ . Since  $-|r|^k \leq r^k \leq |r|^k$ , the claim follows from the squeeze theorem in (ii).

For the sake of convenience, we will often abbreviate the sentence “There exists  $N \in \mathbb{N}$  such that the statement  $P(k)$  is true whenever  $k \geq N$ ” by “ $P(k)$  is true for sufficiently large  $k$ ”. Now that we have some experience with real-valued sequences, let us look at sequences in higher-dimensional Euclidean space.

**Theorem 1.8.2.** A sequence  $\{\mathbf{a}_k\}$  in  $\mathbb{R}^n$  converges to  $\mathbf{a}$  iff  $u^i(\mathbf{a}_k) \rightarrow u^i(\mathbf{a})$  in  $\mathbb{R}$  for  $i = 1, \dots, n$ .

*Proof.* If  $\mathbf{a}_k \rightarrow \mathbf{a}$ , then  $u^i(\mathbf{a}_k) \rightarrow u^i(\mathbf{a})$  since  $|u^i(\mathbf{a}_k) - u^i(\mathbf{a})| \leq |\mathbf{a}_k - \mathbf{a}|$ . Conversely, if  $u^i(\mathbf{a}_k) \rightarrow u^i(\mathbf{a})$  for each  $i = 1, \dots, n$ , then given  $\varepsilon > 0$ ,  $|u^i(\mathbf{a}_k) - u^i(\mathbf{a})| < \varepsilon/\sqrt{n}$  for all  $i$

and sufficiently large  $k$ . But then

$$|\mathbf{a}_k - \mathbf{a}| = \left( \sum_{i=1}^n |u^i(\mathbf{a}_k) - u^i(\mathbf{a})|^2 \right)^{\frac{1}{2}} < \varepsilon$$

for  $k$  large enough, which shows that  $\mathbf{a}_k \rightarrow \mathbf{a}$ .  $\square$

**Corollary 1.8.1.** *If  $\{\mathbf{a}_k\}$  and  $\{\mathbf{b}_k\}$  are sequences in  $\mathbb{R}^n$  that converge to  $\mathbf{a}$  and  $\mathbf{b}$  respectively, then*

- (1)  $\mathbf{a}_k + \mathbf{b}_k \rightarrow \mathbf{a} + \mathbf{b}$ ;
- (2) *If  $c_k \rightarrow c \in \mathbb{R}$ , then  $c_k \mathbf{a}_k \rightarrow c\mathbf{a}$ ;*
- (3)  $\langle \mathbf{a}_k, \mathbf{b}_k \rangle \rightarrow \langle \mathbf{a}, \mathbf{b} \rangle$ ;
- (4)  $|\mathbf{a}_k| \rightarrow |\mathbf{a}|$ .

*Proof.* The first three statements follow immediately from the previous two theorems. The last one is a consequence of the third.  $\square$

Sequences in compact sets have additional properties. Before describing one such, we need some terminology: If  $\{\mathbf{a}_k\}$  is a sequence, and  $k_1 < k_2 < \dots$  is a strictly increasing sequence of positive integers, then the sequence  $\mathbf{a}_{k_1}, \mathbf{a}_{k_2}, \dots$  is called a *subsequence* of  $\mathbf{a}_k$  (more formally, if  $f : \mathbb{N} \rightarrow \mathbb{R}^n$ ,  $f(k) = \mathbf{a}_k$ , is the function defining the sequence, then a subsequence of  $f$  is a function  $f \circ g$ , where  $g : \mathbb{N} \rightarrow \mathbb{N}$  is strictly increasing). Elementary examples are the subsequences  $\{a_{2k}\}$  of even terms and  $\{a_{2k-1}\}$  of odd terms. It is an easy exercise to show that if a sequence converges, then any subsequence converges to the same limit. In fact, we proved this directly for the subsequence  $\{r^{k+1}\}$  of  $\{r^k\}$  in Examples and Remarks 1.8.1 (iii). The general case is similar.

**Theorem 1.8.3.** *If  $K \subset \mathbb{R}^n$  is compact, then any sequence in  $K$  contains a convergent subsequence.*

*Proof.* Let  $A = \{\mathbf{a}_k \mid k \in \mathbb{N}\}$ . If  $A$  is a finite set, then there must be some  $\mathbf{a}$  that equals  $\mathbf{a}_k$  for infinitely many  $k$ . This yields a constant subsequence. We may therefore assume that  $A$  is infinite, and claim that there exists some  $\mathbf{a} \in K$  with the following property: every neighborhood of  $\mathbf{a}$  intersects  $A \setminus \{\mathbf{a}\}$  (a point with this property is called a *limit point* of  $A$ ; notice that any neighborhood of a limit point  $\mathbf{a}$  of  $A$  must contain infinitely many elements of  $A$ : if it contained only finitely many, then we could find  $\mathbf{b} \in A \setminus \{\mathbf{a}\}$  closest to  $\mathbf{a}$ , and the open ball around  $\mathbf{a}$  with radius  $|\mathbf{a} - \mathbf{b}|$  would no longer intersect  $A \setminus \{\mathbf{a}\}$ ). To establish the claim, we argue by contradiction: suppose that every  $\mathbf{a}$  in  $K$  has a neighborhood that does not intersect  $A \setminus \{\mathbf{a}\}$ . Apply this first to points in the complement of  $A$ : every point outside  $A$  (be it in  $K$  or outside of  $K$ ) has a neighborhood disjoint from  $A$ , so that the complement of  $A$  is open, and  $A$  is closed. By Theorem 1.7.1,  $A$  is compact. Next, apply it to points inside  $A$ : every  $\mathbf{a} \in A$  has an open neighborhood  $U_{\mathbf{a}}$  such that  $U_{\mathbf{a}} \cap A = \{\mathbf{a}\}$ . Since  $A$  is infinite, this means that  $\{U_{\mathbf{a}} \mid \mathbf{a} \in A\}$  is an open cover of  $A$  with no finite subcover, contradicting compactness.

Now that we have established the existence of a limit point  $\mathbf{a}$ , we can construct inductively a subsequence that converges to it. By assumption, there is an integer  $k_1$  such that  $\mathbf{a}_{k_1} \in B_1(\mathbf{a})$ . If  $k_j$  is a given integer, then the ball  $B_{1/(j+1)}(\mathbf{a})$  contains infinitely many elements of  $A$ , and we may therefore choose an integer  $k_{j+1} > k_j$  such that  $\mathbf{a}_{k_{j+1}}$  belongs to it. The subsequence  $\{\mathbf{a}_{k_j}\}$  thus constructed satisfies  $|\mathbf{a}_{k_j} - \mathbf{a}| < 1/j$  for all  $j$ , and therefore converges to  $\mathbf{a}$ .  $\square$

**Corollary 1.8.2.** *A bounded sequence in  $\mathbb{R}^k$  has a convergent subsequence.*

*Proof.* By assumption, the sequence lies inside some ball, and this ball has compact closure, so Theorem 1.8.3 applies.  $\square$

**Example 1.8.1.** Let  $r > 0$ ,  $\alpha \in \mathbb{R}$ , and  $\mathbf{a}_k = [r^k \cos(k\alpha) \quad r^k \sin(k\alpha)]^T \in \mathbb{R}^2$ . If  $r > 1$ , then  $\{\mathbf{a}_k\}$  diverges, because the sequence is unbounded:  $|\mathbf{a}_k| = r^k$ . Suppose next that  $r < 1$ . Since

$$0 \leq |u^1(\mathbf{a}_k)|, |u^2(\mathbf{a}_k)| \leq r^k,$$

and  $r^k \rightarrow 0$ ,  $|u^i(\mathbf{a}_k)| \rightarrow 0$  for  $i = 1, 2$  by the squeeze theorem. But then  $u^i(\mathbf{a}_k) \rightarrow 0$  because

$$-|u^i(\mathbf{a}_k)| \leq u^i(\mathbf{a}_k) \leq |u^i(\mathbf{a}_k)|.$$

Finally, when  $r = 1$ , the sequence will, in general, diverge, but must admit a convergent subsequence by Corollary 1.8.2 because  $|\mathbf{a}_k| = 1$ .

There is an alternative characterization of convergent sequences, one that is often useful because it does not explicitly involve a limit:

**Definition 1.8.2.**  $\{\mathbf{a}_k\}$  is said to be a *Cauchy sequence* if for any  $\varepsilon > 0$ , there exists an integer  $N$  such that  $|\mathbf{a}_k - \mathbf{a}_l| < \varepsilon$  for all  $k, l \geq N$ .

If  $\mathbf{a}_k \rightarrow \mathbf{a}$ , then  $\{\mathbf{a}_k\}$  is Cauchy: given  $\varepsilon > 0$  choose  $N$  so that  $|\mathbf{a}_k - \mathbf{a}| < \varepsilon/2$  whenever  $k \geq N$ . If  $k, l \geq N$ , then

$$|\mathbf{a}_k - \mathbf{a}_l| = |(\mathbf{a}_k - \mathbf{a}) + (\mathbf{a} - \mathbf{a}_l)| \leq |\mathbf{a}_k - \mathbf{a}| + |\mathbf{a}_l - \mathbf{a}| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.$$

On the other hand, if  $\{\mathbf{a}_k\}$  is Cauchy, then it is bounded: there exists a positive integer  $N$  such that  $|\mathbf{a}_k - \mathbf{a}_l| < 1$  whenever  $k, l \geq N$ . But then the sequence is contained inside the bounded set

$$B_1(\mathbf{a}_N) \cup \{\mathbf{a}_1, \dots, \mathbf{a}_{N-1}\}.$$

This, together with Corollary 1.8.2, implies:

**Theorem 1.8.4.** *A sequence in  $\mathbb{R}^n$  converges if and only if it is a Cauchy sequence.*

*Proof.* We've already established that a convergent sequence is Cauchy, and that conversely, a Cauchy sequence  $\{\mathbf{a}_k\}$  in  $\mathbb{R}^n$  has a convergent subsequence, say,  $\mathbf{a}_{k_j} \rightarrow \mathbf{a}$ . But then the sequence itself must converge to  $\mathbf{a}$ : given  $\varepsilon > 0$ , choose  $N_1$  such that



$|\mathbf{a}_k - \mathbf{a}_l| < \varepsilon/2$  whenever  $k, l \geq N_1$ , and choose  $N_2$  such that  $|\mathbf{a} - \mathbf{a}_{k_j}| < \varepsilon/2$  if  $j \geq N_2$ . Let  $N$  denote the larger of  $N_1$  and  $N_2$ . Observing that  $k_N \geq N$ , we have

$$|\mathbf{a} - \mathbf{a}_k| \leq |\mathbf{a} - \mathbf{a}_{k_N}| + |\mathbf{a}_{k_N} - \mathbf{a}_k| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} \text{ if } k \geq N,$$

which establishes the claim.  $\square$

There are metric spaces in which not every Cauchy sequence converges: consider for example the set  $\mathbb{Q}$  of rational numbers with the distance function inherited as a subset of the real numbers. If  $x_n$  is the number representing the first  $n$  digits in the decimal expansion of  $\sqrt{2}$ , then  $x_n \rightarrow \sqrt{2}$  in  $\mathbb{R}$ , and is therefore Cauchy. Each  $x_n$  is rational, and the sequence does not converge in  $\mathbb{Q}$ . A metric space in which every Cauchy sequence converges is said to be *complete*. Notice that every closed subset of a complete metric space is complete, and  $\mathbb{Q}$  is of course not closed in  $\mathbb{R}$ .

Given a sequence  $\{\mathbf{a}_m\}$  in  $\mathbb{R}^n$ , consider the sequence  $\{\mathbf{s}_k\}$  of *partial sums*, where

$$\mathbf{s}_k = \mathbf{a}_1 + \cdots + \mathbf{a}_k.$$

We say the *infinite series*  $\sum_{m=1}^{\infty} \mathbf{a}_m$  converges to  $L$  (and write  $\sum_{m=1}^{\infty} \mathbf{a}_m = L$ ) if  $\lim_{k \rightarrow \infty} \mathbf{s}_k = L$ . The series is said to *converge absolutely* if the series of real numbers  $\sum_m |\mathbf{a}_m|$  converges.

Completeness implies the following:

**Theorem 1.8.5.** *An absolutely convergent series converges.*

*Proof.* Let  $\mathbf{s}_k = \mathbf{a}_1 + \cdots + \mathbf{a}_k$ ,  $s_k = |\mathbf{a}_1| + \cdots + |\mathbf{a}_k|$ . The claim will follow once we show that  $\{\mathbf{s}_k\}$  is a Cauchy sequence. So let  $\varepsilon > 0$ . Since  $\{s_k\}$  converges, it is Cauchy, and there exists an integer  $N$  such that  $|s_m - s_l| < \varepsilon$  whenever  $m > l \geq N$ . By the triangle inequality, if  $m > l \geq N$ ,

$$|\mathbf{s}_m - \mathbf{s}_l| = |\mathbf{a}_{l+1} + \cdots + \mathbf{a}_m| \leq |\mathbf{a}_{l+1}| + \cdots + |\mathbf{a}_m| = |s_m - s_l| < \varepsilon. \quad \square$$

## 1.9 Limits and continuity

Recall from Section 1.8 that  $\mathbf{a}$  is said to be a *limit point* of a set  $A \subset \mathbb{R}^n$  if every neighborhood of  $\mathbf{a}$  intersects  $A \setminus \{\mathbf{a}\}$ . Maps from  $A$  to  $\mathbb{R}^m$  will be denoted in bold font when  $m > 1$  (provided they are not linear). If  $\mathbf{f} : A \rightarrow \mathbb{R}^m$  is a map, its *component functions* are the maps  $f^i = u^i \circ \mathbf{f} : A \rightarrow \mathbb{R}$ ,  $i = 1, \dots, m$ . The *graph* of  $\mathbf{f}$  is the subset of  $\mathbb{R}^{n+m}$  that consists of all  $(\mathbf{a}, \mathbf{f}(\mathbf{a}))$  as  $\mathbf{a}$  ranges over  $A$ . If  $B \subset A$ , define the *image* of  $B$  as  $\mathbf{f}(B) = \{\mathbf{f}(\mathbf{b}) \mid \mathbf{b} \in B\}$ . Finally, if  $C \subset \mathbb{R}^m$ , the *pre-image* of  $C$  is  $\mathbf{f}^{-1}(C) = \{\mathbf{a} \in A \mid \mathbf{f}(\mathbf{a}) \in C\}$ .

**Definition 1.9.1.** Let  $A \subset \mathbb{R}^n$ ,  $\mathbf{f} : A \rightarrow \mathbb{R}^m$  a map from  $A$  to  $\mathbb{R}^m$ , and  $\mathbf{a}$  a limit point of  $A$ . We say the *limit of  $\mathbf{f}$  at  $\mathbf{a}$  equals  $\mathbf{b}$*   $\in \mathbb{R}^m$ , and write

$$\lim_{\mathbf{p} \rightarrow \mathbf{a}} \mathbf{f}(\mathbf{p}) = \mathbf{b}, \quad \text{or} \quad \mathbf{f}(\mathbf{p}) \rightarrow \mathbf{b} \text{ as } \mathbf{p} \rightarrow \mathbf{a},$$

if for any  $\varepsilon > 0$  there exists  $\delta > 0$  such that

$$|\mathbf{f}(\mathbf{p}) - \mathbf{b}| < \varepsilon \quad \text{whenever} \quad 0 < |\mathbf{p} - \mathbf{a}| < \delta, \quad \mathbf{p} \in A.$$

Notice that  $\mathbf{a}$  need not belong to the domain  $A$  of  $\mathbf{f}$ . An alternative characterization is that for any  $\varepsilon > 0$  there exists  $\delta > 0$  such that

$$\mathbf{f}(B_\delta(\mathbf{a}) \setminus \{\mathbf{a}\}) \subset B_\varepsilon(\mathbf{b}).$$

A useful way of determining limits is by means of the following:

**Theorem 1.9.1.** *With notation as in Definition 1.9.1, the limit of  $\mathbf{f}$  at  $\mathbf{a}$  equals  $\mathbf{b}$  if and only if for every sequence  $\{\mathbf{a}_k\}$  in  $A$ ,  $\mathbf{a}_k \neq \mathbf{a}$ , that converges to  $\mathbf{a}$ , the sequence  $\{\mathbf{f}(\mathbf{a}_k)\}$  converges to  $\mathbf{b}$ .*

*Proof.* Suppose  $\lim_{\mathbf{p} \rightarrow \mathbf{a}} \mathbf{f}(\mathbf{p}) = \mathbf{b}$ , and consider a sequence  $\mathbf{a}_k \rightarrow \mathbf{a}$ ,  $\mathbf{a}_k \neq \mathbf{a}$ . Given  $\varepsilon > 0$ , we must show that there exists some integer  $N$  such that  $|\mathbf{f}(\mathbf{a}_k) - \mathbf{b}| < \varepsilon$  for  $k \geq N$ . By hypothesis, for this  $\varepsilon$  there exists  $\delta > 0$  such that  $|\mathbf{f}(\mathbf{p}) - \mathbf{b}| < \varepsilon$  whenever  $0 < |\mathbf{p} - \mathbf{a}| < \delta$ ,  $\mathbf{p} \in A$ . Since  $\mathbf{a}_k \rightarrow \mathbf{a}$ , there exists a positive integer  $N$  such that  $0 < |\mathbf{a}_k - \mathbf{a}| < \delta$  when  $k \geq N$ . Thus, for  $k \geq N$ ,  $|\mathbf{f}(\mathbf{a}_k) - \mathbf{b}| < \varepsilon$ , and  $\mathbf{f}(\mathbf{a}_k) \rightarrow \mathbf{b}$ .

Conversely, if the limit of  $\mathbf{f}$  at  $\mathbf{a}$  does not equal  $\mathbf{b}$ , then there exists an  $\varepsilon > 0$  such that for every  $\delta > 0$ , there exists some  $\mathbf{p} \in A$  within distance less than  $\delta$  from  $\mathbf{a}$  such that  $|\mathbf{f}(\mathbf{p}) - \mathbf{b}| \geq \varepsilon$ . Choose some such point  $\mathbf{a}_k$  for each  $\delta = 1/k$ ,  $k \in \mathbb{N}$ . Then  $\{\mathbf{a}_k\}$  is a sequence that converges to  $\mathbf{a}$ , and  $\lim_{k \rightarrow \infty} \mathbf{f}(\mathbf{a}_k) \neq \mathbf{b}$ .  $\square$

Since sequences have at most one limit, the limit of a function, if it exists, is also unique. Furthermore, the properties of sequences imply the following:

**Theorem 1.9.2.** *Let  $A \subset \mathbb{R}^n$ ,  $f, g, h : A \rightarrow \mathbb{R}$ , and  $\mathbf{a}$  a limit point of  $A$ .*

- (a) *If  $\lim_{\mathbf{p} \rightarrow \mathbf{a}} f(\mathbf{p}) = L$ , and  $\lim_{\mathbf{p} \rightarrow \mathbf{a}} g(\mathbf{p}) = M$ , then  $\lim_{\mathbf{p} \rightarrow \mathbf{a}} (f + g)(\mathbf{p}) = L + M$ . A similar property holds for  $fg$  and for  $f/g$  (the latter provided  $M \neq 0$ ).*
- (b) *Let  $U$  be a neighborhood of  $\mathbf{a}$ . If  $f(\mathbf{p}) \leq g(\mathbf{p}) \leq h(\mathbf{p})$  for all  $\mathbf{p}$  in  $U \setminus \{\mathbf{a}\}$ , and if  $\lim_{\mathbf{p} \rightarrow \mathbf{a}} f(\mathbf{p}) = \lim_{\mathbf{p} \rightarrow \mathbf{a}} h(\mathbf{p}) = L$ , then  $\lim_{\mathbf{p} \rightarrow \mathbf{a}} g(\mathbf{p}) = L$ .*

**Examples 1.9.1.** (i) A polynomial of degree  $k$  on  $\mathbb{R}^n$  is a function  $f$  of the form

$$f = a_0 + \sum_{j=1}^k \sum_{1 \leq i_1, \dots, i_j \leq n} a_{i_1 \dots i_j} u^{i_1} \cdots u^{i_j}, \quad a_0, a_{i_1 \dots i_j} \in \mathbb{R}.$$

By Theorem 1.9.2, if  $f$  is a polynomial, then  $\lim_{\mathbf{p} \rightarrow \mathbf{a}} f(\mathbf{p}) = f(\mathbf{a})$  for any  $\mathbf{a}$ .

- (ii) Let  $f : \mathbb{R}^2 \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}$  be given by  $f(x, y) = xy/(x^2 + y^2)$ . Then  $\mathbf{0}$  is a limit point of the domain of  $f$ . If  $\mathbf{a}_n = (1/n, 0)$  then  $\mathbf{a}_n \rightarrow \mathbf{0}$  and  $f(\mathbf{a}_n) = 0$ . If  $\mathbf{b}_n = (1/n, 1/n)$ , then  $\mathbf{b}_n \rightarrow \mathbf{0}$  but  $f(\mathbf{b}_n) = 1/2$  for all  $n$ . Thus,  $f$  has no limit at  $\mathbf{0}$ . The graph of  $f$  in  $\mathbb{R}^3$  intersects each plane  $ax + by = 0$  containing the  $z$ -axis in two half-lines at constant height  $-ab/(a^2 + b^2)$  above the plane  $\mathbb{R}^2 \times \{0\}$ .

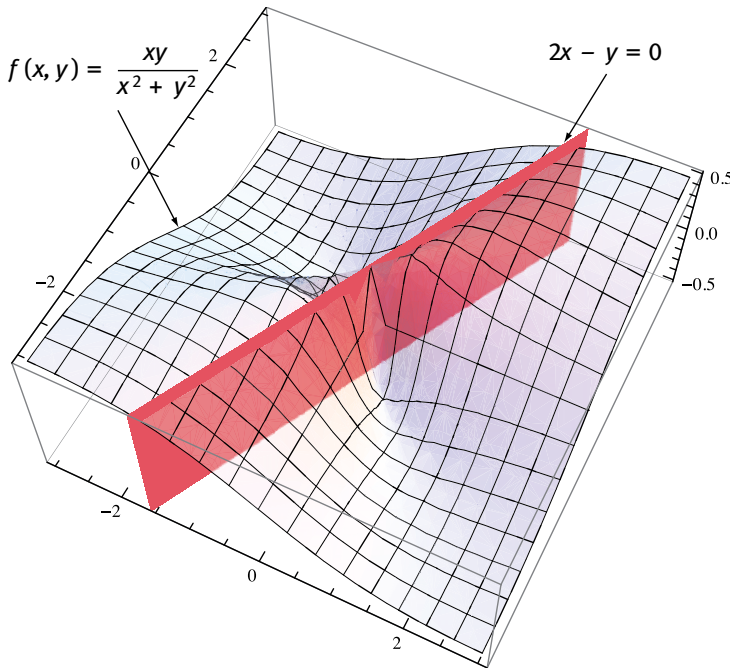


Fig. 1.2: The graph of the function from Examples 1.9.1 (ii)

Recall that the *composition*  $f \circ g$  of maps  $f$  and  $g$  with appropriate domain and range is defined by  $(f \circ g)(a) = f(g(a))$ .

**Corollary 1.9.1.** Suppose  $a$  is a limit point of  $A \subset \mathbb{R}^n$ ,  $f, g : A \rightarrow \mathbb{R}^m$ . If  $\lim_{p \rightarrow a} f(p) = u$  and  $\lim_{p \rightarrow a} g(p) = v$ , then for any  $c \in \mathbb{R}$ ,

$$\lim_{p \rightarrow a} (f + g)(p) = u + v, \quad \lim_{p \rightarrow a} c f(p) = c u, \quad \lim_{p \rightarrow a} \langle f(p), g(p) \rangle = \langle u, v \rangle.$$

Furthermore, if  $h : B \subset \mathbb{R}^m \rightarrow \mathbb{R}^l$  has limit  $r$  at  $u \in B$ , then

$$\lim_{p \rightarrow a} (h \circ f) = r.$$

*Proof.* The first statement follows from Theorem 1.9.1 and Corollary 1.8.1. The second one is immediate from the definition of limit: Let  $a_i \rightarrow a$ . Since  $\lim_{p \rightarrow a} f(p) = u$ ,  $f(a_i) \rightarrow u$ . But  $\lim_{p \rightarrow u} h(p) = r$ , so  $(h \circ f)(a_i) \rightarrow r$ , which implies the claimed limit.  $\square$

The following fundamental concept is defined here in the context of Euclidean spaces, since these remain our main focus. The reader should be aware, though, that the same definition can be, and is used in arbitrary metric spaces. Furthermore, all the results mentioned below that do not rely on the algebraic or order structures of  $\mathbb{R}^n$  hold in these metric spaces.

**Definition 1.9.2.** Let  $a \in A \subset \mathbb{R}^n$ . A map  $f : A \rightarrow \mathbb{R}^m$  is said to be *continuous at*  $a$  if for every  $\varepsilon > 0$ , there exists some  $\delta > 0$  such that

$$f(B_\delta(a) \cap A) \subset B_\varepsilon(f(a)).$$

$f$  is said to be *continuous on*  $A$  if it is continuous at every  $a \in A$ .

If  $\mathbf{a}$  is a limit point of  $A$ , then  $\mathbf{f}$  is continuous at  $\mathbf{a}$  if and only if  $\lim_{\mathbf{p} \rightarrow \mathbf{a}} \mathbf{f}(\mathbf{p}) = \mathbf{f}(\mathbf{a})$ . As an application, any polynomial is continuous on all Euclidean space by Examples 1.9.1 (i). If  $\mathbf{a}$  is not a limit point of  $A$ , then  $\mathbf{f}$  is automatically continuous at  $\mathbf{a}$ . For example, every sequence  $\mathbf{f} : \mathbb{N} \subset \mathbb{R} \rightarrow \mathbb{R}^n$  is continuous on its domain: no matter what  $\mathbf{f}(k)$  is,  $B_{1/2}(k) \cap \mathbb{N} = \{k\}$ , and so its image is contained in  $B_\varepsilon(\mathbf{f}(k))$  for any  $\varepsilon$ .

**Proposition 1.9.1.**  $\mathbf{f} : A \rightarrow \mathbb{R}^m$  is continuous at  $\mathbf{a} \in A$  if and only if  $\mathbf{f}(\mathbf{a}_k) \rightarrow \mathbf{f}(\mathbf{a})$  for every sequence  $\mathbf{a}_k \rightarrow \mathbf{a}$ .

*Proof.* If  $\mathbf{a}$  is a limit point of  $A$ , the claim follows from Theorem 1.9.1. Otherwise, we have seen that the map is necessarily continuous at  $\mathbf{a}$ ; on the other hand, if  $\mathbf{a}_k \rightarrow \mathbf{a}$ , then  $\mathbf{a}_k$  must equal  $\mathbf{a}$  for  $k$  sufficiently large, so certainly  $\mathbf{f}(\mathbf{a}_k) \rightarrow \mathbf{f}(\mathbf{a})$ .  $\square$

The above proposition, together with previous results, implies that continuity is well behaved with respect to common operations on maps.

**Theorem 1.9.3.** Suppose  $\mathbf{f}, \mathbf{g} : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathbf{h} : B \subset \mathbb{R}^m \rightarrow \mathbb{R}^l$ ,  $c \in \mathbb{R}$ . If  $\mathbf{f}$  and  $\mathbf{g}$  are continuous at  $\mathbf{a} \in A$ , then so are  $\mathbf{f} + \mathbf{g}$ ,  $c\mathbf{f}$ , and  $\langle \mathbf{f}, \mathbf{g} \rangle$ . If  $\mathbf{h}$  is continuous at  $\mathbf{f}(\mathbf{a})$ , then  $\mathbf{h} \circ \mathbf{f}$  is continuous at  $\mathbf{a}$ .

*Proof.* This follows, as noted above, from corresponding properties of sequences. Alternatively, one can use Corollary 1.9.1 when  $\mathbf{a}$  is a limit point of  $A$ . If  $\mathbf{a}$  is not a limit point, then there is nothing to prove.  $\square$

**Examples 1.9.2.** (i) A *rational function* is a quotient of two polynomials. Such a function is continuous on its domain.

(ii) Consider the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by

$$f(x, y) = \begin{cases} \frac{xy^2}{x^2+y^2}, & \text{if } (x, y) \neq \mathbf{0}, \\ 0, & \text{if } (x, y) = \mathbf{0}. \end{cases}$$

By (i),  $f$  is continuous everywhere except perhaps at the origin. But  $|x| \leq (x^2 + y^2)^{1/2}$  and  $y^2 \leq x^2 + y^2$ , so that  $0 \leq f(x, y) \leq (x^2 + y^2)^{1/2}$ . Since the term on the right goes to 0 as  $(x, y) \rightarrow \mathbf{0}$ , so does  $f$ , and  $f$  is continuous everywhere.

(iii) Since  $\det : M_{n,n} \cong \mathbb{R}^{n^2} \rightarrow \mathbb{R}$  is a polynomial, it is continuous everywhere. It follows that the *general linear group*  $GL(n)$ , which by definition is the collection of all invertible  $n \times n$  matrices, is open; this is because  $GL(n) = \det^{-1}(\mathbb{R} \setminus \{0\})$ .

Recall that for  $\mathbf{f} : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ , the *pre-image*  $\mathbf{f}^{-1}(U)$  of a set  $U \subset \mathbb{R}^m$  is the collection of all points  $\mathbf{a}$  in  $A$  such that  $\mathbf{f}(\mathbf{a}) \in U$ . One often useful characterization of continuity is the following:

**Theorem 1.9.4.** A map  $\mathbf{f} : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  is continuous if and only if the preimage  $\mathbf{f}^{-1}(U)$  of any open set  $U$  in  $\mathbb{R}^m$  is open in  $A$ .

*Proof.* Suppose  $\mathbf{f}$  is continuous, and  $U$  is open in  $\mathbb{R}^m$ . It must be shown that if  $\mathbf{a} \in \mathbf{f}^{-1}(U)$ , then there exists an open set  $V$  in  $\mathbb{R}^n$  such that  $\mathbf{a} \in V \cap A \subset \mathbf{f}^{-1}(U)$ . Since  $U$  is

open, there exists an  $\varepsilon > 0$  such that the open ball of radius  $\varepsilon$  about  $\mathbf{f}(\mathbf{a})$  is contained in  $U$ . Continuity of  $\mathbf{f}$  means that there exists  $\delta > 0$  such that  $\mathbf{f}(B_\delta(\mathbf{a}) \cap A) \subset B_\varepsilon(\mathbf{f}(\mathbf{a}))$ . In other words,  $B_\delta(\mathbf{a}) \cap A \subset \mathbf{f}^{-1}(B_\varepsilon(\mathbf{f}(\mathbf{a}))) \subset \mathbf{f}^{-1}(U)$ , so that  $\mathbf{f}^{-1}(U)$  is open in  $A$ . Conversely, suppose open sets have open pre-images in  $A$ . Given  $\mathbf{a} \in A$  and  $\varepsilon > 0$ ,  $\mathbf{f}^{-1}(B_\varepsilon(\mathbf{f}(\mathbf{a})))$  is then open in  $A$ , and thus equals  $U \cap A$  for some open set  $U$  in  $\mathbb{R}^n$ . Since  $\mathbf{a} \in U$ , there exists  $\delta > 0$  such that the ball of radius  $\delta$  around  $\mathbf{a}$  is contained inside  $U$ . But then  $B_\delta(\mathbf{a}) \cap A \subset U \cap A = \mathbf{f}^{-1}(B_\varepsilon(\mathbf{f}(\mathbf{a})))$ , and  $\mathbf{f}(B_\delta(\mathbf{a}) \cap A) \subset B_\varepsilon(\mathbf{f}(\mathbf{a}))$ . This shows that  $\mathbf{f}$  is continuous at every  $\mathbf{a} \in A$ .  $\square$

Theorem 1.9.4 provides an easy proof of the fact that continuous maps send compact sets to compact sets:

**Theorem 1.9.5.** *If  $\mathbf{f} : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  is continuous, then  $\mathbf{f}(K)$  is compact for every compact set  $K$  contained in  $A$ .*

*Proof.* Let  $\{U_\alpha\}$  be an open cover of  $\mathbf{f}(K)$ . By Theorem 1.9.4, each  $\mathbf{f}^{-1}(U_\alpha)$  equals  $V_\alpha \cap A$  for some open set  $V_\alpha$  in  $\mathbb{R}^n$ . Since  $K$  is compact, there exist finitely many indices  $\alpha_1, \dots, \alpha_k$  such that  $K \subset V_{\alpha_1} \cup \dots \cup V_{\alpha_k}$ . But then,  $\mathbf{f}(K) \subset U_{\alpha_1} \cup \dots \cup U_{\alpha_k}$ , and  $\mathbf{f}(K)$  is compact. Notice that we have used the fact, easily verified, that for any sets  $U$  and  $V$ ,  $\mathbf{f}(\mathbf{f}^{-1}(U) \cup \mathbf{f}^{-1}(V)) \subset U \cup V$ .  $\square$

Theorem 1.9.5 has the following immediate application, usually referred to as the *extreme value theorem*:

**Theorem 1.9.6.** *Let  $f : K \subset \mathbb{R}^n \rightarrow \mathbb{R}$  be continuous. If  $K$  is compact, then there exist  $\mathbf{a}, \mathbf{b} \in K$  such that  $f(\mathbf{a}) \leq f(\mathbf{p}) \leq f(\mathbf{b})$  for every  $\mathbf{p} \in K$ . (The numbers  $f(\mathbf{a}), f(\mathbf{b})$  are called the minimum and maximum values, respectively, of  $f$  on  $K$ .)*

*Proof.* Since  $f(K)$  is compact by Theorem 1.9.5, it is closed and bounded by Theorem 1.7.4. Being bounded, it has a least upper bound  $\alpha$ , and being closed, it contains  $\alpha$ , so that  $\alpha = f(\mathbf{b})$  for some  $\mathbf{b} \in K$ . This implies the second inequality in the theorem. The first one is proved in a similar way. Alternatively, one can apply the above argument to  $-f$ , since the maximum value of  $-f$  is the negative of the minimum value of  $f$ .  $\square$

There is a stronger version of continuity that plays an important role in Calculus:

**Definition 1.9.3.** A map  $\mathbf{f} : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  is said to be *uniformly continuous* if for any  $\varepsilon > 0$ , there exists a  $\delta > 0$  such that

$$\mathbf{f}(B_\delta(\mathbf{a}) \cap A) \subset B_\varepsilon(\mathbf{f}(\mathbf{a})) \quad (1.9.1)$$

for any  $\mathbf{a} \in A$ . Alternatively, for any  $\varepsilon > 0$ , there exists  $\delta > 0$  such that

$$|\mathbf{f}(\mathbf{a}) - \mathbf{f}(\mathbf{b})| < \varepsilon, \text{ whenever } |\mathbf{a} - \mathbf{b}| < \delta, \quad \mathbf{a}, \mathbf{b} \in A.$$

For  $\mathbf{f}$  to be continuous, one needs to find for each  $\varepsilon > 0$  and each  $\mathbf{a} \in A$  some  $\delta > 0$  (which usually depends on *both*  $\varepsilon$  and  $\mathbf{a}$ ) such (1.9.1) holds. When  $\mathbf{f}$  is uniformly continuous, for a given  $\varepsilon > 0$  there exists some  $\delta > 0$  which works for *all*  $\mathbf{a} \in A$ . In particular,

a uniformly continuous map is continuous at every point. The converse is not true: for example, the function  $f : (0, \infty) \rightarrow (0, \infty)$ , where  $f(x) = 1/x$ , is continuous, but not uniformly continuous. Intuitively, given a fixed  $\varepsilon$ , the corresponding  $\delta$  becomes smaller and smaller as the point  $a$  approaches 0. Specifically, let  $\varepsilon = 1$ . We claim there is no  $\delta > 0$  such that  $|1/a - 1/b| < 1$  for all  $a, b$  at distance less than  $\delta$  from each other. In fact, if  $\delta \geq 1$ , take  $a = 1/4, b = 1/2$ .  $a$  and  $b$  are less than  $\delta$  apart, but  $1/a - 1/b = 2$ . If  $\delta < 1$ , let  $a = \delta/4, b = \delta/2$ . Then  $|a - b| < \delta$ , but  $1/a - 1/b = 2/\delta > 2$ .

The obstruction in the above example is due to the fact that the domain of  $f$  is not compact:

**Theorem 1.9.7.** *If  $K \subset \mathbb{R}^n$  is compact, then any continuous map  $f : K \rightarrow \mathbb{R}^m$  is uniformly continuous.*

*Proof.* Given  $\varepsilon > 0$ , there exists, for each  $\mathbf{a} \in K$ , some  $\delta(\mathbf{a}) > 0$  such that

$$f(B_{\delta(\mathbf{a})}(\mathbf{a})) \subset B_{\varepsilon/2}(f(\mathbf{a})),$$

by continuity of  $f$ . The collection of all  $U_{\mathbf{a}} := B_{\delta(\mathbf{a})/2}(\mathbf{a})$ , as  $\mathbf{a}$  ranges over  $K$ , is an open cover of  $K$ , and therefore admits a finite subcover  $U_{\mathbf{a}_1}, \dots, U_{\mathbf{a}_k}$ . Set  $\delta := \min\{\delta(\mathbf{a}_1)/2, \dots, \delta(\mathbf{a}_k)/2\} > 0$ . We claim that if  $\mathbf{a}$  and  $\mathbf{b}$  are points in  $K$  at distance less than  $\delta$  from each other, then  $|f(\mathbf{a}) - f(\mathbf{b})| < \varepsilon$ . To see this, observe that by assumption  $\mathbf{a}$  belongs to some  $U_{\mathbf{a}_i}$ , and therefore

$$|f(\mathbf{a}) - f(\mathbf{a}_i)| < \frac{\varepsilon}{2}. \quad (1.9.2)$$

On the other hand,

$$|\mathbf{b} - \mathbf{a}_i| \leq |\mathbf{b} - \mathbf{a}| + |\mathbf{a} - \mathbf{a}_i| < \delta + \frac{\delta(\mathbf{a}_i)}{2} \leq \delta(\mathbf{a}_i),$$

so that

$$|f(\mathbf{b}) - f(\mathbf{a}_i)| < \frac{\varepsilon}{2}. \quad (1.9.3)$$

(1.9.2) and (1.9.3) together with the triangle inequality then yield the claim.  $\square$

**Examples and Remarks 1.9.1.** (i) Let  $N$  be a norm on  $\mathbb{R}^n$ , not necessarily the standard one. We claim  $N$  is uniformly continuous. To see this, observe that if  $M = \max\{N(\mathbf{e}_i) \mid 1 \leq i \leq n\}$ , then

$$\begin{aligned} N(\mathbf{a}) &= N\left(\sum_i a_i \mathbf{e}_i\right) \leq \sum_i |a_i| N(\mathbf{e}_i) \leq M \left(\sum_i |a_i|\right) \leq nM \max_i \{|a_i|\} \\ &\leq nM |\mathbf{a}|. \end{aligned}$$

By the triangle inequality,

$$|N(\mathbf{a}) - N(\mathbf{b})| \leq N(\mathbf{a} - \mathbf{b}) \leq nM |\mathbf{a} - \mathbf{b}|,$$

which establishes the claim. One consequence is the following equivalence of norms property for Euclidean space: given any norm  $N$  on  $\mathbb{R}^n$ , there exist  $\alpha, \beta > 0$

such that

$$\alpha|\mathbf{a}| \leq N(\mathbf{a}) \leq \beta|\mathbf{a}|, \quad \mathbf{a} \in \mathbb{R}^n. \quad (1.9.4)$$

The second inequality has already been proved. For the first one, we may assume that  $\mathbf{a} \neq \mathbf{0}$ . If  $\alpha$  denotes the minimum value of  $N$  on the compact unit sphere, then  $N(\mathbf{a}/|\mathbf{a}|) \geq \alpha$ , so  $N(\mathbf{a})$  is indeed no less than  $\alpha|\mathbf{a}|$ .

(1.9.4) in turn has the following important consequence: Euclidean space is complete with respect to any norm: the first inequality in (1.9.4) implies that a sequence which is Cauchy in the  $N$ -norm is Cauchy in the Euclidean one, hence converges in the Euclidean one. The second inequality then guarantees convergence in the  $N$ -norm.

- (ii) A linear transformation from a vector space to itself is also called an *operator*. Consider the space of operators on  $\mathbb{R}^n$ , which we also identify with  $M_{n,n}$  in the usual way. Let us denote the norm operator by  $N : M_{n,n} \rightarrow \mathbb{R}$  for the moment, to distinguish it from the Euclidean norm on  $M_{n,n} = \mathbb{R}^{n^2}$ . It was shown in (i) that the space of operators on  $\mathbb{R}^n$  is complete in the operator norm. We now revert to the old notation; i.e.,  $|L|$  will denote the operator norm of  $L$ . By the proof of Theorem 1.8.5, any absolutely convergent series converges. We may use this property to introduce the exponential of a linear  $L : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . This concept will be useful in the next chapter, as it provides further insight into a large class of vector fields. The reader is assumed to be familiar with the fact that for a real number  $\alpha$ ,  $e^\alpha = \sum_{k=0}^{\infty} \alpha^k/k!$ . The discussion will mostly be framed in the context of  $n \times n$  matrices in view of the isomorphism  $L_A \leftrightarrow A$  between the two spaces together with the fact that  $L_A \circ L_B = L_{AB}$  for  $A, B \in M_{n,n}$ . First of all, notice that  $|AB| \leq |A||B|$ , where  $|A|$  is defined to be  $|L_A|$ : indeed, if  $\mathbf{u} \in \mathbb{R}^n$  has unit norm, then by (1.4.1)

$$|(AB)\mathbf{u}| = |A(B\mathbf{u})| \leq |A||B\mathbf{u}| \leq |A||B|.$$

Induction now implies that  $|A^k| \leq |A|^k$ . In particular, the *exponential series*

$$\exp(A) = e^A = \sum_{k=0}^{\infty} \frac{A^k}{k!}, \quad A \in M_{n,n},$$

converges absolutely (to the real number  $e^{|A|}$ ), and therefore also converges. We call  $e^A$  the *exponential* of  $A$ . When  $n = 1$ , this is the usual exponential.

Absolute convergence also implies that

$$|e^A| \leq e^{|A|}.$$

We emphasize again that if  $L^k$  denotes the composition of the operator  $L$  with itself  $k$  times, then, as noted earlier, the series

$$\exp(L) = e^L = \sum_{k=0}^{\infty} \frac{L^k}{k!}$$

converges absolutely; if  $L$  has matrix  $A$  in the standard basis, then  $e^L$  has matrix  $e^A$ .

(iii) As a concrete example, let  $b \in \mathbb{R}$ , and

$$B = \begin{bmatrix} 0 & -b \\ b & 0 \end{bmatrix}.$$

We assume the reader is familiar with the series

$$\cos x = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}, \quad \sin x = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!},$$

which converge for any real number  $x$ . A direct computation yields  $B^2 = -b^2 I_2$ , and an easy induction implies that  $B^{2k} = (-1)^k b^{2k} I_2$ . Thus,

$$B^{2k} = \begin{bmatrix} (-1)^k b^{2k} & 0 \\ 0 & (-1)^k b^{2k} \end{bmatrix},$$

and multiplying the above expression by  $B$  yields

$$B^{2k+1} = (-1)^k b^{2k} B = \begin{bmatrix} 0 & -(-1)^k b^{2k+1} \\ (-1)^k b^{2k+1} & 0 \end{bmatrix}.$$

This means that

$$\sum_{i=0}^k \frac{B^i}{i!} = \begin{bmatrix} s_k & -t_k \\ t_k & s_k \end{bmatrix},$$

where  $s_k$  and  $t_k$  are the  $k$ -th partial sums for the above  $\sin b$  and  $\cos b$  series respectively. In other words,

$$e^B = \begin{bmatrix} \cos b & -\sin b \\ \sin b & \cos b \end{bmatrix}.$$

## 1.10 Exercises

**1.1.** Let  $a_i, i = 1, \dots, n$  be real numbers. Determine which, if any, of the following sets are subspaces of  $\mathbb{R}^n$ :

- (i)  $\{(x_1, \dots, x_n) \mid \sum_{i=1}^n a_i x_i = 0\}$ ;
- (ii)  $\{(x_1, \dots, x_n) \mid \sum_{i=1}^n a_i x_i^2 = 0\}$ ;
- (iii)  $\{(x_1, \dots, x_n) \mid \sum_{i=1}^n a_i x_i = 1\}$ ;
- (iv)  $\{(x_1, \dots, x_n) \mid (\sum_{i=1}^n a_i^2) x_1 + a_n x_2 = 0\}$ .

**1.2.** (a) Show that the collection  $P_n$  of all polynomials (with real coefficients) of degree no larger than  $n$  is a subspace of the vector space of all real-valued functions.

(b) Prove that  $P_n$  has  $\{1, x, x^2, \dots, x^n\}$  as basis, and therefore has dimension  $n + 1$ .

(c) Show that  $\{1, 1 + x, 1 + x + x^2, \dots, 1 + x + \dots + x^n\}$  is also a basis of  $P_n$ . Find the coordinate vector of  $a + bx + cx^2$  with respect to this basis.



**1.3.** Prove that the vector space of all polynomials is infinite-dimensional, and exhibit a basis for it.

**1.4.** Prove that any subset of a vector space that contains the zero vector is linearly dependent.

**1.5.** Let  $V$  denote an  $n$ -dimensional vector space, and  $E \subset V$  a subset consisting of  $n$  elements. Prove that  $E$  is linearly independent if and only if  $E$  spans  $V$ , and in this case a basis of  $V$ .

**1.6.** Determine whether the following sets form a basis of  $\mathbb{R}^3$ :

(i)  $\{(1, 0, 1), (2, 1, 0), (4, 1, 2)\}$ ;

(ii)  $\{(1, 0, 1), (2, 1, 0)\}$ ;

(iii)  $\{(1, 0, 1), (2, 1, 0), (4, 1, 2), (1, 1, 1)\}$ ;

(iv)  $\{(1, 0, 1), (2, 1, 0), (5, 1, 2)\}$ .

**1.7.** Let  $L : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  be given by

$$L \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x + 2y + 4z \\ y + z \\ x + 2z \end{bmatrix}.$$

Find bases for the kernel and for the image of  $L$ .

**1.8.** Let  $P_n$  denote the space of polynomials of degree  $\leq n$  with its standard basis  $\{1, x, \dots, x^n\}$ . Find the matrix of the derivative operator  $D : P_n \rightarrow P_n$ ,  $Dp(x) := p'(x)$ , with respect to this basis.

**1.9.** A set  $A \subset \mathbb{R}^n$  is said to be *convex* if it contains the line segment joining any two points of  $A$ ; i.e.,  $\mathbf{a} + t(\mathbf{b} - \mathbf{a}) \in A$  for all  $\mathbf{a}, \mathbf{b} \in A$  and  $t \in [0, 1]$ . Show that if  $A$  is convex and  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is linear, then  $L(A)$  is also convex.

**1.10.** Let  $V$  be a vector space with basis  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ . Show that for any  $n$  elements  $\mathbf{w}_1, \dots, \mathbf{w}_n$  in some vector space  $W$ , there exists one and only one linear transformation  $L : V \rightarrow W$  such that  $L\mathbf{v}_i = \mathbf{w}_i$  for  $i = 1, \dots, n$ . Prove that if the set  $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$  is linearly independent, then  $L$  is one-to-one, and if it is a basis, then  $L$  is an isomorphism.

**1.11.** (a) Show that if  $L : V \rightarrow W$  is a linear transformation, then the kernel of  $L$  is a subspace of  $V$  and the image of  $L$  a subspace of  $W$ .

(b) Prove that  $L$  is one-to-one if and only if its kernel consists of the zero vector only.

**1.12.** Suppose  $L : V \rightarrow W$  is a linear transformation between vector spaces of the same dimension. Use Theorem 1.2.2 to show that the following statements are equivalent:

(i)  $L$  is one-to-one;

(ii)  $L$  is onto;

(iii)  $L$  is an isomorphism.

**1.13.** A linear operator  $L$  on an inner product space  $V$  is said to be a *linear isometry* or an *orthogonal transformation* if it preserves the inner product; i.e., if  $\langle L\mathbf{v}, L\mathbf{w} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle$  for all  $\mathbf{v}, \mathbf{w} \in V$ .

- (a) Show that  $L$  is a linear isometry if and only if  $L$  preserves norms; i.e., iff  $|L\mathbf{v}| = |\mathbf{v}|$  for all  $\mathbf{v} \in V$ . *Hint:*  $\langle \mathbf{v}, \mathbf{w} \rangle = \frac{1}{2}(|\mathbf{v} + \mathbf{w}|^2 - |\mathbf{v}|^2 - |\mathbf{w}|^2)$ .
- (b) Prove that a linear isometry is an isomorphism.
- (c) Show that a linear operator on an  $n$ -dimensional inner product space  $V$  is a linear isometry if and only if its matrix  $A$  with respect to any orthonormal basis satisfies  $AA^T = I_n$ .

**1.14.** Recall that the *trace* of an  $n \times n$  matrix  $A$  is the sum  $\text{tr } A = \sum_i a_{ii}$  of its diagonal elements.

- (a) Show that  $\text{tr}(AB) = \text{tr}(BA)$ . *Hint:* Write the  $(i, i)$ -th element of  $AB$  as  $\sum_k a_{ik}b_{ki}$ . Next write the  $(k, k)$ -th element of  $BA$  as  $\sum_i b_{ki}a_{ik}$ .
- (b) Use Example 1.2.5 to show that similar matrices have the same trace. One may therefore define the trace of a linear transformation  $L : V \rightarrow V$  to be the trace of its matrix with respect to any basis of  $V$ .

**1.15.** Let  $V$  be a (finite-dimensional) inner product space, and  $A$  a nonempty subset of  $V$ . Prove that  $(A^\perp)^\perp = \text{span } A$ .

**1.16.** The collection  $\mathcal{L}(V, W)$  of all linear transformations  $L : V \rightarrow W$  is a vector space with the operations  $(L_1 + L_2)\mathbf{v} = L_1\mathbf{v} + L_2\mathbf{v}$ ,  $(aL)\mathbf{v} = aL\mathbf{v}$  for  $L_i \in \mathcal{L}(V, W)$ ,  $a \in \mathbb{R}$ ,  $\mathbf{v} \in V$ . Show that if  $V$  is  $n$ -dimensional with basis  $\mathcal{B}$  and  $W$   $m$ -dimensional with basis  $\mathcal{C}$ , then the map

$$\begin{aligned} \mathcal{L}(V, W) &\rightarrow M_{m,n}, \\ L &\mapsto [L]_{[\mathcal{B}], [\mathcal{C}]} \end{aligned}$$

which assigns to each linear transformation its matrix with respect to the given bases is an isomorphism. This shows in particular that  $\mathcal{L}(V, W)$  has dimension  $(\dim V)(\dim W)$ .

**1.17.** Find the determinant of the matrix

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{bmatrix}$$

whose elements below the diagonal are all zero.

**1.18.** Prove that adding a multiple of a column to another column (or a multiple of a row to another row) of a matrix does not change its determinant. This allows us to replace the matrix by one that has only one nonzero entry in some row or column and

expand along that row or column. Use this to compute

$$\det \begin{bmatrix} 1 & 2 & 3 \\ 1 & 3 & 4 \\ 2 & 1 & 2 \end{bmatrix}.$$

**1.19.** Prove (1.3.6). Show that the set  $U \subset \mathbb{R}^{n^2}$  of invertible matrices is open, and that the map  $U \rightarrow U$  which sends a matrix to its inverse is continuous.

**1.20.** Show that equality holds in the Cauchy-Schwarz inequality from Theorem 1.4.1 if and only if the vectors are linearly dependent.

**1.21.** Let  $V_1, V_2$  be subspaces of  $V$ . Prove that  $\dim(V_1 + V_2) = \dim(V_1) + \dim(V_2) - \dim(V_1 \cap V_2)$ . *Hint:* Start out with a basis of  $V_1 \cap V_2$ , and extend it first to a basis of  $V_1$ , then to a basis of  $V_2$ .

**1.22.** (a) Show that the interior  $A^0$  of a subset  $A$  of a metric space is equal to the union of all open sets that are contained in  $A$ . In other words,  $A^0$  is the largest open set contained in  $A$ .

(b) Prove that the closure  $\bar{A}$  of  $A$  is equal to the intersection of all closed sets that contain  $A$ ; i.e.,  $\bar{A}$  is the smallest closed set that contains  $A$ .

**1.23.** Let  $A, B \subset \mathbb{R}^n$ .

(a) Show that  $(A \cap B)^0 = A^0 \cap B^0$ , but that in general,  $(A \cup B)^0 \neq A^0 \cup B^0$ .

(b) Prove that  $\overline{A \cup B} = \bar{A} \cup \bar{B}$ , but in general  $\overline{A \cap B} \neq \bar{A} \cap \bar{B}$ .

**1.24.** Determine the interior, boundary, and closure of the following subsets of  $\mathbb{R}^2$ :

(a)  $\{(x, y) \mid y \geq x^2\}$ ;

(b)  $\mathbb{R} \times \{0\}$ ;

(c)  $\mathbb{R} \times \mathbb{Q}$ ;

(d)  $\mathbb{Q} \times \mathbb{Q}$ .

(You may use the fact that any open interval in  $\mathbb{R}$  contains a rational number).

**1.25.** (a) Show that if  $\mathbf{a}$  is a boundary point of a set  $A$ , then there exists a sequence in  $A$  that converges to  $\mathbf{a}$ .

(b) Prove that a set  $C$  is closed if and only if for any convergent sequence contained in  $C$ , the limit of the sequence belongs to  $C$ .

**1.26.** Prove that a set  $A \subset \mathbb{R}^n$  is open if and only if  $A \cap \partial A = \emptyset$ , and closed if and only if  $\partial A \subset A$ .

**1.27.** Show that a subset  $C$  of  $A$  is closed in  $A$  if and only if for any sequence in  $C$  that converges in  $A$ , the limit of the sequence belongs to  $C$ .

**1.28.** For a subset  $A$  of a metric space  $X$ , denote by  $A^c$  the complement  $X \setminus A$  of  $A$ . Prove that  $\partial A = \bar{A} \cap \overline{A^c}$ .

**1.29.** Let  $U$  be open in  $\mathbb{R}^n$ ,  $f : U \rightarrow \mathbb{R}^n$  a continuous map, and  $A$  a set whose closure lies in  $U$ .

(a) Show that  $f(\overline{A}) \subset \overline{f(A)}$ .

(b) Show that if  $f$  is a *homeomorphism* (i.e.,  $f$  is one-to-one and its inverse is also continuous), then  $f(\partial A) = \partial f(A)$ .

**1.30.** Let  $A$  be a nonempty subset of  $\mathbb{R}^n$ . Prove or disprove:

(i) Any boundary point of  $A$  is a limit point of  $A$ .

(ii) Any limit point of  $A$  is a boundary point of  $A$ .

(iii) Any boundary point of  $A$  that does not belong to  $A$  is a limit point of  $A$ .

(iv) Any limit point of  $A$  that does not belong to  $A$  is a boundary point of  $A$ .

(v) If  $A$  is closed, then it contains all its limit points.

(vi) If  $A$  contains all its limit points, then it is closed.

**1.31.** Give one or more examples of a set  $A \subset \mathbb{R}^n$  that

(i) is neither open nor closed;

(ii) is both open and closed;

(iii) has empty boundary;

(iv) has all of  $\mathbb{R}^n$  as boundary.

**1.32.** (a) Let  $f : [a, b] \rightarrow \mathbb{R}$  be a continuous function. Show that the graph  $\{(x, f(x)) \mid a \leq x \leq b\}$  is a closed subset of  $\mathbb{R}^2$ . Give an example of a noncontinuous function  $f : [a, b] \rightarrow \mathbb{R}$  whose graph is not closed.

(b) Is the set  $A = \{(x, \sin(1/x)) \mid 0 < x < 1\} \subset \mathbb{R}^2$  closed? If not, find its closure. ( $\overline{A}$  is called the *topologist's sine curve*).

**1.33.** Suppose  $\{\mathbf{a}_k\}$  is a convergent sequence in  $\mathbb{R}^n$ . Prove that any subsequence converges to the same limit.

**1.34.** Determine whether the sequence

$$\sqrt{2}, \sqrt{2 + \sqrt{2}}, \sqrt{2 + \sqrt{2 + \sqrt{2}}}, \dots$$

converges. If it does converge, determine its limit.

**1.35.** Suppose  $\{\mathbf{a}_k\}$  is a sequence in  $\mathbb{R}^n$  that satisfies

$$|\mathbf{a}_{k+1} - \mathbf{a}_k| \leq \frac{1}{2} |\mathbf{a}_k - \mathbf{a}_{k-1}|, \quad k > 1.$$

Prove that the sequence converges.

**1.36.** Suppose  $C_k \subset \mathbb{R}^n$ ,  $k = 1, 2, \dots$ , is a sequence of nonempty compact sets with  $C_k \supset C_{k+1}$  for all  $k$ .

(a) Show that  $\bigcap_{k=1}^{\infty} C_k \neq \emptyset$ .

(b) Prove, by means of an example, that the statement is in general false if the  $C_k$  are not compact.

**1.37.** The *diameter* of a bounded set  $A \subset \mathbb{R}^n$  is

$$\text{diam } A = \sup\{|\mathbf{a} - \mathbf{b}| \mid \mathbf{a}, \mathbf{b} \in A\}.$$

Prove that if  $A$  is compact, then there exist  $\mathbf{a}, \mathbf{b} \in A$  such that  $|\mathbf{a} - \mathbf{b}| = \text{diam } A$ .

**1.38.** Given nonempty subsets  $A, B$  of  $\mathbb{R}^n$ , define the *distance* between  $A$  and  $B$  to be the number  $d(A, B) = \inf\{|\mathbf{a} - \mathbf{b}| \mid \mathbf{a} \in A, \mathbf{b} \in B\}$ .

- (a) Give an example of two disjoint closed sets  $A$  and  $B$  with  $d(A, B) = 0$ .  
 (b) Prove that if  $A$  is compact and  $B$  is closed, then there exist  $\mathbf{a} \in A$  and  $\mathbf{b} \in B$  such that  $d(A, B) = |\mathbf{a} - \mathbf{b}|$ .

**1.39.** Let  $U \subset \mathbb{R}^n$ ,  $f : U \rightarrow \mathbb{R}$  be uniformly continuous. Show that if  $U$  is bounded, then so is  $f$ ; i.e.,  $|f(\mathbf{a})| \leq M$  for some  $M > 0$ .

**1.40.** Suppose  $U \subset \mathbb{R}^n$ , and  $\mathbf{f} : U \rightarrow \mathbb{R}^m$  is uniformly continuous. Show that if  $\{\mathbf{a}_k\}$  is a Cauchy sequence in  $U$ , then  $\{\mathbf{f}(\mathbf{a}_k)\}$  is also Cauchy. Show by means of an example that the conclusion is not necessarily true if  $\mathbf{f}$  is merely continuous.

**1.41.** A subset  $A \subset \mathbb{R}^n$  is said to be *complete* if any Cauchy sequence in  $A$  converges to some point in  $A$ . Prove that if  $A$  is closed, then it is complete. Give examples that show the above conclusion is not necessarily true if  $A$  is not closed.

**1.42.** Let  $U$  denote an open set in  $\mathbb{R}^2$ , and  $f : U \rightarrow \mathbb{R}$  a continuous function. Prove that  $f$  cannot be one-to-one.

**1.43.** Determine whether or not the following maps are continuous:

(i)  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , where

$$f(x, y) = \begin{cases} \frac{x^2 y^{1/3}}{x^2 + y^2} & \text{if } (x, y) \neq (0, 0); \\ 0 & \text{if } (x, y) = (0, 0). \end{cases}$$

(ii)  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , where  $f(x, y) = \max\{x, y\}$ .

(iii)  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , where  $f(\mathbf{a}) = |\mathbf{a}|$ .

(iv)  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , where

$$\mathbf{f}(\mathbf{a}) = \begin{cases} \mathbf{a}/\sqrt{|\mathbf{a}|} & \text{if } \mathbf{a} \neq \mathbf{0}; \\ \mathbf{0} & \text{if } \mathbf{a} = \mathbf{0}. \end{cases}$$

**1.44.** Show that any linear transformation  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is absolutely continuous on  $\mathbb{R}^n$ .

**1.45.** A *separation* of a metric space  $X$  is a pair  $U, V$  of nonempty disjoint open sets whose union equals  $X$ .  $X$  is said to be *connected* if there is no separation of  $X$ .

- (a) Prove that  $X$  is connected if and only if the only subsets of  $X$  that are both open and closed are  $\emptyset$  and  $X$ .  
 (b) Show that a pair  $U, V$  of nonempty subsets of  $X$  whose union equals  $X$  form a separation of  $X$  if and only if  $\overline{U} \cap V = U \cap \overline{V} = \emptyset$ .

- 1.46.** (a) With the terminology from Exercise 1.45, suppose  $U$  and  $V$  form a separation of  $X$ . Show that any connected subset of  $X$  is entirely contained in either  $U$  or  $V$ .  
 (b) Prove that if a subset of  $X$  is connected, then so is its closure.

**1.47.** (a) Let  $A \subset \mathbb{R}^n$ ,  $f : A \rightarrow \mathbb{R}^m$  a continuous map. Prove that  $f(A)$  is connected if  $A$  is.

(b)  $A \subset \mathbb{R}^n$  is said to be *path connected* if given any  $\mathbf{p}, \mathbf{q} \in A$ , there exists a continuous map  $\mathbf{c} : [a, b] \rightarrow A$  such that  $\mathbf{c}(a) = \mathbf{p}$  and  $\mathbf{c}(b) = \mathbf{q}$ . Show that if  $A$  is path connected, then it is connected.

(c) Prove that the topologist's sine curve  $E$  from Exercise 1.32 is connected but not path connected. *Hint:* for connectivity, use part (a) and Exercise 1.46 (b). For path connectivity, consider a point  $\mathbf{p}$  in  $E$  that lies on the  $y$ -axis. Show that the image of any continuous  $\mathbf{c} : [a, b] \rightarrow E$  with  $\mathbf{c}(a) = \mathbf{p}$  is contained inside the  $y$ -axis.

**1.48.** Let  $A \subset \mathbb{R}^n$  be connected,  $f : A \rightarrow \mathbb{R}^m$  a continuous map. Prove the intermediate value theorem: If  $\mathbf{a}, \mathbf{b} \in A$  and  $c$  is a number between  $f(\mathbf{a})$  and  $f(\mathbf{b})$ , then there exists some  $\mathbf{x} \in A$  such that  $f(\mathbf{x}) = c$ .

**1.49.** Suppose that  $\sum_{k=0}^{\infty} A_k = A$ ,  $\sum_{k=0}^{\infty} B_k = B$  are absolutely convergent series of  $n \times n$  matrices (or operators on  $\mathbb{R}^n$ ). Define

$$C_k = \sum_{i=0}^k A_i B_{k-i}.$$

The goal of this exercise is to establish that  $\sum_{k=0}^{\infty} C_k = AB$ .

(a) Let  $a_k, b_k$ , and  $c_k$  denote the  $k$ -th partial sums of the series  $\sum A_i, \sum B_i$ , and  $\sum C_i$  respectively. Show that  $c_{2k} - a_k b_k = s_k + t_k$ , where

$$s_k = \sum_{\substack{0 \leq i \leq k \\ k+1 \leq j \leq 2k \\ i+j \leq 2k}} A_i B_j, \quad t_k = \sum_{\substack{0 \leq i \leq k \\ k+1 \leq j \leq 2k \\ i+j \leq 2k}} A_j B_i.$$

(b) Prove that

$$|s_k| \leq \left( \sum_{i=0}^{\infty} |A_i| \right) \left( \sum_{j=k+1}^{2k} |B_j| \right),$$

and deduce that  $|s_k| \rightarrow 0$ .

(c) Prove a similar result for  $|t_k|$ , and conclude that  $\sum_{k=0}^{\infty} C_k = AB$ .

**1.50.** (a) Use Exercise 1.49 to prove if two  $n \times n$  matrices  $A$  and  $B$  commute (i.e.,  $AB = BA$ ), then  $e^{A+B} = e^A e^B = e^B e^A$ .

(b) Show that for  $A, B \in M_n$ , if  $B$  is invertible, then  $Be^A B^{-1} = e^{BAB^{-1}}$ .

**1.51.** (a) Prove that for any  $n \times n$  matrix  $A$ ,  $e^A$  is invertible, and  $(e^A)^{-1} = e^{-A}$ .

(b) Show that

$$\exp \begin{bmatrix} a & -b \\ b & a \end{bmatrix} = e^a \begin{bmatrix} \cos b & -\sin b \\ \sin b & \cos b \end{bmatrix}.$$

**1.52.** An  $n \times n$  matrix  $A$  is said to be *skew-symmetric* if  $A + A^T = \mathbf{0}$ , and *orthogonal* if  $AA^T = I_n$ , see also Exercise 1.13. Show that the exponential of a skew-symmetric matrix is orthogonal.

**1.53.** Let  $A \subset \mathbb{R}^n$ ,  $f : A \rightarrow \mathbb{R}$  a bounded function. Given  $B \subset A$ , define

$$m_B(f) = \inf\{f(\mathbf{b}) \mid \mathbf{b} \in B\}, \quad M_B(f) = \sup\{f(\mathbf{b}) \mid \mathbf{b} \in B\}.$$

Denote by  $D$  the set of points in  $A$  where  $f$  is discontinuous.

(a) Show that for any  $\varepsilon > 0$ , the set

$$D_\varepsilon = \{\mathbf{a} \in A \mid M_U(f) - m_U(f) \geq \varepsilon \text{ for any neighborhood } U \text{ of } \mathbf{a}\}$$

is contained inside  $D$ .

(b) Prove that  $D = \bigcup_{k=1}^{\infty} D_{1/k}$ .





## 2 Differentiation

In calculus of one variable, the derivative of a function  $f : A \subset \mathbb{R} \rightarrow \mathbb{R}$  at an interior point  $a$  of  $A$  is defined to be the number

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h},$$

provided the limit exists. It is geometrically interpreted as the slope of the line that best approximates the graph of  $f$  at  $(a, f(a))$ . This line, when parallel translated so that it passes through the origin is the graph of the linear transformation  $L : \mathbb{R} \rightarrow \mathbb{R}$  given by  $Lh = f'(a)h$ . The above equation can be rewritten in terms of  $L$ :

$$\lim_{h \rightarrow 0} \frac{f(a+h) - f(a) - Lh}{h} = 0.$$

If we now replace numerator and denominator by their absolute values or norms, then this expression makes sense even when  $f$  is a map between Euclidean spaces of dimension higher than one. This is the approach we will adopt.

### 2.1 The derivative

**Definition 2.1.1.** Let  $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a map.  $f$  is said to be *differentiable* at an interior point  $\mathbf{a}$  of  $A$  if there exists a linear transformation  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  such that

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|f(\mathbf{a} + \mathbf{h}) - f(\mathbf{a}) - L\mathbf{h}\|}{\|\mathbf{h}\|} = 0. \quad (2.1.1)$$

It is worth remarking that if such an  $L$  exists, it is unique: Indeed, if  $M$  is another linear transformation with the same property, then for any unit vector  $\mathbf{u} \in \mathbb{R}^n$ , and  $0 \neq t \in \mathbb{R}$ ,

$$\begin{aligned} \|(L - M)\mathbf{u}\| &= \frac{\|(L - M)(t\mathbf{u})\|}{|t\mathbf{u}|} \\ &= \frac{\|(L(t\mathbf{u}) - f(\mathbf{a} + t\mathbf{u}) + f(\mathbf{a})) + (f(\mathbf{a} + t\mathbf{u}) - f(\mathbf{a}) - M(t\mathbf{u}))\|}{|t\mathbf{u}|} \\ &\leq \frac{\|f(\mathbf{a} + t\mathbf{u}) - f(\mathbf{a}) - Lt(\mathbf{u})\|}{|t\mathbf{u}|} + \frac{\|f(\mathbf{a} + t\mathbf{u}) - f(\mathbf{a}) - M(t\mathbf{u})\|}{|t\mathbf{u}|}. \end{aligned}$$

Letting  $t \rightarrow 0$ , we see that  $\|(L - M)\mathbf{u}\|$  is less than any positive number, so that  $L\mathbf{u} = M\mathbf{u}$ . Since this is true for arbitrary unit  $\mathbf{u}$ ,  $L = M$  (recall that a linear transformation is entirely determined by the image of basis vectors).

**Definition 2.1.2.** If  $f$  is differentiable at  $\mathbf{a}$ , the *derivative*  $Df(\mathbf{a})$  of  $f$  at  $\mathbf{a}$  is defined to be the linear transformation  $L$  from (2.1.1).

**Examples 2.1.1.** (i) If  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a linear transformation, then  $L$  is differentiable at any  $\mathbf{a} \in \mathbb{R}^n$ , and  $DL(\mathbf{a}) = L$ . This is clear, since when substituting in (2.1.1), the numerator  $L(\mathbf{a} + \mathbf{h}) - L(\mathbf{a}) - L(\mathbf{h})$  is identically zero.

(ii) A curve in  $\mathbb{R}^n$  is a continuous map  $\mathbf{c} : I \rightarrow \mathbb{R}^n$ , where  $I$  is an interval in the real line. Let  $c^i = \mathbf{u}^i \circ \mathbf{c}$  denote the  $i$ -th component function of  $\mathbf{c}$ . It is a real-valued function of one variable. Now,  $\mathbf{c}$  is differentiable at an interior point  $t$  of  $I$  if and only if there exists a linear transformation  $D\mathbf{c}(t) : \mathbb{R} \rightarrow \mathbb{R}^n$  such that

$$\lim_{h \rightarrow 0} \frac{|\mathbf{c}(t+h) - \mathbf{c}(t) - D\mathbf{c}(t)h|}{h} = 0.$$

This is equivalent to

$$\lim_{h \rightarrow 0} \frac{c^i(t+h) - c^i(t) - (\mathbf{u}^i \circ D\mathbf{c}(t))h}{h} = 0, \quad i = 1, \dots, n.$$

But  $c^i$  is an ordinary function, so that  $\mathbf{c}$  is differentiable if and only if each component function is differentiable in the usual sense, and in this case, the matrix of  $D\mathbf{c}(t)$  is commonly referred to as the *velocity vector*

$$\mathbf{c}'(t) = \begin{bmatrix} c^{1'}(t) \\ \vdots \\ c^{n'}(t) \end{bmatrix}$$

at  $t$  of the curve. Its norm is called the *speed* of the curve.

(iii) Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f(\mathbf{u}) = |\mathbf{u}|^2$ , denote the square of the distance from the origin function. For any  $\mathbf{a} \in \mathbb{R}^n$ , the map  $L_{\mathbf{a}} : \mathbb{R}^n \rightarrow \mathbb{R}$ , where  $L_{\mathbf{a}}\mathbf{h} = 2\langle \mathbf{a}, \mathbf{h} \rangle$ , is linear. We claim that  $Df(\mathbf{a}) = L_{\mathbf{a}}$ . In fact,

$$\frac{1}{|\mathbf{h}|}(f(\mathbf{a} + \mathbf{h}) - f(\mathbf{a}) - L_{\mathbf{a}}\mathbf{h}) = \frac{1}{|\mathbf{h}|}(|\mathbf{a} + \mathbf{h}|^2 - |\mathbf{a}|^2 - 2\langle \mathbf{a}, \mathbf{h} \rangle) = |\mathbf{h}|$$

clearly goes to zero as  $\mathbf{h} \rightarrow \mathbf{0}$ . This establishes the claim.

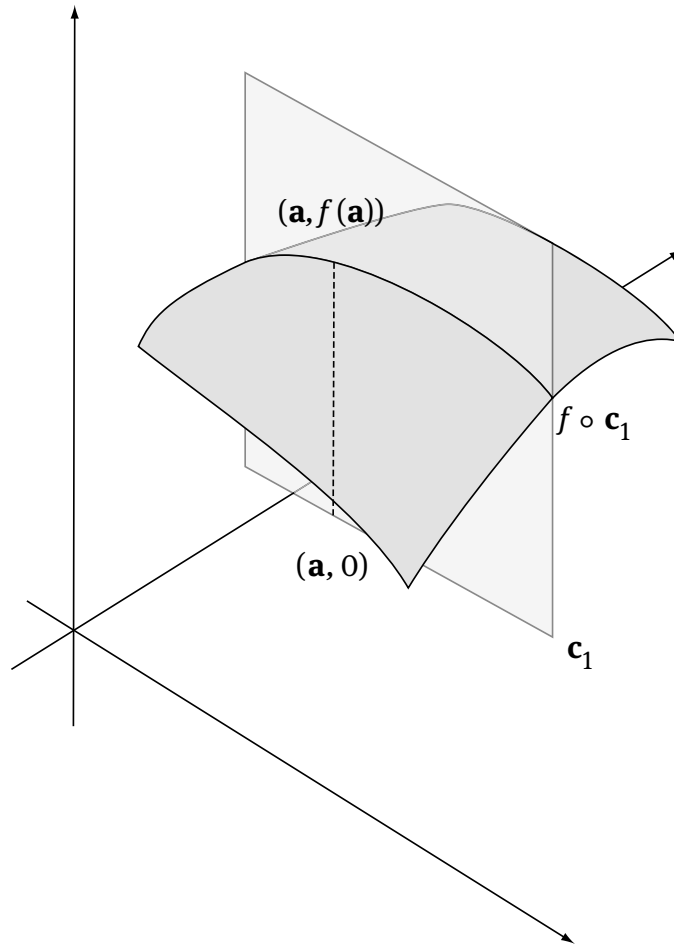
In order to see what the matrix of the linear transformation  $Df(\mathbf{a})$  looks like in general, we need some terminology:

**Definition 2.1.3.** Let  $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}$  be a real-valued function. The  $i$ -th *partial derivative*  $D_i f(\mathbf{a})$  of  $f$  at an interior point  $\mathbf{a}$  of  $A$  is defined to be the limit

$$D_i f(\mathbf{a}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{a} + t\mathbf{e}_i) - f(\mathbf{a})}{t}, \quad i = 1, \dots, n,$$

if it exists.

Notice that  $D_i f(\mathbf{a})$  is just the derivative at 0 of the real-valued function of one variable  $f \circ \mathbf{c}_i$ , where  $t \mapsto \mathbf{c}_i(t) = \mathbf{a} + t\mathbf{e}_i$  is a parametrization of the line through  $\mathbf{a}$  parallel to the  $i$ -th axis. The graph of  $f \circ \mathbf{c}_i$  is the intersection in  $\mathbb{R}^{n+1}$  of the graph of  $f$  with the the plane through  $(\mathbf{a}, 0)$  that is parallel to the  $i$ -th and the  $(n+1)$ -th coordinate axes. We will sometimes use the classical notation  $\partial f / \partial x^i$  for  $D_i f$ .



The matrix of  $Df(\mathbf{a})$  with respect to the standard bases is called the *Jacobian matrix of  $f$  at  $\mathbf{a}$* , and will be denoted  $[Df(\mathbf{a})]$ .

**Theorem 2.1.1.** Suppose  $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  is differentiable at  $\mathbf{a}$ . If  $f^i = u^i \circ f$ ,  $i = 1, \dots, m$ , denote the component functions of  $f$ , then the partial derivatives  $D_j f^i(\mathbf{a})$  exist, and the Jacobian matrix of  $f$  at  $\mathbf{a}$  is

$$[Df(\mathbf{a})] = \begin{bmatrix} D_1 f^1(\mathbf{a}) & D_2 f^1(\mathbf{a}) & \dots & D_n f^1(\mathbf{a}) \\ D_1 f^2(\mathbf{a}) & D_2 f^2(\mathbf{a}) & \dots & D_n f^2(\mathbf{a}) \\ \vdots & \vdots & \ddots & \vdots \\ D_1 f^m(\mathbf{a}) & D_2 f^m(\mathbf{a}) & \dots & D_n f^m(\mathbf{a}) \end{bmatrix}.$$

Thus, the  $i$ -th row of  $[Df(\mathbf{a})]$  is the Jacobian matrix of  $f^i$  at  $\mathbf{a}$ .

*Proof.* For  $j = 1, \dots, n$ , define a curve  $\mathbf{r}_j$  in a neighborhood of 0 by

$$\mathbf{r}_j(t) = \mathbf{f}(\mathbf{a} + t\mathbf{e}_j) - \mathbf{f}(\mathbf{a}) - Df(\mathbf{a})(t\mathbf{e}_j).$$

Then  $|\mathbf{r}_j(t)|/t \rightarrow 0$  as  $t \rightarrow 0$ , and by linearity of  $Df(\mathbf{a})$ ,

$$\frac{\mathbf{f}(\mathbf{a} + t\mathbf{e}_j) - \mathbf{f}(\mathbf{a})}{t} = Df(\mathbf{a})(\mathbf{e}_j) + \frac{\mathbf{r}_j(t)}{t}.$$

Taking the  $i$ -th component on both sides and letting  $t \rightarrow 0$ , we obtain

$$D_j f^i(\mathbf{a}) = u^i \circ D\mathbf{f}(\mathbf{a})(\mathbf{e}_j).$$

The right side of this identity is the  $(i, j)$ -th entry of  $[D\mathbf{f}(\mathbf{a})]$  by Definition 1.2.2 and Theorem 1.2.1. This establishes the claim.  $\square$

It is worth noting that the matrix of partial derivatives may exist without  $\mathbf{f}$  being differentiable: Consider the function  $f(x, y) = xy/(x^2 + y^2)$  from Example 1.9.1, and define it to equal zero at the origin. Since  $f(x, 0) \equiv 0$ ,  $D_1 f(0, 0) = 0$ , and similarly  $D_2 f(0, 0) = 0$ .  $f$ , however, is not differentiable at  $\mathbf{0}$ . This follows from the fact that it has no limit at that point (and is in particular discontinuous there), together with:

**Theorem 2.1.2.** *If  $\mathbf{f} : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  is differentiable at an interior point  $\mathbf{a}$  of  $A$ , then it is continuous at that point.*

*Proof.* Define a map  $\mathbf{r}$  from a neighborhood of  $\mathbf{0} \in \mathbb{R}^n$  to  $\mathbb{R}^m$  by

$$\mathbf{r}(\mathbf{h}) = \mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a}) - D\mathbf{f}(\mathbf{a})\mathbf{h}.$$

Then  $|\mathbf{r}(\mathbf{h})|/|\mathbf{h}| \rightarrow 0$  as  $\mathbf{h} \rightarrow \mathbf{0}$ , and

$$|\mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a})| = |D\mathbf{f}(\mathbf{a})\mathbf{h} + \mathbf{r}(\mathbf{h})| \leq |D\mathbf{f}(\mathbf{a})\mathbf{h}| + |\mathbf{r}(\mathbf{h})| \leq |D\mathbf{f}(\mathbf{a})||\mathbf{h}| + |\mathbf{r}(\mathbf{h})|.$$

Since the last term on the right approaches 0 as  $\mathbf{h} \rightarrow \mathbf{0}$ ,  $\mathbf{f}$  is continuous at  $\mathbf{a}$ .  $\square$

It turns out that if the partial derivatives exist and are continuous, then  $\mathbf{f}$  is differentiable. Before arguing this, we introduce some terminology:

**Definition 2.1.4.** (1)  $\mathbf{f} : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  is said to be *differentiable on  $A$*  if there exists an open set  $U$  containing  $A$  and a map  $\mathbf{g} : U \rightarrow \mathbb{R}^m$  that is differentiable at every point of  $U$  and whose restriction  $\mathbf{g}|_A$  to  $A$  equals  $\mathbf{f}$ .

(2) Suppose  $\mathbf{f}$  is differentiable on an open set  $U$ . If for every  $\mathbf{u} \in U$  and every  $\varepsilon > 0$ , there exists  $\delta > 0$  such that  $|D\mathbf{f}(\mathbf{p}) - D\mathbf{f}(\mathbf{u})| < \varepsilon$  whenever  $\mathbf{p} \in U$  and  $|\mathbf{p} - \mathbf{u}| < \delta$ , then  $\mathbf{f}$  is said to be *continuously differentiable* on  $U$ , and we write  $\mathbf{f} \in \mathcal{C}^1(U)$ .

For example, according to the first part of the above definition, the absolute value function  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(t) = |t|$ , is differentiable on  $[0, \infty)$ . The following fact sheds some light on the meaning of the second part:

**Lemma 2.1.1.** *Suppose  $\mathbf{f}$  is differentiable on an open set  $U$ . Then  $\mathbf{f}$  is continuously differentiable on  $U$  if and only if all partial derivatives  $D_j f^i$  are continuous on  $U$ .*

*Proof.* The statement easily follows from the string of inequalities:

$$\begin{aligned} |D_j f^i(\mathbf{p}) - D_j f^i(\mathbf{u})| &\leq |D\mathbf{f}(\mathbf{p}) - D\mathbf{f}(\mathbf{u})| \\ &\leq \sum_i \left( \sum_j (D_j f^i(\mathbf{p}) - D_j f^i(\mathbf{u}))^2 \right)^{\frac{1}{2}}. \end{aligned}$$

The first inequality holds because

$$\begin{aligned} |D_j f^i(\mathbf{p}) - D_j f^i(\mathbf{u})| &= |\langle (D\mathbf{f}(\mathbf{p}) - D\mathbf{f}(\mathbf{u}))\mathbf{e}_j, \mathbf{e}_i \rangle| \leq |(D\mathbf{f}(\mathbf{p}) - D\mathbf{f}(\mathbf{u}))\mathbf{e}_j| \\ &\leq |D\mathbf{f}(\mathbf{p}) - D\mathbf{f}(\mathbf{u})|. \end{aligned}$$

For the second inequality, recall that if  $L$  is any linear transformation, then its norm  $|L|$  satisfies  $|L| \leq \sum_i |L\mathbf{e}_i|$ , cf. the paragraph following Definition 1.4.1. Thus,

$$\begin{aligned} |D\mathbf{f}(\mathbf{p}) - D\mathbf{f}(\mathbf{u})| &\leq \sum_i |(D\mathbf{f}(\mathbf{p}) - D\mathbf{f}(\mathbf{u}))\mathbf{e}_i| \\ &= \sum_i \left( \sum_j \langle (D\mathbf{f}(\mathbf{p}) - D\mathbf{f}(\mathbf{u}))\mathbf{e}_i, \mathbf{e}_j \rangle^2 \right)^{\frac{1}{2}} \\ &= \sum_i \left( \sum_j (D_j f^i(\mathbf{p}) - D_j f^i(\mathbf{u}))^2 \right)^{\frac{1}{2}}. \quad \square \end{aligned}$$

Next, we drop the assumption of differentiability in the lemma:

**Theorem 2.1.3.** *Let  $\mathbf{f}$  map an open set  $U \subset \mathbb{R}^n$  to  $\mathbb{R}^m$ . Then  $\mathbf{f}$  is continuously differentiable on  $U$  if and only if the partial derivatives  $D_j f^i$  exist and are continuous on  $U$ .*

*Proof.* In light of Lemma 2.1.1, it only remains to show that if the partial derivatives exist and are continuous on  $U$ , then  $\mathbf{f}$  is differentiable on  $U$ . We begin by considering the case when  $\mathbf{f}$  is a real-valued function  $f$ ; i.e., when  $m = 1$ . So let  $\mathbf{a} \in U$ ,  $r > 0$  small enough that  $B_r(\mathbf{a}) \subset U$ . For  $\mathbf{h} \in \mathbb{R}^n$  with  $|\mathbf{h}| < r$ , set  $h_i = u^i(\mathbf{h})$ , so that  $\mathbf{h} = \sum h_i \mathbf{e}_i$ . Define  $\mathbf{u}_0 = \mathbf{0}$ , and  $\mathbf{u}_j = h_1 \mathbf{e}_1 + \cdots + h_j \mathbf{e}_j$  for  $j = 1, \dots, n$ . Then

$$f(\mathbf{a} + \mathbf{h}) - f(\mathbf{a}) = \sum_{i=1}^n (f(\mathbf{a} + \mathbf{u}_i) - f(\mathbf{a} + \mathbf{u}_{i-1})) = \sum_{i=1}^n (g_i(h_i) - g_i(0)),$$

where  $g_i(t) = f(\mathbf{a} + \mathbf{u}_{i-1} + t\mathbf{e}_i)$ . Notice that the line segment connecting  $\mathbf{a} + \mathbf{u}_{i-1}$  with  $\mathbf{a} + \mathbf{u}_i$  lies inside  $U$ , since its endpoints lie in  $B_r(\mathbf{a})$  and the latter ball is convex. We may therefore apply the ordinary mean-value theorem to  $g_i$ , and conclude that there exists  $t_i \in (0, h_i)$  such that

$$g_i(h_i) - g_i(0) = h_i g_i'(t_i) = h_i D_i f(\mathbf{a} + \mathbf{u}_{i-1} + t_i \mathbf{e}_i).$$

Let  $L : \mathbb{R}^n \rightarrow \mathbb{R}$  denote the linear transformation whose matrix in the standard basis has the partial derivatives of  $f$  at  $\mathbf{a}$  as entries. Then

$$\begin{aligned} \frac{|f(\mathbf{a} + \mathbf{h}) - f(\mathbf{a}) - L\mathbf{h}|}{|\mathbf{h}|} &\leq \sum_{i=1}^n \frac{|D_i f(\mathbf{a} + \mathbf{u}_{i-1} + t_i \mathbf{e}_i) - D_i f(\mathbf{a})| |h_i|}{|\mathbf{h}|} \\ &\leq \sum_i |D_i f(\mathbf{a} + \mathbf{u}_{i-1} + t_i \mathbf{e}_i) - D_i f(\mathbf{a})|. \end{aligned}$$

Since the right side of the inequality approaches 0 as  $\mathbf{h} \rightarrow \mathbf{0}$  by continuity of the partial derivatives, this proves the claim when  $m = 1$ . For general  $m$ , let  $L$  denote the linear

transformation whose matrix in the standard basis has  $D_j f^i(\mathbf{a})$  as  $(i, j)$ -th entry. Then

$$\frac{|\mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a}) - L\mathbf{h}|}{|\mathbf{h}|} = \frac{\left(\sum_{i=1}^n (f^i(\mathbf{a} + \mathbf{h}) - f^i(\mathbf{a}) - Df^i(\mathbf{a})\mathbf{h})^2\right)^{\frac{1}{2}}}{|\mathbf{h}|}$$

approaches 0 as  $\mathbf{h} \rightarrow \mathbf{0}$  by the one-dimensional case.  $\square$

**Remark 2.1.1.** Continuity of the partial derivatives of a function is not a necessary condition for the function to be differentiable. In the next section, it will be shown that sums and compositions of differentiable functions are again differentiable. Assuming this for now, consider first the function  $f : \mathbb{R} \rightarrow \mathbb{R}$  given by

$$f(t) = \begin{cases} t^2 \sin \frac{1}{t} & \text{if } t \neq 0, \\ 0 & \text{if } t = 0. \end{cases}$$

$f$  is differentiable everywhere:  $f'(0) = \lim_{h \rightarrow 0} h \sin(1/h) = 0$ , since  $0 \leq |h \sin(1/h)| \leq |h|$ , and  $f'(t) = 2t \sin(1/t) - \cos(1/t)$  if  $t \neq 0$ . However the limit of  $f'$  does not exist at 0. Consider now  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ , where  $g = f \circ u^1 + f \circ u^2$ . As a sum of compositions of differentiable functions,  $g$  is differentiable everywhere. Nevertheless,  $D_1 g = f' \circ u^1$  and  $D_2 g = f' \circ u^2$  are not continuous at the origin.

## 2.2 Basic properties of the derivative

Many of the derivation techniques for real-valued functions of one variable carry over to maps defined on Euclidean spaces. We begin with the following:

**Theorem 2.2.1.** Let  $U$  be open in  $\mathbb{R}^n$ ,  $c \in \mathbb{R}$ .

– If  $\mathbf{f}, \mathbf{g} : U \rightarrow \mathbb{R}^m$  are differentiable at  $\mathbf{a} \in U$ , then so are  $\mathbf{f} + \mathbf{g}$ ,  $c\mathbf{f}$ , and

$$D(\mathbf{f} + \mathbf{g})(\mathbf{a}) = D\mathbf{f}(\mathbf{a}) + D\mathbf{g}(\mathbf{a}), \quad D(c\mathbf{f})(\mathbf{a}) = cD\mathbf{f}(\mathbf{a}).$$

– If  $f, g : U \rightarrow \mathbb{R}$  are differentiable at  $\mathbf{a} \in U$ , then so are  $f \cdot g$ ,  $f/g$  (the latter provided  $g(\mathbf{a}) \neq 0$ ), and

$$D(f \cdot g)(\mathbf{a}) = f(\mathbf{a})Dg(\mathbf{a}) + g(\mathbf{a})Df(\mathbf{a}),$$

$$D(f/g)(\mathbf{a}) = \frac{g(\mathbf{a})Df(\mathbf{a}) - f(\mathbf{a})Dg(\mathbf{a})}{g^2(\mathbf{a})}.$$

*Proof.* Notice that if the maps are known to be continuously differentiable, then the statements follow immediately from the analogous ones for functions of one variable, since by Theorem 2.1.3, it suffices to look at partial derivatives, which, as noted before, are ordinary derivatives of functions of one variable. We will prove the first statement in each item, and leave the others to the reader.

For  $\mathbf{f} + \mathbf{g}$ , it must be shown that

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{|(\mathbf{f} + \mathbf{g})(\mathbf{a} + \mathbf{h}) - (\mathbf{f} + \mathbf{g})(\mathbf{a}) - (D(\mathbf{f}(\mathbf{a})) + D(\mathbf{g}(\mathbf{a})))\mathbf{h}|}{|\mathbf{h}|} = 0.$$

But this follows from the fact that the expression inside the limit is nonnegative and bounded above by

$$\frac{|\mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a}) - D\mathbf{f}(\mathbf{a})\mathbf{h}|}{|\mathbf{h}|} + \frac{|\mathbf{g}(\mathbf{a} + \mathbf{h}) - \mathbf{g}(\mathbf{a}) - D\mathbf{g}(\mathbf{a})\mathbf{h}|}{|\mathbf{h}|},$$

which goes to 0 as  $\mathbf{h} \rightarrow \mathbf{0}$ .

For  $\mathbf{f} \cdot \mathbf{g}$ , we must show that

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{|f\mathbf{g}(\mathbf{a} + \mathbf{h}) - f\mathbf{g}(\mathbf{a}) - (f(\mathbf{a})D\mathbf{g}(\mathbf{a}) + g(\mathbf{a})D\mathbf{f}(\mathbf{a}))\mathbf{h}|}{|\mathbf{h}|} = 0. \quad (2.2.1)$$

Let

$$r_f(\mathbf{h}) = f(\mathbf{a} + \mathbf{h}) - f(\mathbf{a}) - D\mathbf{f}(\mathbf{a})\mathbf{h}, \quad r_g(\mathbf{h}) = g(\mathbf{a} + \mathbf{h}) - g(\mathbf{a}) - D\mathbf{g}(\mathbf{a})\mathbf{h}.$$

By assumption,  $|r_f(\mathbf{h})|/|\mathbf{h}|, |r_g(\mathbf{h})|/|\mathbf{h}| \rightarrow 0$  as  $\mathbf{h} \rightarrow \mathbf{0}$ . The expression inside the limit in (2.2.1) is nonnegative and can be rewritten

$$\frac{1}{|\mathbf{h}|} |f(\mathbf{a} + \mathbf{h})r_g(\mathbf{h}) + (f(\mathbf{a} + \mathbf{h}) - f(\mathbf{a}))D\mathbf{g}(\mathbf{a})\mathbf{h} + g(\mathbf{a})r_f(\mathbf{h})|,$$

which is no larger than

$$|f(\mathbf{a} + \mathbf{h})| \frac{|r_g(\mathbf{h})|}{|\mathbf{h}|} + |f(\mathbf{a} + \mathbf{h}) - f(\mathbf{a})| |D\mathbf{g}(\mathbf{a})| + |g(\mathbf{a})| \frac{|r_f(\mathbf{h})|}{|\mathbf{h}|}.$$

The first and last term in this sum go to 0 as  $\mathbf{h} \rightarrow \mathbf{0}$  as noted earlier, and the middle term also by continuity of  $f$  at  $\mathbf{a}$  (notice that we also need continuity at  $\mathbf{a}$  to ensure that the first term goes to 0). This establishes (2.2.1).  $\square$

**Theorem 2.2.2** (The chain rule). *Let  $\mathbf{f} : U \rightarrow \mathbb{R}^m$ , where  $U$  is an open subset of  $\mathbb{R}^n$ , be differentiable at  $\mathbf{a} \in U$ . Let  $V$  be a neighborhood of  $\mathbf{f}(\mathbf{a})$  in  $\mathbb{R}^m$ , and suppose  $\mathbf{g} : V \rightarrow \mathbb{R}^k$  is differentiable at  $\mathbf{f}(\mathbf{a})$ . Then  $\mathbf{g} \circ \mathbf{f}$  is differentiable at  $\mathbf{a}$ , and*

$$D(\mathbf{g} \circ \mathbf{f})(\mathbf{a}) = D\mathbf{g}(\mathbf{f}(\mathbf{a})) \circ D\mathbf{f}(\mathbf{a}).$$

*Proof.* Set  $T = D\mathbf{f}(\mathbf{a})$ ,  $S = D\mathbf{g}(\mathbf{f}(\mathbf{a}))$ . Begin by rewriting

$$\mathbf{g}(\mathbf{f}(\mathbf{a} + \mathbf{h})) - \mathbf{g}(\mathbf{f}(\mathbf{a})) - (S \circ T)\mathbf{h}$$

as

$$S(\mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a}) - T\mathbf{h}) + \mathbf{g}(\mathbf{f}(\mathbf{a}) + \mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a})) - \mathbf{g}(\mathbf{f}(\mathbf{a})) - S(\mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a})), \quad (2.2.2)$$

and let

$$\begin{aligned} r_f(\mathbf{h}) &= \mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a}) - T\mathbf{h}, \\ r_g(\mathbf{h}) &= \mathbf{g}(\mathbf{f}(\mathbf{a}) + \mathbf{h}) - \mathbf{g}(\mathbf{f}(\mathbf{a})) - S\mathbf{h}, \\ \rho(\mathbf{h}) &= \frac{|r_g(\mathbf{h})|}{|\mathbf{h}|} \text{ if } \mathbf{h} \neq \mathbf{0}, \quad \rho(\mathbf{0}) = 0. \end{aligned}$$

By assumption,  $|r_f(\mathbf{h})|/|\mathbf{h}|, \rho(\mathbf{h}) \rightarrow 0$  as  $\mathbf{h} \rightarrow \mathbf{0}$ . Finally, let

$$\mathbf{k}(\mathbf{h}) = \mathbf{f}(\mathbf{a} + \mathbf{h}) - \mathbf{f}(\mathbf{a}) = r_f(\mathbf{h}) + T\mathbf{h},$$

which also goes to 0 as  $\mathbf{h} \rightarrow \mathbf{0}$ . Now, by (2.2.2),

$$\begin{aligned} \frac{|\mathbf{g}(\mathbf{f}(\mathbf{a} + \mathbf{h})) - \mathbf{g}(\mathbf{f}(\mathbf{a})) - ST\mathbf{h}|}{|\mathbf{h}|} &= \frac{|S(r_f(\mathbf{h})) + r_g(\mathbf{k}(\mathbf{h}))|}{|\mathbf{h}|} \\ &\leq \frac{|S(r_f(\mathbf{h}))|}{|\mathbf{h}|} + \frac{|r_g(\mathbf{k}(\mathbf{h}))|}{|\mathbf{h}|}, \end{aligned}$$

so it suffices to check that each of the last two terms goes to 0 as  $\mathbf{h} \rightarrow \mathbf{0}$ . But

$$\frac{|S(r_f(\mathbf{h}))|}{|\mathbf{h}|} \leq |S| \frac{|r_f(\mathbf{h})|}{|\mathbf{h}|} \rightarrow 0 \text{ as } \mathbf{h} \rightarrow \mathbf{0},$$

whereas

$$\begin{aligned} \frac{|r_g(\mathbf{k}(\mathbf{h}))|}{|\mathbf{h}|} &= \rho(\mathbf{k}(\mathbf{h})) \frac{|\mathbf{k}(\mathbf{h})|}{|\mathbf{h}|} = \rho(\mathbf{k}(\mathbf{h})) \frac{|r_f(\mathbf{h}) + T\mathbf{h}|}{|\mathbf{h}|} \\ &\leq \rho(\mathbf{k}(\mathbf{h})) \left( \frac{|r_f(\mathbf{h})|}{|\mathbf{h}|} + |T| \right) \rightarrow 0 \text{ as } \mathbf{h} \rightarrow \mathbf{0}. \quad \square \end{aligned}$$

Recall from Examples 2.1.1 (i) that the derivative of a linear transformation is the transformation itself. We now generalize this result.

**Definition 2.2.1.** Let  $V, V_i$  be vector spaces,  $1 \leq i \leq k$ . A map  $M : V_1 \times \cdots \times V_k \rightarrow V$  is said to be *multilinear* if for every  $i = 1, \dots, k$ , and any choice of  $\mathbf{v}_j \in V_j, j \neq i$ , the map

$$\begin{aligned} V_i &\rightarrow V, \\ \mathbf{v} &\mapsto M(\mathbf{v}_1, \dots, \mathbf{v}_{i-1}, \mathbf{v}, \mathbf{v}_{i+1}, \dots, \mathbf{v}_k) \end{aligned}$$

is linear.

Thus, the map obtained by fixing all but one arbitrary variable is linear.

**Theorem 2.2.3.** *If  $M : \mathbb{R}^{n_1} \times \cdots \times \mathbb{R}^{n_k} \rightarrow \mathbb{R}^m$  is multilinear, then  $M$  is continuously differentiable on its domain, and*

$$DM(\mathbf{a}_1, \dots, \mathbf{a}_k)(\mathbf{b}_1, \dots, \mathbf{b}_k) = \sum_{i=1}^k M(\mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \mathbf{b}_i, \mathbf{a}_{i+1}, \dots, \mathbf{a}_k).$$



*Proof.* For each  $l = 1, \dots, k$ , define

$$\begin{aligned} \iota_l : \mathbb{R}^{n_l} &\rightarrow \mathbb{R}^{n_1} \times \dots \times \mathbb{R}^{n_k}, \\ \mathbf{v} &\mapsto (\mathbf{a}_1, \dots, \mathbf{a}_{l-1}, \mathbf{v}, \mathbf{a}_{l+1}, \dots, \mathbf{a}_k). \end{aligned}$$

We may assume without loss of generality that  $m = 1$ . If  $i \leq n_l$ , then the  $(n_1 + \dots + n_{l-1} + i)$ -th partial derivative of  $M$  at  $(\mathbf{a}_1, \dots, \mathbf{a}_k)$  equals the  $i$ -th partial derivative of  $M \circ \iota_l$  at  $\mathbf{a}_l$ . But  $M \circ \iota_l$  is linear, so this derivative equals the  $i$ -th entry of (the matrix of)  $M \circ \iota_l$ . In particular, it is continuous, and  $M$  is continuously differentiable by Theorem 2.1.3. Furthermore,

$$DM(\mathbf{a}_1, \dots, \mathbf{a}_k) = (M \circ \iota_1, \dots, M \circ \iota_k),$$

and

$$\begin{aligned} DM(\mathbf{a}_1, \dots, \mathbf{a}_k)(\mathbf{b}_1, \dots, \mathbf{b}_k) &= \sum_{l=1}^k (M \circ \iota_l) \mathbf{b}_l \\ &= \sum_{l=1}^k M(\mathbf{a}_1, \dots, \mathbf{a}_{l-1}, \mathbf{b}_l, \mathbf{a}_{l+1}, \dots, \mathbf{a}_k). \quad \square \end{aligned}$$

**Corollary 2.2.1.** *The derivative of the inner product map  $\langle, \rangle : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  is given by*

$$D\langle, \rangle(\mathbf{a}_1, \mathbf{a}_2)(\mathbf{b}_1, \mathbf{b}_2) = \langle \mathbf{a}_1, \mathbf{b}_2 \rangle + \langle \mathbf{b}_1, \mathbf{a}_2 \rangle.$$

*In particular, if  $\mathbf{c}_i : I \rightarrow \mathbb{R}^n$ ,  $i = 1, 2$ , are differentiable curves defined on a common interval  $I$ , then the real-valued function  $\langle \mathbf{c}_1, \mathbf{c}_2 \rangle$  has derivative*

$$\langle \mathbf{c}_1, \mathbf{c}_2 \rangle'(t) = \langle \mathbf{c}'_1(t), \mathbf{c}_2(t) \rangle + \langle \mathbf{c}_1(t), \mathbf{c}'_2(t) \rangle.$$

*Proof.* The first identity is an immediate consequence of Theorem 2.2.3. The second one follows from the first together with the chain rule: if  $\mathbf{c}$  is the curve  $(\mathbf{c}_1, \mathbf{c}_2)$  in  $\mathbb{R}^{2n}$ , then  $\langle \mathbf{c}_1, \mathbf{c}_2 \rangle = \langle, \rangle \circ \mathbf{c}$ , so that

$$\langle \mathbf{c}_1, \mathbf{c}_2 \rangle'(t) = D\langle, \rangle(\mathbf{c}_1(t), \mathbf{c}_2(t))(\mathbf{c}'_1(t), \mathbf{c}'_2(t)) = \langle \mathbf{c}'_1(t), \mathbf{c}_2(t) \rangle + \langle \mathbf{c}_1(t), \mathbf{c}'_2(t) \rangle. \quad \square$$

If  $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable on an open set  $U$ , its partial derivatives  $D_j f$  are again real-valued functions on  $U$ . When the  $D_j f$  are themselves differentiable, the *second partial derivatives* of  $f$  are defined by

$$D_{ij} f = D_j(D_i f), \quad 1 \leq i, j \leq n.$$

This process can of course be iterated. We say  $f$  is of class  $C^k$ , and write  $f \in C^k(U)$ , if all partial derivatives of order  $k$  of  $f$  exist and are continuous on  $U$ . If this holds for any  $k$ , we say  $f$  is *smooth*, and write  $f \in C^\infty(U)$ . Our next theorem, which is usually interpreted as saying that if  $f$  is of class  $C^2$ , then  $D_{ij} f = D_{ji} f$ , can be stated somewhat more generally:

**Theorem 2.2.4.** Suppose  $f$  is a real-valued function defined on an open set  $U \subset \mathbb{R}^n$ . If  $f$ ,  $D_i f$ ,  $D_j f$ , and  $D_{ij} f$  exist and are continuous on  $U$ , then  $D_{ji} f$  exists on  $U$ , and equals  $D_{ij} f$ .

*Proof.* We may, without loss of generality, assume that  $n = 2$ : for if  $\mathbf{a} = (a_1, \dots, a_n) \in U$  and  $i < j$ , define

$$\begin{aligned} \iota : \mathbb{R}^2 &\rightarrow \mathbb{R}^n, \\ (x, y) &\mapsto (a_1, \dots, a_{i-1}, x, a_{i+1}, \dots, a_{j-1}, y, a_{j+1}, \dots, a_n). \end{aligned}$$

Then  $D_i f(\mathbf{a}) = D_1(f \circ \iota)(a_i, a_j)$ , and  $D_j(f \circ \iota)(\mathbf{a}) = D_2 f(a_i, a_j)$ . So let  $(a, b) \in U$ , and define a real-valued function  $\varphi$  on a neighborhood of  $\mathbf{0} \in \mathbb{R}^2$  by

$$\varphi(s, t) = [f(a + s, b + t) - f(a + s, b)] - [f(a, b + t) - f(a, b)].$$

Notice that if  $g(x) = f(x, b + t) - f(x, b)$ , then  $\varphi(s, t) = g(a + s) - g(a)$ . Applying the mean value theorem to  $g$  on  $[a, a + s]$ , there is an  $\varepsilon_1 \in (0, 1)$  such that

$$\varphi(s, t) = sg'(a + \varepsilon_1 s) = s[D_1 f(a + \varepsilon_1 s, b + t) - D_1 f(a + \varepsilon_1 s, b)].$$

Next, apply the mean value theorem to  $y \mapsto D_1 f(a + \varepsilon_1 s, y)$  on  $[b, b + t]$  to conclude that

$$\varphi(s, t) = stD_{12} f(a + \varepsilon_1 s, b + \varepsilon_2 t)$$

for some  $0 < \varepsilon_1, \varepsilon_2 < 1$ . Thus, if  $t \neq 0$ ,

$$\frac{f(a + s, b + t) - f(a + s, b)}{t} - \frac{f(a, b + t) - f(a, b)}{t} = sD_{12} f(a + \varepsilon_1 s, b + \varepsilon_2 t).$$

Taking the limit as  $t \rightarrow 0$  and using continuity of  $D_{12} f$ , we obtain

$$D_2 f(a + s, b) - D_2 f(a, b) = sD_{12} f(a + \varepsilon_1 s, b).$$

Dividing this equation by  $s$  and letting  $s \rightarrow 0$  yields  $D_{21} f(a, b) = D_{12} f(a, b)$ , again by continuity of  $D_{12} f$ . This establishes the claim.  $\square$

We end this section with the construction of a “bump” function that will be needed in later chapters. Given  $r > 0$ , denote by  $C_r$  the open “cube”  $(-r, r)^n$  in  $\mathbb{R}^n$ .

**Lemma 2.2.1.** For any  $0 < r < R$ , there exists a differentiable function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$  with the following properties:

- (1)  $\varphi \equiv 1$  on  $\bar{C}_r$ ;
- (2)  $0 < \varphi < 1$  on  $C_R - \bar{C}_r$ , and
- (3)  $\varphi \equiv 0$  on  $\mathbb{R}^n \setminus C_R$ .

*Proof.* Let  $h : \mathbb{R} \rightarrow \mathbb{R}$  be given by

$$h(x) = \begin{cases} e^{-1/x}, & \text{if } x > 0, \\ 0, & \text{otherwise.} \end{cases}$$

$h$  is  $C^\infty$  by l'Hospital's rule. Define

$$f(x) = \frac{h(R+x)h(R-x)}{h(R+x)h(R-x) + h(x-r) + h(-x-r)}.$$

This expression makes sense because  $h(x-r) + h(-x-r)$  is nonnegative, and equals 0 only when  $|x| \leq r$ , in which case  $h(R+x)h(R-x) > 0$ . Furthermore,  $f(x) = 1$  if  $|x| \leq r$ ,  $0 < f(x) < 1$  if  $r < |x| < R$ , and  $f(x) = 0$  if  $|x| \geq R$ . Now let  $\varphi(a_1, \dots, a_n) = \prod_{i=1}^n f(a_i)$ .  $\square$

The *support*  $\text{supp } \varphi$  of a function  $\varphi : U \subset \mathbb{R}^n \rightarrow \mathbb{R}$  is defined to be the closure of the set  $\{\mathbf{p} \in U \mid \varphi(\mathbf{p}) \neq 0\}$ . The function  $\varphi$  from the Lemma has its support in the closure of  $C_R$ , but since  $R$  is arbitrary, it is also true that there exists a function satisfying (1) and (2), but with support in  $C_R$  itself.

**Theorem 2.2.5.** *Given any open set  $U$  in  $\mathbb{R}^n$  and any compact subset  $K$  of  $U$ , there exists a differentiable function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$  such that*

- (1)  $0 \leq \varphi \leq 1$ ;
- (2)  $\varphi \equiv 1$  on  $K$ , and
- (3)  $\text{supp } \varphi \subset U$ .

*Proof.* For each  $\mathbf{a} \in K$ , choose  $R(\mathbf{a}) > 0$  such that the cube  $C_R(\mathbf{a}) := \prod_{i=1}^n (a_i - R, a_i + R)$  of radius  $R$  centered at  $\mathbf{a}$  is contained in  $U$ . Since the collection of all cubes of the form  $C_{R/2}(\mathbf{a}_i)$ , with  $\mathbf{a}_i \in K$ , covers  $K$ , there exist  $\mathbf{a}_1, \dots, \mathbf{a}_k \in K$  and  $R_1, \dots, R_k > 0$  such that

$$K \subset \bigcup_{i=1}^k C_{R_i/2}(\mathbf{a}_i).$$

Lemma 2.2.1 guarantees for any  $R > 0$  the existence of a differentiable function  $\varphi_R$  on  $\mathbb{R}^n$  that equals 1 on  $C_{R/2}$ , has support in  $C_R$ , and takes values between 0 and 1. For each  $i$  between 1 and  $k$ , define a function  $\varphi_i$  by

$$\varphi_i(\mathbf{p}) = \varphi_{R_i}(\mathbf{p} - \mathbf{a}_i).$$

Then  $\varphi_i$  is identically 1 on  $C_{R_i/2}(\mathbf{a}_i)$ , has support in  $C_{R_i}(\mathbf{a}_i) \subset U$ , and takes values between 0 and 1. Thus,

$$\varphi := \frac{1}{k} \sum_{i=1}^k \varphi_i$$

satisfies the conditions of the theorem.  $\square$

We will later see that the compactness assumption on  $K$  may be dropped.

## 2.3 Differentiation of integrals

The sole purpose of this section is to establish a result that will be needed in a later chapter: suppose  $f : [a, b] \times [c, d] \rightarrow \mathbb{R}$  is a continuous function. Integrating  $f$  with respect to the first variable  $x$  over  $[a, b]$  defines a function  $\varphi$  of  $y$  only,  $\varphi(y) = \int_a^b f(x, y) dx$ ,

$c \leq y \leq d$ . We wish to determine when  $\varphi$  is differentiable, and what its derivative  $\frac{d}{dy} \int_a^b f(x, y) dx$  is. It turns out that, loosely speaking,

$$\frac{d}{dy} \left( \int_a^b f(x, y) dx \right) = \int_a^b \frac{\partial f}{\partial y}(x, y) dx$$

whenever the partial derivative  $\partial f / \partial y = D_2 f$  is continuous. The latter condition is actually sufficient but not necessary. Nevertheless, we will adopt it since it is not too restrictive for our purposes.

**Theorem 2.3.1.** *If  $f$  and  $D_2 f$  are continuous on  $[a, b] \times [c, d]$ , then the function*

$$\begin{aligned} \varphi : [c, d] &\rightarrow \mathbb{R}, \\ y &\mapsto \int_a^b f(x, y) dx \end{aligned}$$

is differentiable on  $(c, d)$ , and

$$\varphi'(y_0) = \int_a^b D_2 f(x, y_0) dx, \quad y_0 \in (c, d).$$

*Proof.* Fix  $y_0 \in (c, d)$ , and set

$$g(x, y) = \frac{f(x, y) - f(x, y_0)}{y - y_0}.$$

Then

$$\frac{\varphi(y) - \varphi(y_0)}{y - y_0} = \int_a^b g(x, y) dx,$$

so that the theorem will follow once we establish that

$$\lim_{y \rightarrow y_0} \int_a^b g(x, y) dx = \int_a^b D_2 f(x, y_0) dx. \quad (2.3.1)$$

To prove this, let  $\varepsilon > 0$ . Using the fact that  $D_2 f$  is continuous, for each  $x_0 \in [a, b]$  there exists a  $\delta_{x_0} > 0$  such that

$$|D_2 f(x_0, y) - D_2 f(x_0, y_0)| < \frac{1}{2} \frac{\varepsilon}{b - a}, \quad \text{provided } |y - y_0| < \delta_{x_0}.$$

Again by continuity, there exists, for each  $x_0 \in [a, b]$ , a neighborhood  $U_{x_0}$  of  $x_0$  such that

$$|D_2 f(x, y) - D_2 f(x, y_0)| < \frac{\varepsilon}{b - a}, \quad x \in U_{x_0}, \quad |y - y_0| < \delta_{x_0}.$$

Next, choose some finite subcover  $\{U_{x_1}, \dots, U_{x_k}\}$  of the cover  $\{U_x \mid x \in [a, b]\}$  of  $[a, b]$ , and set  $\delta = \min\{\delta_{x_1}, \dots, \delta_{x_k}\}$ . It follows that

$$|D_2 f(x, y) - D_2 f(x, y_0)| < \frac{\varepsilon}{b-a}, \quad x \in [a, b], \quad |y - y_0| < \delta. \quad (2.3.2)$$

By the mean value theorem, for any  $(x, y) \in [a, b] \times [c, d]$ , there is some  $s$  between  $y$  and  $y_0$  such that  $g(x, y) = D_2 f(x, s)$ . (2.3.2) then implies

$$|g(x, y) - D_2 f(x, y_0)| < \frac{\varepsilon}{b-a}, \quad x \in [a, b], \quad |y - y_0| < \delta,$$

and therefore

$$\left| \int_a^b g(x, y) \, dx - \int_a^b D_2 f(x, y_0) \, dx \right| < \varepsilon$$

whenever  $|y - y_0| < \delta$ . This proves (2.3.1).  $\square$

## 2.4 Curves

In Examples 2.1.1 (ii), we introduced curves as continuous maps from an interval to Euclidean space. In geometry, however, many concepts require smoothness. We will therefore restrict ourselves to differentiable curves, or more precisely to the following somewhat larger class, that we still simply call curves for convenience:

**Definition 2.4.1.** A map  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^n$  is said to be a *curve* if it is piecewise smooth; i.e., if there exists a partition  $P : t_0 = a < t_1 < \dots < t_k = b$  of  $[a, b]$  such that the restriction of  $\mathbf{c}$  to each subinterval  $(t_{i-1}, t_i)$  is smooth.

If  $c^i = u^i \circ \mathbf{c}$ , so that  $\mathbf{c} = [c^1 \ \dots \ c^n]^T$ , we define

$$\int_a^b \mathbf{c} := \left[ \int_a^b c^1 \ \dots \ \int_a^b c^n \right]^T.$$

Thus, by the fundamental theorem of Calculus,

$$\int_a^b \mathbf{c}' = \sum_{i=1}^n \int_{t_{i-1}}^{t_i} \mathbf{c}' = \sum_{i=1}^n \mathbf{c}(t_i) - \mathbf{c}(t_{i-1}) = \mathbf{c}(b) - \mathbf{c}(a).$$

**Lemma 2.4.1.** If  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^n$  is a curve, then

$$\left| \int_a^b \mathbf{c}' \right| \leq \int_a^b |\mathbf{c}'|.$$

*Proof.* Define  $\alpha_i = \int_a^b c^{i'}$ ,  $i = 1, \dots, n$ . Then

$$\left| \int_a^b \mathbf{c}' \right|^2 = \sum_i \left( \int_a^b c^{i'} \right)^2 = \sum \alpha_i \int_a^b c^{i'} = \int_a^b \sum \alpha_i c^{i'}$$

By the Cauchy-Schwarz inequality,

$$\sum \alpha_i c^{i'} \leq (\sum \alpha_i^2)^{1/2} (\sum c^{i'2})^{1/2} = \left( \left| \int_a^b \mathbf{c}' \right| \right) |\mathbf{c}'|$$

Thus,

$$\left| \int_a^b \mathbf{c}' \right|^2 \leq \left| \int_a^b \mathbf{c}' \right| \cdot \int_a^b |\mathbf{c}'|,$$

and the claim now follows. □

Our next goal is to define the length of a piecewise smooth curve  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^n$ . To this end, associate to each partition (not necessarily the one from Definition 2.4.1)  $P : t_0 = a < t_1 < \dots < t_k = b$  of  $[a, b]$  the number

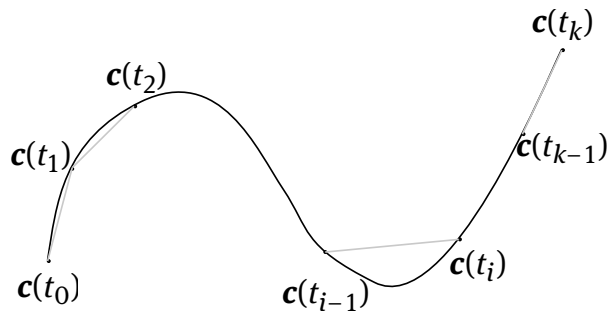
$$\ell(P, \mathbf{c}) = \sum_{i=1}^k |\mathbf{c}(t_i) - \mathbf{c}(t_{i-1})|,$$

which represents the sum of the distances between consecutive points  $\mathbf{c}(t_0), \mathbf{c}(t_1), \dots, \mathbf{c}(t_k)$ . This number should be no larger than the length of the curve, and approaches it as the partition becomes finer, thereby motivating the following:

**Definition 2.4.2.** The *length* of a curve  $c : [a, b] \rightarrow \mathbb{R}^n$  is

$$\ell(\mathbf{c}) = \sup\{\ell(P, \mathbf{c}) \mid P \text{ is a partition of } [a, b]\},$$

provided this number exists, and  $\infty$  otherwise.



**Theorem 2.4.1.** Any curve  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^n$  has finite length  $\ell(\mathbf{c}) = \int_a^b |\mathbf{c}'|$ .

*Proof.* Assume first that  $\mathbf{c}$  is smooth. The result will follow once we establish that (1)  $\ell(P, \mathbf{c}) \leq \int_a^b |\mathbf{c}'|$  for any partition  $P$  of  $[a, b]$ , and (2) given any  $\varepsilon > 0$ , there exists a partition  $P$  such that  $\int_a^b |\mathbf{c}'| \leq \ell(P, \mathbf{c}) + \varepsilon$ .

For the first assertion, consider a partition  $P : a = t_0 < t_1 < \cdots < t_k = b$ . Then by Lemma 2.4.1,

$$\ell(P, \mathbf{c}) = \sum_i |\mathbf{c}(t_i) - \mathbf{c}(t_{i-1})| = \sum_i \left| \int_{t_{i-1}}^{t_i} \mathbf{c}' \right| \leq \sum_i \int_{t_{i-1}}^{t_i} |\mathbf{c}'| = \int_a^b |\mathbf{c}'|.$$

For the second assertion, observe that  $\mathbf{c}'$  is uniformly continuous by compactness of  $[a, b]$ . Given  $\varepsilon > 0$ , choose  $\delta > 0$  so that  $|\mathbf{c}'(x) - \mathbf{c}'(y)| < \varepsilon/2(b-a)$  whenever  $|x - y| < \delta$ , and consider any partition  $P : a = t_0 < t_1 < \cdots < t_k = b$  with  $t_i - t_{i-1} < \delta$  for all  $i$ . Then

$$|\mathbf{c}'(t) - \mathbf{c}'(t_i)| < \frac{\varepsilon}{2(b-a)} \text{ and } |\mathbf{c}'(t)| < |\mathbf{c}'(t_i)| + \frac{\varepsilon}{2(b-a)}$$

for all  $t \in [t_{i-1}, t_i]$ . It follows that

$$\begin{aligned} \int_{t_{i-1}}^{t_i} |\mathbf{c}'| &< \left( |\mathbf{c}'(t_i)| + \frac{\varepsilon}{2(b-a)} \right) (t_i - t_{i-1}) \\ &= \left| \int_{t_{i-1}}^{t_i} \mathbf{c}'(t_i) - \mathbf{c}' + \mathbf{c}' \right| + \frac{\varepsilon}{2(b-a)} (t_i - t_{i-1}) \\ &\leq \left| \int_{t_{i-1}}^{t_i} \mathbf{c}'(t_i) - \mathbf{c}' \right| + \left| \int_{t_{i-1}}^{t_i} \mathbf{c}' \right| + \frac{\varepsilon}{2(b-a)} (t_i - t_{i-1}) \\ &\leq \left( \int_{t_{i-1}}^{t_i} |\mathbf{c}'(t_i) - \mathbf{c}'| \right) + |\mathbf{c}(t_i) - \mathbf{c}(t_{i-1})| + \frac{\varepsilon}{2(b-a)} (t_i - t_{i-1}) \\ &\leq |\mathbf{c}(t_i) - \mathbf{c}(t_{i-1})| + \frac{\varepsilon}{(b-a)} (t_i - t_{i-1}). \end{aligned}$$

Adding these inequalities for each  $i$  then yields the second claim.

If  $\mathbf{c}$  is only piecewise smooth, then the above argument holds over every sub-interval on which  $\mathbf{c}$  is smooth, and the result follows from the next lemma.  $\square$

**Lemma 2.4.2.** *If  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^n$  is a curve, then for any  $t_0 \in [a, b]$ ,  $\ell(\mathbf{c}) = \ell(\mathbf{c}|_{[a, t_0]}) + \ell(\mathbf{c}|_{[t_0, b]})$ .*

*Proof.* If  $P_1$  is a partition of the first sub-interval, and  $P_2$  one of the second, then  $P := P_1 \cup P_2$  is a partition of  $[a, b]$ , so that

$$\ell(P_1, \mathbf{c}|_{[a, t_0]}) + \ell(P_2, \mathbf{c}|_{[t_0, b]}) = \ell(P, \mathbf{c}) \leq \ell(\mathbf{c}).$$

Since this is true for any partitions,  $\ell(\mathbf{c}) \geq \ell(\mathbf{c}|_{[a,t_0]}) + \ell(\mathbf{c}|_{[t_0,b]})$ . On the other hand, if  $P$  is a partition of  $[a, b]$ , then the partition  $P'$  obtained by adding  $t_0$  to  $P$  is finer than  $P$ , and is a union of a partition  $P_1$  of  $[a, t_0]$  and a partition  $P_2$  of  $[t_0, b]$ . Hence

$$\ell(P, \mathbf{c}) \leq \ell(P', \mathbf{c}) = \ell(P_1, \mathbf{c}|_{[a,t_0]}) + \ell(P_2, \mathbf{c}|_{[t_0,b]}) \leq \ell(\mathbf{c}|_{[a,t_0]}) + \ell(\mathbf{c}|_{[t_0,b]}),$$

and since  $P$  was arbitrary,  $\ell(\mathbf{c}) \leq \ell(\mathbf{c}|_{[a,t_0]}) + \ell(\mathbf{c}|_{[t_0,b]})$ . This establishes the result.  $\square$

**Definition 2.4.3.** Let  $I$  and  $\tilde{I}$  denote compact intervals in  $\mathbb{R}$ . A curve  $\tilde{\mathbf{c}} : \tilde{I} \rightarrow \mathbb{R}^n$  is called a *reparametrization* of the curve  $\mathbf{c} : I \rightarrow \mathbb{R}^n$  if  $\tilde{\mathbf{c}} = \mathbf{c} \circ f$ , where  $f : \tilde{I} \rightarrow I$  is a strictly monotone (i.e., either increasing or decreasing) bijection.

The reader can easily verify that the relation  $\mathbf{c} \sim \tilde{\mathbf{c}}$  if  $\tilde{\mathbf{c}}$  is a reparametrization of  $\mathbf{c}$  is an equivalence relation: i.e.,  $\mathbf{c} \sim \mathbf{c}$  for any  $\mathbf{c}$ ,  $\mathbf{c} \sim \tilde{\mathbf{c}}$  implies  $\tilde{\mathbf{c}} \sim \mathbf{c}$ , and if  $\mathbf{c}_1 \sim \mathbf{c}_2$ ,  $\mathbf{c}_2 \sim \mathbf{c}_3$ , then  $\mathbf{c}_1 \sim \mathbf{c}_3$ . This divides the collection of all such curves into disjoint subsets called equivalence classes, with two curves belonging to the same equivalence class if and only if one is a reparametrization of the other. Two curves in the same class have the same image, and in particular the same length. It turns out that any piecewise-smooth curve admits a smooth reparametrization. In order to show this, we will need the following multi-purpose lemma:

**Lemma 2.4.3.** *Given  $a < b$ , there exists a smooth function  $f : \mathbb{R} \rightarrow \mathbb{R}$  such that:*

- (1)  $f(t) = 0$  if  $t \leq a$ ,  $0 < f(t) < 1$  if  $t \in (a, b)$ , and  $f(t) = 1$  if  $t \geq b$ ;
- (2)  $f$  is strictly increasing on  $(a, b)$ .

*Proof.* The result is actually an easy consequence of Theorem 2.2.5. Alternatively, define  $h : \mathbb{R} \rightarrow \mathbb{R}$  by

$$h(t) = \begin{cases} e^{-\frac{1}{t-a} + \frac{1}{t-b}} & \text{if } t \in (a, b), \\ 0 & \text{otherwise.} \end{cases}$$

By l'Hospital's rule, the left and right derivatives of  $h$  of any order exist and equal 0 at  $a$  and  $b$ . Thus,  $h$  is differentiable on  $\mathbb{R}$  and strictly positive on  $(a, b)$ . Now let  $f$  be given by:

$$f(t) = \frac{\int_a^t h}{\int_a^b h}. \quad \square$$

**Proposition 2.4.1.** *Any curve  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^n$  admits a smooth reparametrization  $\tilde{\mathbf{c}} : [a, b] \rightarrow \mathbb{R}^n$ .*

*Proof.* By assumption, there exist numbers  $a = t_0 < t_1 < \dots < t_k = b$  such that the restriction of  $\mathbf{c}$  to each  $(t_{i-1}, t_i)$ ,  $i = 1, \dots, k$ , is smooth. By Lemma 2.4.3, there exist functions  $f_i : [a, b] \rightarrow \mathbb{R}$  such that  $f_i(t) = 0$  if  $t \leq t_{i-1}$ ,  $f_i(t) = 1$  if  $t \geq t_i$ , and  $f_i$  is strictly increasing on  $[t_{i-1}, t_i]$ . Define a function  $f$  on  $[a, b]$  by

$$f(t) = t_0 + \sum_{i=1}^k (t_i - t_{i-1})f_i(t).$$



$f$  is a differentiable, strictly increasing function with  $f(t_i) = t_i$  for each  $i = 1, \dots, k$ . The curve  $\tilde{\mathbf{c}} := \mathbf{c} \circ f$  is then a reparametrization of  $\mathbf{c}$  which is smooth, since  $u^i \circ \tilde{\mathbf{c}}$  has vanishing left and right derivative at any  $t_j$  for  $i = 1, \dots, n, j = 1, \dots, k$ .  $\square$

**Definition 2.4.4.** A curve  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^n$  is said to be *regular* if  $\mathbf{c}'(t) \neq \mathbf{0}$  for  $t \in [a, b]$ . A regular curve is said to be *normal* or *parametrized by arc length* if  $|\mathbf{c}'(t)| = 1$  for all  $t$ .

Since a normal curve has unit speed, its length equals the length of the interval on which it is defined. It follows from the proof of Proposition 2.4.1 that the smooth reparametrization of a piecewise smooth curve is in general not regular. Regularity, though, has one nice advantage:

**Theorem 2.4.2.** *Every regular curve admits a reparametrization by arc length.*

*Proof.* Let  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^n$  be a curve, and denote by  $L$  its length  $\ell(\mathbf{c})$ . Define a function  $\ell_{\mathbf{c}} : [a, b] \rightarrow [0, L]$  by letting  $\ell_{\mathbf{c}}(t)$  equal to the length  $\ell(\mathbf{c}|_{[a,t]})$  of the restriction of  $\mathbf{c}$  to  $[a, t]$ . Since  $\mathbf{c}$  is regular,  $\ell'_{\mathbf{c}}(t) = |\mathbf{c}'(t)| > 0$ , so that  $\ell_{\mathbf{c}}$  admits a differentiable strictly increasing inverse  $\ell_{\mathbf{c}}^{-1}$ . We claim that  $\tilde{\mathbf{c}} := \mathbf{c} \circ \ell_{\mathbf{c}}^{-1} : [0, L] \rightarrow \mathbb{R}^n$  is normal. To see this, notice that if  $f : [\tilde{a}, \tilde{b}] \rightarrow [a, b]$  is an increasing differentiable function, then

$$\ell_{\mathbf{c} \circ f}(t) = \int_{\tilde{a}}^t |(\mathbf{c} \circ f)'| = \int_{\tilde{a}}^t |\mathbf{c}' \circ f| f' = \int_a^{f(t)} |\mathbf{c}'| = \ell_{\mathbf{c}} \circ f(t).$$

Thus,  $\ell_{\tilde{\mathbf{c}}} = \ell_{\mathbf{c} \circ \ell_{\mathbf{c}}^{-1}} = \ell_{\mathbf{c}} \circ \ell_{\mathbf{c}}^{-1}$  is the identity, and  $|\tilde{\mathbf{c}}'(t)| = \ell'_{\tilde{\mathbf{c}}}(t) = 1$ , as claimed.  $\square$

We end this section with a couple of applications of differentiation that involve the use of curves. In order to state them, we need the following:

**Definition 2.4.5.** A subset  $E$  of  $\mathbb{R}^n$  is said to be *path connected* if any two points of  $E$  can be joined by a curve that lies in  $E$ . It is said to be *convex* if it contains the line segment joining any two of its points.

It is easy to see that a path connected set  $E$  is also connected in the sense of Exercise 1.45: otherwise there would exist a separation  $X = U \cup V$  of  $X$ , and for any curve  $\mathbf{c} : [a, b] \rightarrow E$ , the connected set  $\mathbf{c}[a, b]$  would have to lie entirely in  $U$  or in  $V$ . Thus, there would be no curve joining a point in  $U$  to one in  $V$ , contradicting the assumption that  $E$  is path connected. The converse is not true, however. The topologist's sine curve from Exercise 1.32 is connected but not path connected. For our purposes, however, only path connectivity matters, so from now on, for the sake of brevity, a connected set will mean a path connected one.

**Lemma 2.4.4.** *Let  $\mathbf{c} : [0, 1] \rightarrow \mathbb{R}^n$  be a curve, and suppose its speed  $|\mathbf{c}'|$  does not exceed  $\alpha > 0$ . Then the distance  $|\mathbf{c}(1) - \mathbf{c}(0)|$  between its endpoints is less than or equal to  $\alpha$ .*

*Proof.* This is an easy consequence of Lemma 2.4.1:

$$|\mathbf{c}(1) - \mathbf{c}(0)| = \left| \int_0^1 \mathbf{c}' \right| \leq \int_0^1 |\mathbf{c}'| \leq \int_0^1 \alpha = \alpha. \quad \square$$

**Theorem 2.4.3.** Let  $\mathbf{f} : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a differentiable map.

- (1) If  $U$  is connected and  $D\mathbf{f} = 0$  on  $U$ , then  $\mathbf{f}$  is constant.
- (2) If  $U$  is convex and  $|D\mathbf{f}| \leq M$  on  $U$ , then

$$|\mathbf{f}(\mathbf{b}) - \mathbf{f}(\mathbf{a})| \leq M|\mathbf{b} - \mathbf{a}| \text{ for all } \mathbf{a}, \mathbf{b} \in U.$$

*Proof.* We begin with (1): Given  $\mathbf{a}, \mathbf{b} \in U$ , we must show that  $\mathbf{f}(\mathbf{a}) = \mathbf{f}(\mathbf{b})$ . Let  $\mathbf{c} : [0, 1] \rightarrow U$  be a curve that joins  $\mathbf{a}$  to  $\mathbf{b}$ , and assume for now that  $\mathbf{c}$  is smooth. The curve  $\mathbf{f} \circ \mathbf{c}$  in  $\mathbb{R}^m$  has zero velocity everywhere, since by the chain rule,  $(\mathbf{f} \circ \mathbf{c})' = D(\mathbf{f} \circ \mathbf{c}) = ((D\mathbf{f}) \circ \mathbf{c}) \circ D\mathbf{c}$ . But a curve that has zero velocity is constant, since its components are real-valued functions with vanishing derivative on an interval. Thus  $\mathbf{f}(\mathbf{a}) = (\mathbf{f} \circ \mathbf{c})(0) = (\mathbf{f} \circ \mathbf{c})(1) = \mathbf{f}(\mathbf{b})$ . If  $\mathbf{c}$  is merely piecewise-smooth, the same argument shows that  $\mathbf{f} \circ \mathbf{c}$  is constant on each subinterval where  $\mathbf{c}$  is smooth, and again we conclude that  $\mathbf{f}(\mathbf{a}) = \mathbf{f}(\mathbf{b})$ .

For (2), the curve  $\mathbf{c} : [0, 1] \rightarrow \mathbb{R}^n$ , where  $\mathbf{c}(t) = \mathbf{a} + t(\mathbf{b} - \mathbf{a})$ , lies in  $U$ , since it parametrizes the line segment between its endpoints. Thus, the speed

$$|(\mathbf{f} \circ \mathbf{c})'| = |((D\mathbf{f}) \circ \mathbf{c}) \circ D\mathbf{c}| \leq |(D\mathbf{f}) \circ \mathbf{c}| |D\mathbf{c}| = |(D\mathbf{f}) \circ \mathbf{c}| |\mathbf{b} - \mathbf{a}|$$

of the curve  $\mathbf{f} \circ \mathbf{c}$  never exceeds  $M|\mathbf{b} - \mathbf{a}|$ , and the claim follows from Lemma 2.4.4. For the inequality above, we used the fact that if  $L, M$  are linear transformations for which the composition  $L \circ M$  is defined, then  $|L \circ M| \leq |L||M|$ . This is because

$$|(L \circ M)\mathbf{u}| = |L(M(\mathbf{u}))| \leq |L||M\mathbf{u}| \leq |L||M||\mathbf{u}|$$

for any  $\mathbf{u}$  by (1.4.1). □

Finally, we mention a property of connected sets that will often be useful in the future:

**Theorem 2.4.4.** Let  $E$  be a connected set. If  $A$  is a nonempty subset that is both open and closed in  $E$ , then  $A = E$ .

*Proof.* This result is actually the object of Exercise 1.45. For the sake of completeness, we provide an independent proof in the context of path connected sets. Fix some  $\mathbf{a} \in A$ , and let  $\mathbf{b} \in E$  be arbitrary. It must be shown that  $\mathbf{b} \in A$ . Consider a curve  $\mathbf{c} : [0, 1] \rightarrow E$  with  $\mathbf{c}(0) = \mathbf{a}$  and  $\mathbf{c}(1) = \mathbf{b}$ , and define  $I = \{t \in [0, 1] \mid \mathbf{c}(t) \in A\}$ .  $I$  is nonempty because it contains 0, and is open by continuity of  $\mathbf{c}$ : in fact, if  $t_0 \in I$ , then  $A$  is a neighborhood of  $\mathbf{c}(t_0)$ , and there exists some  $\varepsilon > 0$  such that  $\mathbf{c}(t) \in A$  for all  $t$  satisfying  $|t - t_0| < \varepsilon$ . We claim that  $I$  is also closed in  $[0, 1]$ : suppose  $t_0$  is a boundary point of  $I$ . Then for any natural number  $i$ , there exists  $t_i \in I$  such that  $|t_i - t_0| < 1/i$ . Thus,

$t_i \rightarrow t_0$ , and  $\mathbf{c}(t_i) \rightarrow \mathbf{c}(t_0)$  by continuity of  $\mathbf{c}$ . This implies that any neighborhood of  $\mathbf{c}(t_0)$  contains points of  $A$ , so that  $\mathbf{c}(t_0)$  belongs to  $A$  or is a boundary point of  $A$ . But  $A$ , being closed in  $E$ , contains all its boundary points that lie in  $E$ , and so  $\mathbf{c}(t_0) \in A$ . This shows that  $I$  is also closed. By Proposition 1.7.1,  $I = [0, 1]$ , and  $\mathbf{b} \in A$ .  $\square$

## 2.5 The inverse and implicit function theorems

**Definition 2.5.1.** A map  $\mathbf{f} : U \subset \mathbb{R}^n \rightarrow \mathbf{f}(U) \subset \mathbb{R}^n$  is said to be a *diffeomorphism* of class  $C^k$  ( $k \geq 1$ ) if it is of class  $C^k$  on  $U$ , and admits an inverse  $\mathbf{f}^{-1}$  of class  $C^k$  on  $\mathbf{f}(U)$ .

The following theorem has many applications; it is yet another indication that the behavior of a map is locally similar to that of its derivative:

**Theorem 2.5.1** (Inverse function theorem). *Suppose  $\mathbf{f} : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$  is continuously differentiable. If  $D\mathbf{f}(\mathbf{a})$  is invertible at some  $\mathbf{a} \in U$ , then there exists a neighborhood  $V$  of  $\mathbf{a}$  such that the restriction of  $\mathbf{f}$  to  $V$  is a diffeomorphism of class  $C^1$ . Furthermore, given  $\mathbf{p} \in V$ , and  $\mathbf{q} = \mathbf{f}(\mathbf{p})$ ,  $D(\mathbf{f}^{-1})(\mathbf{q}) = (D\mathbf{f}(\mathbf{p}))^{-1}$ .*

The argument will use a property of complete metric spaces (recall that a metric space is complete if every Cauchy sequence converges). It will only be applied in a very specific context, namely for a closed subset  $C$  of Euclidean space, but its proof is the same in the general setting. A map  $f : X \rightarrow X$  from a metric space  $(X, d)$  into itself is said to be a *contraction* if there exists  $C \in [0, 1)$  such that

$$d(f(p), f(q)) \leq C d(p, q), \quad p, q \in X.$$

Notice that a contraction is always continuous; it is, in fact, uniformly continuous: given  $\varepsilon > 0$ , take  $\delta = \varepsilon/C$  if  $C > 0$  (if  $C = 0$ , any  $\delta$  will do).

**Lemma 2.5.1.** *If  $f : X \rightarrow X$  is a contraction of a complete metric space, then  $f$  has a unique fixed point; i.e., a point  $p$  such that  $f(p) = p$ .*

*Proof of lemma.* Uniqueness is clear: if  $p$  and  $q$  are fixed points of  $f$ , then  $d(p, q) = d(f(p), f(q)) \leq C d(p, q)$ , which can only happen if  $d(p, q) = 0$ ; i.e., if  $p = q$ . For existence, let  $q$  be any point of  $X$ , and define a sequence recursively by  $x_1 = q$ ,  $x_{k+1} = f(x_k)$ . Then  $d(x_k, x_{k+1}) \leq C d(x_{k-1}, x_k)$ , and using induction,

$$d(x_k, x_{k+1}) \leq C^{k-1} d(x_1, x_2).$$

It follows that for  $m > n$ ,

$$\begin{aligned} d(x_n, x_m) &\leq d(x_n, x_{n+1}) + d(x_{n+1}, x_{n+2}) + \cdots + d(x_{m-1}, x_m) \\ &\leq (C^{n-1} + C^n + \cdots + C^{m-2}) d(x_1, x_2) \\ &= C^{n-1} (1 + C + \cdots + C^{m-n-1}) d(x_1, x_2) \\ &\leq C^{n-1} \left( \sum_{k=0}^{\infty} C^k \right) d(x_1, x_2) = C^{n-1} \frac{d(x_1, x_2)}{1 - C}. \end{aligned}$$

Since  $C < 1$ ,  $C^n \rightarrow 0$ , and  $\{x_n\}$  is a Cauchy sequence, which must therefore converge to some  $x$ . But  $f$  is continuous, so that

$$f(x) = f\left(\lim_{n \rightarrow \infty} x_n\right) = \lim_{n \rightarrow \infty} f(x_n) = \lim_{n \rightarrow \infty} x_{n+1} = x,$$

which proves that  $x$  is a fixed point of  $f$ .  $\square$

*Proof of Theorem 2.5.1.* To simplify notation, it may be assumed that  $\mathbf{a} = \mathbf{f}(\mathbf{a}) = \mathbf{0}$ : for the map  $\mathbf{F}$ , where  $\mathbf{F}(\mathbf{p}) = \mathbf{f}(\mathbf{p} + \mathbf{a}) - \mathbf{f}(\mathbf{a})$ , satisfies these conditions, and if  $\mathbf{F}$  is a diffeomorphism in a neighborhood of  $\mathbf{0}$ , then  $\mathbf{p} \mapsto \mathbf{f}(\mathbf{p}) = \mathbf{F}(\mathbf{p} - \mathbf{a}) + \mathbf{f}(\mathbf{a})$ , being a composition of local diffeomorphisms, is one too. Similarly,  $D\mathbf{f}(\mathbf{0})$  may be assumed to be the identity  $I$  (replacing  $\mathbf{f}$ , if need be, by  $\mathbf{p} \mapsto (D\mathbf{f}(\mathbf{0}))^{-1} \circ \mathbf{f}(\mathbf{p})$ , which by the chain rule has the identity as derivative at  $\mathbf{0}$ ). Thus, the associated map

$$\begin{aligned} \mathbf{g} : U &\rightarrow \mathbb{R}^n, \\ \mathbf{p} &\mapsto \mathbf{p} - \mathbf{f}(\mathbf{p}) \end{aligned}$$

satisfies  $\mathbf{g}(\mathbf{0}) = \mathbf{0}$  and  $D\mathbf{g}(\mathbf{0}) = 0$ .  $B_r$  will denote the closed ball of radius  $r$  centered at  $\mathbf{0}$ . Recalling that  $D\mathbf{f}$  is continuous and that  $D\mathbf{f}(\mathbf{0}) = I$ , there exists  $r > 0$  small enough that

$$|D\mathbf{g}(\mathbf{p})| = |I - D\mathbf{f}(\mathbf{p})| < \frac{1}{2}, \quad \mathbf{p} \in B_r.$$

Together with Theorem 2.4.3, this means that

$$|\mathbf{g}(\mathbf{p}) - \mathbf{g}(\mathbf{q})| \leq \frac{1}{2}|\mathbf{p} - \mathbf{q}|, \quad \mathbf{p}, \mathbf{q} \in B_r, \quad (2.5.1)$$

and consequently,

$$|\mathbf{f}(\mathbf{p}) - \mathbf{f}(\mathbf{q})| \geq |\mathbf{p} - \mathbf{q}| - |\mathbf{g}(\mathbf{p}) - \mathbf{g}(\mathbf{q})| \geq \frac{1}{2}|\mathbf{p} - \mathbf{q}|, \quad \mathbf{p}, \mathbf{q} \in B_r. \quad (2.5.2)$$

Notice that (2.5.2) already implies that  $\mathbf{f}$  is one-to-one on  $B_r$ . Next, we claim that  $\mathbf{f}(B_r)$  contains  $B_{r/2}$ . To see why, consider any  $\mathbf{q} \in B_{r/2}$ . In order to find some  $\mathbf{p} \in B_r$  that gets mapped to  $\mathbf{q}$ , observe that  $\mathbf{g}_q$ , where  $\mathbf{g}_q(\mathbf{p}) = \mathbf{g}(\mathbf{p}) + \mathbf{q}$ , sends  $B_r$  to itself. But  $B_r$ , being closed, is complete, and by (2.5.1),  $\mathbf{g}_q$  is a contraction, since

$$|\mathbf{g}_q(\mathbf{p}) - \mathbf{g}_q(\mathbf{r})| = |\mathbf{g}(\mathbf{p}) - \mathbf{g}(\mathbf{r})| \leq \frac{1}{2}|\mathbf{p} - \mathbf{r}|.$$

By Lemma 2.5.1,  $\mathbf{g}_q$  has a unique fixed point  $\mathbf{p}$ ; i.e.,  $\mathbf{g}_q(\mathbf{p}) = \mathbf{p} - \mathbf{f}(\mathbf{p}) + \mathbf{q} = \mathbf{p}$ , or  $\mathbf{q} = \mathbf{f}(\mathbf{p})$ , as claimed.

Set  $W = \text{int}(B_{r/2})$ ,  $V = \mathbf{f}^{-1}(W)$ . Then the restriction  $\mathbf{f} : V \rightarrow W$  admits an inverse  $\mathbf{f}^{-1} : W \rightarrow V$ . This inverse is certainly continuous, since by (2.5.2),

$$|\mathbf{f}^{-1}(\mathbf{p}) - \mathbf{f}^{-1}(\mathbf{q})| \leq 2|\mathbf{p} - \mathbf{q}|.$$

We claim it is actually differentiable at any  $\mathbf{q}_0 \in W$ . To see this, let  $\mathbf{p}_0 = \mathbf{f}^{-1}(\mathbf{q}_0)$ ,  $L = D\mathbf{f}(\mathbf{p}_0)$ , and  $r_f$  the map defined by

$$r_f(\mathbf{h}) = \mathbf{f}(\mathbf{p}_0 + \mathbf{h}) - \mathbf{f}(\mathbf{p}_0) - L\mathbf{h},$$

so that as usual,

$$\mathbf{f}(\mathbf{p}_0 + \mathbf{h}) = \mathbf{f}(\mathbf{p}_0) + L\mathbf{h} + r_f(\mathbf{h}), \text{ where } \lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{|r_f(\mathbf{h})|}{|\mathbf{h}|} = 0.$$

Apply  $L^{-1}$  to both sides of the above identity to conclude that

$$\mathbf{h} = L^{-1}(\mathbf{f}(\mathbf{p}_0 + \mathbf{h}) - \mathbf{f}(\mathbf{p}_0)) - L^{-1}(r_f(\mathbf{h})). \quad (2.5.3)$$

Set  $\mathbf{k} = \mathbf{f}(\mathbf{p}_0 + \mathbf{h}) - \mathbf{f}(\mathbf{p}_0)$ . Then

$$\begin{aligned} \mathbf{f}^{-1}(\mathbf{q}_0 + \mathbf{k}) - \mathbf{f}^{-1}(\mathbf{q}_0) &= \mathbf{f}^{-1}(\mathbf{q}_0 + \mathbf{f}(\mathbf{p}_0 + \mathbf{h}) - \mathbf{f}(\mathbf{p}_0)) - \mathbf{f}^{-1}(\mathbf{q}_0) \\ &= \mathbf{f}^{-1}(\mathbf{f}(\mathbf{p}_0 + \mathbf{h})) - \mathbf{p}_0 \\ &= \mathbf{h}, \end{aligned}$$

so that by (2.5.3) and the definition of  $\mathbf{k}$ ,

$$\mathbf{f}^{-1}(\mathbf{q}_0 + \mathbf{k}) - \mathbf{f}^{-1}(\mathbf{q}_0) - L^{-1}\mathbf{k} = -L^{-1} \circ r_f(\mathbf{h}).$$

This means that  $\mathbf{f}$  is differentiable at  $\mathbf{q}_0$  with derivative  $L^{-1}$ , provided

$$\lim_{\mathbf{k} \rightarrow \mathbf{0}} \frac{|L^{-1} \circ r_f(\mathbf{h})|}{|\mathbf{k}|} = 0. \quad (2.5.4)$$

To establish this identity, notice that the definition of  $\mathbf{k}$  together with (2.5.2) implies that  $|\mathbf{k}| \geq |\mathbf{h}|/2$ , so that

$$\frac{|L^{-1} \circ r_f(\mathbf{h})|}{|\mathbf{k}|} \leq |L^{-1}| \frac{|r_f(\mathbf{h})|}{|\mathbf{h}|} \frac{|\mathbf{h}|}{|\mathbf{k}|} \leq 2|L^{-1}| \frac{|r_f(\mathbf{h})|}{|\mathbf{h}|}.$$

By continuity of  $\mathbf{f}^{-1}$ ,  $\mathbf{h} \rightarrow \mathbf{0}$  when  $\mathbf{k} \rightarrow \mathbf{0}$ . This proves (2.5.4) and the differentiability of  $\mathbf{f}^{-1}$ . It only remains to show that  $\mathbf{f}^{-1}$  is continuously differentiable. But this is an immediate consequence of the fact that the matrix of  $D(\mathbf{f}^{-1})$  is the inverse of the matrix of  $D\mathbf{f}$ , and inversion of matrices is a continuous map, see Exercise 1.19.  $\square$

**Examples 2.5.1.** (i) The *polar coordinates* of a point in the plane  $\mathbb{R}^2$  are  $(r, \theta)$ , where  $r$  equals the distance from the point to the origin, and  $\theta$  is the angle between the vector representing the point and  $\mathbf{e}_1$ . They are well defined for any point different from the origin; in fact, in terms of Cartesian coordinates,

$$(r, \theta)(x, y) = \mathbf{f}(x, y) = (\sqrt{x^2 + y^2}, \arctan \frac{y}{x} + c), \quad (2.5.5)$$

if  $x \neq 0$ . Here,  $c = 0, \pi$ , or  $2\pi$  depending on which quadrant  $(x, y)$  lies, cf. Section 4.6.1. When  $x = 0$ ,  $\theta = \pi/2$  if  $y > 0$ , and  $3\pi/2$  if  $y < 0$ . The Jacobian matrix of  $\mathbf{f}$  is

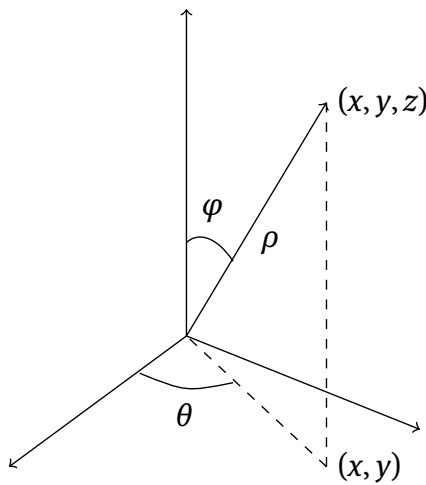
$$D\mathbf{f}(x, y) = \begin{bmatrix} \frac{x}{\sqrt{x^2 + y^2}} & \frac{y}{\sqrt{x^2 + y^2}} \\ \frac{-y}{x^2 + y^2} & \frac{x}{x^2 + y^2} \end{bmatrix},$$

which has determinant  $(x^2 + y^2)^{-1/2}$ . By Theorem 2.5.1,  $f$  is invertible in a neighborhood of any point  $\mathbf{p}$  in its domain, with

$$(x, y) = f^{-1}(r, \theta) = (r \cos \theta, r \sin \theta).$$

The domain of the inverse can be taken to be  $(0, \infty) \times [a, a + 2\pi)$ ,  $a \in \mathbb{R}$ , with the value of  $a$  depending on the point  $\mathbf{p}$ . When  $\mathbf{p}$  does not lie on the positive  $x$ -axis, it is customary to take  $a = 0$ .

- (ii) The *spherical coordinates* of a point  $(x, y, z)$  in  $\mathbb{R}^3$  are  $(\rho, \theta, \varphi)$ , where  $\rho$  is the distance from the point to the origin,  $\theta$  is the polar coordinate angle of the projection  $(x, y)$  of the point in the plane, and  $\varphi$  is the angle between the vector representing the point and  $\mathbf{e}_3$ .



Spherical coordinates  $(\rho, \theta, \varphi)$  of the point  $(x, y, z)$  in  $\mathbb{R}^3$

$\rho$ ,  $\theta$ , and  $\varphi$  are the components of the map  $f$  given by

$$f(x, y, z) = \left( (x^2 + y^2 + z^2)^{1/2}, \arctan \frac{y}{x} + c, \arccos \frac{z}{(x^2 + y^2 + z^2)^{1/2}} \right)$$

(with  $c$  as above), and its inverse is

$$f^{-1}(\rho, \theta, \varphi) = (\rho \sin \varphi \cos \theta, \rho \sin \varphi \sin \theta, \rho \cos \varphi).$$

One easily computes that the Jacobian of  $f^{-1}$  has determinant  $-\rho^2 \sin \varphi$ , see Section 4.6.2.

- (iii) Although the inverse function theorem asserts that a map from  $\mathbb{R}^n$  to itself is one-to-one in a neighborhood of any point where the derivative is one-to-one, the converse is not true already for  $n = 1$ . For example,  $f : \mathbb{R} \rightarrow \mathbb{R}$ , where  $f(x) = x^3$ , is globally one-to-one even though its derivative is zero at the origin. Of course, even though a map may still be one-to-one with vanishing derivative at a point, it cannot in that case be a diffeomorphism, since diffeomorphisms always have invertible derivative.

The inverse function theorem, which deals with maps from  $\mathbb{R}^n$  to  $\mathbb{R}^n$ , admits generalizations to maps from  $\mathbb{R}^n$  to  $\mathbb{R}^k$ , when  $k$  is larger or smaller than  $n$ . They are known as the implicit function theorems. Before stating them, we need some terminology. If  $\mathbf{f}$  is differentiable, its *rank* at  $\mathbf{p}$  is defined to be the rank of the linear map  $D\mathbf{f}(\mathbf{p})$ ; i.e., the dimension of the image  $D\mathbf{f}(\mathbf{p})\mathbb{R}^n$  of the derivative at  $\mathbf{p}$ . Since this image is spanned by the columns  $D\mathbf{f}(\mathbf{p})\mathbf{e}_1, \dots, D\mathbf{f}(\mathbf{p})\mathbf{e}_n$  of the Jacobian matrix, the rank equals the number of linearly independent columns of the Jacobian. It turns out that the dimension of the space spanned by the columns of an  $m \times n$  matrix  $A$  equals that of the space spanned by the rows of  $A$ : indeed, apply Theorem 1.2.2 to the linear transformation  $L_A : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to deduce that the former equals  $n - \dim \ker L_A$ . The latter, on the other hand, equals the dimension of the space spanned by  $\mathbf{a}_1, \dots, \mathbf{a}_m$ , where  $\mathbf{a}_i$  denotes the transpose (so it can be viewed as a vector) of the  $i$ -th row of the matrix. But for  $\mathbf{p} \in \mathbb{R}^n$ ,

$$A\mathbf{p} = \begin{bmatrix} \langle \mathbf{a}_1, \mathbf{p} \rangle \\ \vdots \\ \langle \mathbf{a}_m, \mathbf{p} \rangle \end{bmatrix},$$

so that the kernel of  $L_A$  coincides with the orthogonal complement of the row space. By Proposition 1.4.1, the row space also has dimension  $n - \dim \ker L_A$ , as claimed.

**Theorem 2.5.2** (Implicit Function Theorem). *Let  $U$  be a neighborhood of  $\mathbf{0}$  in  $\mathbb{R}^n$ ,  $\mathbf{f} : U \rightarrow \mathbb{R}^k$  a continuously differentiable map with  $\mathbf{f}(\mathbf{0}) = \mathbf{0}$ . For  $n \leq k$ , let  $\iota : \mathbb{R}^n \rightarrow \mathbb{R}^k$  denote the inclusion*

$$\iota(a_1, \dots, a_n) = (a_1, \dots, a_n, 0, \dots, 0),$$

and for  $n \geq k$ , let  $\pi : \mathbb{R}^n \rightarrow \mathbb{R}^k$  denote the projection

$$\pi(a_1, \dots, a_k, \dots, a_n) = (a_1, \dots, a_k).$$

- (1) *If  $n \leq k$  and  $\mathbf{f}$  has maximal rank ( $= n$ ) at  $\mathbf{0}$ , then there exists a diffeomorphism  $\mathbf{g}$  of a neighborhood of  $\mathbf{0}$  in  $\mathbb{R}^k$  such that  $\mathbf{g} \circ \mathbf{f} = \iota$  in a neighborhood of  $\mathbf{0} \in \mathbb{R}^n$ .*
- (2) *If  $n \geq k$  and  $\mathbf{f}$  has maximal rank ( $= k$ ) at  $\mathbf{0}$ , then there exists a diffeomorphism  $\mathbf{h}$  of a neighborhood of  $\mathbf{0}$  in  $\mathbb{R}^n$  such that  $\mathbf{f} \circ \mathbf{h} = \pi$ .*

*Proof.* As usual, we denote the component functions  $u^i \circ \mathbf{f}$  of  $\mathbf{f}$  by  $f^i$ . In order to prove (1), observe that the  $k \times n$  matrix  $[D_j f^i(\mathbf{0})]$  has rank  $n$ . By rearranging the component functions  $f^i$  of  $\mathbf{f}$  if necessary (which amounts to composing  $\mathbf{f}$  with an invertible transformation, hence a diffeomorphism of  $\mathbb{R}^k$ ), we may assume that the  $n \times n$  submatrix  $[D_j f^i(\mathbf{0})]_{1 \leq i, j \leq n}$  is invertible. Define  $\mathbf{F} : U \times \mathbb{R}^{k-n} \rightarrow \mathbb{R}^k$  by

$$\mathbf{F}(a_1, \dots, a_n, a_{n+1}, \dots, a_k) := \mathbf{f}(a_1, \dots, a_n) + (0, \dots, 0, a_{n+1}, \dots, a_k).$$

Then  $\mathbf{F} \circ \iota = \mathbf{f}$ , and the Jacobian matrix of  $\mathbf{F}$  at  $\mathbf{0}$  is

$$\begin{bmatrix} [D_j f^i(\mathbf{0})]_{1 \leq i, j \leq n} & \mathbf{0} \\ [D_j f^i(\mathbf{0})]_{n+1 \leq i \leq k} & \mathbf{1}_{\mathbb{R}^{k-n}} \end{bmatrix},$$

which has nonzero determinant. Consequently,  $\mathbf{F}$  has a local inverse  $\mathbf{g}$ , and  $\mathbf{g} \circ \mathbf{f} = \mathbf{g} \circ \mathbf{F} \circ \iota = \iota$ . This establishes (1). Similarly, in (2), we may assume that the  $k \times k$  submatrix  $[D_j f^i(\mathbf{0})]_{1 \leq i, j \leq k}$  is invertible. Define  $\mathbf{F} : U \rightarrow \mathbb{R}^n$  by

$$\mathbf{F}(a_1, \dots, a_n) := (\mathbf{f}(a_1, \dots, a_n), a_{k+1}, \dots, a_n).$$

Then  $\mathbf{f} = \pi \circ \mathbf{F}$ , and the Jacobian of  $\mathbf{F}$  at  $\mathbf{0}$  is

$$\begin{bmatrix} [D_j f^i(\mathbf{0})]_{1 \leq j \leq k} & [D_j f^i(\mathbf{0})]_{k+1 \leq j \leq n} \\ 0 & 1_{\mathbb{R}^{n-k}} \end{bmatrix},$$

which is invertible. Thus,  $\mathbf{F}$  has a local inverse  $\mathbf{h}$ , and  $\mathbf{f} \circ \mathbf{h} = \pi \circ \mathbf{F} \circ \mathbf{h} = \pi$ .  $\square$

The reason the above is called the implicit function theorem is that under certain circumstances, given a differentiable map  $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ , the equation

$$\mathbf{f}(\mathbf{p}, \mathbf{q}) = \mathbf{0}, \quad \mathbf{p} \in \mathbb{R}^n, \quad \mathbf{q} \in \mathbb{R}^m,$$

implicitly defines  $\mathbf{p}$  as a function  $\mathbf{g}(\mathbf{q})$  of  $\mathbf{q}$ ; i.e., for each  $\mathbf{q}$ , there exists a unique  $\mathbf{g}(\mathbf{q})$  such that  $\mathbf{f}(\mathbf{g}(\mathbf{q}), \mathbf{q}) = \mathbf{0}$ , and the map  $\mathbf{g}$  is differentiable; more precisely, we have:

**Corollary 2.5.1** (Classical implicit function theorem). *Let  $\mathbf{f}$  be a continuously differentiable map from an open neighborhood of  $(\mathbf{0}, \mathbf{0}) \in \mathbb{R}^n \times \mathbb{R}^m$  to  $\mathbb{R}^n$  such that  $\mathbf{f}(\mathbf{0}, \mathbf{0}) = \mathbf{0}$ .*

*If the matrix  $[D_j f^i(\mathbf{0}, \mathbf{0})]_{1 \leq i, j \leq n}$  is invertible, then there exist open neighborhoods  $V$  of  $\mathbf{0} \in \mathbb{R}^m$ ,  $U$  of  $(\mathbf{0}, \mathbf{0}) \in \mathbb{R}^n \times \mathbb{R}^m$  such that for every  $\mathbf{q} \in V$ , there exists a unique  $\mathbf{p}$  satisfying  $(\mathbf{p}, \mathbf{q}) \in U$  and  $\mathbf{f}(\mathbf{p}, \mathbf{q}) = \mathbf{0}$ . Furthermore, the map  $\mathbf{g} : V \rightarrow \mathbb{R}^n$  that assigns to each  $\mathbf{q}$  the unique  $\mathbf{p} = \mathbf{g}(\mathbf{q})$  such that  $\mathbf{f}(\mathbf{g}(\mathbf{q}), \mathbf{q}) = \mathbf{0}$  is continuously differentiable.*

*Proof.* By Theorem 2.5.2(2), there exists a local diffeomorphism  $\mathbf{h}$  such that  $\mathbf{f} \circ \mathbf{h}$  equals the projection  $\pi_1 : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$  onto the first factor. If  $\pi_2$  is the projection onto the second factor, then the inverse  $\mathbf{F}$  of  $\mathbf{h}$  was shown to satisfy  $\pi_2 \circ \mathbf{F} = \pi_2$ . This implies that if  $\mathbf{f}(\mathbf{p}, \mathbf{q}) = \mathbf{0}$ , then

$$\mathbf{0} = \mathbf{f}(\mathbf{p}, \mathbf{q}) = \mathbf{f} \circ \mathbf{h} \circ \mathbf{F}(\mathbf{p}, \mathbf{q}) = \pi_1 \circ \mathbf{F}(\mathbf{p}, \mathbf{q}),$$

so that  $\mathbf{F}(\mathbf{p}, \mathbf{q}) = (\mathbf{0}, \mathbf{q})$ ; i.e.,  $\mathbf{p} = \pi_1 \circ \mathbf{F}^{-1}(\mathbf{0}, \mathbf{q}) = \pi_1 \circ \mathbf{h}(\mathbf{0}, \mathbf{q})$ . This shows not only uniqueness of  $\mathbf{p}$ , but also existence and smoothness of  $\mathbf{g}$ : let

$$\begin{aligned} \iota : V &\mapsto \mathbb{R}^n \times \mathbb{R}^m, \\ \mathbf{q} &\mapsto (\mathbf{0}, \mathbf{q}). \end{aligned}$$

Then the map  $\mathbf{g} := \pi_1 \circ \mathbf{h} \circ \iota$ , being a composition of  $\mathcal{C}^1$  maps, is continuously differentiable, and

$$\begin{aligned} \mathbf{f}(\mathbf{g}(\mathbf{q}), \mathbf{q}) &= \mathbf{f}(\pi_1(\mathbf{h}(\mathbf{0}, \mathbf{q})), \mathbf{q}) = \mathbf{f}(\pi_1(\mathbf{h}(\mathbf{0}, \mathbf{q})), \pi_2(\mathbf{h}(\mathbf{0}, \mathbf{q}))) \\ &= (\mathbf{f} \circ \mathbf{h})(\mathbf{0}, \mathbf{q}) = \pi_1(\mathbf{0}, \mathbf{q}) \\ &= \mathbf{0}. \end{aligned} \quad \square$$



**Examples and Remarks 2.5.1.** (i) Define

$$f : \mathbb{R} \times \mathbb{R}^2 \rightarrow \mathbb{R},$$

$$(x, (y, z)) \mapsto (x - r)^2 + y^2 + z^2 - r^2,$$

with  $r > 0$ . The set of points  $(x, y, z) \in \mathbb{R}^3$  satisfying  $f(x, y, z) = 0$  is a sphere of radius  $r$  that passes through the origin. Since  $D_1 f(0, 0) = -2r \neq 0$ , a part of the sphere containing the origin may be represented by the graph of a function  $x = g(y, z)$ , even though the whole sphere cannot. In fact, solving  $f(x, y, z) = 0$  in terms of  $x$  yields  $x = g(y, z) = r - \sqrt{r^2 - y^2 - z^2}$ .

(ii) Although the classical implicit function theorem was stated at the origin, it is straightforward to obtain a version of it that is valid at any point: specifically, if  $\mathbf{f}$  is a continuously differentiable map from a neighborhood of  $(\mathbf{p}_0, \mathbf{q}_0) \in \mathbb{R}^n \times \mathbb{R}^m$  to  $\mathbb{R}^n$  such that  $\mathbf{f}(\mathbf{p}_0, \mathbf{q}_0) = \mathbf{r}_0 \in \mathbb{R}^n$ , and if the matrix  $[D_j f^i(\mathbf{p}_0, \mathbf{q}_0)]_{1 \leq i, j \leq n}$  is invertible, then there exists an open neighborhood  $V$  of  $\mathbf{q}_0 \in \mathbb{R}^m$ , and a continuously differentiable map  $\mathbf{g} : V \rightarrow \mathbb{R}^n$  such that  $\mathbf{f}(\mathbf{g}(\mathbf{q}), \mathbf{q}) = \mathbf{r}_0$ .

To see this, observe that the map  $\mathbf{h}$ , where  $\mathbf{h}(\mathbf{p}, \mathbf{q}) := \mathbf{f}(\mathbf{p} + \mathbf{p}_0, \mathbf{q} + \mathbf{q}_0) - \mathbf{r}_0$ , satisfies the hypotheses of the implicit function theorem, so that there exists a continuously differentiable map  $\mathbf{k} : V_0 \rightarrow \mathbb{R}^n$  defined on a neighborhood  $V_0$  of the origin in  $\mathbb{R}^m$  such that  $\mathbf{h}(\mathbf{k}(\mathbf{q}), \mathbf{q}) = \mathbf{0}$ . Then  $V := \{\mathbf{q} + \mathbf{q}_0 \mid \mathbf{q} \in V_0\}$  is a neighborhood of  $\mathbf{q}_0$ , and  $\mathbf{g} : V \rightarrow \mathbb{R}^n$ , where  $\mathbf{g}(\mathbf{q}) := \mathbf{k}(\mathbf{q} - \mathbf{q}_0) + \mathbf{p}_0$ , is the map we are looking for, since

$$\mathbf{f}(\mathbf{g}(\mathbf{q}), \mathbf{q}) = \mathbf{f}(\mathbf{k}(\mathbf{q} - \mathbf{q}_0) + \mathbf{p}_0, \mathbf{q}) = \mathbf{h}(\mathbf{k}(\mathbf{q} - \mathbf{q}_0), \mathbf{q} - \mathbf{q}_0) + \mathbf{r}_0 = \mathbf{r}_0.$$

(iii) To illustrate (ii) above, suppose we are asked to show that there is a neighborhood  $U$  of  $(1, 0)$  in  $\mathbb{R}^2$ , and a differentiable function  $f : U \rightarrow \mathbb{R}$  satisfying

$$xf^2(x, y) + yf(x, y) - x^2f^5(x, y) + y = 0, \quad f(1, 0) = 1.$$

Finding an explicit formula for  $f$  seems pretty much hopeless, so let us define  $F : \mathbb{R}^3 \rightarrow \mathbb{R}$  by  $F(a_1, a_2, a_3) = a_2 a_1^2 + a_3 a_1 - a_2^2 a_1^5 + a_3$ . Then  $F(1, 1, 0) = 0$ , and  $D_1 F(1, 1, 0) = -3 \neq 0$ , so there exists a function  $f$  on a neighborhood of  $(1, 0)$  such that

$$F(f(a_2, a_3), a_2, a_3) = a_2 f^2(a_2, a_3) + a_3 f(a_2, a_3) - a_2^2 f^5(a_2, a_3) + a_3 = 0.$$

Substituting  $x$  for  $a_2$  and  $y$  for  $a_3$  shows this is the function we are looking for.

## 2.6 The spectral theorem and scalar products

In this section, we take a closer look at symmetric matrices since they pop up frequently in calculus and differential geometry: for instance, to each point of a surface,

one can associate a symmetric matrix that measures how curved the surface is at that point. As another example, in the next section, we will use symmetric matrices to classify extreme values of real-valued functions. The discussion will be couched in terms of linear maps rather than matrices.

Recall that a linear map from a vector space  $V$  to itself is also called an *operator* on  $V$ . If  $V$  is an inner product space, there are several types of operators on  $V$  that play a special role; we now introduce two of these:

**Definition 2.6.1.** A linear operator  $L$  on an inner product space  $V$  is said to be *self-adjoint* if  $\langle L\mathbf{u}, \mathbf{v} \rangle = \langle \mathbf{u}, L\mathbf{v} \rangle$  for all  $\mathbf{u}, \mathbf{v} \in V$ .  $L$  is said to be *skew-adjoint* if  $\langle L\mathbf{u}, \mathbf{v} \rangle = -\langle \mathbf{u}, L\mathbf{v} \rangle$  for all  $\mathbf{u}, \mathbf{v} \in V$ .

**Proposition 2.6.1.** Let  $V$  be an inner product space with ordered orthonormal basis  $\mathcal{B} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$ .

- (1) If  $L$  is a linear operator on  $V$ , then the  $(i, j)$ -th entry of the matrix  $[L]_{\mathcal{B}, \mathcal{B}}$  of  $L$  in the basis  $\mathcal{B}$  is  $\langle \mathbf{u}_i, L\mathbf{u}_j \rangle$ .
- (2)  $L$  is self-adjoint if and only if  $[L]_{\mathcal{B}, \mathcal{B}}$  is symmetric.
- (3)  $L$  is skew-adjoint if and only if  $[L]_{\mathcal{B}, \mathcal{B}}$  is skew-symmetric.

*Proof.* For (1), recall that the  $j$ -th column of the matrix of  $L$  in the basis  $\mathcal{B}$  is the coordinate vector  $[L\mathbf{u}_j]_{\mathcal{B}}$  of  $L\mathbf{u}_j$  in the basis  $\mathcal{B}$ . By Theorem 1.4.3, the  $i$ -th entry of that vector is  $\langle \mathbf{u}_i, L\mathbf{u}_j \rangle$ . For (2), if  $L$  is self-adjoint, then the matrix of  $L$  in  $\mathcal{B}$  is symmetric by (1). Conversely, if the matrix of  $L$  is symmetric, then by (1),  $\langle L\mathbf{u}_i, \mathbf{u}_j \rangle = \langle \mathbf{u}_i, L\mathbf{u}_j \rangle$  for all  $i$  and  $j$ , and given  $\mathbf{u}, \mathbf{v} \in V$ ,

$$\begin{aligned} \langle L\mathbf{u}, \mathbf{v} \rangle &= \left\langle L \left( \sum_i \langle \mathbf{u}, \mathbf{u}_i \rangle \mathbf{u}_i \right), \sum_j \langle \mathbf{v}, \mathbf{u}_j \rangle \mathbf{u}_j \right\rangle \\ &= \sum_{i,j} \langle \mathbf{u}, \mathbf{u}_i \rangle \langle \mathbf{v}, \mathbf{u}_j \rangle \langle L\mathbf{u}_i, \mathbf{u}_j \rangle = \sum_{i,j} \langle \mathbf{u}, \mathbf{u}_i \rangle \langle \mathbf{v}, \mathbf{u}_j \rangle \langle \mathbf{u}_i, L\mathbf{u}_j \rangle \\ &= \left\langle \sum_i \langle \mathbf{u}, \mathbf{u}_i \rangle \mathbf{u}_i, L \left( \sum_j \langle \mathbf{v}, \mathbf{u}_j \rangle \mathbf{u}_j \right) \right\rangle = \langle \mathbf{u}, L\mathbf{v} \rangle, \end{aligned}$$

so that  $L$  is self-adjoint. The proof of (3) is similar.  $\square$

A third important type of operator that occurs in the presence of an inner product can be formulated more generally in terms of a linear transformation between two (possibly) different inner product spaces:

**Definition 2.6.2.** A linear transformation  $L : (V_1, \langle \cdot, \cdot \rangle_1) \rightarrow (V_2, \langle \cdot, \cdot \rangle_2)$  between two inner product spaces of the same dimension is said to be a *linear isometry* if it preserves the inner product; i.e., if

$$\langle L\mathbf{u}, L\mathbf{v} \rangle_2 = \langle \mathbf{u}, \mathbf{v} \rangle_1, \quad \mathbf{u}, \mathbf{v} \in V_1.$$

Equivalently,  $L$  is a linear isometry if it preserves the norm of vectors. It follows that a linear isometry has trivial kernel (for if  $L\mathbf{u} = \mathbf{0}$ , then  $|\mathbf{u}| = |L\mathbf{u}| = 0$ , so  $\mathbf{u} = \mathbf{0}$ ), and is

therefore an isomorphism. Two inner product spaces with a linear isometry between them are said to be *isometric*. Isometric spaces are essentially the same both from a vector space perspective and from an inner product one. For example, it is easy to check that if  $T : V \rightarrow W$  is a linear isometry and  $L$  is a self-adjoint (respectively skew-adjoint) operator on  $W$ , then  $T^{-1} \circ L \circ T$  is a self-adjoint (resp. skew-adjoint) operator on  $V$ .

**Proposition 2.6.2.** *Let  $V$  be an inner product space with ordered orthonormal basis  $\mathcal{B} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$ , and  $L$  an operator on  $V$ . Then  $L$  is an isometry if and only if its matrix  $P$  with respect to  $\mathcal{B}$  is orthogonal, meaning that  $P^{-1} = P^T$ .*

*Proof.* By Exercise 2.16, the isomorphism  $V \rightarrow \mathbb{R}^n$  that maps a vector to its coordinate vector in the basis  $\mathcal{B}$  is a linear isometry, so that

$$\langle \mathbf{u}, \mathbf{v} \rangle = [\mathbf{u}]_{\mathcal{B}}^T [\mathbf{v}]_{\mathcal{B}}, \quad \mathbf{u}, \mathbf{v} \in V.$$

Recall that the *Kronecker delta*  $\delta_{ij}$  is defined to be 1 when  $i = j$  and 0 otherwise. If  $L$  is an isometry, then

$$\delta_{ij} = \langle \mathbf{u}_i, \mathbf{u}_j \rangle = \langle L\mathbf{u}_i, L\mathbf{u}_j \rangle = [L\mathbf{u}_i]_{\mathcal{B}}^T [L\mathbf{u}_j]_{\mathcal{B}}.$$

The last term in the above identity is the  $(i, j)$ -th entry of  $P^T P$ , which shows that  $P$  is orthogonal.

Conversely, if  $P^T P = I_n$ , then

$$\delta_{ij} = [L\mathbf{u}_i]_{\mathcal{B}}^T [L\mathbf{u}_j]_{\mathcal{B}} = \langle L\mathbf{u}_i, L\mathbf{u}_j \rangle,$$

so that  $\langle \mathbf{u}_i, \mathbf{u}_j \rangle = \langle L\mathbf{u}_i, L\mathbf{u}_j \rangle$ . Thus, given  $\mathbf{v}, \mathbf{w} \in V$ ,

$$\begin{aligned} \langle L\mathbf{v}, L\mathbf{w} \rangle &= \left\langle L \left( \sum_i \langle \mathbf{v}, \mathbf{u}_i \rangle \mathbf{u}_i \right), L \left( \sum_j \langle \mathbf{w}, \mathbf{u}_j \rangle \mathbf{u}_j \right) \right\rangle \\ &= \sum_{i,j} \langle \mathbf{v}, \mathbf{u}_i \rangle \langle \mathbf{w}, \mathbf{u}_j \rangle \langle L\mathbf{u}_i, L\mathbf{u}_j \rangle = \sum_{i,j} \langle \mathbf{v}, \mathbf{u}_i \rangle \langle \mathbf{w}, \mathbf{u}_j \rangle \delta_{ij} \\ &= \sum_i \langle \mathbf{v}, \mathbf{u}_i \rangle \langle \mathbf{w}, \mathbf{u}_i \rangle = [\mathbf{v}]_{\mathcal{B}}^T [\mathbf{w}]_{\mathcal{B}} \\ &= \langle \mathbf{v}, \mathbf{w} \rangle, \end{aligned}$$

which shows that  $L$  is an isometry. □

For example, the operator on  $\mathbb{R}^2$  consisting of counterclockwise rotation by angle  $\theta$  about the origin is a linear isometry. This can be checked by either noticing that such a rotation leaves the norm of vectors unchanged, or by recalling that the matrix of this operator in the standard basis is

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix},$$

which is clearly orthogonal.

**Definition 2.6.3.** A real number  $\lambda$  is said to be an *eigenvalue* of a linear operator  $L$  on  $V$  if there exists some  $\mathbf{v} \neq \mathbf{0}$  in  $V$  such that  $L\mathbf{v} = \lambda\mathbf{v}$ . In this case,  $\mathbf{v}$  is called an *eigenvector* of  $L$ . When  $L$  is the operator on  $\mathbb{R}^n$  given by left multiplication  $L_A$  by some  $n \times n$  matrix  $A$ , its eigenvalues and eigenvectors are also often called the eigenvalues and eigenvectors of  $A$ .

The reason why one requires  $\mathbf{v} \neq \mathbf{0}$  in the above definition is that otherwise any  $\lambda$  would always be an eigenvalue of every operator  $L$ , since  $L\mathbf{0} = \lambda\mathbf{0}$ . With this restriction, on the other hand, 0 is an eigenvalue of  $L$  if and only if  $L$  has nontrivial kernel, or equivalently (at least if  $V$  is finite-dimensional) if and only if  $L$  is not an isomorphism. Notice also that if  $\lambda$  is an eigenvalue of  $L$ , then the collection of corresponding eigenvectors becomes a subspace of  $V$ , once we add the vector  $\mathbf{0}$  to it. It is called the  $\lambda$ -*eigenspace* of  $L$ .

In general, an operator need not admit any eigenvalues. One such is a rotation by, say,  $\pi/4$  in the plane. The situation is entirely different when  $L$  is self-adjoint, however:

**Theorem 2.6.1.** *Any self-adjoint operator on an inner product space  $V$  admits an eigenvalue.*

*Proof.* Let us assume for now that  $V$  is  $\mathbb{R}^n$  with the standard inner product, and let  $L$  denote the operator. The function

$$\begin{aligned} f : \mathbb{R}^n &\rightarrow \mathbb{R}, \\ \mathbf{v} &\mapsto \langle L\mathbf{v}, \mathbf{v} \rangle \end{aligned}$$

is continuous (in fact, differentiable, being a composition of differentiable maps), and therefore admits a maximum when restricted to the unit sphere  $S^{n-1}$  by compactness of the sphere. Denote by  $\mathbf{u}$  a point of the sphere at which  $f$  takes on this maximum value. We claim that  $\mathbf{u}$  is an eigenvalue of  $L$ . To see this, it is enough to show that  $L\mathbf{u}$  is a multiple of  $\mathbf{u}$ , or equivalently, that  $L\mathbf{u}$  is orthogonal to  $\mathbf{v}$  for any  $\mathbf{v}$  orthogonal to  $\mathbf{u}$ . Notice that we only need to check this for any  $\mathbf{v}$  of norm 1 orthogonal to  $\mathbf{u}$ . Now, the curve  $\mathbf{c} : \mathbb{R} \rightarrow S^{n-1}$ , where  $\mathbf{c}(t) = \cos t\mathbf{u} + \sin t\mathbf{v}$ , has as image the great circle through  $\mathbf{u}$  and  $\mathbf{v}$ , and passes through  $\mathbf{u}$  at 0. Since  $f \circ \mathbf{c}$  is differentiable and has a maximum at 0, its derivative at 0 vanishes. But by the chain rule together with Corollary 2.2.1 and the fact that the derivative of an operator is the operator itself,

$$\begin{aligned} (f \circ \mathbf{c})'(t) &= \langle L \circ \mathbf{c}, \mathbf{c} \rangle'(t) = \langle (L \circ \mathbf{c})'(t), \mathbf{c}(t) \rangle + \langle L(\mathbf{c}(t)), \mathbf{c}'(t) \rangle \\ &= \langle L(\mathbf{c}'(t)), \mathbf{c}(t) \rangle + \langle L(\mathbf{c}(t)), \mathbf{c}'(t) \rangle = 2\langle L(\mathbf{c}(t)), \mathbf{c}'(t) \rangle. \end{aligned}$$

Thus,

$$0 = (f \circ \mathbf{c})'(0) = 2\langle L(\mathbf{c}(0)), \mathbf{c}'(0) \rangle = 2\langle L\mathbf{u}, \mathbf{v} \rangle,$$

which establishes the theorem in the case when  $V$  is Euclidean space. In the general case, consider an orthonormal basis  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of  $V$ . The map  $T : V \rightarrow \mathbb{R}^n$  that maps a vector to its coordinate vector in this basis is then a linear isometry, and as remarked

earlier,  $T \circ L \circ T^{-1}$  is now a self-adjoint operator on Euclidean space. If  $\mathbf{u} \in \mathbb{R}^n$  is a  $\lambda$ -eigenvector of  $T \circ L \circ T^{-1}$ , then  $\mathbf{v} = T^{-1}\mathbf{u}$  is a  $\lambda$ -eigenvector of  $L$ , because

$$L\mathbf{v} = L \circ T^{-1}\mathbf{u} = T^{-1}(T \circ L \circ T^{-1}\mathbf{u}) = T^{-1}(\lambda\mathbf{u}) = \lambda T^{-1}\mathbf{u} = \lambda\mathbf{v}. \quad \square$$

**Theorem 2.6.2** (Spectral theorem). *Let  $L$  be a self-adjoint operator on an inner product space  $V$ . Then  $V$  has an orthonormal basis consisting of eigenvectors of  $L$ . The matrix of  $L$  with respect to this basis is diagonal, with the diagonal entries being the eigenvalues of  $L$ .*

*Proof.* The argument will be by induction on the dimension  $n$  of  $V$ . If  $n = 1$ , then any unit vector in  $V$  is an eigenvector and forms an orthonormal basis. Assume the statement holds for any space of dimension smaller than  $n$ , and consider an  $n$ -dimensional inner product space  $V$  with self-adjoint operator  $L$ . By Theorem 2.6.1,  $L$  admits an eigenvalue  $\lambda$  with corresponding eigenvector  $\mathbf{u}$ , which may be assumed to have norm 1. Observe that  $L$  restricts to an operator on the subspace  $\mathbf{u}^\perp$  of  $V$  consisting of all vectors orthogonal to  $\mathbf{u}$ ; i.e.,  $L(\mathbf{u}^\perp) \subset \mathbf{u}^\perp$ : indeed, if  $\mathbf{v} \in \mathbf{u}^\perp$ , then

$$\langle L\mathbf{v}, \mathbf{u} \rangle = \langle \mathbf{v}, L\mathbf{u} \rangle = \langle \mathbf{v}, \lambda\mathbf{u} \rangle = \lambda \langle \mathbf{v}, \mathbf{u} \rangle = 0,$$

so that  $L\mathbf{v} \in \mathbf{u}^\perp$  as claimed. Since  $L$  remains self-adjoint when restricted to  $\mathbf{u}^\perp$ , the induction hypothesis asserts that this subspace admits an orthonormal basis of eigenvectors. Adding  $\mathbf{u}$  to this set yields an orthonormal basis of  $V$  that consists of eigenvectors.

The last statement of the theorem is clear: if  $\mathcal{B} = \mathbf{v}_1, \dots, \mathbf{v}_n$  is a basis of  $V$  with  $L\mathbf{v}_i = \lambda_i\mathbf{v}_i$ , then  $[L\mathbf{v}_i]_{\mathcal{B}}$  is the vector in  $\mathbb{R}^n$  with  $\lambda_i$  in the  $i$ -th entry and 0 elsewhere, and by definition, the matrix of  $L$  in  $\mathcal{B}$  is

$$\begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}. \quad \square$$

The above theorem does not provide a constructive way of finding eigenvalues and eigenvectors of a self-adjoint operator  $L$ . However,  $\mathbf{u}$  is a  $\lambda$ -eigenvector of  $L$  iff  $L\mathbf{u} = \lambda\mathbf{u}$ ; this is equivalent to the equation

$$(L - \lambda 1_V)\mathbf{u} = \mathbf{0}.$$

In other words, to determine the eigenvalues of  $L$ , it suffices to determine those values of  $\lambda$  for which the operator  $L - \lambda 1_V$  has nontrivial kernel. The nonzero elements in that kernel then correspond to the eigenvectors. Notice that  $L - \lambda 1_V$  has nontrivial kernel if and only if its determinant is zero. We illustrate the procedure in the following:

**Example 2.6.1.** Consider the operator  $L$  on  $\mathbb{R}^3$ , where

$$L \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x + z \\ y + z \\ x + y + z \end{bmatrix}.$$

Thus,  $L$  is left multiplication by the matrix

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix},$$

and since  $A$  is symmetric,  $L$  is self-adjoint by Proposition 2.6.1. Now, expanding along the first row, we obtain

$$\begin{aligned} \det(L - \lambda 1_{\mathbb{R}^3}) &= \det(A - \lambda I_3) = \det \begin{bmatrix} 1 - \lambda & 0 & 1 \\ 0 & 1 - \lambda & 1 \\ 1 & 1 & 1 - \lambda \end{bmatrix} \\ &= (1 - \lambda) \left( (1 - \lambda)^2 - 1 \right) - (1 - \lambda) = (1 - \lambda) \left( (1 - \lambda)^2 - 2 \right) \\ &= (1 - \lambda)(1 - \lambda - \sqrt{2})(1 - \lambda + \sqrt{2}). \end{aligned}$$

Setting this determinant equal to zero implies that  $L$  has  $1$ ,  $1 - \sqrt{2}$ , and  $1 + \sqrt{2}$  as eigenvalues. To find an eigenvector corresponding to the eigenvalue  $1$ , we must find a nontrivial vector in the kernel of  $L - \lambda 1_{\mathbb{R}^3}$  with  $\lambda = 1$ . The equation  $(L - 1_{\mathbb{R}^3})\mathbf{u} = \mathbf{0}$  is equivalent to  $(A - I_3)\mathbf{u} = \mathbf{0}$ ; i.e., to

$$(A - I_3)\mathbf{u} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

Thus,  $x + y = z = 0$ , and the  $1$ -eigenvectors are the nonzero multiples of  $[1 \ -1 \ 0]^T$ .

The same procedure works for  $\lambda = 1 - \sqrt{2}$ :

$$(A - (1 - \sqrt{2})I_3)\mathbf{u} = \begin{bmatrix} \sqrt{2} & 0 & 1 \\ 0 & \sqrt{2} & 1 \\ 1 & 1 & \sqrt{2} \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix},$$

which yields  $\sqrt{2}x + z = 0$ ,  $\sqrt{2}y + z = 0$ , and  $x + y + \sqrt{2}z = 0$ . If one multiplies the first two equations by  $1/\sqrt{2}$  and adds them, the resulting equation is the third one, so the latter may be discarded. The solutions to the first two are all multiples of  $[1 \ 1 \ -\sqrt{2}]^T$ . Repeating this procedure with  $\lambda = 1 + \sqrt{2}$  yields all nontrivial multiples of  $[1 \ 1 \ \sqrt{2}]^T$  as eigenvectors.

The polynomial  $\det(A - \lambda I_n)$  whose roots are the eigenvalues of  $A$  is called the *characteristic polynomial* of  $A$ .

**Corollary 2.6.1.** Any symmetric matrix  $A$  is conjugate to a diagonal matrix  $D$ ; i.e., there exists an invertible matrix  $P$  such that  $D = P^{-1}AP$ . In fact,  $P$  may be chosen to be orthogonal.

*Proof.* As remarked already in the previous example, Proposition 2.6.1 applied to  $\mathbb{R}^n$  with its standard basis  $\mathcal{S}$  and left multiplication  $L_A$  by  $A$  implies that the operator  $L_A$  is self-adjoint. By the spectral theorem, there exists an orthonormal basis  $\mathcal{B}$  of eigenvectors of  $L_A$ , and in that basis, the matrix  $[L_A]_{\mathcal{B}}$  of  $L_A$  is some diagonal matrix  $D$ . (1.2.1) now says that the matrices of  $L_A$  in the two bases satisfy

$$D = [L_A]_{\mathcal{B}} = P^{-1}[L_A]_{\mathcal{S}}P = P^{-1}AP,$$

where  $P = [1_{\mathbb{R}^n}]_{\mathcal{B},\mathcal{S}}$  is the change of basis matrix from  $\mathcal{B}$  to  $\mathcal{S}$ . This means that  $P$  has as columns the vectors  $\mathbf{u}_1, \dots, \mathbf{u}_n$  of  $\mathcal{B}$ , and thus, the  $(i, j)$ -th entry of  $P^T P$  is  $\mathbf{u}_i^T \mathbf{u}_j = \langle \mathbf{u}_i, \mathbf{u}_j \rangle = \delta_{ij}$ . In other words  $P$  is orthogonal.  $\square$

**Example 2.6.2.** Consider the matrix  $A$  from Example 2.6.1. Normalizing the eigenvectors of  $L_A$  so they have length 1, the argument from the above corollary implies that

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 - \sqrt{2} & 0 \\ 0 & 0 & 1 + \sqrt{2} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} & 0 \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{-\sqrt{2}} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{\sqrt{2}} \end{bmatrix} A \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{2} & \frac{1}{2} \\ \frac{-1}{\sqrt{2}} & \frac{1}{2} & \frac{1}{2} \\ 0 & \frac{1}{-\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}.$$

Observe that the order in which the eigenvalues appear in the diagonal matrix is the same as the order in which the eigenvectors appear in the last matrix on the right.

**Definition 2.6.4.** A symmetric bilinear form or scalar product on a vector space  $V$  is a map  $b : V \times V \rightarrow \mathbb{R}$  such that  $b(\mathbf{u}, \mathbf{v}) = b(\mathbf{v}, \mathbf{u})$  and  $b(a\mathbf{u} + \mathbf{v}, \mathbf{w}) = a b(\mathbf{u}, \mathbf{w}) + b(\mathbf{v}, \mathbf{w})$  for all  $\mathbf{u}, \mathbf{v}, \mathbf{w} \in V$  and  $a \in \mathbb{R}$ .

Notice that an inner product is just a symmetric bilinear form that is *positive definite*; i.e.,  $b(\mathbf{u}, \mathbf{u}) > 0$  if  $\mathbf{u} \neq \mathbf{0}$ .  $b$  is said to be *negative definite* if  $-b$  is positive definite. A bilinear form that is either positive definite or negative definite is said to be *definite*.

It is quite easy to construct examples of scalar products:

**Example 2.6.3.** Let  $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  denote a basis of a vector space  $V$ . Any symmetric matrix  $A$  induces a symmetric bilinear form  $b$  on  $V$  by defining

$$b(\mathbf{u}, \mathbf{v}) = [\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}}, \quad \mathbf{u}, \mathbf{v} \in V,$$

where as usual we identify the  $1 \times 1$  matrix  $[\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}}$  on the right with its single entry. The fact that  $b$  as defined above is bilinear follows immediately from properties of matrix operations. Symmetry in turn follows from the fact that a  $1 \times 1$  matrix is always symmetric, so that

$$b(\mathbf{u}, \mathbf{v}) = [\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}} = ([\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}})^T = [\mathbf{v}]_{\mathcal{B}}^T A^T [\mathbf{u}]_{\mathcal{B}}^{TT} = [\mathbf{v}]_{\mathcal{B}}^T A [\mathbf{u}]_{\mathcal{B}} = b(\mathbf{v}, \mathbf{u}).$$

It turns out that the above is essentially the only example possible:

**Theorem 2.6.3.** *Let  $b$  denote a symmetric bilinear form on  $V$ . Given any basis  $\mathcal{B}$  of  $V$ , there exists a unique symmetric matrix  $A$  such that  $b(\mathbf{u}, \mathbf{v}) = [\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}}$  for all  $\mathbf{u}, \mathbf{v} \in V$ .*

*Proof.* Uniqueness is clear, because if  $\mathbf{v}_1, \dots, \mathbf{v}_n$  are the (ordered) basis vectors, then the above formula implies that

$$b(\mathbf{v}_i, \mathbf{v}_j) = [\mathbf{v}_i]_{\mathcal{B}}^T A [\mathbf{v}_j]_{\mathcal{B}} = \mathbf{e}_i^T A \mathbf{e}_j = a_{ij},$$

which also shows symmetry of  $A$ . To check existence, define the  $(i, j)$ -th entry of  $A$  to be  $b(\mathbf{v}_i, \mathbf{v}_j)$ . By linearity of  $b$  in each component, given  $\mathbf{u} = \sum_i x_i \mathbf{v}_i$  and  $\mathbf{v} = \sum_j y_j \mathbf{v}_j$ ,

$$\begin{aligned} b(\mathbf{u}, \mathbf{v}) &= b\left(\sum_i x_i \mathbf{v}_i, \sum_j y_j \mathbf{v}_j\right) = \sum_{ij} x_i b(\mathbf{v}_i, \mathbf{v}_j) y_j = \sum_{ij} x_i a_{ij} y_j \\ &= [\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}}. \end{aligned} \quad \square$$

The theorem above can also be proved in a more elegant way, by observing that there is a one-to-one correspondence between scalar products and self-adjoint operators on an inner product space: If  $L$  is self-adjoint, then the formula

$$b(\mathbf{u}, \mathbf{v}) = \langle L\mathbf{u}, \mathbf{v} \rangle \tag{2.6.1}$$

defines a scalar product. Conversely, given a symmetric bilinear form  $b$  on an inner product space  $V$ , define  $L : V \rightarrow V$  as follows: for  $\mathbf{u} \in V$ , the map  $\mathbf{v} \mapsto b(\mathbf{u}, \mathbf{v})$  defines an element of the dual space  $V^*$ , so that by Corollary 1.4.2, there is unique element  $L\mathbf{u} \in V$  satisfying (2.6.1). The fact that  $L$  is linear and self-adjoint is immediate. The theorem now follows by observing that

$$[\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}} = \langle \mathbf{u}, L\mathbf{v} \rangle.$$

The spectral theorem yields a convenient way of expressing scalar products: Select any basis  $\mathcal{B}$  of  $V$  to obtain a matrix  $A$  such that  $b(\mathbf{u}, \mathbf{v}) = [\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}}$  as in the theorem. Since  $A$  is symmetric, there exists, by Corollary 2.6.1, a diagonal matrix  $D$  and an orthogonal matrix  $Q$  (set  $Q$  equal to  $P^T$  in the corollary) such that  $A = Q^T D Q$ . Since  $Q$  is invertible, the vectors  $\mathbf{u}_i = \sum_j q_{ij} \mathbf{v}_j$ ,  $i = 1, \dots, n$ , form a basis  $\mathcal{C}$  of  $V$ , and by the change of basis formula,

$$[\mathbf{u}]_{\mathcal{C}} = Q[\mathbf{u}]_{\mathcal{B}}, \quad \mathbf{u} \in V.$$

Thus, if  $\mathbf{u} = \sum_i x_i \mathbf{u}_i$ ,  $\mathbf{v} = \sum_j y_j \mathbf{u}_j$ , and  $D$  is the diagonal matrix with entries  $\lambda_1, \dots, \lambda_n$ , then

$$\begin{aligned} b(\mathbf{u}, \mathbf{v}) &= [\mathbf{u}]_{\mathcal{B}}^T A [\mathbf{v}]_{\mathcal{B}} = [\mathbf{u}]_{\mathcal{B}}^T Q^T D Q [\mathbf{v}]_{\mathcal{B}} = (Q[\mathbf{u}]_{\mathcal{B}})^T D (Q[\mathbf{v}]_{\mathcal{B}}) = [\mathbf{u}]_{\mathcal{C}}^T D [\mathbf{v}]_{\mathcal{C}} \\ &= \sum_i \lambda_i x_i y_i. \end{aligned}$$



In particular,

$$b(\mathbf{u}, \mathbf{u}) = \sum_i \lambda_i x_i^2,$$

so that  $b$  is positive-(respectively negative-)definite if and only if all the eigenvalues of  $A$  are positive (respectively) negative.

**Example 2.6.4.** Consider the symmetric bilinear form on  $\mathbb{R}^2$  given by

$$b\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) = 3x_1y_1 - x_1y_2 - x_2y_1 + 3x_2y_2.$$

The matrix of  $b$  in the standard basis is

$$A = \begin{bmatrix} 3 & -1 \\ -1 & 3 \end{bmatrix}.$$

Its characteristic polynomial is

$$\det(A - \lambda I_2) = (3 - \lambda)^2 - 1 = (4 - \lambda)(2 - \lambda),$$

which has roots 4 and 2. Thus,  $b$  is positive definite. Using the same procedure as that in the previous example, one easily computes that  $\mathbf{v}_1 = [1/\sqrt{2} \quad 1/\sqrt{2}]^T$  and  $\mathbf{v}_2 = [-1/\sqrt{2} \quad 1/\sqrt{2}]^T$  are unit eigenvectors corresponding to the eigenvalues 2 and 4 respectively. If  $\mathbf{u} = x_1\mathbf{v}_1 + x_2\mathbf{v}_2$  and  $\mathbf{v} = y_1\mathbf{v}_1 + y_2\mathbf{v}_2$ , then  $b(\mathbf{u}, \mathbf{v}) = 2x_1y_1 + 4x_2y_2$ .

## 2.7 Taylor polynomials and extreme values

A point  $\mathbf{a} \in U \subset \mathbb{R}^n$  is said to be a *critical point* of a map  $\mathbf{f} : U \rightarrow \mathbb{R}^m$  if  $D\mathbf{f}(\mathbf{a}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is either not onto or does not exist. In this case,  $\mathbf{f}(\mathbf{a})$  is called a *critical value* of  $\mathbf{f}$ . Notice that if  $m = 1$ , then  $\mathbf{a}$  is critical if and only if  $Df(\mathbf{a}) = 0$  or does not exist. If in addition,  $n = 1$ , we recover the classical notion of critical point from one-variable Calculus. In the case of a real-valued function  $f : U \rightarrow \mathbb{R}$ ,  $f$  is said to have a local *maximum* (resp. *minimum*) at  $\mathbf{a}$  if there exists a neighborhood  $V$  of  $\mathbf{a}$  such that  $f(\mathbf{p}) \leq$  (resp.  $\geq$ )  $f(\mathbf{a})$  for all  $\mathbf{p} \in V$ . A point  $\mathbf{a}$  where  $f$  has a local maximum or minimum is called an *extremum*, and  $f(\mathbf{a})$  is then said to be an *extreme value*.

**Theorem 2.7.1.** *If  $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}$  has an extremum at  $\mathbf{a}$ , then  $\mathbf{a}$  is a critical point of  $f$ .*

*Proof.* Let  $\mathbf{a}$  be an extremum of  $f$ . We may assume that  $f$  is differentiable at  $\mathbf{a}$ , for otherwise the conclusion certainly holds. In particular,  $f$  is defined on some open neighborhood of  $\mathbf{a}$ . We wish to show that for any  $\mathbf{h} \in \mathbb{R}^n$ ,  $Df(\mathbf{a})\mathbf{h} = 0$ . Choose a small enough interval  $I$  around 0 such that the curve  $\mathbf{c} : I \rightarrow \mathbb{R}^n$ , where  $\mathbf{c}(t) = \mathbf{a} + t\mathbf{h}$ , has its image inside  $U$ . Then the function  $g := f \circ \mathbf{c} : I \rightarrow \mathbb{R}$  of one variable has an extremum at 0, and being differentiable there,

$$0 = g'(0) = (f \circ \mathbf{c})'(0) = Df(\mathbf{c}(0))\mathbf{c}'(0) = Df(\mathbf{a})\mathbf{h}. \quad \square$$

The converse is already false when  $n = 1$ , as the function  $x \mapsto x^3$  shows: 0 is a critical point of  $f$ , but not an extremum. Recall that for functions of one variable, if  $f'(a) = 0$ , then  $a$  is a local minimum (resp. maximum) if  $f''(a) > 0$  (resp.  $< 0$ ). Now, a function defined on  $\mathbb{R}^n$  does not have one second derivative if  $n > 1$ , but it does have an  $n \times n$  matrix  $[D_{ij}f(\mathbf{a})]$  of second partial derivatives, which is symmetric if  $f$  is  $C^2$ . We will see that if the corresponding symmetric bilinear form is positive (resp. negative) definite, then  $f$  has a local minimum (resp. maximum) at  $\mathbf{a}$ . In order to show this, we first generalize Taylor polynomials to functions defined on Euclidean space.

**Definition 2.7.1.** Let  $f : U \rightarrow \mathbb{R}$  be a function of class  $C^k$  on an open set  $U \subset \mathbb{R}^n$ . The *Taylor polynomial of degree  $k$*  of  $f$  at  $\mathbf{a} \in U$  is the function  $P_{f,\mathbf{a},k} : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$P_{f,\mathbf{a},k} = \sum_{j=0}^k \frac{1}{j!} \sum_{1 \leq i_1, \dots, i_j \leq n} D_{i_1 \dots i_j} f(\mathbf{a}) u^{i_1} u^{i_2} \cdots u^{i_j}.$$

Thus, for  $\mathbf{p} = (x_1, \dots, x_n) \in \mathbb{R}^n$ ,

$$P_{f,\mathbf{a},k}(\mathbf{p}) = \sum_{j=0}^k \frac{1}{j!} \sum_{1 \leq i_1, \dots, i_j \leq n} D_{i_1 \dots i_j} f(\mathbf{a}) x_{i_1} x_{i_2} \cdots x_{i_j},$$

with  $D_0 f = f$ .

**Example 2.7.1.** If  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  is defined by  $f(x, y) = e^{xy}$ , then  $D_1 f(x, y) = ye^{xy}$ ,  $D_2 f(x, y) = xe^{xy}$ ,  $D_{12} f(x, y) = D_{21} f(x, y) = (1 + xy)e^{xy}$ ,  $D_{11} f(x, y) = y^2 e^{xy}$ , and  $D_{22} f(x, y) = x^2 e^{xy}$ . Thus the Taylor polynomial of degree 2 of  $f$  at  $\mathbf{0}$  is

$$P_{f,\mathbf{0},2}(x, y) = 1 + xy.$$

**Theorem 2.7.2** (Taylor's theorem). Suppose  $f : U \rightarrow \mathbb{R}$  is of class  $C^{k+1}$  on the open set  $U \subset \mathbb{R}^n$ , and let  $\mathbf{a} \in U$ ,  $\varepsilon > 0$  small enough so that the closed ball of radius  $\varepsilon$  around  $\mathbf{a}$  is contained in  $U$ . Define a function  $r$  on this ball by setting

$$r(\mathbf{p}) = f(\mathbf{a} + \mathbf{p}) - P_{f,\mathbf{a},k}(\mathbf{p}).$$

Then the remainder  $r$  satisfies

$$\lim_{\mathbf{p} \rightarrow \mathbf{0}} \frac{r(\mathbf{p})}{|\mathbf{p}|^k} = 0.$$

*Proof.* Fix any  $\mathbf{p} = (x_1, \dots, x_n) \in \mathbb{R}^n$  of norm less than  $\varepsilon$ . The curve  $\mathbf{c} : [0, 1] \rightarrow \mathbb{R}^n$ ,  $\mathbf{c}(t) = \mathbf{a} + t\mathbf{p}$ , has its image in  $U$ , so we have a well-defined function  $g := f \circ \mathbf{c}$  of class  $C^{k+1}$ . By the chain rule,

$$g'(t) = Df(\mathbf{c}(t))\mathbf{c}'(t) = \sum_{i=1}^n D_i f(\mathbf{c}(t))x_i,$$

and arguing inductively, we obtain

$$g^{(j)}(t) = \sum_{1 \leq i_1, \dots, i_j \leq n} D_{i_1 \dots i_j} f(\mathbf{c}(t))x_{i_1} \cdots x_{i_j}$$

for  $j = 1, \dots, k + 1$ . Taylor's theorem for functions of one variable asserts that there exists some  $t \in (0, 1)$  such that

$$g(1) = \sum_{j=0}^k \frac{g^{(j)}(0)}{j!} + \frac{g^{(k+1)}(t)}{(k+1)!};$$

i.e.,

$$\begin{aligned} f(\mathbf{a} + \mathbf{p}) &= \sum_{j=0}^k \frac{1}{j!} \sum_{1 \leq i_1, \dots, i_j \leq n} D_{i_1 \dots i_j} f(\mathbf{a}) x_{i_1} \cdots x_{i_j} + \frac{g^{(k+1)}(t)}{(k+1)!} \\ &= P_{f, \mathbf{a}, k}(\mathbf{p}) + r(\mathbf{p}), \end{aligned}$$

where

$$r(\mathbf{p}) = \frac{1}{(k+1)!} \sum D_{i_1 \dots i_{k+1}} f(\mathbf{a} + t\mathbf{p}) x_{i_1} \cdots x_{i_{k+1}},$$

and  $t = t(\mathbf{p})$  depends on  $\mathbf{p}$ . It only remains to show that  $r(\mathbf{p})/|\mathbf{p}|^k \rightarrow 0$  as  $\mathbf{p} \rightarrow \mathbf{0}$ . But  $f$  is of class  $C^{k+1}$ , so that each  $D_{i_1 \dots i_{k+1}} f(\mathbf{a} + t\mathbf{p})$  converges to the value of that function at  $\mathbf{a}$ . Furthermore, each  $|x_i| \leq |\mathbf{p}|$ , so that

$$\frac{|x_{i_1} \cdots x_{i_{k+1}}|}{|\mathbf{p}|^k} \leq \min\{|x_{i_1}|, \dots, |x_{i_{k+1}}|\}, \quad (2.7.1)$$

and the right side goes to 0 as  $\mathbf{p} \rightarrow \mathbf{0}$ . This proves the assertion.  $\square$

**Remark 2.7.1.** Since the right side of the inequality 2.7.1 is no larger than  $|\mathbf{p}|$ , for each sufficiently small neighborhood  $V$  of  $\mathbf{a}$ , there exists a constant  $\alpha$  that depends on  $f$  and  $V$  such that  $|r(\mathbf{p})| < \alpha |\mathbf{p}|^{k+1}$  for all  $\mathbf{p}$  such that  $\mathbf{a} + \mathbf{p} \in V$ . Furthermore, if  $C$  is any compact set in  $U$ , uniform continuity of the derivatives of order  $k + 1$  of  $f$  on  $C$  implies that there exists  $\beta > 0$  that depends on  $C$  and  $f$  such that

$$|r(\mathbf{p}, \mathbf{h})| := |f(\mathbf{p} + \mathbf{h}) - P_{f, \mathbf{p}, k}(\mathbf{h})| \leq \beta |\mathbf{h}|^{k+1} \text{ for all } \mathbf{p}, \mathbf{p} + \mathbf{h} \in C.$$

**Definition 2.7.2.** Let  $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}$  be of class  $C^2$ , and  $\mathbf{a}$  an interior point of  $U$ . The Hessian  $H_f(\mathbf{a})$  of  $f$  at  $\mathbf{a}$  is the symmetric operator on  $\mathbb{R}^n$  given by left multiplication by the matrix  $[D_{ij}f(\mathbf{a})]$  of second derivatives of  $f$  at  $\mathbf{a}$ . The Hessian form  $h_f(\mathbf{a})$  of  $f$  at  $\mathbf{a}$  is the associated symmetric bilinear form

$$h_f(\mathbf{a})(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \begin{bmatrix} D_{11}f(\mathbf{a}) & \cdots & D_{1n}f(\mathbf{a}) \\ \vdots & \cdots & \vdots \\ D_{n1}f(\mathbf{a}) & \cdots & D_{nn}f(\mathbf{a}) \end{bmatrix} \mathbf{v}, \quad \mathbf{u}, \mathbf{v} \in \mathbb{R}^n.$$

Observe that  $h_f(\mathbf{a})(\mathbf{u}, \mathbf{u})$  equals two times the term corresponding to  $j = 2$  in the Taylor polynomial of  $f$  at  $\mathbf{a}$ . In particular, if a  $C^2$  function  $f$  has a critical point at  $\mathbf{a}$ , then its Taylor polynomial of order 2 at  $\mathbf{a}$  is

$$P_{f, \mathbf{a}, 2}(\mathbf{u}) = f(\mathbf{a}) + \frac{1}{2} h_f(\mathbf{a})(\mathbf{u}, \mathbf{u}).$$

This fact is used to prove the following:

**Theorem 2.7.3.** Suppose  $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}$  is of class  $\mathcal{C}^2$ ,  $\mathbf{a} \in U$ , and  $Df(\mathbf{a}) = 0$ .

- (1) If the Hessian form  $h_f(\mathbf{a})$  of  $f$  at  $\mathbf{a}$  is positive (resp. negative) definite, then  $f$  has a local minimum (resp. maximum) at  $\mathbf{a}$ .
- (2) If the Hessian  $H_f(\mathbf{a})$  has at least one positive and one negative eigenvalues, then  $f$  does not have a local maximum nor a local minimum at  $\mathbf{a}$ .  $f$  is then said to have a saddle point at  $\mathbf{a}$ .

*Proof.* (1): We will consider the case when the Hessian is positive definite. The other case follows by looking at  $-f$ . By assumption, the function  $\mathbf{u} \mapsto h_f(\mathbf{a})(\mathbf{u}, \mathbf{u})$  assumes a positive minimum value  $2\alpha$  when restricted to the unit sphere  $S^{n-1}$ , since the latter is compact. Let  $0 < \varepsilon < \alpha$ . By Taylor's theorem, there is  $\delta > 0$  such that for  $0 < |\mathbf{u}| < \delta$ ,

$$\frac{|r(\mathbf{u})|}{|\mathbf{u}|^2} = \left| \frac{f(\mathbf{a} + \mathbf{u}) - f(\mathbf{a})}{|\mathbf{u}|^2} - \frac{1}{2}h_f(\mathbf{a})\left(\frac{\mathbf{u}}{|\mathbf{u}|}, \frac{\mathbf{u}}{|\mathbf{u}|}\right) \right| < \varepsilon. \quad (2.7.2)$$

Thus, for sufficiently small  $|\mathbf{u}| > 0$ ,

$$\frac{f(\mathbf{a} + \mathbf{u}) - f(\mathbf{a})}{|\mathbf{u}|^2} > \frac{1}{2}h_f(\mathbf{a})\left(\frac{\mathbf{u}}{|\mathbf{u}|}, \frac{\mathbf{u}}{|\mathbf{u}|}\right) - \varepsilon \geq \alpha - \varepsilon > 0,$$

and  $f(\mathbf{a} + \mathbf{u}) > f(\mathbf{a})$  as claimed.

(2): By assumption,  $H_f(\mathbf{a})$  has some positive eigenvalue  $2\alpha$ . We will show that if  $\mathbf{w}$  is a corresponding eigenvector with small enough norm, then  $f(\mathbf{a} + \mathbf{w}) > f(\mathbf{a})$ . It then follows that  $f$  does not have a local maximum at  $\mathbf{a}$ . So choose some  $0 < \varepsilon < \alpha$ . By Taylor's theorem, there exists  $\delta > 0$  such that (2.7.2) holds for any  $\mathbf{u}$  such that  $0 < |\mathbf{u}| < \delta$ . Now, if  $\mathbf{w}$  is a  $2\alpha$ -eigenvector of the Hessian, then

$$h_f(\mathbf{a})\left(\frac{\mathbf{w}}{|\mathbf{w}|}, \frac{\mathbf{w}}{|\mathbf{w}|}\right) = \frac{1}{|\mathbf{w}|^2} \mathbf{w}^T H_f(\mathbf{a}) \mathbf{w} = \frac{1}{|\mathbf{w}|^2} \mathbf{w}^T 2\alpha \mathbf{w} = 2\alpha,$$

so that for  $0 < |\mathbf{w}| < \delta$ ,

$$\frac{f(\mathbf{a} + \mathbf{w}) - f(\mathbf{a})}{|\mathbf{w}|^2} > \frac{1}{2}h_f(\mathbf{a})\left(\frac{\mathbf{w}}{|\mathbf{w}|}, \frac{\mathbf{w}}{|\mathbf{w}|}\right) - \varepsilon = \alpha - \varepsilon > 0.$$

The same argument applied to  $-f$  and any negative eigenvalue of  $H_f(\mathbf{a})$  implies that  $f$  does not have a local minimum at  $\mathbf{a}$ . □

**Examples 2.7.2.** (i) Suppose we are asked to find and classify the critical points of  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f(x, y) = x^2 + y^2 - 2x^2y$ . The Jacobian matrix of  $f$  at  $(x, y)$  is

$$[Df](x, y) = [2x - 4xy \quad 2y - 2x^2] = 2 [x(1 - 2y) \quad y - x^2],$$

and  $(x, y)$  is a critical point of  $f$  if and only if  $x = 0$  or  $y = 1/2$ , and  $y = x^2$ . Thus, there are 3 critical points: the origin, and  $(\pm 1/\sqrt{2}, 1/2)$ . The Hessian of  $f$  has as matrix

$$\begin{bmatrix} 2 - 4y & -4x \\ -4x & 2 \end{bmatrix}.$$

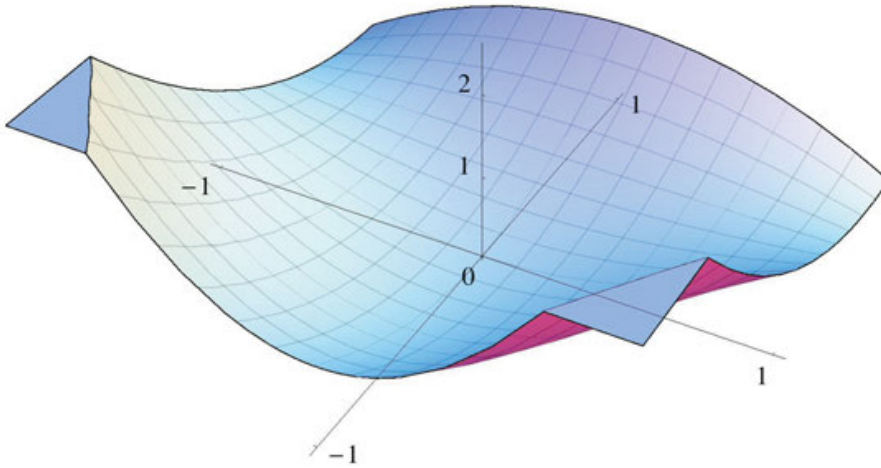


Fig. 2.1:  $z = x^2 + y^2 - 2x^2y$

At the origin, the matrix equals  $2I_2$ , so  $f$  has a local minimum there. At the other two points, the Hessian becomes

$$\begin{bmatrix} 0 & \frac{4}{\sqrt{2}} \\ \frac{4}{\sqrt{2}} & 2 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} 0 & -\frac{4}{\sqrt{2}} \\ -\frac{4}{\sqrt{2}} & 2 \end{bmatrix}.$$

The eigenvalue-finding procedure outlined in the previous section easily yields the same eigenvalues 4 and  $-2$  for both matrices. Thus, both points are saddle points.

- (ii) The classification of critical points in the previous example could also have been done without explicitly finding the eigenvalues of the Hessian. This is because for a  $2 \times 2$  symmetric matrix  $A$ , the sign of the eigenvalues is easily determined: since there is a diagonal matrix of eigenvalues that is similar to  $A$ , the determinant of  $A$  is the product of the eigenvalues, and the trace equals the sum (this is of course true for matrices of any size). Thus, if  $\det A < 0$ , then the eigenvalues have different signs, corresponding to a saddle point. If the determinant is positive, both eigenvalues have the same sign, which coincides with the sign of the  $(1, 1)$  element of the matrix: the reason being that both diagonal terms must have the same sign (otherwise the determinant would not be positive) and their sum equals the sum of the eigenvalues.

Summarizing, if  $f$  is a function of two variables which has  $\mathbf{a}$  as a critical point, and if

$$A = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

denotes the Hessian matrix at  $\mathbf{a}$ , then

- $f$  has a saddle point at  $\mathbf{a}$  if  $\det A < 0$ ;
- $f$  has a local minimum at  $\mathbf{a}$  if  $\det A > 0$  and  $a > 0$ ;
- $f$  has a local maximum at  $\mathbf{a}$  if  $\det A > 0$  and  $a < 0$ ;

Notice that if the determinant is positive, then  $a$  cannot vanish. This means that the above covers all but one possibility, namely the determinant being zero. If this is the case, then no conclusion can be made: for example, the function  $f$  given by  $f(x, y) = x^2y^2$  has a global minimum at the origin but its Hessian there is zero. Taking its negative yields a function with a maximum at the origin with zero Hessian. Finally,  $f$ , where  $f(x, y) = x^3$ , has the whole  $y$ -axis as critical points with vanishing Hessian, but none of them are maxima or minima.

## 2.8 Vector fields

From now on, differentiable maps will be assumed to be smooth (i.e.,  $C^\infty$ ) unless specified otherwise.

One tends to visualize vectors in Euclidean space as arrows whose base point can be assigned arbitrarily. When dealing with vector fields, it is more efficient to have fixed base points. In order to specify a vector with base point at, say,  $\mathbf{p}$ , we use the following concept:

### Definition 2.8.1.

The *tangent space* of  $\mathbb{R}^n$  at  $\mathbf{p} \in \mathbb{R}^n$  is the collection

$$\mathbb{R}_{\mathbf{p}}^n = \{\mathbf{p}\} \times \mathbb{R}^n = \{(\mathbf{p}, \mathbf{u}) \mid \mathbf{u} \in \mathbb{R}^n\}.$$

An element of  $\mathbb{R}_{\mathbf{p}}^n$  is called a *tangent vector* at  $\mathbf{p}$ .  $\mathbb{R}_{\mathbf{p}}^n$  is a vector space with the operations

$$(\mathbf{p}, \mathbf{u}) + (\mathbf{p}, \mathbf{q}) = (\mathbf{p}, \mathbf{u} + \mathbf{v}), \quad a(\mathbf{p}, \mathbf{u}) = (\mathbf{p}, a\mathbf{u}), \quad a \in \mathbb{R}, \mathbf{u}, \mathbf{v} \in \mathbb{R}^n.$$

There is a canonical isomorphism  $\mathcal{I}_{\mathbf{p}} : \mathbb{R}^n \rightarrow \mathbb{R}_{\mathbf{p}}^n$  that maps  $\mathbf{u}$  to  $(\mathbf{p}, \mathbf{u})$ ,  $\mathbf{u} \in \mathbb{R}^n$ . For any  $\mathbf{p}, \mathbf{q} \in \mathbb{R}^n$ , the map

$$\begin{aligned} \mathcal{I}_{\mathbf{q}} \circ \mathcal{I}_{\mathbf{p}}^{-1} : \mathbb{R}_{\mathbf{p}}^n &\rightarrow \mathbb{R}_{\mathbf{q}}^n, \\ (\mathbf{p}, \mathbf{u}) &\mapsto (\mathbf{q}, \mathbf{u}) \end{aligned}$$

is an isomorphism between tangent spaces, called *parallel translation*.

Derivatives can be reformulated in terms of tangent spaces:

**Definition 2.8.2.** If  $\mathbf{f} : U \rightarrow \mathbb{R}^m$  is a differentiable map on an open set  $U \subseteq \mathbb{R}^n$ , the *derivative of  $\mathbf{f}$  at  $\mathbf{p} \in U$*  is the linear transformation

$$\begin{aligned} \mathbf{f}_{*\mathbf{p}} : \mathbb{R}_{\mathbf{p}}^n &\rightarrow \mathbb{R}_{\mathbf{f}(\mathbf{p})}^m \\ (\mathbf{p}, \mathbf{u}) &\mapsto (\mathbf{f}(\mathbf{p}), D\mathbf{f}(\mathbf{p})\mathbf{u}). \end{aligned} \tag{2.8.1}$$

It is worth emphasizing that the only difference between  $\mathbf{f}_*$  and  $D\mathbf{f}$  is one of notation: the former includes the base point of the vector, the latter doesn't. In terms of

the canonical isomorphism between Euclidean space and its tangent space at  $\mathbf{p}$ , the following diagram commutes:

$$\begin{array}{ccc} \mathbb{R}_{\mathbf{p}}^n & \xrightarrow{f_{*\mathbf{p}}} & \mathbb{R}_{f(\mathbf{p})}^m \\ \mathcal{I}_{\mathbf{p}} \uparrow & & \uparrow \mathcal{I}_{f(\mathbf{p})} \\ U & \xrightarrow{Df(\mathbf{p})} & \mathbb{R}^m. \end{array}$$

In order to avoid cumbersome notation, we will often denote  $(\mathbf{p}, \mathbf{e}_i)$  by  $\mathbf{D}_i(\mathbf{p})$ , and  $(t, 1) \in \mathbb{R}_t$  by  $D(t)$  when  $n = 1$ . Thus, the tangent vector  $(\mathbf{p}, \sum a_i \mathbf{e}_i)$  can also be written as  $\sum a_i \mathbf{D}_i(\mathbf{p})$ . Furthermore, when it is clear that  $\mathbf{u}$  belongs to the tangent space of, say,  $\mathbf{p}$ , we often write  $\mathbf{f}_* \mathbf{u}$  instead of  $\mathbf{f}_{*\mathbf{p}} \mathbf{u}$  for brevity.

**Definition 2.8.3.** A vector field on an open set  $U \subset \mathbb{R}^n$  is a map  $\mathbf{X}$  that assigns to each  $\mathbf{p} \in U$  a tangent vector  $\mathbf{X}(\mathbf{p}) \in \mathbb{R}_{\mathbf{p}}^n$  at  $\mathbf{p}$ .

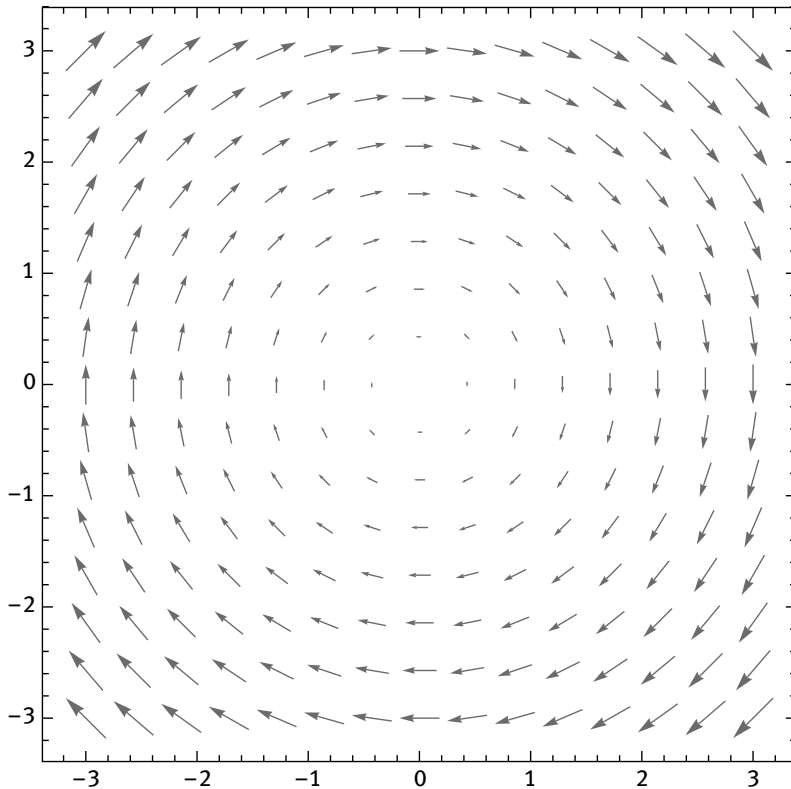


Fig. 2.2: The vector field  $\mathbf{X} = u^2 \mathbf{D}_1 - u^1 \mathbf{D}_2$  in the plane

By definition, vector fields on  $U$  are in one-one correspondence with maps  $\mathbf{f} : U \rightarrow \mathbb{R}^n$ ; for each  $\mathbf{p} \in U$ , the value of  $\mathbf{X}$  at  $\mathbf{p}$  can be written as  $(\mathbf{p}, \mathbf{f}(\mathbf{p}))$  for a unique map  $\mathbf{f} : U \rightarrow \mathbb{R}^n$ , and conversely, any such map defines a vector field  $\mathbf{X}$  by the formula  $\mathbf{X}(\mathbf{p}) = \mathcal{I}_{\mathbf{p}} \mathbf{f}(\mathbf{p})$ . We say the vector field is differentiable or smooth if  $\mathbf{f}$  has that property.

If the vector field  $\mathbf{X}$  is represented by the map  $\mathbf{f}$ , i.e., if  $\mathbf{X}(\mathbf{p}) = \mathcal{I}_{\mathbf{p}}\mathbf{f}(\mathbf{p})$  for all  $\mathbf{p}$  in the domain, then

$$\mathbf{X} = \sum_i f^i \mathbf{D}_i, \quad f^i = u^i \circ \mathbf{f}.$$

For example, the *position vector field*  $\mathbf{P}$  on  $\mathbb{R}^n$ , which is defined by  $\mathbf{P}(\mathbf{p}) = \mathcal{I}_{\mathbf{p}}\mathbf{p}$ , can be written  $\mathbf{P} = \sum_i u^i \mathbf{D}_i$ .

When  $\mathbf{f}$  is a constant map, the vector field is said to be *parallel*; equivalently, such a vector field is obtained by parallel translating its value at any one point to every other point.

It is useful to introduce a more general concept:

**Definition 2.8.4.** Let  $\mathbf{f} : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a map. A *vector field along  $\mathbf{f}$*  is a map  $\mathbf{X}$  that assigns to each  $\mathbf{p} \in U$  an element  $\mathbf{X}(\mathbf{p}) \in \mathbb{R}_{\mathbf{f}(\mathbf{p})}^m$  of the tangent space of  $\mathbb{R}^m$  at  $\mathbf{f}(\mathbf{p})$ . As above, such a field  $\mathbf{X}$  is represented by a unique map  $\mathbf{g} : U \rightarrow \mathbb{R}^m$ , where  $\mathbf{X}(\mathbf{p}) = (\mathbf{f}(\mathbf{p}), \mathbf{g}(\mathbf{p}))$ , and  $\mathbf{X}$  is said to be smooth if  $\mathbf{g}$  is.

A recurring example of such a vector field is the *velocity field* or *tangent field* of a curve  $\mathbf{c} : I \rightarrow \mathbb{R}^n$ : by definition it is the vector field  $\dot{\mathbf{c}}$  along  $\mathbf{c}$  given by

$$\dot{\mathbf{c}}(t) = \mathcal{I}_{\mathbf{c}(t)}\mathbf{c}'(t), \quad t \in I.$$

Alternatively,

$$\dot{\mathbf{c}}(t) = \mathbf{c}_{*t}D(t).$$

For example, if  $\mathbf{c}(t) = (\cos t, -\sin t)$ , then

$$\dot{\mathbf{c}}(t) = -\sin t \mathbf{D}_1 \circ \mathbf{c}(t) - \cos t \mathbf{D}_2 \circ \mathbf{c}(t) = (u^2 \mathbf{D}_1 - u^1 \mathbf{D}_2) \circ \mathbf{c}(t).$$

Notice that if  $\mathbf{X}$  is the vector field from Figure 2.2, then  $\dot{\mathbf{c}} = \mathbf{X} \circ \mathbf{c}$ ; i.e., the “restriction” of  $\mathbf{X}$  to the image of  $\mathbf{c}$  is the velocity field of the curve.

**Definition 2.8.5.** An *integral curve* of a vector field  $\mathbf{X}$  is a curve  $\mathbf{c}$  that satisfies  $\dot{\mathbf{c}} = \mathbf{X} \circ \mathbf{c}$ .

It turns out that for any one point in the domain of a given vector field, there exists an integral curve of that field that passes through the point. This remarkable fact is a consequence of a basic theorem from the theory of ordinary differential equations, a proof of which can be found for example in [9]. The relation between integral curves and differential equations stems from the fact that if  $\mathbf{X}$  is represented by  $\mathbf{f} : U \rightarrow \mathbb{R}^n$ , then  $\mathbf{c}$  is an integral curve of  $\mathbf{X}$  if and only if  $\mathbf{c}' = \mathbf{f} \circ \mathbf{c}$ ; in terms of the components  $x_i = u^i \circ \mathbf{c}$  and  $f^i = u^i \circ \mathbf{f}$  of  $\mathbf{c}$  and  $\mathbf{f}$  respectively, this is equivalent to the system of ordinary differential equations (ODEs)

$$\begin{aligned} x_1'(t) &= f^1(x_1(t), \dots, x_n(t)) \\ &\vdots \\ x_n'(t) &= f^n(x_1(t), \dots, x_n(t)) \end{aligned}$$



We state without proof the existence and uniqueness theorems for integral curves of vector fields:

**Theorem 2.8.1** (Existence of Solutions). *Let  $f : U \rightarrow \mathbb{R}^n$  be a differentiable map, where  $U$  is open in  $\mathbb{R}^n$ . For any  $\mathbf{a} \in U$ , there exists a neighborhood  $W$  of  $\mathbf{a}$ , an interval  $I$  around  $0 \in \mathbb{R}$ , and a differentiable map  $\Psi : I \times W \rightarrow U$  such that*

- (1)  $\Psi(0, \mathbf{u}) = \mathbf{u}$ , and
- (2)  $D\Psi(t, \mathbf{u})\mathbf{e}_1 = f \circ \Psi(t, \mathbf{u})$

for  $t \in I$  and  $\mathbf{u} \in W$ .

Recalling that maps  $f : U \rightarrow \mathbb{R}^n$  are in bijective correspondence with vector fields  $\mathbf{X}$  defined on  $U$ , Theorem 2.8.1 asserts that integral curves  $t \mapsto \mathbf{c}_{\mathbf{u}}(t) := \Psi(t, \mathbf{u})$  exist for arbitrary “initial conditions”  $\mathbf{c}_{\mathbf{u}}(0) = \mathbf{u}$ , that they depend smoothly on the initial conditions, and that at least locally, they can be defined on a fixed common interval. A map  $\Psi$  as in the above theorem is called a *local flow* of  $\mathbf{X}$ .

In the same way, uniqueness of solutions of the above system of differential equations implies uniqueness of integral curves:

**Theorem 2.8.2** (Uniqueness of Solutions). *If  $\mathbf{c}, \tilde{\mathbf{c}} : I \rightarrow U$  are two integral curves of a vector field  $\mathbf{X}$  with  $\mathbf{c}(t_0) = \tilde{\mathbf{c}}(t_0)$  for some  $t_0 \in I$ , then  $\mathbf{c} = \tilde{\mathbf{c}}$ .*

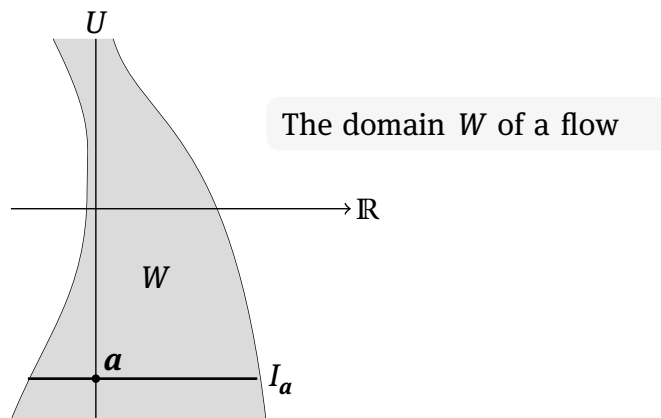
Our next aim is to group all local flows into a single one so that the corresponding integral curves are defined on a maximal interval: For each  $\mathbf{a} \in U$ , let  $I_{\mathbf{a}}$  denote the maximal open interval around 0 on which the (unique by Theorem 2.8.2) integral curve  $\Psi_{\mathbf{a}}$  of  $\mathbf{X}$  that equals  $\mathbf{a}$  at 0 is defined.

**Theorem 2.8.3.** *Given any vector field  $\mathbf{X}$  on  $U \subset \mathbb{R}^n$ , there exists a unique open set  $W \subset \mathbb{R} \times U$  and a unique differentiable map  $\Psi : W \rightarrow U$  such that*

- (1)  $I_{\mathbf{a}} \times \{\mathbf{a}\} = W \cap (\mathbb{R} \times \{\mathbf{a}\})$  for all  $\mathbf{a} \in U$ , and
- (2)  $\Psi(t, \mathbf{a}) = \Psi_{\mathbf{a}}(t)$  if  $(t, \mathbf{a}) \in W$ .

$\Psi$  is called the *flow* of  $\mathbf{X}$ . By (2),  $\{0\} \times U \subset W$ , and (1), (2) of Theorem 2.8.1 are satisfied.

*Proof.* (1) determines  $W$  uniquely, while (2) does the same for  $\Psi$ . It thus remains to show that  $W$  is open, and that  $\Psi$  is differentiable.



Fix  $\mathbf{a} \in U$ , and let  $I$  denote the set of all  $t \in I_{\mathbf{a}}$  for which there exists a neighborhood of  $(t, \mathbf{a})$  contained in  $W$  on which  $\Psi$  is differentiable. We will establish that  $I$  is nonempty, open and closed in  $I_{\mathbf{a}}$ , so that  $I = I_{\mathbf{a}}$  by Proposition 1.7.1:  $I$  is nonempty because  $0 \in I$  by Theorem 2.8.1, and is open by definition. To see that it is closed, consider a point  $t_0$  in the closure  $\bar{I}$  of  $I$ ; by 2.8.1, there exists a local flow  $\Psi' : I' \times V' \rightarrow U$  with  $0 \in I'$  and  $\Psi'_{\mathbf{a}}(t_0) \in V'$ . Let  $t_1 \in I$  be close enough to  $t_0$  that  $t_0 - t_1 \in I'$  (recall that  $t_0$  belongs to the closure of  $I$ ) and  $\Psi'_{\mathbf{a}}(t_1) \in V'$  (by continuity of  $\Psi'_{\mathbf{a}}$ ). Choose an interval  $I_0$  around  $t_0$  such that  $t - t_1 \in I'$  for  $t \in I_0$ . Finally, by continuity of  $\Psi$  at  $(t_1, \mathbf{a})$ , there exists a neighborhood  $V$  of  $\mathbf{a}$  such that  $\Psi(t_1 \times V) \subset V'$ .

We claim that  $\Psi$  is defined and differentiable on  $I_0 \times V$ , so that  $t_0 \in I$ : Indeed, if  $t \in I_0$  and  $\mathbf{u} \in V$ , then by definition of  $I_0$  and  $V$ ,  $t - t_1 \in I'$  and  $\Psi(t_1, \mathbf{u}) \in V'$ , so that  $\Psi'(t - t_1, \Psi(t_1, \mathbf{u}))$  is defined. The curve  $s \mapsto \Psi'(s - t_1, \Psi(t_1, \mathbf{u}))$  is an integral curve of  $\mathbf{X}$  which equals  $\Psi(t_1, \mathbf{u})$  at  $t_1$ . By uniqueness,  $\Psi(t, \mathbf{u}) = \Psi'(t - t_1, \Psi(t_1, \mathbf{u}))$  is defined, and  $\Psi$  is therefore differentiable at  $(t, \mathbf{u})$ .  $\square$

A vector field on  $U$  is said to be *complete* if its flow has domain  $\mathbb{R} \times U$ ; i.e., if its integral curves are defined for all time. The vector field on  $\mathbb{R}^2$  from Figure 2.2 is complete. The vector field  $\mathbf{X}$  on  $\mathbb{R}$ , with  $X(t) = -t^2 D(t)$ , is not: for  $a \neq 0$ , the maximal integral curve  $c_a$  of  $\mathbf{X}$  with  $c_a(0) = a$  is given by  $c_a(t) = 1/(t + 1/a)$ , since

$$\dot{c}_a(t) = -\left(\frac{1}{t + \frac{1}{a}}\right)^2 D\left(\frac{1}{t + \frac{1}{a}}\right) = X\left(\frac{1}{t + \frac{1}{a}}\right) = (X \circ c_a)(t).$$

**Example 2.8.1.** Consider a vector field  $\mathbf{X}$  on  $\mathbb{R}^n$  of the form

$$\mathbf{X} = \sum_{i=1}^n \left( \sum_{j=1}^n a_{ij} u^j \right) \mathbf{D}_i, \quad a_{ij} \in \mathbb{R}.$$

Notice that  $\mathbf{X}(\mathbf{u}) = \mathcal{I}_{\mathbf{u}}(A\mathbf{u})$ , where  $A$  is the square matrix  $(a_{ij})$ . Thus,  $\mathbf{c} : I \rightarrow \mathbb{R}^n$  is an integral curve of  $\mathbf{X}$  if and only if  $\mathbf{c}' = A\mathbf{c}$ . Writing out this equation in components yields a so-called *linear system of ordinary differential equations with constant coefficients*

$$\begin{aligned} x_1'(t) &= \sum_{i=1}^n a_{1i} x_i(t) \\ &\vdots \\ x_n'(t) &= \sum_{i=1}^n a_{ni} x_i(t) \end{aligned}$$

where  $x_i = u^i \circ \mathbf{c}$ . We claim that  $\mathbf{X}$  is complete. In fact, the integral curve of  $\mathbf{X}$  that passes through  $\mathbf{u} \in \mathbb{R}^n$  at time 0 is given by

$$\begin{aligned} \mathbf{c} : \mathbb{R} &\rightarrow \mathbb{R}^n, \\ t &\mapsto e^{tA}\mathbf{u}. \end{aligned}$$

Indeed, by Exercise 1.50,

$$\begin{aligned}\mathbf{c}'(t) &= \lim_{h \rightarrow 0} \frac{\mathbf{c}(t+h) - \mathbf{c}(t)}{h} = \lim_{h \rightarrow 0} \frac{e^{(t+h)A}\mathbf{u} - e^{tA}\mathbf{u}}{h} \\ &= \left( \lim_{h \rightarrow 0} \frac{e^{hA} - I_n}{h} \right) e^{tA}\mathbf{u} \\ &= A \cdot \mathbf{c}(t),\end{aligned}$$

where the last line follows from the series definition of the exponential, cf. Examples and Remarks 1.9.1 (ii).

For example, if  $\mathbf{X} = u^2 D_1 - u^1 D_2$  is the vector field from Figure 2.2, the associated matrix is

$$A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix},$$

and by Examples and Remarks 1.9.1 (iii), the integral curve  $\mathbf{c}$  of  $\mathbf{X}$  with  $\mathbf{c}(0) = \mathbf{u}$  is given by

$$\mathbf{c}(t) = \Phi(t, \mathbf{u}) = \begin{bmatrix} \cos t & \sin t \\ -\sin t & \cos t \end{bmatrix} \cdot \mathbf{u},$$

where  $\Phi$  denotes the flow of  $\mathbf{X}$ .

**Definition 2.8.6.** Let  $\mathbf{X}$  denote a vector field along a curve  $\mathbf{c} : I \rightarrow \mathbb{R}^n$ , so that  $\mathbf{X}(t) = \mathcal{I}_{\mathbf{c}(t)}\mathbf{c}'_1(t)$  for some curve  $\mathbf{c}_1$  with the same domain as  $\mathbf{c}$ . The *covariant derivative* of  $\mathbf{X}$  at  $t \in I$  is the tangent vector

$$\mathbf{X}'(t) = \mathcal{I}_{\mathbf{c}(t)}\mathbf{c}'_1(t) \in \mathbb{R}^n_{\mathbf{c}(t)}.$$

Thus, the covariant derivative of a vector field along a curve is again a vector field along the same curve. For example, if  $\mathbf{P}$  is the position vector field on  $\mathbb{R}^n$  and  $\mathbf{c}$  is any curve, then the velocity field of  $\mathbf{c}$  can be written as

$$\dot{\mathbf{c}} = (\mathbf{P} \circ \mathbf{c})'. \quad (2.8.2)$$

This is because  $(\mathbf{P} \circ \mathbf{c})(t) = \mathcal{I}_{\mathbf{c}(t)}\mathbf{c}(t)$ , so that

$$(\mathbf{P} \circ \mathbf{c})'(t) = \mathcal{I}_{\mathbf{c}(t)}\mathbf{c}'(t) = \dot{\mathbf{c}}(t).$$

**Definition 2.8.7.** (1) Let  $\mathbf{X}$  denote a vector field on  $U \subset \mathbb{R}^n$ . For  $\mathbf{p} \in U$ ,  $\mathbf{u} \in \mathbb{R}^n_{\mathbf{p}}$ , the *covariant derivative*  $D_{\mathbf{u}}\mathbf{X}$  of  $\mathbf{X}$  with respect to  $\mathbf{u}$  is defined to be

$$D_{\mathbf{u}}\mathbf{X} := (\mathbf{X} \circ \mathbf{c})'(0),$$

where  $\mathbf{c}$  is any curve with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ .

(2) If  $f$  is a function on  $U$ , the *derivative of  $f$  with respect to  $\mathbf{u}$*  is by definition the number

$$\mathbf{u}(f) := (f \circ \mathbf{c})'(0),$$

with  $\mathbf{c}$  as above.

It must of course be checked that this definition is independent of the particular curve chosen. Let  $\mathbf{f} : U \rightarrow \mathbb{R}^n$  denote the map  $\mathcal{I}^{-1} \circ \mathbf{X}$  representing  $\mathbf{X}$ ; i.e.,  $\mathbf{X}(\mathbf{p}) = \mathcal{I}_{\mathbf{p}}\mathbf{f}(\mathbf{p})$ ,  $\mathbf{p} \in U$ . Then for any curve  $\mathbf{c}$  with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ ,

$$(\mathbf{X} \circ \mathbf{c})'(0) = \mathcal{I}_{\mathbf{c}(0)}(\mathbf{f} \circ \mathbf{c})'(0) = \mathcal{I}_{\mathbf{c}(0)}D\mathbf{f}(\mathbf{p})\mathbf{c}'(0) = \mathcal{I}_{\mathbf{p}}D\mathbf{f}(\mathbf{p})\mathcal{I}_{\mathbf{p}}^{-1}\mathbf{u},$$

and the derivative is indeed independent of the curve. The argument for the directional derivative of a function is similar. Recalling that  $\mathcal{I}_{\mathbf{f}(\mathbf{p})} \circ D\mathbf{f}(\mathbf{p}) = \mathbf{f}_{*\mathbf{p}} \circ \mathcal{I}_{\mathbf{p}}$ , the right side of the above equality may be rewritten as  $\mathcal{I}_{\mathbf{p}}\mathcal{I}_{\mathbf{f}(\mathbf{p})}^{-1}\mathbf{f}_{*\mathbf{p}}\mathbf{u}$ . The identity

$$D_{\mathbf{u}}\mathbf{X} = \mathcal{I}_{\mathbf{p}}\mathcal{I}_{\mathbf{f}(\mathbf{p})}^{-1}\mathbf{f}_{*\mathbf{p}}\mathbf{u}, \quad \mathbf{f} = \mathcal{I}^{-1} \circ \mathbf{X}, \quad (2.8.3)$$

says that the covariant derivative of a vector field represented by a map  $\mathbf{f}$  with respect to a given tangent vector  $\mathbf{u}$  at a point  $\mathbf{p}$  is essentially the ordinary derivative  $\mathbf{f}_*\mathbf{u}$ . The latter, however, is a vector based at  $\mathbf{f}(\mathbf{p})$ , and must therefore be parallel translated back to  $\mathbf{p}$ . Similarly, if the tangent vector  $\mathbf{u}$  equals  $\mathcal{I}_{\mathbf{p}}\mathbf{v}$  with  $\mathbf{v} \in \mathbb{R}^n$ , then  $\mathbf{u}(\mathbf{f}) = D\mathbf{f}(\mathbf{p})\mathbf{v}$ . The motivation behind the somewhat strange notation  $\mathbf{u}(\mathbf{f})$  is that when  $\mathbf{u} = \mathbf{D}_i(\mathbf{p})$ , then  $\mathbf{u}(\mathbf{f})$  equals the  $i$ -th partial derivative of  $\mathbf{f}$  at  $\mathbf{p}$ ; i.e.,  $\mathbf{D}_i(\mathbf{p})(\mathbf{f}) = D_i\mathbf{f}(\mathbf{p})$ .

More generally, one can define covariant derivatives of vector fields along maps by using the same equation as in Definition 2.8.7:

**Definition 2.8.8.** Let  $\mathbf{f} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^m$ ,  $\mathbf{X}$  a vector field along  $\mathbf{f}$ . Given  $\mathbf{p} \in U$ ,  $\mathbf{u} \in \mathbb{R}_{\mathbf{p}}^n$ , define the *covariant derivative of  $\mathbf{X}$  with respect to  $\mathbf{u}$*  to be  $D_{\mathbf{u}}\mathbf{X} = (\mathbf{X} \circ \mathbf{c})'(0)$ , where  $\mathbf{c}$  is a curve in  $U$  with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ .

Of course, if  $\mathbf{f}$  is the identity map, one recovers the original concept from Definition 2.8.7. On the other hand, if  $\mathbf{X}$  is a vector field on  $\mathbb{R}^m$  and  $\mathbf{f} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^m$ , then  $\mathbf{X} \circ \mathbf{f}$  is a vector field along  $\mathbf{f}$ , and by the above definition,  $D_{\mathbf{u}}(\mathbf{X} \circ \mathbf{f}) = (\mathbf{X} \circ \mathbf{f} \circ \mathbf{c})'(0)$ , where  $\mathbf{c}$  is a curve with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ . But  $\mathbf{f}'\dot{\mathbf{c}}(0) = \mathbf{f}_*\mathbf{u}$ , so that

$$D_{\mathbf{u}}(\mathbf{X} \circ \mathbf{f}) = D_{\mathbf{f}_*\mathbf{u}}\mathbf{X}. \quad (2.8.4)$$

Similarly, if  $\varphi : \mathbb{R}^m \rightarrow \mathbb{R}$ , then

$$\mathbf{u}(\varphi \circ \mathbf{f}) = \mathbf{f}_*\mathbf{u}(\varphi). \quad (2.8.5)$$

**Examples 2.8.2.** (i) The “coordinate vector field”  $\mathbf{D}_i$  is parallel, meaning that  $D_{\mathbf{u}}\mathbf{D}_i = 0$  for any  $\mathbf{u}$ . More generally, a vector field  $\mathbf{X}$  defined on  $U \subset \mathbb{R}^n$  is parallel if and only if it is of the form

$$\mathbf{X} = a_1\mathbf{D}_1 + \cdots + a_n\mathbf{D}_n, \quad a_1, \dots, a_n \in \mathbb{R}.$$

(ii) The position vector field  $\mathbf{P}$  is represented by the identity map. Since the derivative of the identity is the identity,

$$D_{\mathbf{u}}\mathbf{P} = \mathbf{u}$$

for any vector  $\mathbf{u}$ . Alternatively, this also follows directly from (2.8.2) together with the definition of covariant derivative.

- (iii) The *acceleration* of a curve  $\mathbf{c} : I \rightarrow \mathbb{R}^n$  is the vector field  $\dot{\mathbf{c}}'$  along  $\mathbf{c}$ . It follows from the definition that  $\dot{\mathbf{c}}' = \mathcal{I}_{\mathbf{c}}\mathbf{c}''$ . In particular,  $\mathbf{c}$  has zero acceleration if and only if  $\mathbf{c}'$  is a constant vector  $\mathbf{v}$ , and in this case  $\mathbf{c}(t) = \mathbf{c}(0) + t\mathbf{v}$ ; the image of  $\mathbf{c}$  is a straight line parallel to  $\mathbf{v}$ .
- (iv) The *gradient* of a function  $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}$  is the vector field  $\nabla f$  on  $U$  given by  $\nabla f(\mathbf{p}) = \mathcal{I}_{\mathbf{p}}[Df(\mathbf{p})]^T$ ,  $\mathbf{p} \in U$ . Alternatively,

$$\nabla f = \sum_i (D_i f) \mathbf{D}_i.$$

The inner product in  $\mathbb{R}^n$  induces one in each tangent space  $\mathbb{R}_{\mathbf{p}}^n$ , by setting

$$\langle \mathcal{I}_{\mathbf{p}}\mathbf{u}, \mathcal{I}_{\mathbf{p}}\mathbf{v} \rangle := \langle \mathbf{u}, \mathbf{v} \rangle, \quad \mathbf{p}, \mathbf{u}, \mathbf{v} \in \mathbb{R}^n.$$

This collection of inner products is called the *standard flat Riemannian metric* on  $\mathbb{R}^n$ .

As a first application, we shall derive a geometric interpretation of the gradient vector field of a function  $f$ . When  $\mathbf{u}$  is a *unit vector* (meaning  $\mathbf{u}$  has norm 1), the derivative  $\mathbf{u}(f)$  is called the *directional derivative of  $f$  in direction  $\mathbf{u}$* . It measures how fast  $f$  is changing along a curve that has  $\mathbf{u}$  as velocity vector. If  $\nabla f(\mathbf{p}) \neq 0$ , we claim that  $(\nabla f/|\nabla f|)(\mathbf{p})$  is the direction of greatest increase of  $f$  at  $\mathbf{p}$ , and  $-(\nabla f/|\nabla f|)(\mathbf{p})$  the direction of greatest decrease: indeed, for any unit  $\mathbf{u} \in \mathbb{R}_{\mathbf{p}}^n$ , if  $\mathbf{c}$  is a curve with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ ,

$$\mathbf{u}(f) = (f \circ \mathbf{c})'(0) = [Df(\mathbf{p})]\mathbf{c}'(0) = \langle [Df(\mathbf{p})]^T, \mathbf{c}'(0) \rangle = \langle \nabla f(\mathbf{p}), \mathbf{u} \rangle,$$

so that if  $\theta$  denotes the angle between  $\mathbf{u}$  and  $\nabla f(\mathbf{p})$ , then

$$\mathbf{u}(f) = \langle \nabla f(\mathbf{p}), \mathbf{u} \rangle = |\nabla f(\mathbf{p})||\mathbf{u}| \cos \theta = |\nabla f(\mathbf{p})| \cos \theta.$$

This expression is maximal when  $\theta = 0$ , i.e., when  $\mathbf{u}$  points in the direction of  $\nabla f$ , and minimal when  $\theta = \pi$ , which corresponds to the opposite direction. The *level sets* of  $f : \mathbb{R}^n \supset U \rightarrow \mathbb{R}$  are the sets  $f^{-1}(c) = \{\mathbf{p} \in U \mid f(\mathbf{p}) = c\}$ ,  $c \in \mathbb{R}$ . It follows from the above that the gradient of  $f$  is always orthogonal to these level sets; i.e., if  $\mathbf{c}$  is a smooth curve contained in a level set, then  $\langle (\nabla f) \circ \mathbf{c}, \dot{\mathbf{c}} \rangle = (f \circ \mathbf{c})' \equiv 0$ . For example, topographic maps of a region usually feature curves connecting points at the same altitude. These curves are the level curves (also called contour lines) of the altitude function, and traveling in a direction perpendicular to them means you are following the path of steepest ascent or descent.

Notice that for vector fields  $\mathbf{X}, \mathbf{Y}$  on  $U$ ,  $\langle \mathbf{X}, \mathbf{Y} \rangle$  is a function on  $U$ , if we let  $\langle \mathbf{X}, \mathbf{Y} \rangle(\mathbf{p}) = \langle \mathbf{X}(\mathbf{p}), \mathbf{Y}(\mathbf{p}) \rangle$ . In much the same way, given  $f : U \rightarrow \mathbb{R}$ , define new vector fields  $f\mathbf{X} + \mathbf{Y}$  and  $D_{\mathbf{X}}\mathbf{Y}$  on  $U$  by

$$(f\mathbf{X} + \mathbf{Y})(\mathbf{p}) := f(\mathbf{p})\mathbf{X}(\mathbf{p}) + \mathbf{Y}(\mathbf{p}), \quad D_{\mathbf{X}}\mathbf{Y}(\mathbf{p}) := D_{\mathbf{X}(\mathbf{p})}\mathbf{Y}, \quad \mathbf{p} \in U.$$

Similarly, if  $\mathbf{Y}$  is a vector field on  $\mathbb{R}^m$ , and  $\mathbf{X}$  a vector field along a map  $\mathbf{f} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^m$ , the formula  $D_{\mathbf{X}}\mathbf{Y}(\mathbf{p}) := D_{\mathbf{X}(\mathbf{p})}\mathbf{Y}$ , for  $\mathbf{p} \in U$ , defines a vector field along  $\mathbf{f}$ .

**Theorem 2.8.4.** Let  $\mathbf{X}, \mathbf{Y}$  denote vector fields on  $U \subset \mathbb{R}^n$ ,  $f : U \rightarrow \mathbb{R}$ ,  $a \in \mathbb{R}$ . Given  $\mathbf{p} \in U$ ,  $\mathbf{u} \in \mathbb{R}_p^n$ ,

$$(1) D_{\mathbf{u}}(a\mathbf{X} + \mathbf{Y}) = aD_{\mathbf{u}}\mathbf{X} + D_{\mathbf{u}}\mathbf{Y};$$

$$(2) D_{a\mathbf{u}+\mathbf{v}}\mathbf{X} = aD_{\mathbf{u}}\mathbf{X} + D_{\mathbf{v}}\mathbf{X};$$

$$(3) D_{\mathbf{u}}(f\mathbf{X}) = (\mathbf{u}(f))\mathbf{X}(\mathbf{p}) + f(\mathbf{p})D_{\mathbf{u}}\mathbf{X};$$

$$(4) D_{f\mathbf{X}}\mathbf{Y} = fD_{\mathbf{X}}\mathbf{Y};$$

(5) Let  $\mathbf{X}, \mathbf{Y}$  denote vector fields along a map  $\mathbf{f} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^m$ ,  $\mathbf{p} \in U$ ,  $\mathbf{u} \in \mathbb{R}_p^n$ . Then

$$D_{\mathbf{u}}\langle \mathbf{X}, \mathbf{Y} \rangle = \langle D_{\mathbf{u}}\mathbf{X}, \mathbf{Y}(\mathbf{p}) \rangle + \langle \mathbf{X}(\mathbf{p}), D_{\mathbf{u}}\mathbf{Y} \rangle.$$

*Proof.* (1), (2), and (4) are left as easy exercises. For (3), let  $\mathbf{g} = \mathcal{I}^{-1}\mathbf{X}$  (as always, in the sense that  $\mathbf{X}(\mathbf{p}) = \mathcal{I}_p\mathbf{g}(\mathbf{p})$ ). If  $\mathbf{c}$  is a curve with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ , then

$$\begin{aligned} D_{\mathbf{u}}f\mathbf{X} &= ((f \circ \mathbf{c})(\mathbf{X} \circ \mathbf{c}))'(0) = \mathcal{I}_p((f \circ \mathbf{c})(\mathbf{g} \circ \mathbf{c}))'(0) \\ &= \mathcal{I}_p(f(\mathbf{p})(\mathbf{g} \circ \mathbf{c})'(0) + (f \circ \mathbf{c})'(0)\mathbf{g}(\mathbf{p})) \\ &= f(\mathbf{p})\mathcal{I}_p(\mathbf{g} \circ \mathbf{c})'(0) + (f \circ \mathbf{c})'(0)\mathcal{I}_p\mathbf{g}(\mathbf{p}) \\ &= f(\mathbf{p})D_{\mathbf{u}}\mathbf{X} + (\mathbf{u}(f))\mathbf{X}(\mathbf{p}). \end{aligned}$$

To prove (5), let  $\mathbf{c} : I \rightarrow U$  denote a curve with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ . Then  $\mathbf{X} \circ \mathbf{c} = \mathcal{I}_c\mathbf{f}$  and  $\mathbf{Y} \circ \mathbf{c} = \mathcal{I}_c\mathbf{g}$  for some curves  $\mathbf{f}, \mathbf{g} : I \rightarrow \mathbb{R}^m$ . Since  $\langle \mathbf{X}, \mathbf{Y} \rangle \circ \mathbf{c} = \langle \mathbf{f}, \mathbf{g} \rangle$ , Corollary 2.2.1 implies

$$\begin{aligned} D_{\mathbf{u}}\langle \mathbf{X}, \mathbf{Y} \rangle &= \langle \mathbf{f}, \mathbf{g} \rangle'(0) = \langle \mathbf{f}', \mathbf{g} \rangle(0) + \langle \mathbf{f}, \mathbf{g}' \rangle(0) \\ &= \langle \mathcal{I}_c\mathbf{f}', \mathcal{I}_c\mathbf{g} \rangle(0) + \langle \mathcal{I}_c\mathbf{f}, \mathcal{I}_c\mathbf{g}' \rangle(0) \\ &= \langle D_{\mathbf{u}}\mathbf{X}, \mathbf{Y}(\mathbf{p}) \rangle + \langle \mathbf{X}(\mathbf{p}), D_{\mathbf{u}}\mathbf{Y} \rangle. \end{aligned} \quad \square$$

The following is an immediate application of the above theorem together with Examples 2.8.2(i):

**Corollary 2.8.1.** Let  $\mathbf{X}$  be a vector field on  $U \subset \mathbb{R}^n$ , and write

$$\mathbf{X} = \sum_i X^i \mathbf{D}_i, \quad X^i = \langle \mathbf{X}, \mathbf{D}_i \rangle, \quad i = 1, \dots, n.$$

Then  $D_{\mathbf{u}}\mathbf{X} = \sum_i (\mathbf{u}(X^i))\mathbf{D}_i(\mathbf{p})$  for any vector  $\mathbf{u} \in \mathbb{R}_p^n$ ,  $\mathbf{p} \in U$ .

**Example 2.8.3.** Consider the vector fields  $\mathbf{X}$  and  $\mathbf{Y}$  in  $\mathbb{R}^2$ , where

$$\mathbf{X} = u^1 \sin u^2 \mathbf{D}_1 + (1 + e^{u^2}) \mathbf{D}_2, \quad \mathbf{Y} = \frac{u^2}{1 + (u^1)^2} \mathbf{D}_1 + 3\mathbf{D}_2.$$

By Theorem 2.8.4 and the fact that the  $\mathbf{D}_i$  are parallel,

$$\begin{aligned} D_{\mathbf{X}}\mathbf{Y} &= u^1 \sin u^2 D_{\mathbf{D}_1}\mathbf{Y} + (1 + e^{u^2}) D_{\mathbf{D}_2}\mathbf{Y} \\ &= u^1 \sin u^2 D_1 \left( \frac{u^2}{1 + (u^1)^2} \right) \mathbf{D}_1 + (1 + e^{u^2}) D_2 \left( \frac{u^2}{1 + (u^1)^2} \right) \mathbf{D}_1 \\ &= \left( -\frac{2(u^1)^2 u^2 \sin u^2}{(1 + (u^1)^2)^2} + \frac{1 + e^{u^2}}{1 + (u^1)^2} \right) \mathbf{D}_1. \end{aligned}$$

## 2.9 Lie brackets

If  $\mathbf{X}$  is a vector field on an open set  $U$  in Euclidean space, one can associate to any function  $f : U \rightarrow \mathbb{R}$  a new function  $\mathbf{X}f$  by the formula

$$(\mathbf{X}f)(\mathbf{p}) = \mathbf{X}(\mathbf{p})(f), \quad \mathbf{p} \in U.$$

Notice that if  $\tilde{\mathbf{X}} : U \rightarrow \mathbb{R}^n$  denotes the map  $\mathcal{I}^{-1} \circ \mathbf{X}$  representing the vector field  $\mathbf{X}$ , then for  $f : U \rightarrow \mathbb{R}$ ,

$$\mathbf{X}f = Df(\tilde{\mathbf{X}}). \quad (2.9.1)$$

Indeed, if  $\mathbf{c}$  denote the maximal integral curve of  $\mathbf{X}$  with  $\mathbf{c}(0) = \mathbf{p} \in U$ , then

$$\mathbf{X}(\mathbf{p})(f) = (f \circ \mathbf{c})'(0) = Df(\mathbf{c}(0))\mathbf{c}'(0) = Df(\mathbf{p})\tilde{\mathbf{X}}(\mathbf{p}).$$

An important special case of (2.9.1) is that the  $i$ -th component  $X^i$  of a vector field  $\mathbf{X} = \sum_i X^i \mathbf{D}_i$  on  $U$  is given by

$$X^i = \langle \mathbf{X}, \mathbf{D}_i \rangle = u^i \circ \tilde{\mathbf{X}} = Du^i(\tilde{\mathbf{X}}) = \mathbf{X}u^i, \quad (2.9.2)$$

since  $u^i$ , being a linear map, equals  $Du^i$ .

Given vector fields  $\mathbf{X}$ ,  $\mathbf{Y}$ , and a function  $f$ , set  $\mathbf{X}\mathbf{Y}f := \mathbf{X}(\mathbf{Y}f)$ .

**Lemma 2.9.1.** *Let  $\mathbf{X}$  and  $\mathbf{Y}$  be vector fields on an open set  $U \subset \mathbb{R}^n$ . Then there exists a unique vector field  $[\mathbf{X}, \mathbf{Y}]$  on  $U$  such that*

$$D_{[\mathbf{X}, \mathbf{Y}]}f = (\mathbf{X}\mathbf{Y} - \mathbf{Y}\mathbf{X})f$$

for any function  $f$  on  $U$ .  $[\mathbf{X}, \mathbf{Y}]$  is called the Lie bracket of  $\mathbf{X}$  and  $\mathbf{Y}$ .

*Proof.* Write  $\mathbf{X} = \sum_i X^i \mathbf{D}_i$ , and similarly for  $\mathbf{Y}$ . Then

$$\begin{aligned} \mathbf{X}\mathbf{Y}f &= \sum_j \mathbf{X}(Y^j D_j f) = \sum_{i,j} X^i D_i(Y^j D_j f) \\ &= \sum_{i,j} X^i Y^j D_{ij} f + X^i D_i Y^j D_j f. \end{aligned}$$

Thus,

$$(\mathbf{X}\mathbf{Y} - \mathbf{Y}\mathbf{X})f = \sum_j \left( \sum_i X^i D_i Y^j - Y^i D_i X^j \right) D_j f.$$

Now, if  $\mathbf{v}_1$  and  $\mathbf{v}_2$  are two vectors in  $\mathbb{R}_p^n$  such that  $\mathbf{v}_1 f = \mathbf{v}_2 f$  for every function  $f$ , then the vectors are equal: indeed, taking  $f = u^i$  in (2.9.2) shows that the  $i$ -th components of both vectors coincide. This proves uniqueness in the definition of the bracket. Existence is also clear: define

$$[\mathbf{X}, \mathbf{Y}] = \sum_j \left( \sum_i X^i D_i Y^j - Y^i D_i X^j \right) \mathbf{D}_j. \quad (2.9.3)$$

□

Together with Theorem 2.2.4, Lemma 2.9.1 immediately implies that  $[D_i, D_j] = \mathbf{0}$ . It turns out that the Lie bracket of two vector fields measures the amount by which their flows fail to commute; specifically, we will see that if  $\mathbf{X}$  has flow  $\Phi$  and  $\mathbf{Y}$  has flow  $\Psi$ , then  $[\mathbf{X}, \mathbf{Y}] \equiv 0$  if and only if  $\Phi_t \circ \Psi_s = \Psi_s \circ \Phi_t$  for all  $s, t$ , where  $\Phi_t(\mathbf{a}) = \Phi(t, \mathbf{a})$ , etc. But first, we state some general properties of the bracket:

**Proposition 2.9.1.**  $[\mathbf{X}, \mathbf{Y}] = D_{\mathbf{X}}\mathbf{Y} - D_{\mathbf{Y}}\mathbf{X}$ .

*Proof.*

$$D_{\mathbf{X}}\mathbf{Y} = \sum_j (\mathbf{X}Y^j)\mathbf{D}_j = \sum_j \left( \sum_i X^i D_i Y^j \right) \mathbf{D}_j = \sum_{ij} X^i D_i Y^j \mathbf{D}_j,$$

and a similar expression holds for  $D_{\mathbf{Y}}\mathbf{X}$ . The claim now follows from (2.9.3).  $\square$

More generally, we have:

**Theorem 2.9.1.** Let  $f : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^m$ . If  $\mathbf{X}, \mathbf{Y}$  are vector fields on  $U$ , then

$$f_*[\mathbf{X}, \mathbf{Y}] = D_{\mathbf{X}}f_*\mathbf{Y} - D_{\mathbf{Y}}f_*\mathbf{X}.$$

*Proof.* Write  $\mathbf{X} = \sum_i X^i \mathbf{D}_i$  as usual, and similarly for  $\mathbf{Y}$ . Then

$$f_*\mathbf{Y} = \sum_j Y^j f_*\mathbf{D}_j = \sum_{ij} (Y^j D_j f^i) \mathbf{D}_i \circ f,$$

so that

$$\begin{aligned} D_{\mathbf{X}}f_*\mathbf{Y} &= \sum_k X^k D_k \left( \sum_{ij} (Y^j D_j f^i) \mathbf{D}_i \circ f \right) \\ &= \sum_{i,j,k} X^k [(D_k Y^j)(D_j f^i) + Y^j (D_k D_j f^i)] \mathbf{D}_i \circ f. \end{aligned}$$

Thus,

$$\begin{aligned} D_{\mathbf{X}}f_*\mathbf{Y} - D_{\mathbf{Y}}f_*\mathbf{X} &= \sum_{i,j,k} [X^k (D_k Y^j) - Y^k (D_k X^j)] (D_j f^i) \mathbf{D}_i \circ f \\ &\quad + \sum_{i,j,k} [X^k Y^j D_{kj} f^i - Y^k X^j D_{kj} f^i] \mathbf{D}_i \circ f. \end{aligned}$$

The second sum in the above identity vanishes, because

$$\begin{aligned} \sum_{j,k} X^k Y^j D_{kj} f^i - Y^k X^j D_{kj} f^i &= \sum_{j,k} X^k Y^j D_{kj} f^i - Y^j X^k D_{jk} f^i \\ &= \sum_{j,k} X^k Y^j (D_{kj} f^i - D_{jk} f^i) = 0. \end{aligned}$$



Therefore

$$\begin{aligned}
 D_X f_* Y - D_Y f_* X &= \sum_{j,k} [X^k (D_k Y^j) - Y^k (D_k X^j)] f_* D_j \\
 &= \sum_j (X(Y^j) - Y(X^j)) f_* D_j \\
 &= f_* \sum_j [D_X(Y^j D_j) - D_Y(X^j D_j)] = f_* (D_X Y - D_Y X) \\
 &= f_* [X, Y]. \quad \square
 \end{aligned}$$

**Proposition 2.9.2.** Let  $X, Y, Z$  denote vector fields on  $U \subset \mathbb{R}^n$ ,  $f, g : U \rightarrow \mathbb{R}$ ,  $c \in \mathbb{R}$ . Then

- (1)  $[cX + Y, Z] = c[X, Z] + [Y, Z]$ ;
- (2)  $[Y, X] = -[X, Y]$ ;
- (3)  $[X, [Y, Z]] + [Y, [Z, X]] + [Z, [X, Y]] = 0$ ;
- (4)  $[fX, Y]g = f[X, Y]g - (Yf)(Xg)$ .

*Proof.* The first two properties are immediate. For (3), write

$$[X, [Y, Z]]f = (XYZ - XZY - YZX + ZYX)f.$$

Writing out the other two remaining terms in the same way and summing, one easily sees that all the terms cancel. To establish the last property, we compute

$$\begin{aligned}
 [fX, Y]g &= fX(Yg) - Y((fX)g) = fXYg - Y(f(Xg)) \\
 &= fXYg - (Yf)(Xg) - fYXg \\
 &= f[X, Y]g - (Yf)(Xg). \quad \square
 \end{aligned}$$

A vector space  $V$  together with an operation  $[\cdot, \cdot] : V \times V \rightarrow \mathbb{R}$  that satisfies the first three properties of Proposition 2.9.2 is called a *Lie algebra*. The third property is known as the *Jacobi identity*.

Let  $U \subset \mathbb{R}^n$ ,  $X$  a vector field on  $U$ , and  $f : U \rightarrow \mathbb{R}^n$ . For each  $p \in U$ ,  $f_* X(p)$  is a vector in  $\mathbb{R}_{f(p)}^n$ , but in general, this formula does not define a vector field on  $f(U)$ , since  $f$  may not be injective. If  $Y$  is a vector field on  $f(U)$ ,  $X$  and  $Y$  are said to be *f-related* if  $f_* X = Y \circ f$ .

**Theorem 2.9.2.** If  $X_i$  is *f-related* to  $Y_i$ ,  $i = 1, 2$ , then  $[X_1, X_2]$  is *f-related* to  $[Y_1, Y_2]$ .

*Proof.* We will show that the  $i$ -th components of  $f_* [X_1, X_2]$  and  $[Y_1, Y_2] \circ f$  coincide. By (2.9.2), that of the latter equals

$$([Y_1, Y_2] \circ f)u^i = (Y_1 \circ f)Y_2 u^i - (Y_2 \circ f)Y_1 u^i,$$

which by (2.8.5) can be written

$$\begin{aligned}
 (f_* X_1)Y_2 u^i - (f_* X_2)Y_1 u^i &= X_1((Y_2 u^i) \circ f) - X_2((Y_1 u^i) \circ f) \\
 &= [X_1, X_2](u^i \circ f) \\
 &= f_* [X_1, X_2]u^i,
 \end{aligned}$$

thereby establishing the claim.  $\square$

Our next goal is to derive an alternative characterization of the bracket, one that involves the flow of the first vector field. For this, we need the following:

**Lemma 2.9.2.** *Let  $I$  be an interval containing 0. Then any function  $f : I \rightarrow \mathbb{R}$  may be written as  $f(t) = f(0) + tg(t)$ , where  $g(0) = f'(0)$ .*

*Proof.* Fix any  $t \in I$ , and define  $\varphi : [0, 1] \rightarrow \mathbb{R}$  by  $\varphi(s) = f(st)$ . Then

$$f(t) - f(0) = \varphi(1) - \varphi(0) = \int_0^1 \varphi'(s) ds = t \int_0^1 f'(st) ds = tg(t),$$

where  $g(t) = \int_0^1 f'(st) ds$ . □

**Theorem 2.9.3.** *Given vector fields  $X, Y$  on an open set  $U \subset \mathbb{R}^n$ , and  $\mathbf{p} \in U$ ,*

$$[X, Y](\mathbf{p}) = \lim_{t \rightarrow 0} \frac{\Phi_{-t*} Y(\Phi_t(\mathbf{p})) - Y(\mathbf{p})}{t},$$

where  $\Phi$  denotes the flow of  $X$ , and  $\Phi_t(\mathbf{q}) := \Phi(t, \mathbf{q})$ ,  $\mathbf{q} \in U$ .

*Proof.* The plan is to show that the two vectors in the above identity have the same components. The  $i$ -th component of the one on the left is

$$([X, Y]u^i)(\mathbf{p}) = (X(\mathbf{p})Y - Y(\mathbf{p})X)(u^i).$$

As usual, denote by  $\tilde{Y} : U \rightarrow \mathbb{R}^n$  the map  $\mathcal{I}^{-1} \circ Y$  representing the vector field  $Y$ . The  $i$ -th component  $Z^i$  of the right side prior to taking limits is

$$\begin{aligned} \frac{1}{t} [u^i \circ D\Phi_{-t}(\Phi_t(\mathbf{p}))(\tilde{Y}(\Phi_t(\mathbf{p}))) - u^i \circ \tilde{Y}(\mathbf{p})] \\ = \frac{1}{t} [D(u^i \circ \Phi_{-t})(\Phi_t(\mathbf{p}))(\tilde{Y}(\Phi_t(\mathbf{p}))) - u^i \circ \tilde{Y}(\mathbf{p})] \end{aligned}$$

by the chain rule and the fact that  $Du^i = u^i$ ,  $u^i$  being linear. Now, if  $h(t, \mathbf{p}) = (u^i \circ \Phi)(t, \mathbf{p})$ , then by the previous lemma,  $h(t, \mathbf{p}) = u^i(\mathbf{p}) + tg(t, \mathbf{p})$ , where  $g(0, \mathbf{p}) = D_1 h(0, \mathbf{p}) = D_1(u^i \circ \Phi)(0, \mathbf{p})$ . By Theorem 2.8.1 (2), the latter equals  $u^i(\tilde{X}(\mathbf{p}))$ , which may be written as  $X(\mathbf{p})(u^i)$ . Setting  $g_t(\mathbf{q}) = g(t, \mathbf{q})$ , we conclude that

$$u^i \circ \Phi_{-t} = u^i - tg_{-t}, \quad g_0 = Xu^i.$$

Thus,

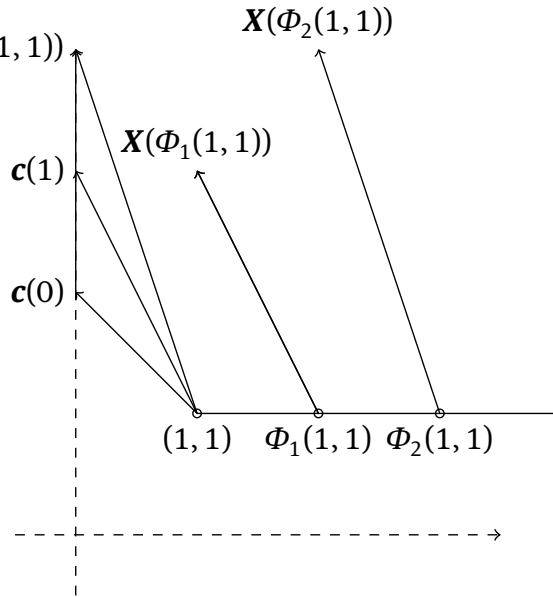
$$\begin{aligned} Z^i &= \frac{1}{t} [Du^i(\Phi_t(\mathbf{p}))(\tilde{Y}(\Phi_t(\mathbf{p}))) - u^i(\tilde{Y}(\mathbf{p}))] - Dg_{-t}(\Phi_t(\mathbf{p}))(\tilde{Y}(\Phi_t(\mathbf{p}))) \\ &= \frac{1}{t} [(u^i \circ \tilde{Y})(\Phi_t(\mathbf{p})) - u^i(\tilde{Y}(\mathbf{p}))] - Dg_{-t}(\Phi_t(\mathbf{p}))(\tilde{Y}(\Phi_t(\mathbf{p}))). \end{aligned}$$

Taking limits as  $t \rightarrow 0$  and recalling that  $t \mapsto \Phi_t(\mathbf{p})$  is the integral curve of  $X$  passing through  $\mathbf{p}$  at time 0, we deduce that the  $i$ -th component of the right side is

$$\begin{aligned} X(\mathbf{p})(u^i \circ \tilde{Y}) - Dg_0(\mathbf{p})(\tilde{Y}(\mathbf{p})) &= X(\mathbf{p})(Yu^i) - D(Xu^i)(\mathbf{p})\tilde{Y}(\mathbf{p}) \\ &= X(\mathbf{p})(Yu^i) - Y(\mathbf{p})(Xu^i), \end{aligned}$$

as claimed. □

For  $X = u^1 D_2 - u^2 D_1$ , the bracket  $[D_1, X](1, 1)$  equals  $c'(0) = D_2(1, 1)$ , where  $c(t) = \Phi_{-t*} X(\Phi_t(1, 1))$ , and  $\Phi$  is the flow of  $D_1$



Our last goal in this section is to use the above theorem and prove that the bracket of two vector fields is identically zero if and only if their flows commute. We begin with the following:

**Lemma 2.9.3.** *Let  $U$  be an open set in  $\mathbb{R}^n$ ,  $f : U \rightarrow f(U) \subset \mathbb{R}^n$  a diffeomorphism, and  $X$  a vector field on  $U$  with flow  $\Phi$ . Then  $X$  is  $f$ -related to itself if and only if  $f \circ \Phi_t = \Phi_t \circ f$ . Here, as always,  $\Phi_t = \Phi(t, \cdot)$ .*

*Proof.* Fix any  $q \in f(U)$ , and consider the curve  $\gamma$ , where

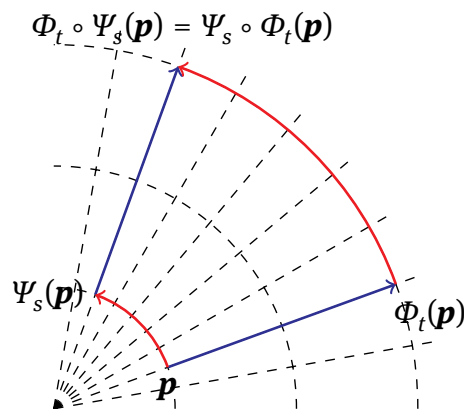
$$\gamma(t) = (f \circ \Phi_t \circ f^{-1})(q).$$

The curve  $t \mapsto c(t) := (\Phi_t \circ f^{-1})(q) = (f^{-1} \circ \gamma)(t)$  is by definition an integral curve of  $X$ , so that

$$\dot{\gamma} = f_*(\dot{c}) = f_* \circ X \circ c = (f_* \circ X \circ f^{-1}) \circ \gamma.$$

Thus, the vector field  $f_* \circ X \circ f^{-1}$  has flow  $\Psi$ , where  $\Psi_t = f \circ \Phi_t \circ f^{-1}$ . The claim follows, since  $X$  is  $f$ -related to itself if and only if  $f_* \circ X \circ f^{-1} = X$ ; i.e., if and only if  $\Phi_t = f \circ \Phi_t \circ f^{-1}$ .  $\square$

The flows  $\Phi$  of  $X = u^1 D_1 + u^2 D_2$  (in blue) and  $\Psi$  of  $Y = u^1 D_2 - u^2 D_1$  (in red) commute



**Theorem 2.9.4.** Let  $\Phi$  and  $\Psi$  denote the flows of  $X$  and  $Y$ . Then  $[X, Y] \equiv \mathbf{0}$  if and only if these flows commute; i.e.,  $\Phi_t \circ \Psi_s = \Psi_s \circ \Phi_t$  for all  $s, t$ .

*Proof.* If  $\Phi_t \circ \Psi_s = \Psi_s \circ \Phi_t$ , then  $Y$  is  $\Phi_t$ -related to itself for every  $t$  by Lemma 2.9.3; i.e.,  $\Phi_{t*} Y = Y \circ \Phi_t$ , and by Theorem 2.9.3,

$$[X, Y](\mathbf{p}) = \lim_{t \rightarrow 0} \frac{\Phi_{-t*} Y(\Phi_t(\mathbf{p})) - Y(\mathbf{p})}{t} = \lim_{t \rightarrow 0} \frac{Y(\mathbf{p}) - Y(\mathbf{p})}{t} = \mathbf{0}$$

for every  $\mathbf{p}$ . Conversely, suppose the vector fields have vanishing bracket, and for any fixed  $\mathbf{p}$ , consider the curve  $\mathbf{c}$  in the tangent space  $\mathbb{R}_\mathbf{p}^n$  of  $\mathbb{R}^n$  at  $\mathbf{p}$  given by  $\mathbf{c}(t) = \Phi_{-t*} Y(\Phi_t(\mathbf{p}))$ . We will be done (again by the lemma) once we show that  $\mathbf{c}(t) = \mathbf{c}(0) = Y(\mathbf{p})$  for all  $t$ ; i.e., once we establish that  $\mathbf{c}' \equiv \mathbf{0}$ . Vanishing of the Lie bracket implies, by Theorem 2.9.3, that  $\mathbf{c}'(0) = \mathbf{0}$ . But if  $\mathbf{q} = \Phi_t(\mathbf{p})$  for any given fixed  $t$ , then using Exercise 2.28, we obtain

$$\begin{aligned} \mathbf{c}'(t) &= \lim_{h \rightarrow 0} \frac{\mathbf{c}(t+h) - \mathbf{c}(t)}{h} \\ &= \lim_{h \rightarrow 0} \frac{1}{h} [\Phi_{-(t+h)*} \circ Y \circ \Phi_{t+h}(\mathbf{p}) - \Phi_{-t*} \circ Y \circ \Phi_t(\mathbf{p})] \\ &= \lim_{h \rightarrow 0} \frac{1}{h} \Phi_{-t*} [(\Phi_{-h*} \circ Y \circ \Phi_h)(\Phi_t(\mathbf{p})) - Y(\Phi_t(\mathbf{p}))] \\ &= \Phi_{-t*} \lim_{h \rightarrow 0} \frac{1}{h} [(\Phi_{-h*} \circ Y \circ \Phi_h)(\mathbf{q}) - Y(\mathbf{q})] = \Phi_{-t*} \mathbf{c}'(0) \\ &= \mathbf{0}. \end{aligned} \quad \square$$

More generally, given  $\mathbf{p} \in \mathbb{R}^n$ , define a curve  $\mathbf{c}_\mathbf{p}$  on an interval around 0 by

$$\mathbf{c}_\mathbf{p}(t) = (\Psi_{-t} \circ \Phi_{-t} \circ \Psi_t \circ \Phi_t)(\mathbf{p}).$$

If the bracket of  $X$  and  $Y$  is identically zero, then  $\mathbf{c}_\mathbf{p}(t) = \mathbf{p}$  for all  $t$ . When the bracket is nonzero, the curve  $\mathbf{c}_\mathbf{p}$  is no longer constant. In the exercises, however, the reader is asked to show that  $\mathbf{c}'_\mathbf{p}(0) = \mathbf{0}$ . It can be shown that the acceleration of  $\mathbf{c}_\mathbf{p}$  at 0 measures the bracket; specifically, for any smooth function  $f : U \rightarrow \mathbb{R}$  defined on a neighborhood  $U$  of  $\mathbf{p}$ ,

$$(f \circ \mathbf{c}_\mathbf{p})''(0) = 2[X, Y](\mathbf{p})(f).$$

## 2.10 Partitions of unity

In this section, we introduce a concept that plays a key role both in the theory of integration and that of manifolds.

Let  $A \subset \mathbb{R}^n$ . Recall that a collection  $\mathcal{C}$  of (not necessarily open) subsets of  $\mathbb{R}^n$  is said to be a *cover* of  $A$  if  $A$  is contained in the union of all the sets in  $\mathcal{C}$ .  $\mathcal{C}$  is said to be *locally finite* if every  $\mathbf{a} \in A$  admits a neighborhood that intersects only finitely many sets in  $\mathcal{C}$ .

**Definition 2.10.1.** A *partition of unity* for  $A \subset \mathbb{R}^n$  is a collection  $\Phi$  of differentiable functions defined on some open set containing  $A$  such that

- (1)  $0 \leq \varphi \leq 1$  for all  $\varphi \in \Phi$ ;
- (2) The collection  $\{\text{supp } \varphi \mid \varphi \in \Phi\}$  is a locally finite cover of  $A$ ;
- (3)  $\sum_{\varphi \in \Phi} \varphi(\mathbf{a}) = 1$  for all  $\mathbf{a} \in A$ .

Even though  $\Phi$  may be an infinite collection, the sum in (3) makes sense because by (2), only finitely many  $\varphi$  are nonzero on some neighborhood of any given point. In fact, only finitely many  $\varphi$  are nonzero on any given compact set  $C \subset A$ : by (2), any  $\mathbf{p} \in C$  has a neighborhood on which only finitely many  $\varphi$  are not identically zero. Since  $C$  is compact, it can be covered by finitely many such sets.

**Theorem 2.10.1.** Any  $A \subset \mathbb{R}^n$  admits a partition of unity  $\Phi$ . Furthermore, if  $\{U_\alpha\}_{\alpha \in J}$  is an open cover of  $A$ , then  $\Phi$  may be chosen to be subordinate to the cover; i.e., for any  $\varphi \in \Phi$ , there exists  $\alpha \in J$  such that  $\text{supp } \varphi \subset U_\alpha$ .

*Proof.* Both statements will be considered at the same time; in other words, we assume some open cover  $\{U_\alpha\}_{\alpha \in J}$  of  $A$  is given (if none is given, we take it to be  $\{\mathbb{R}^n\}$ ). For each  $\mathbf{a} \in A$ , choose some  $U_\alpha$  that contains it, and some bounded open neighborhood  $V_\alpha$  of  $\mathbf{a}$  whose closure is contained in  $U_\alpha$ . The collection  $\{V_\alpha \mid \alpha \in J\}$  is an open cover of  $A$ , which, by Theorem 1.7.5, contains a countable subcover  $\{V_1, V_2, \dots\}$ . By construction, for each  $i \in \mathbb{N}$  there exists some  $\alpha_i \in J$  such that  $\overline{V_i} \subset U_{\alpha_i}$ . We may now appeal to Theorem 2.2.5 to assert the existence of smooth functions  $\tilde{\varphi}_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $i \in \mathbb{N}$ , satisfying

$$0 \leq \tilde{\varphi}_i \leq 1, \quad \tilde{\varphi}_i|_{\overline{V_i}} \equiv 1, \quad \text{supp } \tilde{\varphi}_i \subset U_{\alpha_i}.$$

Set  $\varphi_1 = \tilde{\varphi}_1$ , and

$$\varphi_i = (1 - \tilde{\varphi}_1) \cdots (1 - \tilde{\varphi}_{i-1}) \tilde{\varphi}_i, \quad i > 1.$$

Notice that  $\text{supp } \varphi_i \subset \text{supp } \tilde{\varphi}_i \subset U_{\alpha_i}$ , so that the collection of supports is subordinate to the original cover. Furthermore, it is locally finite: given any  $\mathbf{a} \in A$ , choose some  $V_i$  that contains it. Then  $\tilde{\varphi}_i$  equals 1 on  $V_i$ , and by definition,  $\varphi_k$  vanishes on  $V_i$  for all  $k > i$ . Finally, an easy induction argument shows that

$$\sum_{i=1}^k \varphi_i = 1 - (1 - \tilde{\varphi}_1) \cdots (1 - \tilde{\varphi}_k), \quad k \in \mathbb{N},$$

so that on  $V_k$  (where  $1 - \tilde{\varphi}_k$  is identically zero),

$$\sum_{i=1}^{\infty} \varphi_i|_{V_k} = \sum_{i=1}^k \varphi_i|_{V_k} = 1 - 0 = 1.$$

Since the  $V_k$  cover  $A$ ,  $\sum_i \varphi_i$  is identically 1 on  $A$ . □

**Remarks 2.10.1.** (i) The proof of the Theorem shows that when the cover  $\{U_1, \dots, U_k\}$  of  $A$  is finite, the partition of unity  $\Phi$  may be assumed to consist of the same

number of elements; i.e.  $\Phi = \{\varphi_1, \dots, \varphi_k\}$ , with  $\text{supp } \varphi_i \subset U_i$ . It also shows that in general, the partition of unity is countable.

- (ii) We are now in a position to prove the claim made in Section 2.2, that given any open set  $U \subset \mathbb{R}^n$ , any  $\mathbf{p} \in U$ , and any (not necessarily compact) neighborhood  $V$  of  $\mathbf{p}$  whose closure is contained in  $U$ , there exists a smooth function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $0 \leq \varphi \leq 1$ , that equals 1 on the closure of  $V$  and has its support inside  $U$ : indeed, since  $\{U, \mathbb{R}^n \setminus \overline{V}\}$  is an open cover of  $\mathbb{R}^n$ , there exists a partition of unity  $\{\psi_1, \psi_2\}$  subordinate to this cover, with  $\text{supp } \psi_1 \subset U$  and  $\text{supp } \psi_2 \subset \mathbb{R}^n \setminus \overline{V}$ . Setting  $\varphi = \psi_1$  yields the claim: its support lies in  $U$ , and since  $\psi_2$  vanishes on  $\overline{V}$ ,  $\varphi$  must equal 1 there.

## 2.11 Exercises

**2.1.** Find the derivative of each of the following maps:

- (1)  $f(x, y, z) = (x^y)^z$ ;
- (2)  $g(x, y, z) = x^{(y^z)}$ ;
- (3)  $h(x, y, z) = (x \cos(ye^z), \sin \sqrt{x^2 + z^2 + 1})$ .
- (4)  $\mathbf{k}(x, y, z, w) = ((x - y)e^{(z+w^2)}, \cos(\sin(\log y^4 + z^2 + 5)), 1/(x^2 + y^2 + z^2 + w^2 + 1))$ .

**2.2.** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  be given by  $f(x, y) = x^2y + \cos e^{x+y}$ , and  $\mathbf{c} : \mathbb{R} \rightarrow \mathbb{R}^2$  by  $\mathbf{c}(t) = (t - 1, t^2 + 1)$ .

- (a) Use the chain rule to determine  $(f \circ \mathbf{c})'(t)$  and  $D(\mathbf{c} \circ f)(x, y)$ ;
- (b) Check your answers in (a) by computing  $f \circ \mathbf{c}$ ,  $\mathbf{c} \circ f$  and evaluating their derivatives directly.

**2.3.** Suppose  $g : \mathbb{R}^3 \rightarrow \mathbb{R}$  has Jacobian matrix  $\begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$  at  $(1, 1, 0)$ . Find  $Df(1, 1, 1)$ , if  $f(x, y, z) = g(x^2, yz, x - y)$ .

**2.4.** Suppose  $\mathbf{c} : I \rightarrow \mathbb{R}^n$  is a curve such that the position vector  $\mathbf{c}(t)$  is always orthogonal to the velocity vector  $\mathbf{c}'(t)$ . Show that the image of  $\mathbf{c}$  lies in some sphere centered at the origin.

**2.5.** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , and define  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$  by  $g(r, \theta) = f(r \cos \theta, r \sin \theta)$ . Express the partial derivatives of  $g$  in terms of those of  $f$ .

Many Calculus texts write these partial derivatives in polar coordinates as

$$\frac{\partial f}{\partial r} = \frac{\partial f}{\partial x} \cos \theta + \frac{\partial f}{\partial y} \sin \theta,$$

and another expression for the derivative of  $f$  with respect to  $\theta$ . Explain why, even though the notation is suggestive, the formula is incorrect or at least ambiguous.

**2.6.** Answer the same question as in the previous problem, but for a function  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  of 3 variables, and for spherical coordinates.

**2.7.** If  $E$  is a vector space, and  $f : E \rightarrow \mathbb{R}^m$  is a map, we say  $f$  is *differentiable* if  $f \circ L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is differentiable for some isomorphism  $L : \mathbb{R}^n \rightarrow E$ .

- (a) Prove that this definition makes sense; i.e. it is independent of the choice of  $L$ .  
 (b) Explain why one recovers the usual concept of differentiability when  $E = \mathbb{R}^n$ .  
 (c) Show that any linear transformation  $T : E \rightarrow \mathbb{R}^m$  is differentiable.

**2.8.** Let  $E$  be an inner product space. Prove that the function  $f : E \rightarrow \mathbb{R}$  given by  $f(\mathbf{u}) = |\mathbf{u}|$  is differentiable at any point other than the origin, and is not differentiable at  $\mathbf{0}$  (see Exercise 2.7). Find  $Df(\mathbf{u})$ .

**2.9.** Let  $E$  be an inner product space,  $f : E \rightarrow \mathbb{R}$  a function satisfying  $|f(\mathbf{u})| \leq |\mathbf{u}|^\alpha$  for all  $\mathbf{u} \in E$  and some  $\alpha > 1$ . Show that  $f$  is differentiable at  $\mathbf{0}$  and find  $Df(\mathbf{0})$ . (Notice that the conclusion is false if  $\alpha = 1$  by Exercise 2.8.)

**2.10.** Determine whether  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , where

$$f(\mathbf{a}) = \begin{cases} \frac{\sin |\mathbf{a}|^3}{|\mathbf{a}|} & \text{if } \mathbf{a} \neq \mathbf{0} \\ 0 & \text{if } \mathbf{a} = \mathbf{0}, \end{cases}$$

is differentiable at the origin, and if yes, find  $Df(\mathbf{0})$ .

**2.11.** Suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  satisfies  $f(t\mathbf{a}) = tf(\mathbf{a})$  for all  $\mathbf{a} \in \mathbb{R}^n$  and  $t \in \mathbb{R}$ . Prove that if  $f$  is differentiable at the origin, then  $f$  is linear. *Hint:* It is enough to show that  $f = Df(\mathbf{0})$ .

**2.12.** A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is said to be *homogeneous of degree  $k$*  if  $f(t\mathbf{a}) = t^k f(\mathbf{a})$  for all  $\mathbf{a} \in \mathbb{R}^n$ ,  $t \in \mathbb{R}$ , and some  $k \in \mathbb{N}$ . Show that if  $f$  is homogeneous of degree  $k$  and differentiable, then

$$f = \frac{1}{k} \sum_{i=1}^n u^i D_i f.$$

**2.13.** Show, by means of an example, that in an arbitrary metric space  $X$ , a contraction  $X \rightarrow X$  does not necessarily have a fixed point (in this case,  $X$  cannot of course be complete).

**2.14.** Let  $U$  denote a convex open set in  $\mathbb{R}^n$ ,  $f : U \rightarrow \mathbb{R}^m$  a differentiable map, and  $\mathbf{a}, \mathbf{b} \in U$ . Prove the *mean value theorem*: For any  $\mathbf{u} \in \mathbb{R}^m$ , there exists some  $\mathbf{c}$  on the line segment joining  $\mathbf{a}$  and  $\mathbf{b}$  such that

$$\langle \mathbf{u}, f(\mathbf{b}) - f(\mathbf{a}) \rangle = \langle \mathbf{u}, Df(\mathbf{c})(\mathbf{b} - \mathbf{a}) \rangle.$$

*Hint:* Set  $\mathbf{v} = \mathbf{b} - \mathbf{a}$ , and apply the ordinary mean value theorem to  $g$  on  $[0, 1]$ , where  $g(t) := \langle \mathbf{u}, f(\mathbf{a} + t\mathbf{v}) \rangle$ .

**2.15.** (a) Show that the mean value theorem from Exercise 2.14 does indeed generalize the ordinary mean value theorem.

- (b) Suppose  $U$  is a convex open set in  $\mathbb{R}^n$ ,  $\mathbf{a}, \mathbf{b} \in U$ ,  $f : U \rightarrow \mathbb{R}$  a differentiable function. Prove that there exists some  $\mathbf{c}$  on the line segment joining  $\mathbf{a}$  and  $\mathbf{b}$  such that

$$f(\mathbf{b}) - f(\mathbf{a}) = Df(\mathbf{c})(\mathbf{b} - \mathbf{a}) = \langle \nabla f(\mathbf{c}), \mathbf{b} - \mathbf{a} \rangle.$$

**2.16.** Prove that a linear isometry of an inner product space has determinant  $\pm 1$ .

**2.17.** Show that if  $T : V \rightarrow W$  is a linear isometry and  $L$  is a self-adjoint (respectively skew-adjoint) operator on  $W$ , then  $T^{-1} \circ L \circ T$  is a self-adjoint (resp. skew-adjoint) operator on  $V$ .

**2.18.** The *orthogonal group*  $O(n) \subset M_{n,n}$  is the collection of all  $n \times n$  orthogonal matrices. Suppose  $\mathbf{c} : I \rightarrow O(n)$  is a smooth curve with  $\mathbf{c}(0) = I_n$ . Prove that  $\mathbf{c}'(0)$  is skew-adjoint:  $\mathbf{c}'(0) + \mathbf{c}'^T(0) = \mathbf{0}$ .

**2.19.** (a) Find the derivative of the determinant  $\det : M_{n,n} \cong \mathbb{R}^{n^2} \rightarrow \mathbb{R}$ .

(b) The *special linear group* is the subset  $Sl(n) \subset M_{n,n}$  consisting of all matrices with determinant 1. Suppose  $\mathbf{c} : I \rightarrow Sl(n)$  is a smooth curve with  $\mathbf{c}(0) = I_n$ . Prove that  $\mathbf{c}'(0)$  has trace equal to zero.

**2.20.** Prove that there exists a differentiable function  $f$  defined on some neighborhood of  $(1, 0) \in \mathbb{R}^2$  satisfying  $x \log f(x, y) + yf(x, y) = 0$ .

**2.21.** Find and classify all critical points of  $f$ , if  $f(x, y) = e^{x^2 - y^2 + 1}$ .

**2.22.** Find and classify all critical points of  $f$ , where  $f$  is given by  $f(x, y) = (x^2 - y^2)e^{-(x^2 + y^2)^2}$ .

**2.23.** Show that the function  $f$  given by

$$f(x, y) = \begin{cases} \frac{\sin(xy) - xy}{x^2 + y^2} & \text{if } (x, y) \neq \mathbf{0}, \\ 0 & \text{if } (x, y) = \mathbf{0} \end{cases}$$

is continuous everywhere. *Hint:* Look at the Taylor polynomial of the function  $(x, y) \mapsto \sin(xy)$  at  $\mathbf{0}$ .

**2.24.** Recall from Chapter 1 that a *polynomial of degree  $k$*  on  $\mathbb{R}^n$  is a function  $f$  of the form

$$f = a_0 + \sum_{j=1}^k \sum_{1 \leq i_1, \dots, i_j \leq n} a_{i_1, \dots, i_j} u^{i_1} \cdots u^{i_j}, \quad a_0, a_{i_1, \dots, i_j} \in \mathbb{R}.$$

Prove that the Taylor polynomial of degree  $k$  at any point of such a function  $f$  equals  $f$ .

**2.25.** Write  $2x^2 - y^2 + 3xy - 1$  as a polynomial in  $(x - 1)$  and  $(y + 1)$ .

**2.26.** Let  $b$  denote a scalar product on  $\mathbb{R}^n$ ,  $L$  the associated self-adjoint operator, so that  $b(\mathbf{u}, \mathbf{v}) = \langle L\mathbf{u}, \mathbf{v} \rangle$ . The *quadratic form associated to  $b$*  is the function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  defined by  $f(\mathbf{a}) = b(\mathbf{a}, \mathbf{a})$ .



- (a) Show that the Jacobian matrix  $[Df(\mathbf{a})]$  of  $f$  at  $\mathbf{a}$  equals  $2(L\mathbf{a})^T$ .
- (b) Prove that the Hessian of  $f$  is  $2L$ .
- (c) Part (a) implies that the set of critical points of  $f$  is the kernel of  $L$ , and in particular the origin is always a critical point. Suppose the kernel is trivial, so that the origin is the only critical point. *Without using Theorem 2.73*, show that the origin is an absolute minimum if all eigenvalues of  $L$  are positive, an absolute maximum if they are all negative, and a saddle point in all other cases.

**2.27.** A map  $\mathbf{f} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^{n+k}$  that has rank  $n$  at every point in  $U$  is called an *immersion* of  $U$  into  $\mathbb{R}^{n+k}$ . An injective immersion is called an *imbedding* if  $\mathbf{f}^{-1} : \mathbf{f}(U) \rightarrow U$  is continuous.

- (a) Show that an immersion  $\mathbf{f} : U \rightarrow \mathbb{R}^{n+k}$  is always locally an imbedding; i.e., any  $\mathbf{a} \in U$  admits a neighborhood  $V$  such that the restriction of  $\mathbf{f}$  to  $V$  is an imbedding.
- (b) Let  $\mathbf{c} : (0, 2\pi) \rightarrow \mathbb{R}^2$  be given by  $\mathbf{c}(t) = (\sin t, \sin 2t)$ . The image of  $\mathbf{c}$  is called a *lemniscate*. Show that  $\mathbf{c}$  is an injective immersion but not an imbedding.

**2.28.** If  $\mathbf{X}$  is a vector field with flow  $\Phi$ , define a map  $\Phi_t$  by  $\Phi_t(\mathbf{p}) = \Phi(t, \mathbf{p})$ , for all  $t$  and  $\mathbf{p}$  for which the formula makes sense. Prove that  $\Phi_t \circ \Phi_s = \Phi_s \circ \Phi_t = \Phi_{t+s}$ .

**2.29.** Let  $U$  be open in  $\mathbb{R}^n$ ,  $\mathbf{p} \in U$ ,  $\mathbf{u} = (\mathbf{p}, \mathbf{v}) \in \mathbb{R}_\mathbf{p}^n$ .

- (a) If  $\mathbf{X}$  is a smooth vector field on  $U$  represented by the map  $\mathbf{g}$ , show that  $D_\mathbf{u}\mathbf{X} = (\mathbf{p}, D\mathbf{g}(\mathbf{p})\mathbf{v})$ .
- (b) If  $f : U \rightarrow \mathbb{R}$  is smooth, show that the derivative of  $f$  in direction  $\mathbf{u}$  (as defined in Definition 2.8.7) equals  $Df(\mathbf{p})\mathbf{v}$ .

**2.30.** Determine the velocity and acceleration vector fields of the curve  $\mathbf{c} : I \rightarrow \mathbb{R}^3$ ,  $\mathbf{c}(t) = (a \cos t, a \sin t, bt)$ ,  $a, b > 0$ , which parametrizes a helix.

**2.31.** Given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , extend the Hessian of  $f$  on  $\mathbb{R}^n$  to any tangent space by identifying both spaces in the usual way; i.e., for  $\mathbf{a} \in \mathbb{R}^n$ , define  $\tilde{H}_f : \mathbb{R}_\mathbf{a}^n \rightarrow \mathbb{R}_\mathbf{a}^n$  by  $\tilde{H}_f \mathcal{I}_\mathbf{a} \mathbf{u} = H_f \mathbf{u}$ ,  $\mathbf{u} \in \mathbb{R}^n$ , and extend the Hessian form similarly. For the sake of simplicity, we denote  $\tilde{H}_f$  by  $H_f$  again.

- (a) Show that  $H_f \mathbf{u} = D_\mathbf{u} \nabla f$  for  $\mathbf{u} \in T\mathbb{R}^n$ .
- (b) Given  $\mathbf{a}, \mathbf{u} \in \mathbb{R}^n$ , prove that  $h_f(\mathcal{I}_\mathbf{a} \mathbf{u}, \mathcal{I}_\mathbf{a} \mathbf{u}) = (f \circ \mathbf{c})''(0)$ , where  $\mathbf{c}(t) = \mathbf{a} + t\mathbf{u}$ .

**2.32.** Suppose  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a smooth map such that  $|Df(\mathbf{p})\mathbf{u}| = |\mathbf{u}|$  for all  $\mathbf{p}, \mathbf{u} \in \mathbb{R}^n$ . Prove that  $\mathbf{f}$  is a diffeomorphism which preserves distances; i.e.,  $|\mathbf{f}(\mathbf{p}) - \mathbf{f}(\mathbf{q})| = |\mathbf{p} - \mathbf{q}|$  for any  $\mathbf{p}, \mathbf{q} \in \mathbb{R}^n$ . *Hint:  $\mathbf{f}$  preserves the length of curves, so it cannot increase distances. The same is true for any (a priori only local) inverse, so it cannot decrease them either.*

**2.33.** Determine the flow of the position vector field  $\mathbf{P} = \sum_i u^i D_i$ . Is  $\mathbf{P}$  a complete vector field?

**2.34.** This problem investigates integral curves of vector fields on  $\mathbb{R}^n$  of the form

$$\mathbf{X} = \sum_{i=1}^n \left( \sum_{j=1}^n a_{ij} u^j \right) \mathbf{D}_i, \quad a_{ij} = a_{ji} \in \mathbb{R},$$

without using Example 2.8.1.

- (a) Prove that  $\mathbf{X}(\mathbf{u}) = (\mathbf{u}, A\mathbf{u})$ , where  $A = (a_{ij})$  is a symmetric matrix, and that  $\mathbf{c} : I \rightarrow \mathbb{R}^n$  is an integral curve of  $\mathbf{X}$  if and only if  $\mathbf{c}' = A\mathbf{c}$ .
- (b) By the spectral theorem, there exists an orthogonal matrix  $P$  and a diagonal matrix

$$D = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}$$

such that  $A = PDP^T$ . Show that the integral curve  $\mathbf{c}$  of  $\mathbf{X}$  with  $\mathbf{c}(0) = \mathbf{a} \in \mathbb{R}^n$  is given by

$$\mathbf{c}(t) = P \cdot \begin{bmatrix} e^{\lambda_1 t} & 0 & \dots & 0 \\ 0 & e^{\lambda_2 t} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & e^{\lambda_n t} \end{bmatrix} \cdot P^T \mathbf{a}.$$

In particular,  $\mathbf{X}$  is complete.

**2.35.** Use the results of the previous exercise to find the integral curve  $\mathbf{c}$  of the vector field  $\mathbf{X}$ , where  $\mathbf{X} = (u^2 - u^3)\mathbf{D}_1 + (u^1 - u^3)\mathbf{D}_2 + (u^1 - u^2)\mathbf{D}_3$ , if  $\mathbf{c}(0) = [a_1 \ a_2 \ a_3]^T$ .

**2.36.** Show that  $\mathbb{R}^3$  is a Lie algebra with the cross product.

**2.37.** Let  $\mathbf{X}$  be a parallel vector field on  $\mathbb{R}^n$ .

- (a) Prove that there exists some  $\mathbf{u} \in \mathbb{R}^n$  such that  $\mathbf{X}(\mathbf{p}) = \mathcal{I}_{\mathbf{p}}\mathbf{u}$  for all  $\mathbf{p} \in \mathbb{R}^n$ .
- (b) Show that a vector field  $\mathbf{Y}$  satisfies  $[\mathbf{X}, \mathbf{Y}] \equiv \mathbf{0}$  if and only if  $\mathbf{X}(\mathbf{a} + t\mathbf{u}) = (\mathcal{I}_{\mathbf{a}+t\mathbf{u}} \circ \mathcal{I}_{\mathbf{a}}^{-1})\mathbf{X}(\mathbf{a})$  for all  $\mathbf{a} \in \mathbb{R}^n$ ,  $t \in \mathbb{R}$ . Interpret this geometrically.

**2.38.** Determine the flow  $\Phi$  of the vector field  $\mathbf{Y}$ , where

$$\mathbf{Y} = -u^3 \mathbf{D}_1 - u^4 \mathbf{D}_2 + u^1 \mathbf{D}_3 + u^2 \mathbf{D}_4,$$

see Example 2.8.1. Interpret  $\Phi_t : \mathbb{R}^4 \rightarrow \mathbb{R}^4$  geometrically.

**2.39.** Consider the vector fields  $\mathbf{X}$ ,  $\mathbf{Y}$ ,  $\mathbf{Z}$  on  $\mathbb{R}^4$  given by

$$\mathbf{X} = -u^2 \mathbf{D}_1 + u^1 \mathbf{D}_2 + u^4 \mathbf{D}_3 - u^3 \mathbf{D}_4,$$

$$\mathbf{Y} = -u^3 \mathbf{D}_1 - u^4 \mathbf{D}_2 + u^1 \mathbf{D}_3 + u^2 \mathbf{D}_4,$$

$$\mathbf{Z} = -u^4 \mathbf{D}_1 + u^3 \mathbf{D}_2 - u^2 \mathbf{D}_3 + u^1 \mathbf{D}_4.$$

Prove that the set  $\{a\mathbf{X} + b\mathbf{Y} + c\mathbf{Z} \mid a, b, c \in \mathbb{R}\}$  is a 3-dimensional Lie algebra.

**2.40.** Let  $X, Y$  denote vector fields on  $\mathbb{R}^n$  with flows  $\Phi_t$  and  $\Psi_t$  respectively. Given  $\mathbf{p} \in \mathbb{R}^n$ , define a curve  $\mathbf{c}$  on a neighborhood of  $\mathbf{0}$  by

$$\mathbf{c}(t) = (\Psi_{-t} \circ \Phi_{-t} \circ \Psi_t \circ \Phi_t)(\mathbf{p}).$$

We have seen that if  $[X, Y] \equiv \mathbf{0}$ , then  $\mathbf{c}(t) = \mathbf{p}$  for all  $t$ . The object of this exercise is to show that even when the bracket of  $X$  and  $Y$  is not zero (so that  $\mathbf{c}$  is not a constant curve), its derivative  $\mathbf{c}'(0)$  nevertheless vanishes. To see this, define maps  $F_i : U \rightarrow \mathbb{R}^n$  on a neighborhood  $U \subset \mathbb{R}^2$  of the origin by

$$\begin{aligned} F_1(t, s) &= (\Psi_t \circ \Phi_s)(\mathbf{p}), \\ F_2(t, s) &= (\Phi_{-t} \circ \Psi_s \circ \Phi_s)(\mathbf{p}), \\ F_3(t, s) &= (\Psi_{-t} \circ \Phi_{-s} \circ \Psi_s \circ \Phi_s)(\mathbf{p}). \end{aligned}$$

(a) Prove that

$$F_{1*} \mathbf{D}_1 = Y \circ F_1, \quad F_{2*} \mathbf{D}_1 = -X \circ F_2, \quad F_{3*} \mathbf{D}_1 = -Y \circ F_3,$$

and that  $F_{1*} \mathbf{D}_2(0, s) = (X \circ F_1)(0, s)$ .

(b) Notice that  $\mathbf{c}(t) = F_3(t, t)$ ,  $F_2(0, t) = F_1(t, t)$ , and  $F_3(0, t) = F_2(t, t)$ . Use the chain rule to show that  $\mathbf{c}'(0) = \mathbf{0}$ .

**2.41.** Two Lie algebras  $(V_1, [\cdot, \cdot]_1)$  and  $(V_2, [\cdot, \cdot]_2)$  are said to be isomorphic if there exists a vector space isomorphism  $L : V_1 \rightarrow V_2$  that preserves Lie brackets; i.e.,

$$L[\mathbf{u}, \mathbf{v}]_1 = [L\mathbf{u}, L\mathbf{v}]_2, \quad \mathbf{u}, \mathbf{v} \in V_1.$$

In this case,  $L$  is called a Lie algebra isomorphism. Show that the Lie algebras from Exercises 2.36 and 2.39 are isomorphic. Give an example of a 3-dimensional Lie algebra that is not isomorphic to these. Thus, unlike vector spaces, Lie algebras of the same dimension need not be isomorphic.

**2.42.** (a) Show that if  $L$  is a self-adjoint operator, then so is  $e^L$ .

(b) Show that if  $L$  is a skew-adjoint operator, then  $e^L$  is orthogonal.

**2.43.** This exercise generalizes the discussion from Example 2.8.1 regarding the system of ODEs  $\mathbf{c}' = A\mathbf{c}$  to the case when  $A$  is no longer constant.

(a) Prove that  $\exp : M_{n,n} \cong \mathbb{R}^{n^2} \rightarrow M_{n,n}$  is differentiable with derivative  $D \exp(\mathbf{0})A = \exp(\mathbf{0})A = A$  at the origin.

(b) Show that if  $AB = BA$ , then  $D \exp(B)A = \exp(B)A = A \exp(B)$ .

(c) If  $\mathbf{A} : \mathbb{R} \rightarrow M_{n,n}$  is a smooth curve, prove that  $\mathbf{c} : \mathbb{R} \rightarrow M_{n,n}$ , where  $\mathbf{c}(t) = \exp(\mathbf{A}(t))$ , has derivative  $\mathbf{c}'(t) = \mathbf{A}'(t)\mathbf{c}(t)$ , provided  $\mathbf{A}(t)\mathbf{A}'(t) = \mathbf{A}'(t)\mathbf{A}(t)$  for all  $t$ .

(d) With  $\mathbf{A}$  as in (c), show that the system of linear ODEs  $\mathbf{c}'(t) = \mathbf{A}(t)\mathbf{c}(t)$  has as solution

$$\mathbf{c}(t) = \exp\left(\int_0^t \mathbf{A}\right) \mathbf{c}(0),$$

provided  $\mathbf{A}(s)\mathbf{A}(t) = \mathbf{A}(t)\mathbf{A}(s)$  for all  $s, t$ . Here, the  $(i, j)$ -th entry of  $\int_0^t \mathbf{A}$  is

$$u^{ij} \circ \int_0^t \mathbf{A} = \int_0^t u^{ij} \circ \mathbf{A}.$$

**2.44.** Use the results from the previous exercise to solve the system

$$x'(t) = -x(t) - ty(t)$$

$$y'(t) = tx(t) - y(t)$$

**2.45.** Suppose  $\Phi : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}^n$  is differentiable. For each  $t \in \mathbb{R}$ , define  $\Phi_t : \mathbb{R}^n \rightarrow \mathbb{R}^n$  by  $\Phi_t(\mathbf{a}) = \Phi(t, \mathbf{a})$ ,  $\mathbf{a} \in \mathbb{R}^n$ .  $\{\Phi_t\}_{t \in \mathbb{R}}$  is called a *one-parameter group of diffeomorphisms of  $\mathbb{R}^n$*  if

(1)  $\Phi_0 = 1_{\mathbb{R}^n}$ , and

(2)  $\Phi_{s+t} = \Phi_s \circ \Phi_t$ ,  $s, t \in \mathbb{R}$ .

Prove that one-parameter groups are in 1-1 correspondence with complete vector fields on  $\mathbb{R}^n$ .

**2.46.** Let  $C$  be a closed set in  $\mathbb{R}^n$ , and  $f : C \rightarrow \mathbb{R}$  a smooth function. Show that if  $U$  is any open set containing  $C$ , then there exists a smooth function  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ , such that  $g$  agrees with  $f$  on  $C$ , and the support of  $g$  lies in  $U$ . Show this is no longer necessarily true if  $C$  is not assumed to be closed. In particular, any smooth function that is defined on a closed subset of Euclidean space is extendable to a smooth function on all of Euclidean space.

**2.47.** Let  $\mathbf{X}$  be a smooth vector field on a closed set  $C \subset \mathbb{R}^n$  (as for functions, this means that there exists a smooth vector field  $\tilde{\mathbf{X}}$  on some open set  $U$  containing  $C$  such that the restriction of  $\tilde{\mathbf{X}}$  to  $C$  equals  $\mathbf{X}$ ). Show that  $\mathbf{X}$  is extendable to a smooth vector field on  $\mathbb{R}^n$ ; i.e., there exists a smooth vector field on  $\mathbb{R}^n$  whose restriction to  $C$  equals  $\mathbf{X}$ .

## 3 Manifolds

Differential geometry arguably began with the study of surfaces in three-dimensional Euclidean space. In this chapter, we give a precise meaning to the word “surface”, and generalize this concept to higher dimensions.

### 3.1 Submanifolds of Euclidean space

**Definition 3.1.1.** A subset  $M \subset \mathbb{R}^{n+k}$  is said to be an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+k}$  if each  $\mathbf{p} \in M$  admits an open neighborhood  $U$  in  $\mathbb{R}^{n+k}$  and there exists a one-to-one differentiable map  $\mathbf{h} : \mathbb{R}^n \supset V \rightarrow \mathbb{R}^{n+k}$  defined on some open set  $V$ , such that

- (1)  $\mathbf{h}$  has maximal rank ( $= n$ ) everywhere;
- (2)  $\mathbf{h}(V) = U \cap M$ , and
- (3)  $\mathbf{h}^{-1} : U \cap M \rightarrow V$  is continuous.

Thus, loosely speaking, an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+k}$  is a subset  $M$  of  $\mathbb{R}^{n+k}$  such that each point of  $M$  has an open neighborhood (in  $M$ ) that “looks like” an open set in  $\mathbb{R}^n$ . In the terminology of Exercise 2.27,  $M$  is a submanifold if for any  $\mathbf{p} \in M$  there exists an imbedding of an open set in  $\mathbb{R}^n$  into  $\mathbb{R}^{n+k}$  whose image is an open neighborhood of  $\mathbf{p}$  in  $M$ . Condition (3) in the above definition is equivalent to requiring that open sets in  $V$  get mapped by  $\mathbf{h}$  to open sets of  $M$ . A subset  $M$  for which the first two conditions hold but the third one does not is called an *immersed submanifold*.

The pair  $(V, \mathbf{h})$  in the definition is called a *local parametrization* of  $M$  around  $\mathbf{p}$ . The inverse of a parametrization, or more precisely, the pair  $(\mathbf{h}(V), \mathbf{h}^{-1})$  is called a *chart* of  $M$  around  $\mathbf{p}$ . It is common practice to denote charts by  $(U, \mathbf{x})$ ,  $(V, \mathbf{y})$ , and so on. A collection of charts whose domains form an open cover of  $M$  is called an *atlas* of  $M$ . For the sake of brevity, and when there is no possible confusion about what the “ambient space”  $\mathbb{R}^{n+k}$  is, we sometimes say that  $M$  is an  $n$ -dimensional manifold, or just an  $n$ -manifold. The dimension is often implicitly specified by adding a superscript, as in “ $M^n$ ”.

Notice also that any submanifold of  $\mathbb{R}^{n+k}$  is a submanifold of  $\mathbb{R}^{n+k+l}$  for  $l \in \mathbb{N}$  if one composes every parametrization with the map  $\mathbb{R}^{n+k}$  to  $\mathbb{R}^{n+k+l}$  which sends  $\mathbf{p}$  to  $(\mathbf{p}, \mathbf{0}) \in \mathbb{R}^{n+k} \times \mathbb{R}^l$ . It is customary to choose the ambient space with the smallest dimension.

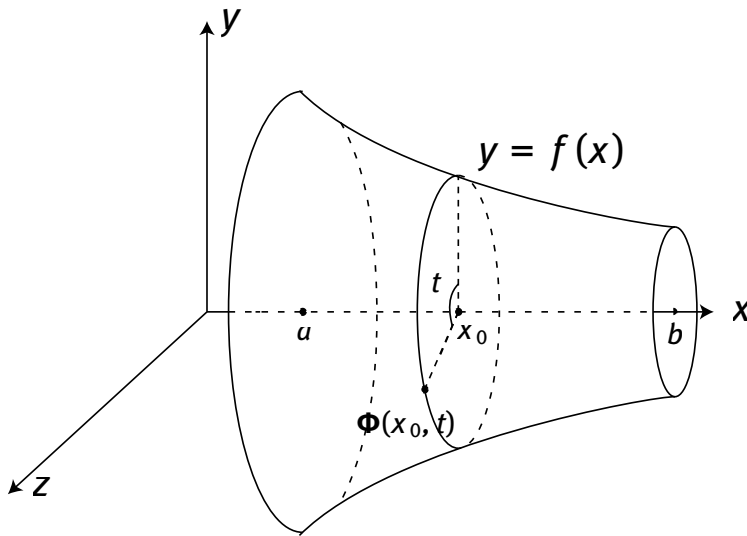
A 2-dimensional submanifold of  $\mathbb{R}^3$  is often called a *surface*. More generally, a *hypersurface* is an  $(n - 1)$ -dimensional submanifold of  $\mathbb{R}^n$  for some  $n \in \mathbb{N}$ .

Any open set  $U$  in  $\mathbb{R}^n$  is a trivial example of a submanifold of  $\mathbb{R}^n$ , with the atlas  $\{(U, \iota)\}$ , where  $\iota : U \rightarrow \mathbb{R}^n$  is the inclusion map. More interesting manifolds are described below:

**Examples 3.1.1.** (i) Let  $f : V \rightarrow \mathbb{R}$  be a differentiable function on an open set  $V$  in  $\mathbb{R}^2$ . Then the graph

$$M = \{(\mathbf{a}, f(\mathbf{a})) \mid \mathbf{a} \in V\}$$

of  $f$  is a 2-dimensional submanifold of  $\mathbb{R}^3$ : The map  $\mathbf{h} : V \rightarrow \mathbb{R}^3$ , where  $\mathbf{h}(\mathbf{a}) = (\mathbf{a}, f(\mathbf{a}))$ , is certainly one-to-one with maximal rank everywhere, and the set  $U$  in the definition may be taken to be  $V \times \mathbb{R}$ . In the same way, if  $V_0$  is open in  $V$ , then  $\mathbf{h}(V_0) = M \cap (V_0 \times \mathbb{R})$  is open in  $M$ , so that  $\mathbf{h}^{-1}$  is continuous. More generally, the graph of a function defined on an open set in  $\mathbb{R}^n$  is an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+1}$ .



**Fig. 3.1:** Parametrizing a surface of revolution

(ii) If  $f : (a, b) \rightarrow \mathbb{R}$  is a positive function, the *surface of revolution* obtained by revolving the graph of  $f$  around the  $x$ -axis is the 2-dimensional submanifold of  $\mathbb{R}^3$  described by two parametrizations  $\mathbf{h} : (a, b) \times (0, 2\pi) \rightarrow \mathbb{R}^3$ , and  $\mathbf{k} : (a, b) \times (\varepsilon, 2\pi + \varepsilon) \rightarrow \mathbb{R}^3$ ,  $0 < \varepsilon < 2\pi$ , both given by the same formula

$$\mathbf{h}(x, t) = \mathbf{k}(x, t) = (x, f(x) \cos t, f(x) \sin t).$$

The first parametrization covers the whole surface except for the original graph, which is why a second one is needed. The Jacobian matrix of either equals

$$\begin{bmatrix} 1 & 0 \\ f'(x) \cos t & -f(x) \sin t \\ f'(x) \sin t & f(x) \cos t \end{bmatrix},$$

which has rank 2, unless  $f$  is zero somewhere, but that was ruled out at the beginning. This establishes the first condition; the other two are straightforward, and left to the reader.

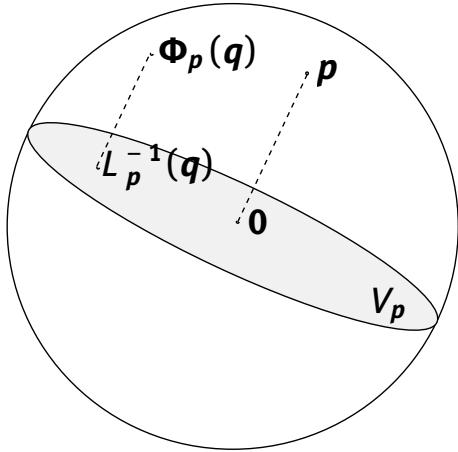


Fig. 3.2: A local parametrization of  $S^n(r)$  around  $\mathbf{p}$

(iii) The sphere  $S^n(r) = \{\mathbf{p} \in \mathbb{R}^{n+1} \mid |\mathbf{p}| = r\}$  of radius  $r > 0$  centered at the origin is an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+1}$ . A local parametrization  $(V, \mathbf{h})$  around the north pole  $r\mathbf{e}_{n+1} = (0, \dots, 0, r)$  is given by

$$\mathbf{h}(\mathbf{a}) = (\mathbf{a}, \sqrt{1 - |\mathbf{a}|^2}), \quad \mathbf{a} \in V := B_r(\mathbf{0}) \subset \mathbb{R}^n.$$

The image of  $\mathbf{h}$  is the open northern hemisphere. This parametrization can be modified to yield one around any  $\mathbf{p} \in S^n(r)$ : Let  $L_{\mathbf{p}} : \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$  be any linear isometry that maps  $\mathbf{p}$  to  $r\mathbf{e}_{n+1}$  (for example, extend  $\mathbf{p}/r$  to an orthonormal basis  $\mathbf{u}_1, \dots, \mathbf{u}_n, \mathbf{p}/r$  of  $\mathbb{R}^{n+1}$ , define  $L_{\mathbf{p}}(\mathbf{u}_i) = \mathbf{e}_i$  for  $i \leq n$ ,  $L_{\mathbf{p}}(\mathbf{p}/r) = \mathbf{e}_{n+1}$ , and extend linearly; i.e.,  $L_{\mathbf{p}}$  maps the vector  $\sum_{i=1}^n a_i \mathbf{u}_i + a_{n+1}(\mathbf{p}/r)$  to  $(a_1, \dots, a_{n+1})$ ). If  $V_{\mathbf{p}}$  denotes the set of all  $\mathbf{q} \in \mathbf{p}^\perp$  with norm less than  $r$ , then  $L_{\mathbf{p}}$  maps  $V_{\mathbf{p}}$  isometrically onto the domain  $V$  of the parametrization  $(V, \mathbf{h})$  above. The pair  $(V, \mathbf{h}_{\mathbf{p}})$  is now a parametrization around  $\mathbf{p}$ , if  $\mathbf{h}_{\mathbf{p}} = L_{\mathbf{p}}^{-1} \circ \mathbf{h}$ .

(iv) Although the above parametrizations are easy to visualize, they are not the most efficient in the sense that a fairly large number of them is needed to cover the whole sphere. It is possible to give an atlas with only two charts: *stereographic projection*  $\chi : S^n(r) \setminus \{r\mathbf{e}_{n+1}\} \rightarrow \mathbb{R}^n$  from the north pole  $r\mathbf{e}_{n+1} = (0, \dots, 0, r)$  is the map that

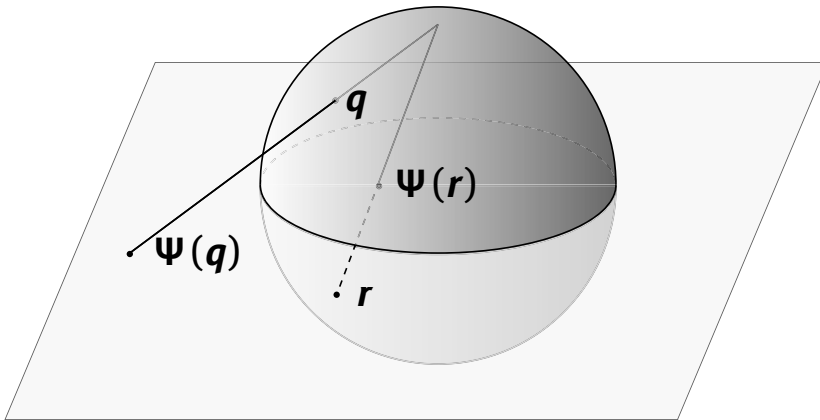


Fig. 3.3: Stereographic projection from the north pole

assigns to each  $\mathbf{q}$  of  $S^n(r)$  distinct from the north pole the unique point where the line from the north pole to  $\mathbf{q}$  intersects the hyperplane  $\mathbb{R}^n \times 0 = \{(\mathbf{p}, 0) \in \mathbb{R}^{n+1} \mid \mathbf{p} \in \mathbb{R}^n\}$ .

To find the point of intersection, notice that the line through the north pole and any other  $\mathbf{q} \in S^n(r)$  can be parametrized by the curve

$$t \mapsto (\mathbf{0}, r) + t(\mathbf{q} - (\mathbf{0}, r)) = t\mathbf{q} + (\mathbf{0}, (1-t)r),$$

with  $\mathbf{0}$  denoting the origin in  $\mathbb{R}^n$ . Its image intersects the equatorial hyperplane when  $t = r/(r - u^{n+1}(\mathbf{q}))$ , so that

$$\mathbf{x}(\mathbf{q}) = \frac{r}{r - u^{n+1}(\mathbf{q})} (u^1(\mathbf{q}), \dots, u^n(\mathbf{q})), \quad \mathbf{q} \in S^n(r) \setminus \{r\mathbf{e}_{n+1}\}.$$

The inverse  $\mathbf{h} : \mathbb{R}^n \rightarrow S^n(r)$  of  $\mathbf{x}$  is the parametrization given by

$$\mathbf{h}(\mathbf{p}) = \frac{1}{|\mathbf{p}|^2 + r^2} (2r^2\mathbf{p}, (|\mathbf{p}|^2 - r^2)r), \quad \mathbf{p} \in \mathbb{R}^n.$$

Combining stereographic projection from the north pole with the one from the south pole now yields an atlas of the sphere. The atlas is minimal in terms of the amount of charts that are needed. A space admitting an atlas with only one chart – as is the case of a sphere with one point deleted – is considered trivial. This perspective will change once curvature is introduced.

- (v) The curve  $\mathbf{c} : (0, 2\pi) \rightarrow \mathbb{R}^2$ , with  $\mathbf{c}(t) = (\sin t, \sin(2t))$ , parametrizes a lemniscate. The image is an immersed submanifold  $M$  of  $\mathbb{R}^2$ : the origin in the plane is the point  $\mathbf{c}(\pi)$ , and the set  $\mathbf{c}(3\pi/4, 5\pi/4)$  is not an open neighborhood of the origin in  $M$ , since any such neighborhood must contain points  $\mathbf{c}(t)$  for values of  $t$  arbitrarily close to 0 and  $2\pi$ .

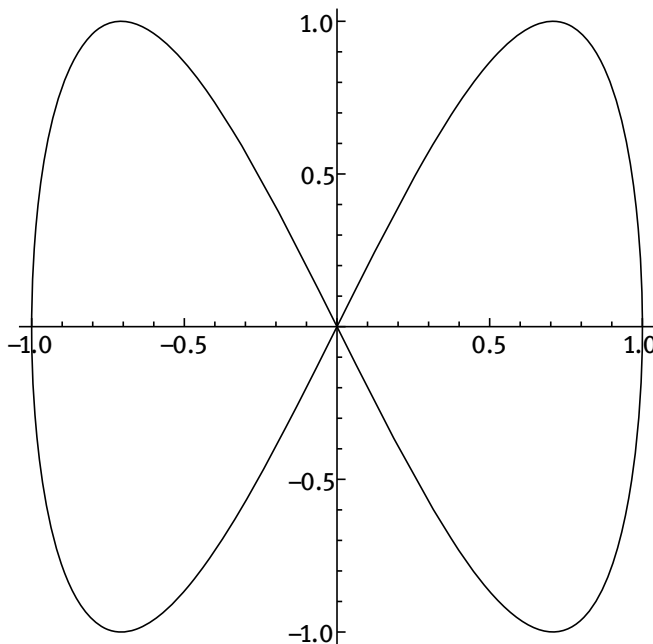


Fig. 3.4: The lemniscate in (v)



(vi) If  $M$  is an  $n$ -manifold, then so is any open subset  $U$  of  $M$ : in fact, any parametrization  $(V, \mathbf{h})$  of  $M$  around some  $\mathbf{p} \in U$ , when restricted to  $V \cap \mathbf{h}^{-1}(U)$ , is a parametrization of  $U$  around  $\mathbf{p}$ .

(vii) If  $M^n \subset \mathbb{R}^k$  and  $N^m \subset \mathbb{R}^l$  are manifolds, then the Cartesian product

$$M \times N = \{(\mathbf{p}, \mathbf{q}) \in \mathbb{R}^k \times \mathbb{R}^l \mid \mathbf{p} \in M, \mathbf{q} \in N\}$$

is a  $(n+m)$ -dimensional submanifold of  $\mathbb{R}^{k+l}$ : any parametrizations  $\mathbf{h} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^k$  of  $M$  and  $\mathbf{k} : \mathbb{R}^m \supset V \rightarrow \mathbb{R}^l$  of  $N$  generate a parametrization  $\mathbf{h} \times \mathbf{k}$  of  $M \times N$ , where  $\mathbf{h} \times \mathbf{k}(\mathbf{a}, \mathbf{b}) = (\mathbf{h}(\mathbf{a}), \mathbf{k}(\mathbf{b}))$ .

There is a particularly useful and simple way of constructing manifolds:  $\mathbf{q} \in \mathbb{R}^k$  is said to be a *regular value* of a map  $\mathbf{f} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^k$  if the derivative  $\mathbf{f}_*$  has maximal rank at every  $\mathbf{p} \in \mathbf{f}^{-1}(\mathbf{q})$ .

**Theorem 3.1.1.** *Let  $U$  be open in  $\mathbb{R}^{n+k}$ , and  $\mathbf{f} : U \rightarrow \mathbb{R}^k$ . If  $\mathbf{a} \in \mathbb{R}^k$  is a regular value of  $\mathbf{f}$  and  $M = \mathbf{f}^{-1}(\mathbf{a})$  is nonempty, then  $M$  is an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+k}$ .*

*Proof.* We shall construct a parametrization  $\mathbf{h}$  around any  $\mathbf{p} \in M$ . It may be assumed first of all that  $\mathbf{a} = \mathbf{0}$ , for if the result holds when  $\mathbf{a} = \mathbf{0}$ , we may, in the general case, apply it to  $\mathbf{g}$ , where  $\mathbf{g}(\mathbf{p}) = \mathbf{f}(\mathbf{p}) - \mathbf{a}$ . Similarly,  $\mathbf{p}$  may be assumed to be  $\mathbf{0}$ .

By the implicit function theorem, there exists a neighborhood  $W$  of  $\mathbf{0}$  in  $\mathbb{R}^{n+k}$  and a diffeomorphism  $\Psi : W \rightarrow \Psi(W) \subset \mathbb{R}^{n+k}$  such that  $\mathbf{f} \circ \Psi$  equals the restriction to  $W$  of the projection  $\pi_2 : \mathbb{R}^{n+k} = \mathbb{R}^n \times \mathbb{R}^k \rightarrow \mathbb{R}^k$  onto the second factor. Let  $\pi_1$  denote the projection onto the first factor  $\mathbb{R}^n$ ,  $V = \pi_1(W)$ ,  $\iota : \mathbb{R}^n \rightarrow \mathbb{R}^n \times \mathbb{R}^k$  the inclusion  $\iota(\mathbf{u}) = (\mathbf{u}, \mathbf{0})$ , and define  $\mathbf{h} : V \rightarrow \mathbb{R}^{n+k}$  to be  $\Psi \circ \iota$ .  $\mathbf{h}$  is one-to-one of maximal rank, being the composition of two maps that enjoy those properties. Furthermore,

$$\mathbf{f} \circ \mathbf{h} = \mathbf{f} \circ \Psi \circ \iota = \pi_2 \circ \iota = \mathbf{0},$$

so that  $\mathbf{h}(V) \subset M$ .

$$\begin{array}{ccc} W & \xrightarrow{\Psi} & \mathbb{R}^{n+k} \\ \pi_1 \downarrow & \nearrow \mathbf{h} & \downarrow \mathbf{f} \\ V & \xrightarrow{\iota} & \mathbb{R}^n \times \mathbb{R}^k \\ & \searrow \pi_2 & \downarrow \mathbf{f} \\ & & \mathbb{R}^k \end{array}$$

We claim that  $\mathbf{h}(V) = M \cap \Psi(W)$ ; more generally, if  $V_0 \subset V$  is open, we claim that

$$\mathbf{h}(V_0) = M \cap (\Psi(\pi_{1|W}^{-1}(V_0))).$$

(If  $V$  is a subset of the domain  $U$  of a map  $f$ , the notation  $f|_U$  refers to the *restriction* of  $f$  to  $U$ ; i.e.,  $f|_U := f \circ j$ , where  $j : U \hookrightarrow V$  is the inclusion map.) Notice that since  $\Psi(\pi_{1|W}^{-1}(V_0))$  is open in  $\mathbb{R}^{n+k}$ , the claim, once established, will imply that  $\mathbf{h}$  maps open sets to open sets in  $M$ , thereby concluding the proof.

One inclusion is already clear:

$$\mathbf{h}(V_0) = \Psi \circ \iota(V_0) \subset \Psi(\pi_{1|W}^{-1}(V_0)),$$

and, of course, it is also contained in  $M$ , since  $\mathbf{h}(V)$  is. Conversely, consider any  $\mathbf{q} \in M \cap (\Psi(\pi_{1|W}^{-1}(V_0)))$ . Then  $\mathbf{q} = \Psi(\mathbf{r})$  for some unique  $\mathbf{r} \in \pi_{1|W}^{-1}(V_0)$ , and  $\pi_2(\mathbf{r}) = (\mathbf{f} \circ \Psi)(\mathbf{r}) = \mathbf{f}(\mathbf{q}) = \mathbf{0}$ . Thus,

$$\mathbf{r} = (\tilde{\mathbf{r}}, \mathbf{0}) = \iota(\tilde{\mathbf{r}})$$

for a unique  $\tilde{\mathbf{r}} = \pi_1(\mathbf{r}) \in V_0$ , and  $\mathbf{q} = \Psi(\mathbf{r}) = \mathbf{h}(\tilde{\mathbf{r}}) \in \mathbf{h}(V_0)$ , as claimed.  $\square$

**Example 3.1.2.** Given  $r > 0$ , the function  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ , with  $f(\mathbf{a}) = |\mathbf{a}|^2 - r^2$ , has 0 as regular value, since  $[Df(\mathbf{p})] = 2\mathbf{p} \neq \mathbf{0}$  if  $\mathbf{p} \in f^{-1}(0) = S^n(r)$ . This yields a shorter proof of the fact that the sphere is a submanifold of Euclidean space.

**Definition 3.1.2.** Let  $M$  be an  $n$ -dimensional submanifold of Euclidean space  $\mathbb{R}^{n+k}$ . The *tangent space*  $M_{\mathbf{p}}$  of  $M$  at  $\mathbf{p} \in M$  is the collection of velocity vectors  $\dot{\mathbf{c}}(0)$  of all curves  $\mathbf{c} : I \rightarrow M$  defined on some open interval  $I$  containing 0 such that  $\mathbf{p} = \mathbf{c}(0)$ .

Thus,  $M_{\mathbf{p}}$  is a subset of  $\mathbb{R}^{n+k}$ . What is not immediately clear is that it is actually a subspace:

**Theorem 3.1.2.** *If  $M$  is an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+k}$ , then  $M_{\mathbf{p}}$  is an  $n$ -dimensional subspace of  $\mathbb{R}^{n+k}$  for every  $\mathbf{p} \in M$ .*

*Proof.* Consider a parametrization  $(V, \mathbf{h})$  around  $\mathbf{p}$ ,  $V \subset \mathbb{R}^n$ . By composing  $\mathbf{h}$  with a translation if necessary, it may be assumed that  $\mathbf{0} \in V$ , and that  $\mathbf{h}(\mathbf{0}) = \mathbf{p}$ . If  $\mathbf{c} : I \rightarrow M$  is a curve passing through  $\mathbf{p}$  at 0, then, after restricting the domain of  $\mathbf{c}$  if necessary,  $\mathbf{h}^{-1} \circ \mathbf{c}$  is a curve in  $V$ , and

$$\dot{\mathbf{c}}(0) = \mathbf{c}_{*0}D(0) = (\mathbf{h} \circ \mathbf{h}^{-1} \circ \mathbf{c})_{*0}D(0) = \mathbf{h}_{*\mathbf{0}}(\mathbf{h}^{-1} \circ \mathbf{c})_{*0}D(0) = \mathbf{h}_{*\mathbf{0}}\mathbf{v},$$

where  $\mathbf{v} = (\mathbf{h}^{-1} \circ \mathbf{c})_{*0}D(0) \in \mathbb{R}_0^n$  is the velocity vector of the curve  $\mathbf{h}^{-1} \circ \mathbf{c}$  in  $V$ . Thus,  $M_{\mathbf{p}} \subset \mathbf{h}_{*\mathbf{0}}\mathbb{R}_0^n$ . Conversely, if  $\mathbf{c}$  is a curve in  $V$  with  $\mathbf{c}(0) = \mathbf{0}$ , then  $\mathbf{h} \circ \mathbf{c}$  is a curve in  $M$  passing through  $\mathbf{p}$  at 0, so that  $\mathbf{h}_*(\dot{\mathbf{c}}(0)) \in M_{\mathbf{p}}$ . We have therefore shown that

$$M_{\mathbf{p}} = \mathbf{h}_{*\mathbf{0}}\mathbb{R}_0^n. \tag{3.1.1}$$

Since  $\mathbf{h}$  has rank  $n$ , this concludes the proof.  $\square$

If instead of parametrizations, the manifold is given by  $\mathbf{f}^{-1}(\mathbf{0})$ , where  $\mathbf{0}$  is a regular value of  $\mathbf{f}$ , the following description of tangent spaces holds:

**Proposition 3.1.1.** *If  $M = \mathbf{f}^{-1}(\mathbf{a})$ , where  $\mathbf{a}$  is a regular value of the map  $\mathbf{f} : \mathbb{R}^{n+k} \supset U \rightarrow \mathbb{R}^k$ , then*

$$M_{\mathbf{p}} = \ker \mathbf{f}_{*\mathbf{p}}$$

for any  $\mathbf{p} \in M$ .

*Proof.* Since both sides of the above identity are vector spaces of the same dimension  $n$ , it suffices to show that one is contained in the other. So consider the velocity vector  $\mathbf{u} = \dot{\mathbf{c}}(0) \in M_p$  of a curve  $\mathbf{c}$  in  $M$  at time 0. Since  $\mathbf{c}(I) \subset M = \mathbf{f}^{-1}(\mathbf{a})$ ,  $\mathbf{f} \circ \mathbf{c} \equiv \mathbf{a}$ . Thus,

$$\mathbf{f}_{*p}\mathbf{u} = (\mathbf{f} \circ \mathbf{c})_{*0}D(0) = \mathbf{0},$$

and  $\mathbf{u} \in \ker \mathbf{f}_{*p}$ . This shows that  $M_p \subset \ker \mathbf{f}_{*p}$ , as claimed.  $\square$

**Examples 3.1.3.** (i) If  $M$  is the graph of a function  $f : \mathbb{R}^n \supset U \rightarrow \mathbb{R}$ , then  $M$  admits the parametrization  $\mathbf{h} : U \rightarrow \mathbb{R}^{n+1}$ ,  $\mathbf{h}(\mathbf{a}) = (\mathbf{a}, f(\mathbf{a}))$ . By (3.1.1), the tangent space of  $M$  at  $\mathbf{h}(\mathbf{a})$  is the image of  $\mathbf{h}_{*\mathbf{a}}$ . Now,

$$[D\mathbf{h}(\mathbf{a})] = \begin{bmatrix} 1_{\mathbb{R}^n} \\ [Df(\mathbf{a})] \end{bmatrix},$$

so that  $D\mathbf{h}(\mathbf{a})\mathbf{e}_i = \mathbf{e}_i + D_i f(\mathbf{a})\mathbf{e}_{n+1}$ . It follows that

$$M_{\mathbf{h}(\mathbf{a})} = \text{span}\{\mathbf{D}_i(\mathbf{h}(\mathbf{a})) + (D_i f)(\mathbf{a})\mathbf{D}_{n+1}(\mathbf{h}(\mathbf{a})) \mid i = 1, \dots, n\}.$$

Another way of describing the tangent space at  $\mathbf{h}(\mathbf{a})$  is to notice that  $M = g^{-1}(0)$ , where  $g : \mathbb{R}^{n+1} \supset U \times \mathbb{R} \rightarrow \mathbb{R}$  is given by  $g(a_1, \dots, a_{n+1}) = a_{n+1} - f(a_1, \dots, a_n)$ . Since

$$Dg(a_1, \dots, a_{n+1}) = \begin{bmatrix} -D_1 f(a_1, \dots, a_n) & \dots & -D_n f(a_1, \dots, a_n) & 1 \end{bmatrix},$$

0 is a regular value of  $g$  and Proposition 3.1.1 says that the tangent space is the orthogonal complement of the gradient of  $g$ : for  $\mathbf{a}, \mathbf{u} \in \mathbb{R}^{n+1}$ ,

$$\langle \nabla g(\mathbf{a}), \mathcal{I}_a \mathbf{u} \rangle = \langle Dg(\mathbf{a})^T, \mathbf{u} \rangle = Dg(\mathbf{a})\mathbf{u},$$

so that  $\mathcal{I}_a \mathbf{u}$  is orthogonal to  $\nabla g(\mathbf{a})$  if and only if  $\mathbf{u}$  belongs to the kernel of  $Dg(\mathbf{a})$ , or equivalently  $\mathcal{I}_a \mathbf{u}$  belongs to the kernel of  $g_{*\mathbf{a}}$ . The reader should check that the two descriptions of the tangent space agree.

(ii) If  $M$  is the sphere of radius  $r$  centered at the origin in  $\mathbb{R}^{n+1}$ , then  $M = f^{-1}(0)$ , where  $f(\mathbf{a}) = |\mathbf{a}|^2 - r^2$ . It was computed earlier that  $[Df(\mathbf{a})] = 2\mathbf{a}^T$ . This means that

$$\ker Df(\mathbf{a}) = \{\mathbf{u} \in \mathbb{R}^{n+1} \mid \mathbf{a}^T \mathbf{u} = \langle \mathbf{a}, \mathbf{u} \rangle = 0\} = \mathbf{a}^\perp,$$

and therefore

$$S^n(r)_a = \mathcal{I}_a \mathbf{a}^\perp.$$

(iii) Suppose  $\mathbf{h} : \mathbb{R}^n \supset U \rightarrow M$  and  $\mathbf{k} : \mathbb{R}^m \supset V \rightarrow N$  are local parametrizations of  $M$  and  $N$  respectively. The canonical isomorphism

$$\begin{aligned} \mathbb{R}_p^n \times \mathbb{R}_q^m &\xrightarrow{\cong} (\mathbb{R}^n \times \mathbb{R}^m)_{(\mathbf{p}, \mathbf{q})} \\ ((\mathbf{p}, \mathbf{u}), (\mathbf{q}, \mathbf{v})) &\longmapsto ((\mathbf{p}, \mathbf{q}), (\mathbf{u}, \mathbf{v})) \end{aligned}$$

induces a canonical isomorphism  $M_{\mathbf{p}} \times N_{\mathbf{q}} \cong (M \times N)_{(\mathbf{p}, \mathbf{q})}$  by means of the commutative diagram

$$\begin{array}{ccc} \mathbb{R}_{\mathbf{p}}^n \times \mathbb{R}_{\mathbf{q}}^m & \xrightarrow{h_{*p} \times k_{*q}} & M_{\mathbf{p}} \times N_{\mathbf{q}} \\ \cong \downarrow & & \downarrow \cong \\ (\mathbb{R}^n \times \mathbb{R}^m)_{(\mathbf{p}, \mathbf{q})} & \xrightarrow{(h, k)_{*(\mathbf{p}, \mathbf{q})}} & (M \times N)_{(\mathbf{p}, \mathbf{q})} \end{array}$$

### 3.2 Differentiable maps on manifolds

**Definition 3.2.1.** Let  $M^n, N^l \subset \mathbb{R}^k$  be manifolds,  $\iota : N \hookrightarrow \mathbb{R}^k$  the inclusion map. A map  $\mathbf{f} : N \rightarrow M$  is said to be *differentiable* or *smooth* if for every  $\mathbf{p} \in N$  and any parametrization  $(V, \mathbf{h})$  of  $N$  around  $\mathbf{p}$ , the map

$$\iota \circ \mathbf{f} \circ \mathbf{h} : V \rightarrow \mathbb{R}^k$$

is differentiable.

**Examples and Remarks 3.2.1.** (i) The identity map  $1_M : M \rightarrow M$  is differentiable, since this is equivalent to requiring that every parametrization be differentiable as a map into the ambient Euclidean space.

(ii) In the above definition, it suffices to check smoothness of  $\iota \circ \mathbf{f} \circ \mathbf{h}$  for *some* parametrization  $(V, \mathbf{h})$  around every  $\mathbf{p}$ : If  $\tilde{\mathbf{h}}$  is another parametrization around  $\mathbf{p}$ , then locally  $\iota \circ \mathbf{f} \circ \tilde{\mathbf{h}} = (\iota \circ \mathbf{f} \circ \mathbf{h}) \circ (\mathbf{h}^{-1} \circ \tilde{\mathbf{h}})$  which is a composition of smooth maps by the following:

**Theorem 3.2.1.** Let  $(V_i, \mathbf{h}_i)$  be two parametrizations around  $\mathbf{p} \in M^n$ . Then

$$\mathbf{h}_2^{-1} \circ \mathbf{h}_1 : \mathbb{R}^n \supset V_1 \cap \mathbf{h}_1^{-1}(\mathbf{h}_2(V_2)) \rightarrow \mathbf{h}_2^{-1}(\mathbf{h}_1(V_1) \cap \mathbf{h}_2(V_2)) \subset \mathbb{R}^n$$

is a diffeomorphism.

*Proof.* The map is bijective, continuous, and has continuous inverse. Once we show it is differentiable, it will follow by symmetry that the inverse is also differentiable, since the inverse is obtained by interchanging the indices. It may be assumed, after composing with (necessarily differentiable) translations if need be, that  $\mathbf{p} = \mathbf{0} = \mathbf{h}_i(\mathbf{0})$ . The implicit function theorem guarantees the existence of diffeomorphisms  $\mathbf{F}_i$  in a neighborhood of  $\mathbf{0}$  such that  $\mathbf{F}_i \circ \mathbf{h}_i = \iota$ , with  $\iota : \mathbb{R}^n \rightarrow \mathbb{R}^{n+k}$  mapping  $\mathbf{a}$  to  $(\mathbf{a}, \mathbf{0})$ . Furthermore, if  $\pi$  is the projection  $\mathbb{R}^{n+k} = \mathbb{R}^n \times \mathbb{R}^k \rightarrow \mathbb{R}^n$  onto the first factor, then

$$(\pi \circ \mathbf{F}_i) \circ \mathbf{h}_i = \pi \circ \iota \tag{3.2.1}$$

equals the identity  $1_{V_i}$  on  $V_i$ , so that  $\pi \circ \mathbf{F}_i$  is a left inverse for  $\mathbf{h}_i$ . Thus,

$$\mathbf{h}_2^{-1} \circ \mathbf{h}_1 = \pi \circ \mathbf{F}_2 \circ \mathbf{h}_1, \tag{3.2.2}$$

which is a composition of differentiable maps.  $\square$

(iii) The above theorem implies that if  $(V, \mathbf{h})$  is a local parametrization of  $M$ , then  $\mathbf{h}^{-1}$  is differentiable as a map from  $\mathbf{h}(V)$  (which is a manifold, since it is open in  $M$ ). In other words, using the terminology of the previous section, charts are differentiable maps.

If, however, we view a local parametrization of  $M$  as a map into Euclidean space, then it admits no differentiable inverse (since the inverse is not defined on an open set of the ambient Euclidean space). Nevertheless, the inverse may be extended to a differentiable map on an open set:

**Proposition 3.2.1.** *If  $(V, \mathbf{h})$  is a local parametrization of  $M^n \subset \mathbb{R}^{n+k}$ , then there exists an open set  $U \subset \mathbb{R}^{n+k}$  and a differentiable map  $\mathbf{G} : U \rightarrow \mathbb{R}^n$  such that  $\mathbf{G} \circ \mathbf{h} = 1_V$ .*

*Proof.* As usual, we may assume that  $V$  contains  $\mathbf{0} \in \mathbb{R}^n$  and that  $\mathbf{h}(\mathbf{0}) = \mathbf{0} \in \mathbb{R}^{n+k}$ . If  $\mathbf{F}$  is the diffeomorphism (guaranteed by the implicit function theorem) that satisfies  $\mathbf{F} \circ \mathbf{h} = \iota$ , set  $\mathbf{G} = \pi \circ \mathbf{F}$  (with  $\iota$  and  $\pi$  denoting the same maps used in the proof of the above theorem). The statement then follows as in (3.2.1).  $\square$

One useful property of manifolds is the existence of “bump functions” in a neighborhood of any point:

**Proposition 3.2.2.** *If  $(U, \mathbf{x})$  is a chart of  $M^n$  around  $\mathbf{p}$ , then there exists a differentiable function  $\psi : M \rightarrow \mathbb{R}$  such that*

- (1)  $0 \leq \psi \leq 1$ ;
- (2)  $\psi \equiv 1$  on some neighborhood of  $\mathbf{p}$ , and
- (3) the support of  $\psi$  is contained in  $U$ .

*Proof.* By Theorem 2.2.5, there exists a smooth  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$  with values in  $[0, 1]$ , support in  $\mathbf{x}(U)$ , that equals 1 in some neighborhood of  $\mathbf{x}(\mathbf{p})$ .  $\psi$  may now be taken to equal  $\varphi \circ \mathbf{x}$  on  $U$  and zero outside  $U$ . This function is smooth because it is so in a neighborhood of any point: if the point is in  $U$ , then the neighborhood may be taken to be  $U$ ; otherwise the point lies in the complement of the support of  $\varphi \circ \mathbf{x}$ , and this is an open set on which the function vanishes.  $\square$

Notice that (3.2.2) also yields a formula for the derivative of  $\mathbf{h}_2^{-1} \circ \mathbf{h}_1$ : even though the inverse of  $\mathbf{h}_2$  is not differentiable, a formula similar to the chain rule holds:

$$D(\mathbf{h}_2^{-1} \circ \mathbf{h}_1) = D(\pi \circ \mathbf{F}_2) \circ D\mathbf{h}_1 = \pi \circ D\mathbf{F}_2 \circ D\mathbf{h}_1.$$

The above considerations are in a sense a special case of the following theorem, which roughly speaking, asserts that an  $n$ -dimensional manifold in  $\mathbb{R}^{n+k}$  looks locally, up to a diffeomorphism of the ambient space, like  $\mathbb{R}^n \times \{\mathbf{0}\} \subset \mathbb{R}^n \times \mathbb{R}^k$ .

**Theorem 3.2.2.**  *$M$  is an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+k}$  if and only if for every  $\mathbf{p} \in M$ , there exists an open neighborhood  $U$  of  $\mathbf{p}$  in  $\mathbb{R}^{n+k}$  and a diffeomorphism  $\mathbf{F} : U \rightarrow \mathbf{F}(U)$  such that*

$$\mathbf{F}(U \cap M) = \mathbf{F}(U) \cap (\mathbb{R}^n \times \{\mathbf{0}\}).$$

*Proof.* Suppose  $M$  satisfies the above condition. We proceed to construct a parametrization  $(V, \mathbf{h})$  around any  $\mathbf{p} \in M$ . Let  $\pi : \mathbb{R}^{n+k} = \mathbb{R}^n \times \mathbb{R}^k \rightarrow \mathbb{R}^n$  denote the projection,  $\iota : \mathbb{R}^n \rightarrow \mathbb{R}^n \times \mathbb{R}^k$  the map sending  $\mathbf{a}$  to  $(\mathbf{a}, \mathbf{0})$ . If  $\mathbf{F}$  is a diffeomorphism as in the statement, define

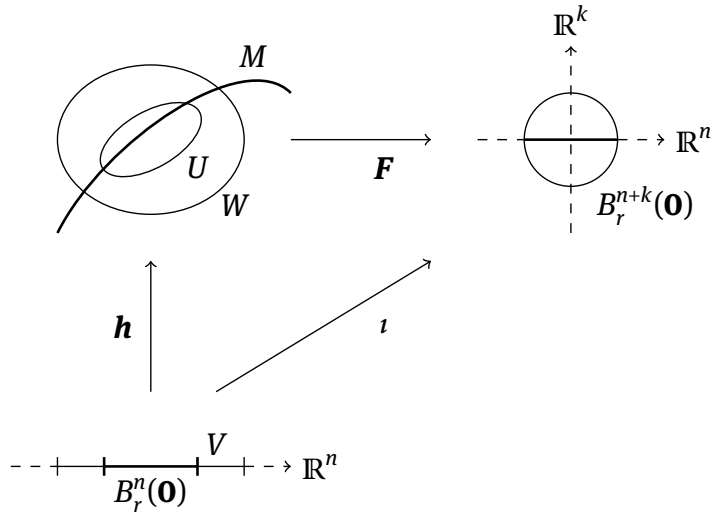
$$V = \pi(\mathbf{F}(U) \cap (\mathbb{R}^n \times \{\mathbf{0}\})), \quad \mathbf{h} = \mathbf{F}^{-1} \circ \iota.$$

$\mathbf{h}$  is one-to-one with maximal rank everywhere because  $\mathbf{F}$  and  $\iota$  have those properties. Its inverse is continuous, for if  $W$  is an open subset of  $V$ , then  $\mathbf{h}(W) = M \cap \mathbf{F}^{-1}(W \times \mathbb{R}^k)$  is open in  $M$ . Furthermore,

$$\begin{aligned} \mathbf{h}(V) &= \mathbf{F}^{-1}(\iota(V)) = \mathbf{F}^{-1}((\iota \circ \pi)(\mathbf{F}(U) \cap \mathbb{R}^n \times \{\mathbf{0}\})) \\ &= \mathbf{F}^{-1}(\mathbf{F}(U) \cap \mathbb{R}^n \times \{\mathbf{0}\}) = \mathbf{F}^{-1}(\mathbf{F}(U \cap M)) \\ &= U \cap M, \end{aligned}$$

and  $M$  is therefore an  $n$ -manifold.

Conversely, suppose  $M$  is a manifold,  $\mathbf{p} \in M$ ,  $(V, \mathbf{h})$  a parametrization around  $\mathbf{p}$ , so that  $\mathbf{h}(V)$  equals the intersection of  $M$  with some open neighborhood  $W$  of  $\mathbf{h}(\mathbf{p})$  in  $\mathbb{R}^{n+k}$ . As usual, we may assume that  $\mathbf{p} = \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{h}(\mathbf{p}) = \mathbf{0} \in \mathbb{R}^{n+k}$ . By the implicit function theorem, there exists a diffeomorphism  $\mathbf{F}$  of a neighborhood of  $\mathbf{0}$  in  $\mathbb{R}^{n+k}$  such that  $\mathbf{F} \circ \mathbf{h} = \iota$ .



For  $r > 0$ ,  $B_r^n(\mathbf{0})$  and  $B_r^{n+k}(\mathbf{0})$  will denote the open balls of radius  $r$  around  $\mathbf{0}$  in  $\mathbb{R}^n$  and  $\mathbb{R}^{n+k}$  respectively. Choose  $r$  small enough that  $B_r^n(\mathbf{0}) \subset V$ , that  $B_r^{n+k}(\mathbf{0})$  is contained in the image of  $\mathbf{F}$ , and that  $U := \mathbf{F}^{-1}(B_r^{n+k}(\mathbf{0}))$  is contained in  $W$ . Then  $B_r^n(\mathbf{0}) = \mathbf{h}^{-1}(U)$ , so that  $\mathbf{h}(B_r^n(\mathbf{0})) = U \cap M$ . We will denote the restriction  $\mathbf{F}|_U$  of  $\mathbf{F}$  to  $U$  by  $\mathbf{F}$  for simplicity. Then  $\mathbf{F} : U \rightarrow \mathbf{F}(U)$  is by definition a diffeomorphism, and

$$\begin{aligned} \mathbf{F}(U \cap M) &= \mathbf{F}(\mathbf{h}(B_r^n(\mathbf{0}))) = \iota(B_r^n(\mathbf{0})) = B_r^{n+k}(\mathbf{0}) \cap (\mathbb{R}^n \times \{\mathbf{0}\}) \\ &= \mathbf{F}(U) \cap (\mathbb{R}^n \times \{\mathbf{0}\}). \end{aligned}$$

□

Recall that Theorem 3.1.1 provides a convenient way of constructing manifolds: if  $f : U \subset \mathbb{R}^{n+k} \rightarrow \mathbb{R}^k$  has  $\mathbf{0} \in \mathbb{R}^k$  as a regular value, then the pre-image of  $\mathbf{0}$ , if nonempty, is an  $n$ -manifold. Theorem 3.2.2 shows that every manifold is at least *locally*, if not globally, obtainable in this way:

**Corollary 3.2.1.**  $M \subset \mathbb{R}^{n+k}$  is an  $n$ -dimensional manifold if and only if for every  $\mathbf{p} \in M$ , there exists a neighborhood  $U$  of  $\mathbf{p}$  in  $\mathbb{R}^{n+k}$ , and a map  $f : U \rightarrow \mathbb{R}^k$  having  $\mathbf{0}$  as a regular value, such that

$$U \cap M = f^{-1}(\mathbf{0}).$$

*Proof.* Suppose  $M$  is a manifold. Given  $\mathbf{p} \in M$ , there exists, by Theorem 3.2.2, a neighborhood  $U$  of  $\mathbf{p}$  in the ambient Euclidean space, and a diffeomorphism  $F : U \rightarrow F(U)$  such that  $U \cap M = U \cap F^{-1}(\mathbb{R}^n \times \{0\})$ . It follows that if  $\pi_2 : \mathbb{R}^n \times \mathbb{R}^k \rightarrow \mathbb{R}^k$  is projection, then  $f := \pi_2 \circ F$  has  $\mathbf{0}$  as regular value (being a composition of maps of maximal rank), and  $U \cap M = f^{-1}(\mathbf{0})$ . The converse is left as an exercise.  $\square$

A smooth map  $f : M \rightarrow N$  that has a differentiable inverse is called a *diffeomorphism*.

**Definition 3.2.2.** Let  $M^n, N^l \subset \mathbb{R}^k$  be manifolds,  $f : M \rightarrow N$  a smooth map. The *derivative* of  $f$  at  $\mathbf{p} \in M$  is the linear transformation  $f_{*\mathbf{p}} : M_{\mathbf{p}} \rightarrow N_{f(\mathbf{p})}$  defined as follows: if  $\mathbf{u} = \dot{\mathbf{c}}(0) \in M_{\mathbf{p}}$  for some curve  $\mathbf{c}$  in  $M$ , then  $f_{*\mathbf{p}}\mathbf{u} = \dot{\gamma}(0)$ , where  $\gamma = f \circ \mathbf{c}$ .

The same argument used when defining the tangent space of  $M$  shows that this does not depend on the particular choice of curve. It is furthermore straightforward to check that if  $g : N \rightarrow P$  is also differentiable, then

$$(g \circ f)_{*\mathbf{p}} = g_{*f(\mathbf{p})} \circ f_{*\mathbf{p}}.$$

In particular, if  $\mathbf{h}$  and  $\mathbf{k}$  are parametrizations of  $M$  and  $N$  mapping  $\mathbf{0}$  to  $\mathbf{p}$  and  $f(\mathbf{p})$  respectively, then the following diagram commutes:

$$\begin{array}{ccc} M_{\mathbf{p}} & \xrightarrow{f_{*\mathbf{p}}} & N_{f(\mathbf{p})} \\ \mathbf{h}_{*\mathbf{0}} \uparrow & & \uparrow \mathbf{k}_{*\mathbf{0}} \\ \mathbb{R}_{\mathbf{0}}^n & \xrightarrow{(k^{-1} \circ f \circ h)_{*\mathbf{0}}} & \mathbb{R}_{\mathbf{0}}^l \end{array}$$

Notice that using the alternative notation  $\{\mathbf{0}\} \times \mathbb{R}^k$  for the tangent space of  $\mathbb{R}^k$  at  $\mathbf{0}$ ,

$$f_{*\mathbf{p}}(\mathbf{p}, Dh(\mathbf{0})\mathbf{u}) = (f(\mathbf{p}), D(f \circ h)(\mathbf{0})\mathbf{u}), \quad \mathbf{u} \in \mathbb{R}^n. \quad (3.2.3)$$

In the special case of a function  $f : M \rightarrow \mathbb{R}$ , there is a further concept which plays an important role:

**Definition 3.2.3.** The *differential* of  $f : M \rightarrow \mathbb{R}$  at  $\mathbf{p} \in M$  is the linear map  $df(\mathbf{p}) : M_{\mathbf{p}} \rightarrow \mathbb{R}$  defined by  $df(\mathbf{p})\mathbf{u} = (f \circ \mathbf{c})'(0)$ , where  $\mathbf{c}$  is any curve in  $M$  with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ .

Thus  $df(\mathbf{p})$  is an element of the dual space of  $M_{\mathbf{p}}$ . When  $M$  is Euclidean space,

$$df(\mathbf{p})(\mathbf{u}) = \mathbf{u}(f) = \langle \nabla f(\mathbf{p}), \mathbf{u} \rangle,$$

where  $\mathbf{u}(f)$  is the derivative of  $f$  with respect to  $\mathbf{u}$  as introduced in Definition 2.8.7. In particular,

$$df(\mathbf{p}) = \sum_i D_i f(\mathbf{p}) du^i(\mathbf{p}),$$

which in classical notation reads

$$df = \sum_i \frac{\partial f}{\partial x^i} dx^i,$$

for  $n > 1$ , and

$$df = f'(x) dx$$

when  $n = 1$ . More generally, in an arbitrary manifold,

$$df(\mathbf{p})\mathbf{u} = \mathcal{I}_{f(\mathbf{p})}^{-1}(f_{*\mathbf{p}}\mathbf{u}),$$

which incidentally also shows that the differential of  $f$  is independent of the curve chosen in Definition 3.2.3.

It is useful to keep the notation for derivative of a function  $f$  in direction  $\mathbf{u}$  in the context of manifolds. Summarizing, we have:

$$df(\mathbf{p})\mathbf{u} = \mathbf{u}(f), \quad f : M \rightarrow \mathbb{R}, \quad \mathbf{p} \in M, \quad \mathbf{u} \in M_{\mathbf{p}}. \quad (3.2.4)$$

We end this section with an application of Corollary 3.2.1 called the *method of Lagrange multipliers*. Consider the following rather trivial problem: find the highest and lowest points on the sphere  $S^2 \subset \mathbb{R}^3$  of radius 1 centered at the origin. These are of course the north and south poles respectively. Now, the height function,  $u^3$ , has gradient  $\mathbf{D}_3$ , and the poles are precisely those points where this gradient is orthogonal to the tangent plane. Alternatively,  $S^2 = f^{-1}(0)$ , where  $f = (u^1)^2 + (u^2)^2 + (u^3)^2 - 1$ , and since  $\nabla f$  spans the orthogonal complement of the tangent space of  $S^2$  at every point,  $\nabla u^3$  is a multiple of  $\nabla f$  at those points where  $f$  has a maximum or minimum: indeed if  $\mathbf{c}$  is any curve in  $S^2$  that passes through, say, the highest point  $\mathbf{p}$  at time 0, then

$$\langle \nabla(u^3)(\mathbf{p}), \dot{\mathbf{c}}(0) \rangle = (u^3 \circ \mathbf{c})'(0) = 0;$$

this says that the gradient of  $u^3$  is orthogonal to the tangent space at  $\mathbf{p}$ . More generally, we have the following:

**Proposition 3.2.3.** *Let  $f : \mathbb{R}^{n+k} \supset U \rightarrow \mathbb{R}^k$  be a map that has  $\mathbf{0}$  as a regular value, so that  $M = f^{-1}(\mathbf{0})$  is an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+k}$ . If a function  $g : U \rightarrow \mathbb{R}$ , when restricted to  $M$ , has a maximum or minimum at  $\mathbf{p} \in M$ , then there exist  $\lambda_i \in \mathbb{R}$ ,  $i = 1, \dots, k$ , such that*

$$\nabla g(\mathbf{p}) = \sum_{i=1}^k \lambda_i \nabla f^i(\mathbf{p}),$$

with  $f^i := u^i \circ f$ .



The numbers  $\lambda_i$  are called *Lagrange multipliers*. When  $k = 1$ , then  $\mathbf{f} = f$  is real-valued, and the identity becomes  $\nabla g = \lambda \nabla f$  as in the example above.

*Proof.* Let  $\mathbf{u} \in M_{\mathbf{p}}$ ,  $\mathbf{c} : I \rightarrow M$  a curve with  $\dot{\mathbf{c}}(0) = \mathbf{u}$ . Then  $g \circ \mathbf{c}$  has a maximum or minimum at 0, and

$$0 = (g \circ \mathbf{c})'(0) = \langle \nabla g(\mathbf{p}), \mathbf{u} \rangle.$$

Since  $\mathbf{u}$  is arbitrary,  $\nabla g(\mathbf{p})$  is orthogonal to  $M_{\mathbf{p}}$ . Now,  $\nabla f^i(\mathbf{p})$  is also orthogonal to  $M_{\mathbf{p}}$  for each  $i$  (recall that  $f^i \circ \mathbf{c} \equiv 0$ ) and they are linearly independent because  $\mathbf{f}$  has maximal rank on  $M$ . Thus,  $\nabla f^i(\mathbf{p})$  span the orthogonal complement of  $M_{\mathbf{p}}$ , and the statement follows.  $\square$

**Example 3.2.1.** We will use Lagrange multipliers to prove that the geometric mean of  $n$  positive numbers is no larger than its arithmetic mean; i.e.,

$$\sqrt[n]{x_1 \cdots x_n} \leq \frac{x_1 + \cdots + x_n}{n}, \quad x_1, \dots, x_n > 0.$$

Consider the first quadrant  $U = \{\mathbf{p} \mid u^i(\mathbf{p}) > 0, \quad i = 1, \dots, n\}$  in  $\mathbb{R}^n$ , and for any fixed  $c > 0$ , the function  $f = \sum_i u^i - c : \mathbb{R}^n \rightarrow \mathbb{R}$ . Let  $g = (u^1 \cdots u^n)^{1/n} : U \rightarrow \mathbb{R}$ . Notice that  $g$  has no minimum on  $U$ , since  $g > 0$  on  $U$ , and  $g(\mathbf{p}) \rightarrow 0$  if  $u^i(\mathbf{p}) \rightarrow 0$ . On the other hand,  $g$  has a maximum on  $U \cap f^{-1}(0)$ : it certainly has one on  $\bar{U} \cap f^{-1}(0)$  which is a closed and bounded set, hence compact. Since  $g$  is zero on the boundary of  $U$ , this maximum lies in the interior of  $U$ . Thus, if the equation  $\nabla g = \lambda \nabla f$  is satisfied at only one point in  $U$ , that point must be the maximum. We compute

$$\nabla g = \frac{1}{n} (u^1 \cdots u^n)^{\frac{1}{n}-1} \sum_{i=1}^n (u^1 \cdots \hat{u}^i \cdots u^n) \mathbf{D}_i,$$

where the accent “ $\hat{\phantom{u}}$ ” indicates the corresponding term is deleted. Similarly,  $\nabla f = \sum_i \mathbf{D}_i$ , so  $\nabla g = \lambda \nabla f$  iff

$$\frac{1}{n} (u^1 \cdots u^n)^{\frac{1}{n}-1} u^1 \cdots \hat{u}^i \cdots u^n = \lambda, \quad i = 1, \dots, n.$$

This can only hold if  $u^i = u^j$  for all  $i$  and  $j$ . In this case, their common value is  $c/n$ , and the corresponding value of  $g$  is  $c/n = (u^1 + \cdots + u^n)/n$ . Since it occurs at a single point, it must be the maximum, and the inequality is proved.

### 3.3 Vector fields on manifolds

**Definition 3.3.1.** A *vector field* on a manifold  $M^n$  is a map  $\mathbf{X}$  which assigns to each  $\mathbf{p} \in M$  an element  $\mathbf{X}(\mathbf{p}) \in M_{\mathbf{p}}$  of the tangent space of  $M$  at  $\mathbf{p}$ , and the map is differentiable in the following sense: for any  $\mathbf{p} \in M$  and any parametrization  $(V, \mathbf{h})$  around  $\mathbf{p}$ , there exists a differentiable (in the usual sense) vector field  $\tilde{\mathbf{X}}$  on  $V$  such that

$$\mathbf{X} \circ \mathbf{h} = \mathbf{h}_* \circ \tilde{\mathbf{X}}; \quad (3.3.1)$$

i.e.,  $X(\mathbf{h}(q)) = \mathbf{h}_* \tilde{X}(q)$  for every  $q \in V$ . Just as in the case of Euclidean space, vector fields satisfying (3.3.1) for some map  $\mathbf{h}$  are said to be  $\mathbf{h}$ -related. If, as is the case for a parametrization,  $\mathbf{h}$  is invertible, we also write  $\mathbf{X} = \mathbf{h}_* \circ \tilde{\mathbf{X}} \circ \mathbf{h}^{-1}$ .

**Examples and Remarks 3.3.1.** (i) (3.3.1) can be rephrased as saying that  $\mathbf{X}$  is differentiable if for any parametrization  $(V, \mathbf{h})$ , the map  $(\mathbf{h}_*)^{-1} \circ \mathbf{X} \circ \mathbf{h}$  is a differentiable vector field on  $V \subset \mathbb{R}^n$ .

(ii) In Definition 3.3.1, it is enough to require that (3.3.1) hold for *some* parametrization  $(V, \mathbf{h})$  around every  $\mathbf{p} \in M$ . For if  $(U, \mathbf{k})$  is another parametrization around  $\mathbf{p}$ , let  $\tilde{\mathbf{X}} = \mathbf{k}_*^{-1} \circ \mathbf{X} \circ \mathbf{k}$ . Then  $\mathbf{X} \circ \mathbf{k} = \mathbf{k}_* \circ \tilde{\mathbf{X}}$ , and  $\tilde{\mathbf{X}}$  is differentiable, since locally

$$\tilde{\mathbf{X}} = (\mathbf{k}^{-1} \circ \mathbf{h})_* \circ \tilde{\mathbf{X}} \circ (\mathbf{h}^{-1} \circ \mathbf{k})$$

is a composition of differentiable maps.

(iii) The *coordinate vector fields* associated to a chart  $(U, \mathbf{x})$  of  $M^n$  are the vector fields

$$\frac{\partial}{\partial x^i} := \mathbf{x}_*^{-1} \circ \mathbf{D}_i \circ \mathbf{x}, \quad i = 1, \dots, n \quad (3.3.2)$$

on  $U$ . By (i) and (ii) they are differentiable, being  $\mathbf{h}$ -related to  $\mathbf{D}_i$ , for the parametrization  $(\mathbf{x}(U), \mathbf{h} = \mathbf{x}^{-1})$ .

(iv) As an illustration of (iii), recall spherical coordinates from Examples 2.5.1, which assign to each  $(x, y, z) \in \mathbb{R}^3$  with  $x \neq 0$  the point  $(\rho, \theta, \varphi)$ , with

$$\begin{aligned} \rho(x, y, z) &= (x^2 + y^2 + z^2)^{1/2}, \\ \theta(x, y, z) &= \arctan \frac{y}{x} + c, \\ \varphi(x, y, z) &= \arccos \frac{z}{(x^2 + y^2 + z^2)^{1/2}}, \end{aligned}$$

where  $c$  depends on the quadrant where  $(x, y) \in \mathbb{R}^2$  lies, cf. Section 4.6.1. (As observed earlier, they are actually defined on  $\mathbb{R}^3 \setminus \{\mathbf{0}\}$ ; we exclude the plane  $x = 0$  only to have a unified formula). The map  $\mathbf{G}$  with  $G^1 = \rho$ ,  $G^2 = \theta$ , and  $G^3 = \varphi$  is a local diffeomorphism of  $\mathbb{R}^3$  with inverse

$$\mathbf{G}^{-1}(a_1, a_2, a_3) = (a_1 \cos a_2 \sin a_3, a_1 \sin a_2 \sin a_3, a_1 \cos a_3).$$

Since for any  $a > 0$ , the sphere  $S^2(a)$  of radius  $a$  centered at the origin equals  $(G^1)^{-1}(a)$ , the proof of Corollary 3.2.1 shows that the map  $\mathbf{x} := (\theta, \varphi)$  is a chart on any such sphere. In order to describe the coordinate vector fields  $\partial/\partial\varphi$ ,  $\partial/\partial\theta$ , notice that the inverse of the chart is the parametrization  $\mathbf{h}$ , where

$$\mathbf{h} = (a \cos u^1 \sin u^2, a \sin u^1 \sin u^2, a \cos u^2);$$

i.e.,  $\mathbf{h}(s, t) = \mathbf{G}^{-1}(a, s, t)$ . Thus, by (3.3.2),

$$\begin{aligned} \frac{\partial}{\partial \theta} \circ \mathbf{h} &= \mathbf{h}_* \mathbf{D}_1 = -a \sin u^1 \sin u^2 (\mathbf{D}_1 \circ \mathbf{h}) + a \cos u^1 \sin u^2 (\mathbf{D}_2 \circ \mathbf{h}) \\ &= -h^2 (\mathbf{D}_1 \circ \mathbf{h}) + h^1 (\mathbf{D}_2 \circ \mathbf{h}). \end{aligned}$$

Since  $h^i \circ \mathbf{h}^{-1} = u^i \circ \mathbf{h} \circ \mathbf{h}^{-1} = u^i$ , composing the above identity with  $\mathbf{h}^{-1}$  yields

$$\frac{\partial}{\partial \theta} = -u^2 \mathbf{D}_1 + u^1 \mathbf{D}_2. \quad (3.3.3)$$

The computations for  $\partial/\partial\varphi$  are similar, only slightly more involved:

$$\begin{aligned} \frac{\partial}{\partial \varphi} \circ \mathbf{h} &= \mathbf{h}_* \mathbf{D}_2 = a \cos u^1 \cos u^2 (\mathbf{D}_1 \circ \mathbf{h}) + a \sin u^1 \cos u^2 (\mathbf{D}_2 \circ \mathbf{h}) \\ &\quad - a \sin u^2 (\mathbf{D}_3 \circ \mathbf{h}). \end{aligned} \quad (3.3.4)$$

Since  $u^2$  takes values in  $(0, \pi)$ , its sine is nonnegative, so that

$$\sin u^2 = |\sin u^2| = ((\cos u^1 \sin u^2)^2 + (\sin u^1 \sin u^2)^2)^{1/2},$$

and the first function on the right side of (3.3.4) may be written

$$\begin{aligned} a \cos u^1 \cos u^2 &= \frac{(a \cos u^1 \sin u^2)(a \cos u^2)}{((a \cos u^1 \sin u^2)^2 + (a \sin u^1 \sin u^2)^2)^{1/2}} \\ &= \frac{h^1 h^3}{((h^1)^2 + (h^2)^2)^{1/2}}. \end{aligned}$$

Similarly,

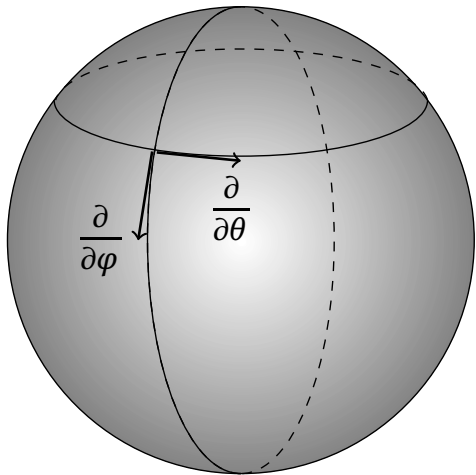
$$a \sin u^1 \cos u^2 = \frac{h^2 h^3}{((h^1)^2 + (h^2)^2)^{1/2}}, \quad -a \sin u^2 = -((h^1)^2 + (h^2)^2)^{1/2},$$

and therefore

$$\begin{aligned} \frac{\partial}{\partial \varphi} \circ \mathbf{h} &= \frac{1}{((h^1)^2 + (h^2)^2)^{1/2}} (h^1 h^3 (\mathbf{D}_1 \circ \mathbf{h}) + h^2 h^3 (\mathbf{D}_2 \circ \mathbf{h}) \\ &\quad - ((h^1)^2 + (h^2)^2) (\mathbf{D}_3 \circ \mathbf{h})), \end{aligned}$$

or equivalently,

$$\frac{\partial}{\partial \varphi} = \frac{1}{((u^1)^2 + (u^2)^2)^{1/2}} (u^1 u^3 \mathbf{D}_1 + u^2 u^3 \mathbf{D}_2 - ((u^1)^2 + (u^2)^2) \mathbf{D}_3). \quad (3.3.5)$$



Coordinate vector fields on a sphere

Notice that  $\partial/\partial\theta$  and  $\partial/\partial\varphi$  are mutually orthogonal everywhere, and are orthonormal along the equator.

- (v) A vector field  $\mathbf{X}$  on  $M$  and a function  $f : M \rightarrow \mathbb{R}$  can be combined to yield a new vector field  $f\mathbf{X}$ , by setting  $(f\mathbf{X})(\mathbf{p}) := f(\mathbf{p})\mathbf{X}(\mathbf{p})$ ,  $\mathbf{p} \in M$ . They can also be combined to form a new function  $\mathbf{X}f$ , where  $(\mathbf{X}f)(\mathbf{p}) := \mathbf{X}(\mathbf{p})(f)$ .

The following proposition generalizes the fact that coordinate vector fields are differentiable.

**Proposition 3.3.1.** *A map  $\mathbf{p} \mapsto \mathbf{X}(\mathbf{p}) \in M_{\mathbf{p}}$  is a vector field (in the sense that  $\mathbf{X}$  is differentiable) if and only if for any  $\mathbf{p} \in M$  there exists a chart  $(U, \mathbf{x})$  of  $M$  around  $\mathbf{p}$  such that*

$$\mathbf{X}|_U = \sum_i f^i \frac{\partial}{\partial x^i}$$

for some differentiable functions  $f^i$  on  $U$ .

*Proof.* The restriction to  $U$  of the map  $\mathbf{X}$  is smooth iff in terms of the above chart, the map  $\mathbf{x}_* \circ \mathbf{X} \circ \mathbf{x}^{-1}$  is differentiable on  $V = \mathbf{x}(U)$ ; i.e., iff it can be written as  $\sum_i f^i \mathbf{D}_i$  for smooth functions  $f^i$  on  $V$ . This is equivalent to requiring that

$$\mathbf{X}|_U \circ \mathbf{x}^{-1} = \mathbf{x}_*^{-1} \left( \sum_i f^i \mathbf{D}_i \right) = \sum_i (f^i \circ \mathbf{x}^{-1}) \mathbf{x}_*^{-1} \mathbf{D}_i;$$

i.e.,

$$\mathbf{X}|_U = \sum_i f^i \frac{\partial}{\partial x^i},$$

by (3.3.2). □

There is another, sometimes more useful way of determining differentiability:

**Theorem 3.3.1.** *Let  $M$  denote an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+k}$ . A map  $\mathbf{X}$  that assigns to each  $\mathbf{p} \in M$  a vector  $\mathbf{X}(\mathbf{p}) \in M_{\mathbf{p}}$  is a vector field (once again, in the sense that the map is differentiable) if and only if for any  $\mathbf{p} \in M$  there exists a vector field  $\mathbf{Y}$  on an open neighborhood  $U$  of  $\mathbf{p}$  in  $\mathbb{R}^{n+k}$  such that the restrictions of both vector fields agree on the induced neighborhood of  $\mathbf{p}$  in  $M$ : i.e.,  $\mathbf{Y}|_{M \cap U} = \mathbf{X}|_{M \cap U}$ .*

*Proof.* Suppose  $\mathbf{X}$  is a vector field on  $M$ , and consider a parametrization  $(V, \mathbf{h})$  around some  $\mathbf{p} \in M$ . As usual, we may assume to simplify matters that  $\mathbf{p} = \mathbf{0} = \mathbf{h}(\mathbf{0})$ . Thus, the map  $\tilde{\mathbf{X}} = (\mathbf{h}^{-1})_* \circ \mathbf{X} \circ \mathbf{h}$  is differentiable on  $V$ . Write  $\tilde{\mathbf{X}} = \sum_i f^i \mathbf{D}_i$ , where  $f^i$  are differentiable functions on  $V$ , and extend  $\tilde{\mathbf{X}}$  to a vector field  $\tilde{\mathbf{Y}}$  on  $V \times \mathbb{R}^k$  by setting

$$\tilde{\mathbf{Y}} = \sum_i (f^i \circ \pi) \mathbf{D}_i,$$

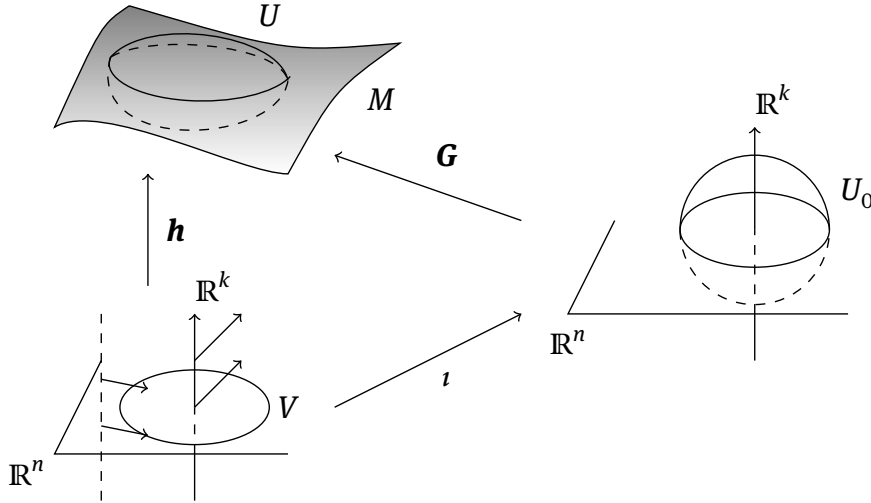
with  $\pi : V \times \mathbb{R}^k \rightarrow V$  denoting projection. Observe that if  $\iota : V \rightarrow V \times \mathbb{R}^k$  maps  $\mathbf{a} \in V$  to  $(\mathbf{a}, \mathbf{0})$ , then

$$\tilde{\mathbf{Y}} \circ \iota = \iota_* \circ \tilde{\mathbf{X}}|_V.$$

By Theorem 3.2.2, there exists a diffeomorphism  $\mathbf{G}$  on a neighborhood  $U_0$  of  $\mathbf{0}$  in  $\mathbb{R}^{n+k}$  such that

$$\begin{aligned} \mathbf{h} &= \mathbf{G} \circ \iota, \\ U_0 \cap (\mathbb{R}^n \times \{\mathbf{0}\}) &= \iota(V), \text{ and} \\ \mathbf{G}^{-1}(U \cap M) &= \mathbf{G}^{-1}(U) \cap (\mathbb{R}^n \times \{\mathbf{0}\}), \end{aligned}$$

with  $U$  denoting  $\mathbf{G}(U_0)$ .



Define a vector field  $\mathbf{Y}$  on  $U$  by setting  $\mathbf{Y} = \mathbf{G}_* \circ \tilde{\mathbf{Y}} \circ \mathbf{G}^{-1}$ . We claim  $\mathbf{Y}$  satisfies the conclusion of the theorem. Indeed,

$$\mathbf{Y} \circ \mathbf{h} = \mathbf{G}_* \circ \tilde{\mathbf{Y}} \circ \mathbf{G}^{-1} \circ \mathbf{h} = \mathbf{G}_* \circ \tilde{\mathbf{Y}} \circ \iota = \mathbf{G}_* \circ \iota_* \circ \tilde{\mathbf{X}} = \mathbf{h}_* \circ \mathbf{X} = \mathbf{X} \circ \mathbf{h}.$$

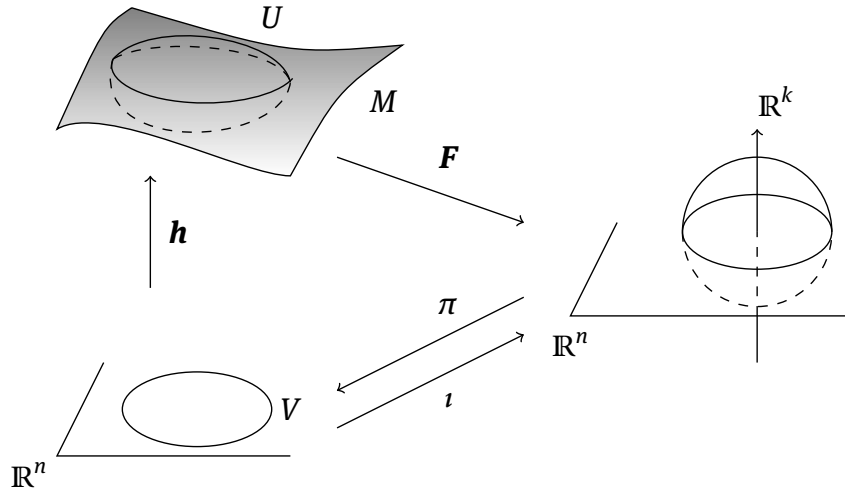
Thus, the restrictions of  $\mathbf{X}$  and  $\mathbf{Y}$  to  $\mathbf{h}(V)$  agree. Since  $\mathbf{h}(V) = (\mathbf{G} \circ \iota)(V) = U \cap M$ , the claim follows.

For the converse, let  $\mathbf{p} \in M$ , which as usual may be assumed to be  $\mathbf{0}$ . By assumption, there exists an open neighborhood  $U$  of  $\mathbf{p}$  in  $\mathbb{R}^{n+k}$ , and a vector field  $\mathbf{Y}$  on  $U$  such that the restrictions of  $\mathbf{X}$  and  $\mathbf{Y}$  to  $M \cap U$  agree. It must be shown that there exists a parametrization  $(V, \mathbf{h})$  of  $M$  with  $\mathbf{p} \in \mathbf{h}(V)$  such that  $\mathbf{h}_*^{-1} \circ \mathbf{X} \circ \mathbf{h}$  is differentiable.

Now, by Theorem 3.2.2, after restricting  $U$  if necessary, there exists a diffeomorphism  $\mathbf{F} : U \rightarrow \mathbf{F}(U)$  with the property that

$$\mathbf{F}(U \cap M) = \mathbf{F}(U) \cap (\mathbb{R}^n \times \{\mathbf{0}\}).$$

Furthermore, if  $\pi : \mathbb{R}^n \times \mathbb{R}^k \rightarrow \mathbb{R}^n$  is projection, if  $\iota : \mathbb{R}^n \rightarrow \mathbb{R}^n \times \mathbb{R}^k$  maps  $\mathbf{a}$  to  $(\mathbf{a}, \mathbf{0})$  and  $V = (\pi \circ \mathbf{F})(U \cap M)$ , then  $(V, \mathbf{h})$ , where  $\mathbf{h} = \mathbf{F}^{-1} \circ \iota$ , is a parametrization of  $M$  around  $\mathbf{p}$ .



Define  $\tilde{X} = F_* \circ Y \circ F^{-1}$ .  $\tilde{X}$  is differentiable and its restriction to  $\mathbb{R}^n \times \{\mathbf{0}\}$  maps by construction each  $\mathbf{p}$  to a vector tangent to  $\mathbb{R}^n$ . Thus, there exists a vector field  $\tilde{X}$  on  $V$  such that  $X \circ \iota = \iota_* \circ \tilde{X}$ : indeed, if  $\tilde{X} = \sum_i g^i \mathbf{D}_i$ , then each  $g_i$  is smooth, and therefore so is the map  $\tilde{X} = \sum (g^i \circ \iota) \mathbf{D}_i$ . This in turn shows that  $X$  is differentiable, since

$$h_* \circ \tilde{X} = (F^{-1} \circ \iota)_* \circ \tilde{X} = F_*^{-1} \circ \tilde{X} \circ \iota = Y \circ F^{-1} \circ \iota = Y \circ h = X \circ h. \quad \square$$

Either the definition of differentiable vector field or the characterization given in Theorem 3.3.1 can now be used to extend the concept of flow to vector fields on manifolds. As in Euclidean space, an *integral curve* of a vector field  $X$  on a manifold  $M$  is a curve  $\mathbf{c} : I \rightarrow M$  such that  $\dot{\mathbf{c}} = X \circ \mathbf{c}$ . Theorem 3.3.1 immediately implies the existence and uniqueness of integral curves, and thus of local flows. If one prefers to use the definition, it suffices to observe that locally, a vector field on a manifold is  $h$ -related to one in Euclidean space for some parametrization  $h$ , and appeal to the following:

**Proposition 3.3.2.** *Let  $X, \tilde{X}$  be  $h$ -related vector fields,  $h_* \circ \tilde{X} = X \circ h$ .*

- (1) *If  $\mathbf{c}$  is an integral curve of  $\tilde{X}$ , then  $h \circ \mathbf{c}$  is an integral curve of  $X$ .*
- (2) *If  $\tilde{\Psi}_t, \Psi_s$  are flows of  $\tilde{X}, X$  respectively, then  $h \circ \tilde{\Psi}_t = \Psi_t \circ h$ .*

*Proof.* If  $\mathbf{c}$  is an integral curve of  $\tilde{X}$ , then

$$h \circ \dot{\mathbf{c}} = h_* \dot{\mathbf{c}} = h_* \circ \tilde{X} \circ \mathbf{c} = (X \circ h) \circ \mathbf{c} = X \circ (h \circ \mathbf{c}),$$

which is the first claim. The second claim is then clear: for any  $\mathbf{p} \in M$ , the curve  $t \mapsto h \circ \tilde{\Psi}_t(\mathbf{p})$  is by (1) an integral curve of  $X$  passing through  $h(\mathbf{p})$  at  $t = 0$ , and so is  $t \mapsto \Psi_t \circ h(\mathbf{p})$ . They are therefore one and the same.  $\square$

The work needed to group local flows into a single global one was already done in the proof of Theorem 2.8.3. Even though it was stated in the setting of Euclidean space, the argument makes no use of that setting, and works equally well in the more general one of manifolds. We restate it for future reference. For  $\mathbf{p} \in M$ , let  $I_{\mathbf{p}}$  denote the largest open interval on which the integral curve  $\Psi_{\mathbf{p}}$  passing through  $\mathbf{p}$  at time 0 is defined.

**Theorem 3.3.2.** *If  $X$  is a vector field on a manifold  $M$ , there exists a unique open set  $W \subset \mathbb{R} \times M$  and a unique differentiable map  $\Psi : W \rightarrow M$  such that*

- (1)  $I_{\mathbf{p}} \times \{\mathbf{p}\} = W \cap (\mathbb{R} \times \{\mathbf{p}\})$  for all  $\mathbf{p} \in M$ , and
- (2)  $\Psi(t, \mathbf{p}) = \Psi_{\mathbf{p}}(t)$  for  $(t, \mathbf{p}) \in W$ .

The map  $\Psi$  in the above theorem is called *the flow* of the vector field  $X$ , as opposed to *a flow*, or *a local flow*, which refers to a restriction of  $\Psi$  to a subset of its domain. As in Euclidean space,  $X$  is said to be *complete* if its flow has domain  $\mathbb{R} \times M$ , or equivalently, if its integral curves are defined on all of  $\mathbb{R}$ . It turns out that any vector field on a compact manifold is complete. In order to show this, we will need the following extension theorem:

**Theorem 3.3.3.** *Let  $c : [a, b) \rightarrow M$  be an integral curve of a vector field  $X$  on  $M$ . If there exists a sequence  $\{t_k\}$  in  $[a, b)$  that converges to  $b$  for which  $\{c(t_k)\}$  converges to some  $\mathbf{p} \in M$ , then  $c$  may be extended to a continuous curve on  $[a, b]$ ; i.e., the curve  $\bar{c} : [a, b] \rightarrow M$ , with  $\bar{c}(t) = c(t)$  if  $t < b$  and  $\bar{c}(b) = \mathbf{p}$ , is continuous. Furthermore, if  $\gamma : I \rightarrow M$  is the maximal integral curve of  $X$  with  $\gamma(b) = \mathbf{p}$ , then  $I$  contains  $[a, b]$  and the restriction of  $\gamma$  to  $[a, b]$  equals  $\bar{c}$ .*

*Proof.* By using a chart around  $\mathbf{p}$  as in Theorem 3.3.1, we may translate the situation to Euclidean space. So assume  $M = \mathbb{R}^n$  and  $\mathbf{p} = \mathbf{0}$ . Suppose the curve  $\bar{c}$  from the statement is not continuous. Then there exists a sequence  $r_k \rightarrow b$  such that  $\{c(r_k)\}$  does not converge to  $\mathbf{0}$ ; i.e., there exists  $\delta > 0$  such that for any  $N \in \mathbb{N}$ ,  $|c(r_k)| > \delta$  for some  $k > N$ . Let  $\varepsilon \in (0, \delta)$ . By passing to subsequences if necessary, we may assume  $|c(t_k)| < \varepsilon$  and  $|c(r_{\tilde{k}_i})| > \delta$  for all  $k$ . The contradiction arises, roughly speaking, from the fact that  $|t_k - r_{\tilde{k}_i}| \rightarrow 0$  but  $|c(t_k) - c(r_{\tilde{k}_i})|$  does not, so that the speed  $|\dot{c}|$  of  $c$  and hence the norm of  $X$  is unbounded on a compact neighborhood of the origin. More precisely, we construct subsequences as follows: set  $k_1 = 1$ , and choose some  $\tilde{k}_1 \in \mathbb{N}$  such that  $r_{\tilde{k}_1} > t_{k_1}$  (such a  $\tilde{k}_1$  exists since otherwise  $\{r_k\}$  would not converge to  $b$ ). For the same reason, there exists some  $k_2 > k_1$  with  $t_{k_2} > r_{\tilde{k}_1}$ . Continue in this fashion to obtain subsequences  $\{t_{k_i}\}$  and  $\{r_{\tilde{k}_i}\}$  that converge to  $b$  and satisfy  $t_{k_i} < r_{\tilde{k}_i} < t_{k_{i+1}}$  for all  $i$ . Now,  $|c(t_{k_i})| < \varepsilon < \delta$ , but  $|c(r_{\tilde{k}_i})| > \delta$ , so that by continuity  $|c(t)|$  must equal  $\delta$  for some  $t \in (t_{k_i}, r_{\tilde{k}_i})$ . If  $s_i$  denotes the supremum of all such  $t$ , then  $|c(s_i)| = \delta$  by continuity again.

The length of the restriction of  $c$  to  $[t_{k_i}, s_i]$  satisfies

$$\int_{t_{k_i}}^{s_i} |\dot{c}| \geq |c(s_i) - c(t_{k_i})| \geq |c(s_i)| - |c(t_{k_i})| \geq \delta - \varepsilon$$

by the triangle inequality. Thus, for any natural number  $m$ ,

$$\sum_{i=1}^m \int_{t_{k_i}}^{s_i} |\dot{c}| \geq m(\delta - \varepsilon). \quad (3.3.6)$$

Since  $\mathbf{c}[t_{k_i}, s_i]$  is contained in the closure  $\overline{B_\delta(\mathbf{0})}$  of the ball of radius  $\delta$  about the origin, this means that the portion of  $\mathbf{c}$  that lies inside it has unbounded length. But this is impossible, because the norm of  $\mathbf{X}$  is bounded on the compact  $\overline{B_\delta(\mathbf{0})}$ : specifically, if  $|\mathbf{X}| \leq M$  on  $\overline{B_\delta(\mathbf{0})}$ , then  $|\dot{\mathbf{c}}(t)| = |\mathbf{X}(\mathbf{c}(t))| \leq M$  for  $t \in [t_{k_i}, s_i]$ , and

$$\sum_{i=1}^m \int_{t_{k_i}}^{s_i} |\dot{\mathbf{c}}| \leq M \sum_i (s_i - t_{k_i}) \leq M(b - a),$$

which contradicts (3.3.6). This shows that  $\mathbf{c}$  is extendable to a continuous curve  $\bar{\mathbf{c}} : [a, b] \rightarrow M$  by setting  $\bar{\mathbf{c}}(b) = \mathbf{0}$ . Furthermore,

$$\lim_{t \rightarrow b^-} \mathbf{X}(\bar{\mathbf{c}}(t)) = \lim_{t \rightarrow b^-} \dot{\bar{\mathbf{c}}}(t) = \mathbf{X}(\mathbf{0}) \tag{3.3.7}$$

by continuity of  $\mathbf{X}$  and  $\bar{\mathbf{c}}$ .

Consider the curve  $\gamma_1 : I \rightarrow M$  that equals  $\bar{\mathbf{c}}$  on  $[a, b]$  and  $\gamma$  on  $I \setminus [a, b]$ . This curve is continuous and differentiable (with  $\dot{\gamma}_1 = \mathbf{X} \circ \gamma_1$ ) everywhere except perhaps at  $b$ . To show smoothness at  $b$  (and thereby conclude the proof of the theorem), let  $\{t_k\}$  be any sequence in  $[a, b)$  that converges to  $b$ . By the mean value theorem, there exists  $s_k \in (t_k, b)$  such that

$$\frac{(u^i \circ \gamma_1)(b) - (u^i \circ \gamma_1)(t_k)}{b - t_k} = (u^i \circ \gamma_1)'(s_k).$$

But  $\dot{\gamma}_1(s_k) = \dot{\bar{\mathbf{c}}}(s_k) \rightarrow \mathbf{X}(\mathbf{0})$  by (3.3.7). Similarly, if  $b < s_k < t_k$  and  $t_k \rightarrow b$ , then  $\dot{\gamma}_1(s_k) = \dot{\gamma}(s_k) \rightarrow \mathbf{X}(\mathbf{0})$ . This shows that  $(u^i \circ \gamma_1)'(b)$  exists. The same argument then applies to derivatives of higher order to show that  $\gamma_1$  is smooth at  $b$ .  $\square$

**Corollary 3.3.1.** *Any vector field on a compact manifold is complete; i.e., its integral curves are defined on all of  $\mathbb{R}$ .*

*Proof.* Let  $\mathbf{X}$  denote a vector field on a compact manifold  $M$ , and  $\mathbf{c}$  the maximal integral curve of  $\mathbf{X}$  with  $\mathbf{c}(0)$  equal to some  $\mathbf{p} \in M$ . It is enough to show that  $[0, \infty)$  is contained in the domain of  $\mathbf{c}$ ; the same argument applied to the integral curve  $t \mapsto \mathbf{c}(-t)$  of  $-\mathbf{X}$  then shows that  $(-\infty, 0]$  is also contained in the domain. With this in mind, let  $I$  denote the set of all  $t \in [0, \infty)$  such that  $\mathbf{c}$  is defined on  $[0, t)$ . Then  $I$  is by definition a nonempty open interval in  $[0, \infty)$ . Theorem 3.3.3 implies  $I$  is also closed. By Proposition 1.7.1,  $I = [0, \infty)$ .  $\square$

The Lie bracket of vector fields on a manifold  $M$  is defined in exactly the same way it was done for Euclidean space, and all relevant theorems proved in Section 2.9 carry over to the manifold setting. In fact, the Lie bracket of  $\mathbf{X}$  and  $\mathbf{Y}$  at a point  $\mathbf{p} \in M$  depends only on the values of these fields in a neighborhood of  $\mathbf{p}$ , and if we take this neighborhood to be the domain of a chart  $(U, \mathbf{x})$ , then the vector fields  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{Y}}$  on  $\mathbf{x}(U)$ , where  $\tilde{\mathbf{X}} = \mathbf{x}_* \circ \mathbf{X}|_U \circ \mathbf{x}^{-1}$  (and similarly for  $\tilde{\mathbf{Y}}$ ) are  $\mathbf{h}$ -related to  $\mathbf{X}$  and  $\mathbf{Y}$ , if  $\mathbf{h} = \mathbf{x}^{-1}$ .



By Theorem 2.9.2, the restriction of the bracket to  $U$  is given by

$$[\mathbf{X}, \mathbf{Y}]|_U = \mathbf{h}_* \circ [\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}] \circ \mathbf{x}.$$

For example, the Lie bracket of coordinate vector fields is always zero, since these fields are related to the standard coordinate vector fields  $\mathbf{D}_i$  in Euclidean space.

### 3.4 Lie groups

The material in this section will not be used in the sequel, and the reader may skip it without loss of continuity. It has been included because it describes a large class of examples of manifolds, which, in addition, possess a rich additional structure. It also illustrates many of the techniques we developed for vector fields.

In modern algebra, a *group* is a pair  $(G, \cdot)$ , where  $G$  is a nonempty set, and  $\cdot : G \times G \rightarrow G$  is a map, called the group product, satisfying the following properties:

- (1)  $a \cdot (b \cdot c) = (a \cdot b) \cdot c$  for all  $a, b, c \in G$ ;
- (2) There exists an element  $e \in G$  such that  $a \cdot e = e \cdot a = a$  for all  $a \in G$ ;
- (3) For any  $a \in G$  there exists an element  $a^{-1} \in G$  such that  $a \cdot a^{-1} = a^{-1} \cdot a = e$ .

$e$  is called the identity element,  $a^{-1}$  the inverse of  $a$ . We will often write  $ab$  instead of  $a \cdot b$ . Examples of groups are plentiful: the set of nonzero reals with the usual multiplication is one. Any vector space is a group with vector addition. So is the set  $GL(n)$  of all  $n \times n$  invertible matrices with the usual matrix multiplication. In the latter case, the identity element is the identity matrix  $I_n$ .

**Definition 3.4.1.** Let  $G$  be a group that admits in addition a manifold structure. If the group product  $G \times G \rightarrow G$  and the inverse map  $G \rightarrow G$  (which sends  $a$  to  $a^{-1}$ ) are differentiable, then  $G$  is called a *Lie group*.

The three examples given earlier are Lie groups: vector addition in a vector space  $E$  is differentiable as a bilinear map on  $E \times E$ , and in the other cases, the group product is differentiable because it is the restriction of a bilinear map to an open subset of a vector space. The inverse maps are also easily seen to be differentiable; in the case of  $GL(n)$  for example, this follows from the alternative formula for the inverse of a matrix given in the proof of Theorem 1.3.6. Another important example is the following subgroup of  $GL(n)$  (a subgroup of a group  $G$  is a subset that is a group in its own right with the restriction of the group product in  $G$ ):

**Proposition 3.4.1.** *The orthogonal group  $O(n) = \{A \in M_n \mid AA^T = I_n\}$  is a Lie group of dimension  $n(n-1)/2$ .*

*Proof.* Denote by  $S_n$  the subspace of  $M_n$  that consists of all symmetric matrices.  $S_n$  has as basis the set  $\{A_{ij} \mid 1 \leq i \leq j \leq n\}$ , where  $A_{ij}$  has entries 1 in the  $(i, j)$  and  $(j, i)$  slots, and zero elsewhere. Its dimension equals the number of elements in the set of all pairs  $(i, j)$  where  $1 \leq i \leq j \leq n$ . This is just  $1 + 2 + \dots + n = n(n+1)/2$ .

Next, consider the map  $\mathbf{f} : GL(n) \rightarrow S_n$  given by  $\mathbf{f}(A) = AA^T$ . We claim it has maximal rank at any  $A \in \mathbf{f}^{-1}(I_n)$ , so that by Theorem 3.1.1,  $O(n) = \mathbf{f}^{-1}(I_n)$  is a manifold of dimension  $n^2 - n(n+1)/2 = n(n-1)/2$ . Now,  $\mathbf{f}$  is the restriction to  $GL(n)$  of the composition  $m \circ \iota$ , where  $\iota : M_n \rightarrow M_n \times M_n$  is given by  $\iota(A) = (A, A^T)$ , and  $m$  denotes the multiplication map  $m : M_n \times M_n \rightarrow M_n$ ,  $m(A, B) = AB$ . But  $\iota$  is linear, and  $m$  is bilinear, so the results from Chapter 2 imply that for  $A \in GL(n)$ ,  $M \in M_n$ ,

$$\begin{aligned} D\mathbf{f}(A)M &= Dm(A, A^T) \circ D\iota(A)M = Dm(A, A^T) \circ \iota(M) \\ &= Dm(A, A^T)(M, M^T) = m(A, M^T) + m(M, A^T) \\ &= AM^T + MA^T. \end{aligned}$$

This means that for  $A \in O(n)$ ,  $D\mathbf{f}(A)$  is onto  $S_n$ : indeed, given a symmetric matrix  $S$ , if  $M = (1/2)SA$ , then

$$D\mathbf{f}(A)M = \frac{1}{2}A(SA)^T + \frac{1}{2}SAA^T = \frac{1}{2}(S^T + S) = S.$$

Thus,  $O(n)$  is a manifold. It is a Lie group because the product and inverse map are the composition of the corresponding smooth operations on  $GL(n)$  with the differentiable inclusion map of  $O(n)$  into the open subset  $GL(n)$  of  $M_n$ .  $\square$

Given an element  $a$  in a Lie group  $G$ , define *left translation* by  $a$  to be the map

$$\begin{aligned} L_a : G &\rightarrow G, \\ b &\mapsto ab. \end{aligned}$$

It is differentiable, being the composition of the group product with the map  $\iota_a : G \rightarrow G \times G$ ,  $\iota_a(b) = (a, b)$ , and is in fact a diffeomorphism of  $G$  with inverse  $L_{a^{-1}}$ . Notice also that  $L_a \circ L_b = L_{ab}$  for  $a, b \in G$ .

**Definition 3.4.2.** A vector field  $\mathbf{X}$  on a Lie group  $G$  is said to be *left-invariant* if it is  $L_g$ -related to itself for any  $g \in G$ ; i.e.,

$$L_{g*}\mathbf{X} = \mathbf{X} \circ L_g, \quad g \in G.$$

Recall that the collection of all vector fields on  $G$  is a (infinite-dimensional) Lie algebra with the Lie bracket of vector fields, cf. Proposition 2.9.2. It follows from the definition that the bracket of two left-invariant vector fields is again left-invariant, so that the collection  $\mathfrak{g}$  of all such vector fields is also a Lie algebra. This one, though, is finite-dimensional. In fact, we have the following:

**Theorem 3.4.1.** *The Lie algebra  $\mathfrak{g}$  of a Lie group  $G$  is a vector space naturally isomorphic to  $G_e$ .*

*Proof.* Define  $\mathbf{h} : \mathfrak{g} \rightarrow G_e$  by  $\mathbf{h}(\mathbf{X}) = \mathbf{X}(e)$  for  $\mathbf{X} \in \mathfrak{g}$ . This map is clearly linear. If  $\mathbf{h}(\mathbf{X}) = \mathbf{X}(e) = \mathbf{0}$ , then for any  $g \in G$ ,  $\mathbf{X}(g) = L_{g*}\mathbf{X}(e) = \mathbf{0}$ , so  $\mathbf{h}$  has trivial kernel. To

see that it is onto, consider an arbitrary  $\mathbf{u} \in G_e$ . Define a vector field  $\mathbf{X}$  on  $G$  by setting  $\mathbf{X}(g) = L_{g*}\mathbf{u}$ .  $\mathbf{X}$  is left-invariant, since for any  $a, b \in G$ ,

$$L_{a*}\mathbf{X}(b) = L_{a*}L_{b*}\mathbf{v} = L_{(ab)*}\mathbf{u} = \mathbf{X}(ab) = \mathbf{X} \circ L_a(b).$$

Furthermore,  $\mathbf{hX} = \mathbf{u}$  by construction. This completes the argument.  $\square$

We next investigate integral curves of left-invariant vector fields. Given groups  $G, H$ , a *group homomorphism* from  $G$  to  $H$  is a map  $\mathbf{f} : G \rightarrow H$  which preserves group products; i.e.,  $\mathbf{f}(ab) = \mathbf{f}(a)\mathbf{f}(b)$  for all  $a, b \in G$ . Any homomorphism maps the identity element  $e \in G$  to the identity  $e \in H$ , because  $\mathbf{f}(e) = \mathbf{f}(e \cdot e) = \mathbf{f}(e) \cdot \mathbf{f}(e)$ ; multiplying both sides by  $\mathbf{f}(e)^{-1}$  then yields  $\mathbf{f}(e) = e$ . This in turn implies that  $\mathbf{f}(a^{-1}) = \mathbf{f}(a)^{-1}$  for any  $a \in G$ .

In the following theorem,  $\mathbb{R}$  is the Lie group with the usual addition.

**Theorem 3.4.2.** *Any smooth Lie group homomorphism  $\mathbf{c} : \mathbb{R} \rightarrow G$  is the integral curve passing through  $e$  at 0 of the left-invariant vector field  $\mathbf{X}$  with  $\mathbf{X}(e) = \dot{\mathbf{c}}(0)$ . The integral curve of  $\mathbf{X}$  that passes through  $g \in G$  at 0 is  $L_g \circ \mathbf{c}$ .*

*Conversely, if  $\mathbf{c}$  is an integral curve of  $\mathbf{X} \in \mathfrak{g}$  with  $\mathbf{c}(0) = e$ , then  $\mathbf{c} : \mathbb{R} \rightarrow G$  is a Lie group homomorphism. In particular, left-invariant vector fields are complete.*

*Proof.* Let  $\mathbf{c} : \mathbb{R} \rightarrow G$  be a homomorphism, and fix any  $t_0 \in \mathbb{R}$ . The curve  $\gamma : \mathbb{R} \rightarrow G$  given by  $\gamma(t) = \mathbf{c}(t + t_0)$  satisfies  $\gamma(t) = \mathbf{c}(t_0)\mathbf{c}(t) = L_{\mathbf{c}(t_0)}\mathbf{c}(t)$ . Thus if  $\mathbf{X} \in \mathfrak{g}$  is the left-invariant vector field that equals  $\dot{\mathbf{c}}(0)$  at  $e$ , then

$$\dot{\mathbf{c}}(t_0) = \dot{\gamma}(0) = L_{\mathbf{c}(t_0)*}\dot{\mathbf{c}}(0) = L_{\mathbf{c}(t_0)*}\mathbf{X}(e) = \mathbf{X}(\mathbf{c}(t_0)).$$

This shows that  $\mathbf{c}$  is indeed the integral curve of  $\mathbf{X}$  that passes through  $e$  at 0. Furthermore, if  $g \in G$  and  $\mathbf{c}_1 = L_g \circ \mathbf{c}$ , then  $\mathbf{c}_1(0) = g$  and

$$\dot{\mathbf{c}}_1(t) = L_{g*}\dot{\mathbf{c}}(t) = L_{g*}\mathbf{X}(\mathbf{c}(t)) = \mathbf{X}(L_g \circ \mathbf{c}(t)) = \mathbf{X}(\mathbf{c}_1(t)).$$

Conversely, let  $\mathbf{c} : (a, b) \rightarrow G$  be the maximal integral curve of some  $\mathbf{X} \in \mathfrak{g}$  passing through  $e$  at 0. Given  $t, \tilde{t} \in \mathbb{R}$  such that  $t, \tilde{t}, t + \tilde{t} \in (a, b)$ , we must have  $\mathbf{c}(t + \tilde{t}) = \mathbf{c}(t)\mathbf{c}(\tilde{t})$  because the curves  $t \mapsto \mathbf{c}(t + \tilde{t})$  and  $t \mapsto \mathbf{c}(t)\mathbf{c}(\tilde{t})$  are both integral curves of  $\mathbf{X}$  and they agree at 0. It only remains to show that  $(a, b) = \mathbb{R}$ . But if, say,  $b < \infty$ , choose  $\tilde{b} \in (0, b)$ ; the curve  $\tilde{\mathbf{c}} : (a + \tilde{b}, b + \tilde{b}) \rightarrow G$  defined by  $\tilde{\mathbf{c}}(t) = \mathbf{c}(\tilde{b})\mathbf{c}(t - \tilde{b})$  is an integral curve of  $\mathbf{X}$  that coincides with  $\mathbf{c}$  at  $\tilde{b}$ , and must therefore coincide with  $\mathbf{c}$  everywhere. Thus, the domain of  $\mathbf{c}$  may be extended beyond  $b$ , contradicting our assumption.  $\square$

**Example 3.4.1.** Let us illustrate these various concepts for the orthogonal group  $G = O(n)$  introduced earlier. For the sake of brevity, the identity matrix  $I_n$  will be denoted by  $e$ . We begin with the Lie algebra of  $G$ , which may be identified with  $G_e$ . By Proposition 3.1.1,  $\iota_*G_e = \ker \mathbf{f}_{*e}$ , where  $\mathbf{f} : GL(n) \rightarrow M_n$  is given by  $\mathbf{f}(M) = MM^T$ , and  $\iota : G \hookrightarrow GL(n)$  denotes inclusion. We already computed that  $D\mathbf{f}(A)M = AM^T + MA^T$ , so that

$$G_e = \ker \mathbf{f}_{*e} = \{\mathcal{I}_e M \mid M + M^T = 0\}$$

is isomorphic to the space of skew-symmetric matrices. The left-invariant field  $\mathbf{X}$  with  $\mathbf{X}(e) = \mathcal{I}_e M$  is given by

$$\mathbf{X}(A) = L_{A*e} \mathbf{X}(e) = L_{A*e} \mathcal{I}_e M = \mathcal{I}_A D L_A(e) M = \mathcal{I}_A (AM), \quad A \in G$$

since  $L_A : M_n \rightarrow M_n$  is linear, so that  $D L_A(e) = L_A$ . The integral curve  $\mathbf{c}$  of  $\mathbf{X}$  that passes through  $e$  at time 0 has as expression  $\mathbf{c}(t) = \exp(tM)$  by the theorem, since  $\mathbf{c} : \mathbb{R} \rightarrow G$  is a homomorphism with  $\dot{\mathbf{c}}(0) = \mathcal{I}_e M$ . Notice that  $\mathbf{c}$  does indeed have its image in  $G$ , because

$$\mathbf{c}(t)\mathbf{c}(t)^T = \exp(tM) \exp(tM^T) = \exp(t(M + M^T)) = \exp(\mathbf{0}) = e$$

by Exercise 1.50. It follows that the integral curve of  $\mathbf{X}$  that passes through  $A$  at time zero is  $t \mapsto A \exp(tM)$ .

Next, we investigate the Lie bracket  $[\mathbf{X}, \mathbf{Y}]$  for  $\mathbf{X}, \mathbf{Y} \in \mathfrak{g}$  with  $\mathbf{X}(e) = \mathcal{I}_e M$  and  $\mathbf{Y}(e) = \mathcal{I}_e N$ . Since  $[\mathbf{X}, \mathbf{Y}] \in \mathfrak{g}$ ,

$$[\mathbf{X}, \mathbf{Y}](A) = \mathcal{I}_A (A \mathcal{I}_e^{-1} [\mathbf{X}, \mathbf{Y}](e)), \quad A \in G,$$

so that we need only identify  $[\mathbf{X}, \mathbf{Y}](e) = D_{\mathbf{X}(e)} \mathbf{Y} - D_{\mathbf{Y}(e)} \mathbf{X}$ . Now, the curve  $\mathbf{c}$ , where  $\mathbf{c}(t) = e + tM$ , has velocity vector  $\mathcal{I}_e M$  at 0, so that  $(D_{\mathbf{X}} \mathbf{Y})(e) = (\mathbf{Y} \circ \dot{\mathbf{c}})(0)$ . Furthermore,  $(\mathbf{Y} \circ \mathbf{c})(t) = \mathcal{I}_{\mathbf{c}(t)}(\mathbf{c}(t)N) = \mathcal{I}_{\mathbf{c}(t)}((e + tM)N)$ , which implies  $(D_{\mathbf{X}} \mathbf{Y})(e) = \mathcal{I}_e(MN)$ . The Lie bracket is therefore given by

$$[\mathbf{X}, \mathbf{Y}](e) = \mathcal{I}_e(MN - NM), \quad \mathbf{X}(e) = \mathcal{I}_e M, \quad \mathbf{Y}(e) = \mathcal{I}_e N.$$

In the computation of the Lie bracket, nowhere did we use the fact that  $\mathbf{X}$  and  $\mathbf{Y}$  belong to the Lie algebra of  $O(n)$ . The same argument shows that the above formula actually holds for the Lie algebra of  $GL(n)$  whose Lie algebra is isomorphic to  $M_n$ . Notice also that  $GL(1)$  is just  $\mathbb{R} \setminus \{0\}$  with ordinary multiplication.

**Remark 3.4.1.** The exponential map on  $M_n$  is extendable to the Lie algebra  $\mathfrak{g}$  of any Lie group  $G$ : given  $\mathbf{X} \in \mathfrak{g}$ , define  $\exp(\mathbf{X}) = \mathbf{c}_{\mathbf{X}}(1)$ , where  $\mathbf{c}_{\mathbf{X}}$  is the integral curve of  $\mathbf{X}$  passing through the identity at time 0. It follows that this integral curve is given by  $t \mapsto \exp(t\mathbf{X})$  for all  $t$ : to see this, let  $t_0 \in \mathbb{R}$ , and consider the curve  $\gamma$ , where  $\gamma(t) = \mathbf{c}_{\mathbf{X}}(t_0 t)$ . Then

$$\dot{\gamma}(t) = t_0 \dot{\mathbf{c}}_{\mathbf{X}}(t_0 t) = (t_0 \mathbf{X} \circ \mathbf{c}_{\mathbf{X}})(t_0 t) = (t_0 \mathbf{X} \circ \gamma)(t),$$

so that  $\gamma$  is the integral curve of  $t_0 \mathbf{X}$  through  $e$  at time 0. By uniqueness,  $\gamma = \mathbf{c}_{t_0 \mathbf{X}}$ . Thus,  $\mathbf{c}_{\mathbf{X}}(t_0) = \gamma(1) = \exp(t_0 \mathbf{X})$ . Since  $t_0$  was arbitrary, the claim follows.

$\exp : \mathfrak{g} \rightarrow G$  shares many of the properties of the matrix exponential map. For example,  $\exp(-\mathbf{X}) = (\exp \mathbf{X})^{-1}$ : the term on the left equals, by the above,  $\mathbf{c}_{\mathbf{X}}(-1)$ , whereas the one on the right is  $\mathbf{c}_{\mathbf{X}}(1)^{-1}$ . But  $\mathbf{c}_{\mathbf{X}}$  is a homomorphism, so  $e = \mathbf{c}_{\mathbf{X}}(0) = \mathbf{c}_{\mathbf{X}}(1 - 1) = \mathbf{c}_{\mathbf{X}}(1) \mathbf{c}_{\mathbf{X}}(-1)$ , which means that  $\mathbf{c}_{\mathbf{X}}(-1)$  is the inverse of  $\mathbf{c}_{\mathbf{X}}(1)$ .

More generally, recall that for  $A, B \in M_n$ , if  $AB = BA$ , then  $\exp(A + B) = \exp A \exp B$ . By the above example, identifying  $M_n$  with the Lie algebra of  $GL(n)$ , this says that for  $\mathbf{X}, \mathbf{Y} \in \mathfrak{g}$ ,  $\exp(\mathbf{X} + \mathbf{Y}) = \exp(\mathbf{X}) \exp(\mathbf{Y})$  whenever  $[\mathbf{X}, \mathbf{Y}] = \mathbf{0}$ . This is also true in arbitrary Lie algebras. A proof is outlined in Exercise 3.22.

### 3.5 The tangent bundle

The main concepts – derivatives and vector fields – introduced in the last three sections involved maps taking values in tangent spaces. These tangent spaces can be grouped together to form a new manifold, and the corresponding maps become differentiable in the process. If  $M^n$  is a submanifold of  $\mathbb{R}^{n+k}$ , define the *tangent bundle*  $TM$  of  $M$  to be

$$TM = \bigcup_{\mathbf{p} \in M} M_{\mathbf{p}}.$$

Each tangent space  $M_{\mathbf{p}}$  is a subset of  $\mathbb{R}_{\mathbf{p}}^{n+k} = \{\mathbf{p}\} \times \mathbb{R}^{n+k} \subset \mathbb{R}^{n+k} \times \mathbb{R}^{n+k}$ , so that the tangent bundle of  $M$  is contained in Euclidean space of dimension  $2(n+k)$ .

**Theorem 3.5.1.** *The tangent bundle  $TM$  of  $M^n$  is a  $2n$ -dimensional manifold.*

*Proof.* Any parametrization  $(U, \mathbf{h})$  of  $M$  induces one of  $TM$ : Define

$$\begin{aligned} \mathbf{F} : U \times \mathbb{R}^n &\rightarrow TM, \\ (\mathbf{q}, \mathbf{u}) &\mapsto (\mathbf{h}(\mathbf{q}), D\mathbf{h}(\mathbf{q})\mathbf{u}). \end{aligned}$$

The first  $n+k$  component functions of  $\mathbf{F}$  are those of  $\mathbf{h}$ , and hence are differentiable. So are the last  $n+k$ , since

$$u^{n+k+j} \circ \mathbf{F}(\mathbf{q}, \mathbf{u}) = \sum_{i=1}^n D_i h^j(\mathbf{q}) u_i.$$

Thus,  $\mathbf{F}$  is differentiable, and has maximal rank because  $\mathbf{h}$  has. It remains to show that  $\mathbf{F}^{-1}$  is continuous. So consider a convergent sequence

$$(\mathbf{h}(\mathbf{p}_n), D\mathbf{h}(\mathbf{p}_n)\mathbf{u}_n) \rightarrow (\mathbf{h}(\mathbf{p}), D\mathbf{h}(\mathbf{p})\mathbf{u})$$

in the image of  $\mathbf{F}$ . It must be shown that  $\mathbf{p}_n \rightarrow \mathbf{p}$  and  $\mathbf{u}_n \rightarrow \mathbf{u}$ . The first sequence converges by continuity of  $\mathbf{h}^{-1}$ . Since  $\mathbf{h}$  is continuously differentiable, this, in turn, implies that

$$D\mathbf{h}(\mathbf{p}_n) \rightarrow D\mathbf{h}(\mathbf{p}), \quad (3.5.1)$$

in the sense that  $|D\mathbf{h}(\mathbf{p}_n) - D\mathbf{h}(\mathbf{p})| \rightarrow 0$  with the norm from Definition 1.4.1. By Proposition 3.2.1, there exists an open set  $V \subset \mathbb{R}^{2(n+k)}$  and a smooth map  $\mathbf{G} : V \rightarrow \mathbb{R}^{n+k}$  such that  $\mathbf{G} \circ \mathbf{h} = 1_U$ . Applying  $D\mathbf{G}(\mathbf{p})$  to (3.5.1), we then obtain

$$D\mathbf{G}(\mathbf{p})D\mathbf{h}(\mathbf{p}_n) \rightarrow 1_{\mathbb{R}^{n+k}}. \quad (3.5.2)$$

Now,

$$|\mathbf{u}_n - \mathbf{u}| \leq |\mathbf{u}_n - D\mathbf{G}(\mathbf{p})D\mathbf{h}(\mathbf{p}_n)\mathbf{u}_n| + |D\mathbf{G}(\mathbf{p})D\mathbf{h}(\mathbf{p}_n)\mathbf{u}_n - \mathbf{u}|.$$

The first term on the right is no larger than  $|1_{\mathbb{R}^{n+k}} - D\mathbf{G}(\mathbf{p})D\mathbf{h}(\mathbf{p}_n)| |\mathbf{u}_n|$ , which goes to zero by (3.5.2) (we are implicitly using here the fact that  $\{\mathbf{u}_n\}$  is bounded, which follows from convergence of  $\{D\mathbf{h}(\mathbf{p}_n)\mathbf{u}_n\}$  and  $\{D\mathbf{h}(\mathbf{p}_n)\}$ ). The second term can be written

$$|D\mathbf{G}(\mathbf{p})D\mathbf{h}(\mathbf{p}_n)\mathbf{u}_n - D\mathbf{G}(\mathbf{p})D\mathbf{h}(\mathbf{p})\mathbf{u}| \leq |D\mathbf{G}(\mathbf{p})| |D\mathbf{h}(\mathbf{p}_n)\mathbf{u}_n - D\mathbf{h}(\mathbf{p})\mathbf{u}|,$$

which also goes to zero since  $Dh(\mathbf{p}_n)\mathbf{u}_n \rightarrow Dh(\mathbf{p})\mathbf{u}$  by hypothesis. This shows that  $\mathbf{F}^{-1}$  is continuous.  $\square$

One could also have argued Theorem 3.5.1 in terms of charts instead of parametrizations. The reader is invited to check that a chart  $(U, \mathbf{x})$  of  $M$  induces a chart  $(\pi^{-1}(U), \bar{\mathbf{x}})$  of  $TM$  where

$$\bar{\mathbf{x}} = (\mathbf{x}, dx^1, \dots, dx^n) \circ \pi. \tag{3.5.3}$$

Given a smooth map  $f : M \rightarrow N$  between manifolds  $M$  and  $N$ , we define a map  $f_* : TM \rightarrow TN$  between their tangent bundles by setting

$$f_*\mathbf{u} = f_{*\mathbf{p}}\mathbf{u}, \quad \text{if } \mathbf{u} \in M_{\mathbf{p}}.$$

**Corollary 3.5.1.** *If  $f : M \rightarrow N$  is differentiable, then so is  $f_* : TM \rightarrow TN$ .*

*Proof.* It must be shown that if  $(U \times \mathbb{R}^n, \mathbf{F})$  is a parametrization of  $TM$  as in Theorem 3.5.1, then  $f_* \circ \mathbf{F}$  is differentiable in the usual sense. But

$$(f_* \circ \mathbf{F})(\mathbf{p}, \mathbf{u}) = f_*(\mathbf{h}(\mathbf{p}), Dh(\mathbf{p})\mathbf{u}) = ((f \circ \mathbf{h})(\mathbf{p}), D(f \circ \mathbf{h})(\mathbf{p})\mathbf{u})$$

by (3.2.3), which establishes the claim.  $\square$

The *bundle projection*  $\pi_M : TM \rightarrow M$  is the map that sends a vector to its base point; i.e., any  $\mathbf{v} \in TM$  belongs to some  $M_{\mathbf{p}}$  for a unique  $\mathbf{p} \in M$ , and  $\pi_M(\mathbf{v})$  is defined to be this  $\mathbf{p}$ . It is an easy exercise to show that  $\pi_M$  is differentiable. When  $M = \mathbb{R}^n$ ,

$$T\mathbb{R}^n = \bigcup_{\mathbf{p} \in \mathbb{R}^n} \mathbb{R}_{\mathbf{p}}^n = \bigcup_{\mathbf{p} \in \mathbb{R}^n} \{\mathbf{p}\} \times \mathbb{R}^n = \mathbb{R}^n \times \mathbb{R}^n,$$

so that the bundle projection is the projection  $\pi_1 : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n$  onto the first factor. If  $\pi_2$  is the projection onto the second factor, then  $\mathbf{u} = \mathcal{I}_{\mathbf{p}}\mathbf{v} \in T\mathbb{R}^n$  if and only if  $\pi_1(\mathbf{u}) = \mathbf{p}$  and  $\pi_2(\mathbf{u}) = \mathbf{v}$ . In general, notice that for  $f : N \rightarrow M$ , the diagram

$$\begin{array}{ccc} TN & \xrightarrow{f_*} & TM \\ \pi_N \downarrow & & \downarrow \pi_M \\ N & \xrightarrow{f} & M \end{array}$$

commutes.

**Corollary 3.5.2.** *A vector field  $\mathbf{X}$  on  $M$  is differentiable as a map  $\mathbf{X} : M \rightarrow TM$ .*

*Proof.*  $\mathbf{X}$  is differentiable if  $\mathbf{X} \circ \mathbf{h}$  is smooth for any local parametrization  $(U, \mathbf{h})$  of  $M$ . But this is immediate from (3.3.1).  $\square$

**Definition 3.5.1.** Let  $N, M$  denote manifolds,  $f : N \rightarrow M \subset \mathbb{R}^{n+k}$  a map. A *vector field along  $f$*  is a differentiable map  $\mathbf{X} : N \rightarrow TM$  that assigns to each  $\mathbf{p} \in N$  a vector

$\mathbf{X}(\mathbf{p}) \in M_{f(\mathbf{p})}$ ; i.e.,  $\mathbf{X}$  is a differentiable map for which the diagram

$$\begin{array}{ccc} & & TM \\ & \nearrow \mathbf{X} & \downarrow \pi_M \\ N & \xrightarrow{f} & M \end{array}$$

commutes.

The above definition clearly generalizes that of vector fields along maps between Euclidean spaces given in Chapter 2.

**Remark 3.5.1.** Any vector field  $\mathbf{X}$  on  $N$  induces a vector field along  $f : N \rightarrow M$ , namely  $\mathbf{p} \mapsto f_{*\mathbf{p}}\mathbf{X}(\mathbf{p})$ . Similarly, any vector field  $\mathbf{Y}$  on  $M$  induces one along  $f$ , namely  $\mathbf{Y} \circ f$ . In general, given  $\mathbf{X}$  as above, there need not exist any vector field  $\mathbf{Y}$  on  $M$  such that  $f_*\mathbf{X} = \mathbf{Y} \circ f$  (i.e., such that  $\mathbf{X}$  and  $\mathbf{Y}$  are  $f$ -related). If  $f_{*\mathbf{p}}$  is one-to-one at  $\mathbf{p} \in N$  however, then  $f$  is one-to-one on a neighborhood of  $\mathbf{p}$  by the inverse function theorem; thus, there exists a neighborhood  $U$  of  $\mathbf{p}$  and a vector field  $\tilde{\mathbf{X}}$  on  $f(U)$  that is  $f$ -related to  $\mathbf{X}$  (namely,  $\tilde{\mathbf{X}} = f_*\mathbf{X} \circ (f|_U)^{-1}$ ). Smoothness of  $\mathbf{X}$  in the definition above is then equivalent to smoothness of  $\tilde{\mathbf{X}}$ .

### 3.6 Covariant differentiation

In Definition 2.8.7, we introduced the covariant derivative  $D_{\mathbf{u}}\mathbf{X}$  of a vector field  $\mathbf{X}$  in Euclidean space  $\mathbb{R}^n$  with respect to a vector  $\mathbf{u} \in \mathbb{R}_p^n$ . We wish to extend this concept to a vector field  $\mathbf{X}$  on a manifold  $M \subset \mathbb{R}^n$ . Merely adopting the same formula does not work, because although  $D_{\mathbf{u}}\mathbf{X} \in \mathbb{R}_p^n$ , it need not belong to the tangent space of  $M$  at  $\mathbf{p}$ . The best we can do is to project this vector back onto  $M_p$ . This can be done using the Riemannian metric on  $\mathbb{R}^n$ : there is a decomposition

$$\mathbb{R}_p^n = M_p \oplus M_p^\perp$$

of  $\mathbb{R}_p^n$  as a direct sum of the tangent space of  $M$  at that point with its orthogonal complement, and any  $\mathbf{u} \in \mathbb{R}_p^n$  has a corresponding unique decomposition

$$\mathbf{u} = \mathbf{u}^\top + \mathbf{u}^\perp \in M_p \oplus M_p^\perp,$$

thus inducing maps  $\top : \mathbb{R}_p^n \rightarrow M_p$ ,  $\perp : \mathbb{R}_p^n \rightarrow M_p^\perp$ , that project  $\mathbf{u}$  onto  $\mathbf{u}^\top$  and  $\mathbf{u}^\perp$  respectively.  $M_p^\perp$  is called the *normal space* of  $M$  at  $\mathbf{p}$ .

**Definition 3.6.1.** If  $M$  is a submanifold of  $\mathbb{R}^n$ ,  $\mathbf{X}$  a vector field on some open set  $U \subset M$ ,  $\mathbf{p} \in U$ , and  $\mathbf{u} \in M_p$ , the *covariant derivative* of  $\mathbf{X}$  with respect to  $\mathbf{u}$  is

$$\nabla_{\mathbf{u}}\mathbf{X} := (D_{\mathbf{u}}\mathbf{X})^\top. \quad (3.6.1)$$

More generally, as in the Euclidean setting, the same formula defines the covariant derivative of a vector field  $\mathbf{X}$  along a map  $\mathbf{f} : N \rightarrow M$  with respect to some  $\mathbf{u} \in N_{\mathbf{p}}$ . Such an  $\mathbf{X}$  will sometimes be called a *vector field in  $M$  along  $\mathbf{f}$*  to emphasize the fact that  $\mathbf{X}$  takes values in the tangent bundle of  $M$  rather than that of the ambient space. In particular, if  $\mathbf{X}$  is a vector field in  $M$  along a curve  $\mathbf{c} : I \rightarrow M$ , its covariant derivative at  $t$  is

$$(\nabla_{\mathbf{D}}\mathbf{X})(t) := \nabla_{\mathbf{D}(t)}\mathbf{X} = (\mathbf{X}'(t))^{\top}. \quad (3.6.2)$$

Likewise, given vector fields  $\mathbf{X}, \mathbf{Y}$  on  $M$ , one obtains a new vector field  $\nabla_{\mathbf{X}}\mathbf{Y}$  by defining  $(\nabla_{\mathbf{X}}\mathbf{Y})(\mathbf{p}) = \nabla_{\mathbf{X}(\mathbf{p})}\mathbf{Y}$ . It must, of course, be checked that  $\nabla_{\mathbf{X}}\mathbf{Y}$  is differentiable. For this, we will use the following:

**Theorem 3.6.1.** *If  $M^n$  is a submanifold of  $\mathbb{R}^{n+k}$ , then for any  $\mathbf{p} \in M$ , there exists an open neighborhood  $U$  of  $\mathbf{p}$  in  $\mathbb{R}^{n+k}$ , and (differentiable) vector fields  $\mathbf{N}_1, \dots, \mathbf{N}_k$  on  $U$  such that  $\mathbf{N}_1(\mathbf{q}), \dots, \mathbf{N}_k(\mathbf{q})$  form an orthonormal basis of the normal space of  $M$  at each  $\mathbf{q} \in U \cap M$ .*

*Proof.* By Corollary 3.2.1, there exists a neighborhood  $U$  of  $\mathbf{p}$  such that  $U \cap M = \mathbf{f}^{-1}(\mathbf{0})$  for some map  $\mathbf{f} : U \rightarrow \mathbb{R}^k$  that has  $\mathbf{0}$  as regular value. Since the tangent space of  $M$  at any point of  $U$  equals the kernel of  $\mathbf{f}_*$  at that point, and since  $\mathbf{u} \in \ker \mathbf{f}_{*\mathbf{q}}$  if and only if  $0 = D_{\mathbf{u}}\mathbf{f}^i = \langle \nabla \mathbf{f}^i(\mathbf{q}), \mathbf{u} \rangle$  for  $i = 1, \dots, k$  (see Examples 2.8.2 (iv)), it follows that the vector fields  $\nabla \mathbf{f}^i$  on  $U$  are linearly independent vector fields that span the normal space of  $M$  at every point of  $U \cap M$ . Apply the Gram-Schmidt orthogonalization process (Theorem 1.4.2), observing that if  $\mathbf{X}$  and  $\mathbf{Y}$  are differentiable, then so is  $\text{proj}_{\mathbf{X}} \mathbf{Y}$ , to obtain differentiable fields  $\mathbf{N}_1, \dots, \mathbf{N}_k$  on  $U$  that form an orthonormal basis of the normal space of  $M$  at each  $\mathbf{q} \in U \cap M$ .  $\square$

To see that  $\nabla_{\mathbf{X}}\mathbf{Y}$  is differentiable if  $\mathbf{X}$  and  $\mathbf{Y}$  are, recall that the latter may be extended to an open neighborhood  $U$  of any  $\mathbf{p}$  (Theorem 3.3.1). This neighborhood may be assumed to be the one in Theorem 3.6.1. Then  $D_{\mathbf{X}}\mathbf{Y}$  is differentiable on  $U$ , and so is

$$D_{\mathbf{X}}\mathbf{Y} - \sum_{i=1}^k \langle D_{\mathbf{X}}\mathbf{Y}, \mathbf{N}_i \rangle \mathbf{N}_i.$$

But by definition this vector field equals  $\nabla_{\mathbf{X}}\mathbf{Y}$  at any point of  $U \cap M$ . This shows that  $\nabla_{\mathbf{X}}\mathbf{Y}$  is smooth.

The exact same argument shows that if  $\mathbf{X}$  is a vector field on  $N$  and  $\mathbf{Y}$  is a vector field along  $\mathbf{f} : N \rightarrow M$ , then the vector field  $\nabla_{\mathbf{X}}\mathbf{Y}$  along  $\mathbf{f}$  is smooth.

The properties of covariant derivatives in Euclidean space that were established in Chapter 2 also hold on manifolds:

**Theorem 3.6.2.** *Let  $\mathbf{X}, \mathbf{Y}$  denote vector fields on an open subset  $U$  of a manifold  $M$ ,  $f : U \rightarrow \mathbb{R}$ ,  $a \in \mathbb{R}$ . Given  $\mathbf{p} \in U$ ,  $\mathbf{u} \in M_{\mathbf{p}}$ ,*

- (1)  $\nabla_{\mathbf{u}}(a\mathbf{X} + \mathbf{Y}) = a\nabla_{\mathbf{u}}\mathbf{X} + \nabla_{\mathbf{u}}\mathbf{Y}$ ;
- (2)  $\nabla_{a\mathbf{u}+\mathbf{v}}\mathbf{X} = a\nabla_{\mathbf{u}}\mathbf{X} + \nabla_{\mathbf{v}}\mathbf{X}$ ;
- (3)  $\nabla_{\mathbf{u}}f\mathbf{X} = \mathbf{u}(f)\mathbf{X}(\mathbf{p}) + f(\mathbf{p})\nabla_{\mathbf{u}}\mathbf{X}$ ;



- (4)  $\nabla_{fX}Y = f\nabla_XY$ ;  
 (5)  $\nabla_XY - \nabla_YX = [X, Y]$ ;  
 (6) Given  $g : N \rightarrow M$ ,  $v \in TN$ ,  $\nabla_v(X \circ g) = \nabla_{g_*v}X$ ;  
 (7) Let  $X, Y$  denote vector fields in  $M$  along a map  $f : N \rightarrow M$ ,  $p \in N$ ,  $u \in N_p$ . Then

$$u\langle X, Y \rangle = \langle \nabla_u X, Y(p) \rangle + \langle X(p), \nabla_u Y \rangle.$$

*Proof.* These properties follow from the corresponding ones in Euclidean space that were established in Theorem 2.8.4, together with linearity of the projection  $\tau$  for the first four. In the same way, (6) is a consequence of (2.8.4). To prove the identity (5), consider any  $p \in M$ , and extend the vector fields  $X$  and  $Y$  to a neighborhood of  $p$  in the ambient Euclidean space. Then  $[X, Y](p) = D_{X(p)}Y - D_{Y(p)}X$ . But the left side of this identity is tangent to  $M$ . Thus,

$$[X, Y](p) = (D_{X(p)}Y - D_{Y(p)}X)^\top = \nabla_{X(p)}Y - \nabla_{Y(p)}X.$$

For the last one, notice that since  $Y(p)$  is tangent to  $M$ ,

$$\begin{aligned} \langle D_u X, Y(p) \rangle &= \langle (D_u X)^\top + (D_u X)^\perp, Y(p) \rangle = \langle (D_u X)^\top, Y(p) \rangle \\ &= \langle \nabla_u Y(p), Y(p) \rangle, \end{aligned}$$

with a similar identity holding for the other term. The claim then follows from Theorem 2.8.4.  $\square$

**Definition 3.6.2.** A vector field  $X$  in  $M$  along a curve  $c : I \rightarrow M$  is said to be *parallel* if

$$\nabla_{\dot{c}}X = \mathbf{0}.$$

Notice that if  $X$  and  $Y$  are parallel along  $c$ , then  $\langle X, Y \rangle$  is a constant function by part (7) of Theorem 3.6.2. In particular, parallel vector fields have constant norm.

**Theorem 3.6.3.** Let  $c : [0, a] \rightarrow M$  be a curve in a manifold  $M^n \subset \mathbb{R}^{n+k}$ . Then

- (1) For any  $u \in M_{c(0)}$ , there exists one and only one parallel field  $X$  along  $c$  such that  $X(0) = u$ .  
 (2) The (well-defined by (1)) map  $M_{c(0)} \rightarrow M_{c(a)}$  which sends  $u \in M_{c(0)}$  to  $X(a)$ , where  $X$  is the parallel field along  $c$  with  $X(0) = u$  is a linear isometry, called parallel translation along  $c$ .

*Proof.* The image of  $c$  can be broken up into finitely many pieces, each of which lie in an open set admitting an orthonormal basis of vector fields normal to  $M$ . Since it suffices to prove the theorem for each such piece, we may assume the whole image lies in such a set, so that there exist vector fields  $N_1, \dots, N_k$  such that  $\{N_i \circ c(t) \mid 1 \leq i \leq k\}$  is an orthonormal basis of  $M_{c(t)}^\perp$  for each  $t \in [0, a]$ . A vector field  $X$  along  $c$  is parallel

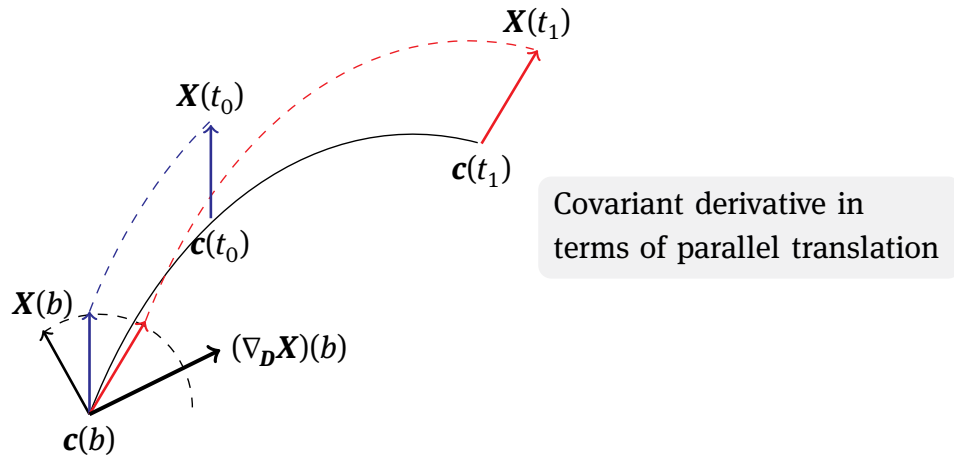
iff

$$\begin{aligned} \mathbf{0} &= \nabla_{\mathbf{D}} \mathbf{X} = \mathbf{X}' - \sum_j \langle \mathbf{X}', \mathbf{N}_j \circ \mathbf{c} \rangle \mathbf{N}_j \circ \mathbf{c} \\ &= \mathbf{X}' - \sum_j \left( \langle \mathbf{X}, \mathbf{N}_j \circ \mathbf{c} \rangle' - \langle \mathbf{X}, (\mathbf{N}_j \circ \mathbf{c})' \rangle \right) \mathbf{N}_j \circ \mathbf{c} \\ &= \mathbf{X}' + \sum_j \langle \mathbf{X}, (\mathbf{N}_j \circ \mathbf{c})' \rangle \mathbf{N}_j \circ \mathbf{c}. \end{aligned}$$

Writing  $\mathbf{X} = \sum X^i \mathbf{D}_i \circ \mathbf{c}$  in components, with  $X^i = \langle \mathbf{X}, \mathbf{D}_i \circ \mathbf{c} \rangle$ , and doing the same for  $\mathbf{N}_j$ , the above equation becomes a system of ordinary differential equations

$$X^{i'} + \sum_{j,l} X^l (N_j^l \circ \mathbf{c}) (N_j^i \circ \mathbf{c})' = 0, \quad i = 1, \dots, n+k.$$

This system is linear, in the sense that a sum of two solutions as well as a scalar multiple of a solution are again solutions. The theory of linear differential equations then guarantees the existence of a unique collection  $X^i$  of solutions defined on  $[0, a]$  with initial conditions  $X^i(0) = u^i$ , where  $u^i = \langle \mathbf{u}, \mathbf{D}_i \circ \mathbf{c}(0) \rangle$ . The second part of Theorem 3.6.3 is an immediate consequence of the remark preceding the statement of the theorem.  $\square$



Conversely, covariant derivatives can be expressed in terms of parallel translation:

**Theorem 3.6.4.** *Let  $c : [0, a] \rightarrow M$  be a curve in a manifold  $M^n \subset \mathbb{R}^{n+k}$ , and  $\mathbf{X}$  a vector field in  $M$  along  $c$ . For  $b \in (0, a)$ , let  $\mathbf{X}^b$  denote the parallel vector field along  $c$  that equals  $\mathbf{X}(b)$  at  $b$ . Then*

$$(\nabla_{\mathbf{D}} \mathbf{X})(b) = \lim_{t \rightarrow b} \frac{\mathbf{X}^t(b) - \mathbf{X}(b)}{t - b}.$$

*Proof.* By Theorem 3.6.3, there exist parallel vector fields  $\mathbf{Y}_i$  in  $M$  along  $c$  such that  $\mathbf{Y}_1(t), \dots, \mathbf{Y}_n(t)$  form a basis of  $M_{c(t)}$  for each  $t \in [0, a]$ . Thus,  $\mathbf{X} = \sum_i f_i \mathbf{Y}_i$  for smooth functions  $f_i : [0, a] \rightarrow \mathbb{R}$ , and the left side of the above identity can be written as

$$(\nabla_{\mathbf{D}} \mathbf{X})(b) = \sum_i f_i'(b) \mathbf{Y}_i(b).$$

By definition of  $\mathbf{X}^t$ ,  $\mathbf{X}^t(b) = \sum_i f_i(t) \mathbf{Y}_i(b)$ , so that the right side reads

$$\lim_{t \rightarrow b} \frac{\mathbf{X}^t(b) - \mathbf{X}(b)}{t - b} = \sum_i \lim_{t \rightarrow b} \frac{f_i(t) - f_i(b)}{t - b} \mathbf{Y}_i(b).$$

Comparing both expressions now implies the claim.  $\square$

The above theorem roughly says that in order to compute the covariant derivative of  $\mathbf{X}$  at  $b$ , one parallel translates for each  $t$  close to  $b$  the vector  $\mathbf{X}(t)$  to the tangent space of  $M$  at  $\mathbf{c}(b)$ . This yields a curve in the vector space  $M_{\mathbf{c}(b)}$ , and its derivative at  $b$  is the desired vector.

**Example 3.6.1.** Let us explore parallel translation along a parametrized circle of latitude in the sphere  $M = S^2(1)$ . The circle at height  $a \in (-1, 1)$  may be parametrized by  $\mathbf{c} : \mathbb{R} \rightarrow M$ , where

$$\mathbf{c}(t) = (b \cos t, b \sin t, a), \quad b := \sqrt{1 - a^2}.$$

The vector fields  $\mathbf{X}$ ,  $\mathbf{Y}$  along  $\mathbf{c}$ , with

$$\begin{aligned} \mathbf{X} &= \frac{1}{|\dot{\mathbf{c}}|} \dot{\mathbf{c}} = -(\sin) \mathbf{D}_1 \circ \mathbf{c} + (\cos) \mathbf{D}_2 \circ \mathbf{c}, \\ \mathbf{Y} &= -a(\cos) \mathbf{D}_1 \circ \mathbf{c} - a(\sin) \mathbf{D}_2 \circ \mathbf{c} + b \mathbf{D}_3 \circ \mathbf{c}, \end{aligned}$$

form an orthonormal basis of the tangent space of the sphere at  $\mathbf{c}(t)$ , which, together with the unit normal vector field

$$\mathbf{N} = b(\cos) \mathbf{D}_1 \circ \mathbf{c} + b(\sin) \mathbf{D}_2 \circ \mathbf{c} + a \mathbf{D}_3 \circ \mathbf{c}$$

yield an orthonormal basis of the tangent space of  $\mathbb{R}^3$  at  $\mathbf{c}(t)$ . Since  $\mathbf{X}' = -(\cos) \mathbf{D}_1 \circ \mathbf{c} - (\sin) \mathbf{D}_2 \circ \mathbf{c}$ , and  $\langle \mathbf{X}', \mathbf{N} \rangle = -b$ , we deduce

$$\begin{aligned} \nabla_{\mathbf{D}} \mathbf{X} &= \mathbf{X}' - \langle \mathbf{X}', \mathbf{N} \rangle \mathbf{N} \\ &= (-\cos + b^2 \cos) \mathbf{D}_1 \circ \mathbf{c} + (-\sin + b^2 \sin) \mathbf{D}_2 \circ \mathbf{c} + ab \mathbf{D}_3 \circ \mathbf{c} \\ &= a \mathbf{Y}, \end{aligned}$$

where the last equality uses the fact that  $b^2 - 1 = -a^2$ . Similarly,

$$\nabla_{\mathbf{D}} \mathbf{Y} = -a \mathbf{X}.$$

Thus, the vector fields  $\mathbf{X}$  and  $\mathbf{Y}$  are parallel only when the circle in question is a great circle (namely, the equator, corresponding to  $a = 0$ ). In this case, parallel translation along  $\mathbf{c}$  from  $\mathbf{c}(0)$  to  $\mathbf{c}(2\pi) = \mathbf{c}(0)$  is the identity. In general, any vector field  $\mathbf{Z}$  in  $M$  along  $\mathbf{c}$  may be expressed as

$$\mathbf{Z} = f \mathbf{X} + g \mathbf{Y}, \quad f := \langle \mathbf{Z}, \mathbf{X} \rangle, \quad g := \langle \mathbf{Z}, \mathbf{Y} \rangle,$$

so that  $\mathbf{Z}$  is parallel when

$$\mathbf{0} = \nabla_{\mathbf{D}}(f\mathbf{X} + g\mathbf{Y}) = f'\mathbf{X} + af\mathbf{Y} + g'\mathbf{Y} - ag\mathbf{X}.$$

This yields the system of linear differential equations

$$f' - ag = 0, \quad g' + af = 0,$$

which may be written as

$$\begin{bmatrix} f' \\ g' \end{bmatrix} = A \begin{bmatrix} f \\ g \end{bmatrix}, \quad A = \begin{bmatrix} 0 & a \\ -a & 0 \end{bmatrix}.$$

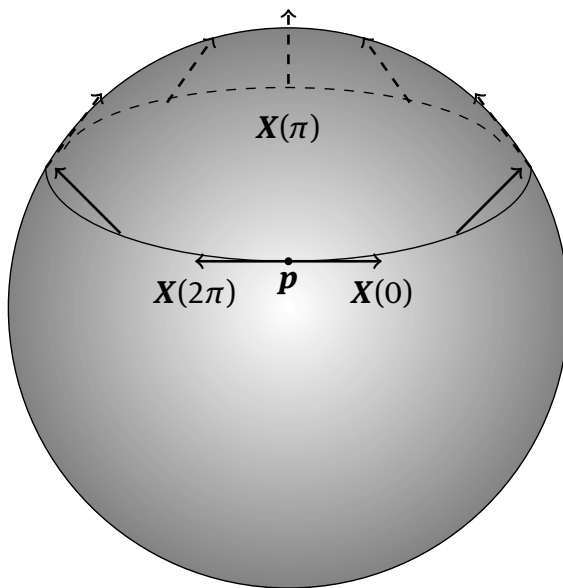
By Example 2.8.1, it admits as general solution

$$\begin{bmatrix} f(t) \\ g(t) \end{bmatrix} = e^{tA} \begin{bmatrix} f(0) \\ g(0) \end{bmatrix} = \begin{bmatrix} \cos(at) & \sin(at) \\ -\sin(at) & \cos(at) \end{bmatrix} \begin{bmatrix} f(0) \\ g(0) \end{bmatrix}.$$

In other words, the parallel field  $\mathbf{Z}$  along  $\mathbf{c}$  with “initial condition”  $\mathbf{Z}(0) = c_1\mathbf{X}(0) + c_2\mathbf{Y}(0)$  is given by

$$\mathbf{Z}(t) = (c_1 \cos(at) + c_2 \sin(at))\mathbf{X}(t) + (c_2 \cos(at) - c_1 \sin(at))\mathbf{Y}(t).$$

Thus, parallel translation along  $\mathbf{c}$  from  $\mathbf{c}(0)$  to  $\mathbf{c}(2\pi) = \mathbf{c}(0)$  consists of rotation by angle  $2\pi a$ . In particular, and as remarked earlier, parallel translation along the equator is the identity. We will see that curvature, a concept to be soon introduced, is responsible for parallel translation along closed curves differing from the identity.



Parallel translation along the circle of latitude through  $\mathbf{p}$  results in a rotation about the origin by angle  $\pi$  in  $S^2_{\mathbf{p}}$  when  $\mathbf{p}$  has height  $1/2$ .

### 3.7 Geodesics

A fundamental concept in differential geometry is that of curves of shortest length between two points in a manifold. Such curves (or rather certain parametrizations of such curves) are called *geodesics*.

**Definition 3.7.1.** A curve  $\mathbf{c} : I \rightarrow M$  is said to be a *geodesic* if its acceleration is everywhere orthogonal to  $M$ ; i.e., if

$$\nabla_{\mathbf{D}} \dot{\mathbf{c}} = (\dot{\mathbf{c}})^\top = \mathbf{0}.$$

Thus, from the surface's point of view, geodesics are those curves that have no acceleration. In particular, they have constant speed  $|\dot{\mathbf{c}}|$ :

$$|\dot{\mathbf{c}}|^{2'} = \langle \dot{\mathbf{c}}, \dot{\mathbf{c}} \rangle' = 2 \langle \nabla_{\mathbf{D}} \dot{\mathbf{c}}, \dot{\mathbf{c}} \rangle = 0.$$

Geodesics are also invariant under linear reparametrizations: if  $\mathbf{c}$  is a geodesic and  $t \mapsto \varphi(t) := at + b$ ,  $a, b \in \mathbb{R}$ , then so is  $\mathbf{c} \circ \varphi$ : in fact,

$$\mathbf{c} \circ \varphi = (\mathbf{c} \circ \varphi)_* \mathbf{D} = \mathbf{c}_* \circ \varphi_* \mathbf{D} = \varphi' \mathbf{c}_* (\mathbf{D} \circ \varphi) = a \dot{\mathbf{c}} \circ \varphi,$$

so that by Theorem 3.6.2,

$$\nabla_{\mathbf{D}} (\mathbf{c} \circ \varphi) = a \nabla_{\mathbf{D}} (\dot{\mathbf{c}} \circ \varphi) = a \nabla_{\varphi_* \mathbf{D}} \dot{\mathbf{c}} = a^2 (\nabla_{\mathbf{D}} \dot{\mathbf{c}}) \circ \varphi \quad (3.71)$$

also vanishes. In particular, any nonconstant geodesic admits a reparametrization by arc length which is again a geodesic. The latter is called a *normal geodesic*. A similar definition can be made for Euclidean space itself: a curve  $\mathbf{c}$  in  $\mathbb{R}^n$  is said to be a geodesic if  $\dot{\mathbf{c}}' \equiv 0$ ; equivalently,  $\mathbf{c}$  has vanishing acceleration. It is easy to see that  $\mathbf{c}$  is a geodesic in  $\mathbb{R}^n$  if and only if  $\mathbf{c}(t) = \mathbf{p} + t\mathbf{u}$  for some  $\mathbf{p}, \mathbf{u} \in \mathbb{R}^n$ .

**Theorem 3.7.1.** For any  $\mathbf{p} \in M$  and any  $\mathbf{v} \in M_{\mathbf{p}}$ , there exists a unique geodesic  $\mathbf{c}_{\mathbf{v}}$  such that  $\mathbf{c}_{\mathbf{v}}(0) = \mathbf{p}$  and  $\dot{\mathbf{c}}_{\mathbf{v}}(0) = \mathbf{v}$ .

*Proof.* Any  $\mathbf{p} \in M$  has a neighborhood  $U$  on which there exists an orthonormal basis  $\mathbf{N}_i$  of vector fields normal to  $M$ . The geodesic equation in  $U$  then reads

$$\dot{\mathbf{c}}' - \sum_i \langle \mathbf{N}_i \circ \mathbf{c}, \dot{\mathbf{c}} \rangle \mathbf{N}_i \circ \mathbf{c} = 0.$$

Since

$$\langle \mathbf{N}_i \circ \mathbf{c}, \dot{\mathbf{c}} \rangle' = \langle \mathbf{N}_i \circ \mathbf{c}, \dot{\mathbf{c}} \rangle' - \langle (\mathbf{N}_i \circ \mathbf{c})', \dot{\mathbf{c}} \rangle = -\langle (\mathbf{N}_i \circ \mathbf{c})', \dot{\mathbf{c}} \rangle,$$

this equation can be rewritten

$$\dot{\mathbf{c}}' + \sum_i \langle (\mathbf{N}_i \circ \mathbf{c})', \dot{\mathbf{c}} \rangle \mathbf{N}_i \circ \mathbf{c} = 0. \quad (3.72)$$

If  $M^n \subset \mathbb{R}^{n+l}$ , let  $c^k = u^k \circ \mathbf{c}$  as usual, and  $N_i^k = \langle \mathbf{N}_i, \mathbf{D}_k \rangle$ . (3.72) is then equivalent to the second order linear system

$$c^{k''} + \sum_{i=1}^l \sum_{j=1}^{n+l} (N_i^j \circ \mathbf{c})' (N_i^k \circ \mathbf{c}) c^{j'} = 0, \quad k = 1, \dots, n+l. \quad (3.73)$$

Standard results from the theory of ordinary differential equations now guarantee the existence of an interval  $I$  containing 0, and solutions  $c^k : I \rightarrow \mathbb{R}$  of (3.73) satisfying arbitrary initial conditions  $c^k(0) = u^k(\mathbf{p})$  and  $c^{k'}(0) = \mathbf{v}(u^k)$ .  $\square$

It should be noted that the system of equations (3.7.3) is not linear, so that the solutions need not be defined on the whole real line.

**Example 3.7.1** (Geodesics on the sphere). Let  $M = S^n(r) \subset \mathbb{R}^{n+1}$ . Example 3.6.1 implies that for  $n = 2$  and  $r = 1$  the equatorial curve  $t \mapsto (\cos t, \sin t, 0)$  is a geodesic. More generally, consider any  $\mathbf{p} \in M$ , and recall that the tangent space of  $M$  at  $\mathbf{p}$  is just  $\mathcal{I}_{\mathbf{p}}\mathbf{p}^\perp$ . Let  $\mathbf{u} \in \mathbb{R}^{n+1}$  be a unit vector orthogonal to  $\mathbf{p}$ , and define  $\mathbf{c}_u : \mathbb{R} \rightarrow M$  by

$$\mathbf{c}_u(t) = \left(\cos \frac{t}{r}\right)\mathbf{p} + \left(\sin \frac{t}{r}\right)r\mathbf{u}.$$

We claim that  $\mathbf{c}_u$  is the geodesic in  $M$  with  $\mathbf{c}_u(0) = \mathbf{p}$ ,  $\dot{\mathbf{c}}_u(0) = \mathcal{I}_{\mathbf{p}}\mathbf{u}$ . Indeed,  $\mathbf{c}_u'' = -(1/r^2)\mathbf{c}_u$ , so that  $\dot{\mathbf{c}}_u' = \mathcal{I}_{\mathbf{c}_u}\mathbf{c}_u''$  is orthogonal to  $M$ . This means that  $\mathbf{c}_u$  is a geodesic in the sphere, and the initial conditions are easily verified.

Implicit in the proof of Theorem 3.7.1 is the fact that a geodesic  $\mathbf{c}$  “depends smoothly” on initial conditions  $\mathbf{c}(0)$  and  $\dot{\mathbf{c}}(0)$ . This is also implied by the following fact, which we next establish: the velocity fields of geodesics are integral curves of a certain vector field on  $TM$ . To see this, recall that if  $M^n \subset \mathbb{R}^{n+l}$ , then  $TM \subset M \times \mathbb{R}^{n+l} \subset \mathbb{R}^{2(n+l)}$  is a  $2n$ -dimensional manifold. If  $\pi : TM \subset M \times \mathbb{R}^{n+l} \rightarrow M$  and  $\pi_2 : TM \rightarrow \mathbb{R}^{n+l}$  denote the projections onto the two factors, then any point  $\mathbf{u} \in TM$  can be written

$$\mathbf{u} = (\pi(\mathbf{u}), \pi_2(\mathbf{u})) = \mathcal{I}_{\pi(\mathbf{u})}\pi_2(\mathbf{u}). \quad (3.7.4)$$

Given a curve  $\mathbf{c} : I \rightarrow M$ , its velocity field  $\dot{\mathbf{c}} : I \rightarrow TM$  is a curve in  $TM$  with

$$\dot{\mathbf{c}}(t) = (\mathbf{c}(t), \mathbf{c}'(t)) = \mathcal{I}_{\mathbf{c}(t)}\mathbf{c}'(t) \in M_{\mathbf{c}(t)}.$$

Repeating the above process yields a curve  $\ddot{\mathbf{c}} : I \rightarrow T(TM)$  given by

$$\ddot{\mathbf{c}}(t) = (\dot{\mathbf{c}}(t), \mathbf{c}'(t), \mathbf{c}''(t)) = \mathcal{I}_{\dot{\mathbf{c}}(t)}(\mathbf{c}'(t), \mathbf{c}''(t)) \in (TM)_{\dot{\mathbf{c}}(t)}.$$

As usual, we denote by  $\mathbf{c}_v$  the geodesic with initial tangent vector  $\mathbf{v}$ . Define a map  $S : TM \rightarrow T(TM)$  by  $S(\mathbf{v}) = \ddot{\mathbf{c}}_v(0)$ .

**Theorem 3.7.2.** (1) *The map  $S$  is differentiable; i.e.,  $S$  is a vector field on  $TM$ , called the geodesic spray of  $M$ .*

(2)  *$\mathbf{c}$  is a geodesic in  $M$  if and only if  $\dot{\mathbf{c}}$  is an integral curve of  $S$ .*

*Proof.* We first show that  $S$  is differentiable at any  $\mathbf{p} \in M$ . Denote by  $U$  a neighborhood of  $\mathbf{p}$  on which there exists an orthonormal basis  $\mathbf{N}_1, \dots, \mathbf{N}_l$  of unit normal fields to  $M$ . We claim that the restriction of  $S$  to  $TU$  is given by

$$\mathbf{v} \mapsto \mathcal{I}_{\mathbf{v}} \left( \mathcal{I}_{\pi(\mathbf{v})}^{-1}\mathbf{v}, -\sum_i \langle \nabla_{\mathbf{v}}\mathbf{N}_i, \mathbf{v} \rangle \mathcal{I}_{\pi(\mathbf{v})}^{-1}\mathbf{N}_i \circ \pi(\mathbf{v}) \right), \quad (3.7.5)$$

so that  $S$  is smooth at  $\mathbf{p}$ . To see this, denote by  $\tilde{S}$  the local vector field determined by (3.7.5). If  $\gamma$  is an integral curve of  $\tilde{S}$ , and  $\mathbf{c} = \pi \circ \gamma$ , then

$$\dot{\gamma} = \tilde{S} \circ \gamma = \mathcal{I}_{\gamma} \left( \mathcal{I}_{\mathbf{c}}^{-1}\dot{\gamma}, -\sum_i \langle \nabla_{\dot{\gamma}}\mathbf{N}_i, \dot{\gamma} \rangle \mathcal{I}_{\mathbf{c}}^{-1}\mathbf{N}_i \circ \mathbf{c} \right). \quad (3.7.6)$$

On the other hand,  $\gamma = (\mathbf{c}, \pi_2 \circ \gamma)$ , so that

$$\dot{\gamma} = (\gamma, \mathbf{c}', (\pi_2 \circ \gamma)') = \mathcal{I}_\gamma(\mathbf{c}', (\pi_2 \circ \gamma)'). \quad (3.7.7)$$

Comparing with the first part of (3.7.6), we deduce that

$$\mathbf{c}' = \mathcal{I}_c^{-1} \dot{\gamma}, \text{ or equivalently, } \gamma = \mathcal{I}_c \mathbf{c}' = \dot{\mathbf{c}}. \quad (3.7.8)$$

The latter just says that  $\gamma$  is the velocity field of the curve  $\mathbf{c}$  in  $M$ . Comparing the second part of the expressions for  $\dot{\gamma}$  in (3.7.6) and (3.7.7) yields

$$(\pi_2 \circ \gamma)' = - \sum_i \langle \nabla_\gamma \mathbf{N}_i, \gamma \rangle \mathcal{I}_c^{-1} \mathbf{N}_i \circ \mathbf{c}.$$

Now, by (3.7.4) and (3.7.8),

$$(\pi_2 \circ \gamma)' = (\mathcal{I}_c^{-1} \circ \gamma)' = (\mathcal{I}_c^{-1} \circ \dot{\mathbf{c}})' = \mathbf{c}'',$$

so that

$$\dot{\mathbf{c}}' = \mathcal{I}_c \mathbf{c}'' = - \sum_i \langle \nabla_\gamma \mathbf{N}_i, \gamma \rangle \mathbf{N}_i \circ \mathbf{c} = - \sum_i \langle \nabla_{\dot{\mathbf{c}}} \mathbf{N}_i, \dot{\mathbf{c}} \rangle \mathbf{N}_i \circ \mathbf{c}.$$

This is just the geodesic equation (3.7.2) for  $\mathbf{c}$ .

Summarizing, integral curves of  $\tilde{S}$  are velocity fields of geodesics of  $M$ . Thus, for any  $\mathbf{v} \in TU$ ,  $\tilde{S}(\mathbf{v}) = \tilde{S}(\dot{\mathbf{c}}_\mathbf{v}(0)) = \ddot{\mathbf{c}}_\mathbf{v}(0)$ , and  $\tilde{S}$  is the restriction of  $S$  to  $TU$ , as claimed. This shows that  $S$  is differentiable. The second statement of the theorem was established in the course of proving the first one.  $\square$

The concept of geodesic spray provides with yet a third way of describing geodesics, one which is both useful and important: we shall construct a differentiable map  $\exp$  from an open subset of  $TM$  to  $M$  that for each  $\mathbf{v} \in TM$  maps the ray  $t \mapsto t\mathbf{v}$  in the tangent space  $M_{\pi(\mathbf{v})}$  to the geodesic  $\mathbf{c}_\mathbf{v}$  in  $M$ . For this, recall that by Theorem 3.3.2, there exists an open subset  $W$  of  $\mathbb{R} \times TM$  containing  $\{0\} \times TM$ , and a differentiable map  $\Psi : W \rightarrow TM$  (namely, the flow of the geodesic spray vector field) such that for any  $\mathbf{u} \in TM$ , if  $\Psi_\mathbf{u}$  denotes the curve  $t \mapsto \Psi(t, \mathbf{u})$ , then the geodesic  $\mathbf{c}_\mathbf{u}$  in direction  $\mathbf{u}$  equals  $\pi \circ \Psi_\mathbf{u}$ . Let

$$\widetilde{TM} = \{\mathbf{u} \in TM \mid (1, \mathbf{u}) \in W\}.$$

Notice that  $\widetilde{TM}$  is open in  $TM$  (because  $W$  is), and contains a neighborhood of the zero section of  $M$  (the zero section is defined as the image of the zero vector field  $\mathbf{Z}$ ,  $\mathbf{Z}(\mathbf{p}) = \mathbf{0} \in M_\mathbf{p}$  for  $\mathbf{p} \in M$ ).

**Definition 3.7.2.** The exponential map of  $M$  is the map  $\exp : \widetilde{TM} \rightarrow M$  given by

$$\exp(\mathbf{u}) = (\pi \circ \Psi)(1, \mathbf{u}) = \mathbf{c}_\mathbf{u}(1).$$

For  $\mathbf{p} \in M$ , let  $\tilde{M}_\mathbf{p} = \widetilde{TM} \cap M_\mathbf{p}$ , and denote by  $\exp_\mathbf{p}$  the restriction of  $\exp$  to  $\tilde{M}_\mathbf{p}$ . Recall also the isomorphism  $\mathcal{I}_\mathbf{u}$  between Euclidean space and its tangent space at any point

$\mathbf{u}$ . It can be written as  $\mathcal{I}_{\mathbf{u}}(\mathbf{v}) = \dot{\gamma}_{\mathbf{v}}(0)$ , where  $\gamma_{\mathbf{v}}(t) = \mathbf{u} + t\mathbf{v}$ . For a manifold  $M$ , given  $\mathbf{p} \in M$ , the tangent space  $M_{\mathbf{p}}$  is both a vector space and a manifold, so  $M_{\mathbf{p}}$  also has a tangent space at any  $\mathbf{u} \in M_{\mathbf{p}}$ . In the same vein, we introduce a canonical isomorphism  $\mathcal{I}_{\mathbf{u}}$  between  $M_{\mathbf{p}}$  and its tangent space at  $\mathbf{u}$  by the formula above. The second statement in the following theorem then says that up to this isomorphism, the derivative  $\exp_{\mathbf{p}*0}$  of  $\exp_{\mathbf{p}}$  at the origin  $\mathbf{Z}(\mathbf{p}) = \mathbf{0}$  in  $M_{\mathbf{p}}$  is just the identity.

**Theorem 3.7.3.** (1)  $\mathbf{c}_{\mathbf{v}}(t) = \exp(t\mathbf{v})$  for any  $\mathbf{v} \in TM$ ;  
 (2)  $\exp$  is differentiable, and for  $\mathbf{p} \in M, \mathbf{v} \in M_{\mathbf{p}}, \exp_{\mathbf{p}*} \mathcal{I}_{\mathbf{Z}(\mathbf{p})}\mathbf{v} = \mathbf{v}$ .

*Proof.* For the first statement, notice that if  $t_0$  belongs to the domain of  $\mathbf{c}_{\mathbf{v}}$ , then the curve  $t \mapsto \gamma(t) := \mathbf{c}_{\mathbf{v}}(t_0 t)$  is a geodesic (being a linear reparametrization of  $\mathbf{c}_{\mathbf{v}}$ ) with initial velocity  $\dot{\gamma}(0) = t_0 \dot{\mathbf{c}}_{\mathbf{v}}(0) = t_0 \mathbf{v}$ , and by uniqueness,  $\gamma = \mathbf{c}_{t_0 \mathbf{v}}$ . Thus,

$$\Psi_{t_0 \mathbf{v}}(1) = \Psi_{\mathbf{v}}(t_0), \quad (t_0, \mathbf{v}) \in W,$$

so that

$$\exp(t_0 \mathbf{v}) = \pi \circ \Psi_{t_0 \mathbf{v}}(1) = \pi \circ \Psi_{\mathbf{v}}(t_0) = \mathbf{c}_{\mathbf{v}}(t_0)$$

for any  $t_0$  in the domain of  $\mathbf{c}_{\mathbf{v}}$ , thereby establishing (1). For (2), let

$$\begin{aligned} \iota_1 : \widetilde{TM} &\hookrightarrow W, \\ \mathbf{u} &\mapsto (1, \mathbf{u}). \end{aligned}$$

Then  $\exp$ , being a composition  $\pi \circ \Psi \circ \iota_1$  of differentiable maps, is also smooth. Finally, if  $\gamma_{\mathbf{v}}$  denotes the ray  $t \mapsto t\mathbf{v}$  in the tangent space  $M_{\pi(\mathbf{v})}$ , and  $\mathbf{0}$  the origin there, then  $\dot{\gamma}_{\mathbf{v}}(0) = \mathcal{I}_{\mathbf{0}}\mathbf{v}$ , so that

$$\mathbf{v} = \dot{\mathbf{c}}_{\mathbf{v}}(0) = \exp_{\mathbf{p}*}(\dot{\gamma}_{\mathbf{v}}(0)) = \exp_{\mathbf{p}*} \mathcal{I}_{\mathbf{0}}\mathbf{v}. \quad \square$$

In particular, for every  $\mathbf{p} \in M$ ,  $\exp_{\mathbf{p}}$  has maximal rank at the origin in  $M_{\mathbf{p}}$ . The inverse function theorem then implies:

**Corollary 3.7.1.** For any  $\mathbf{p} \in M$ , there exists a neighborhood  $U$  of the origin in  $M_{\mathbf{p}}$  such that  $\exp_{\mathbf{p}}$  maps  $U$  diffeomorphically onto a neighborhood of  $\mathbf{p}$  in  $M$ .

**Remark 3.7.1.** There is a stronger version of Corollary 3.7.1 that will be used later on. Consider the map  $(\pi, \exp) : \widetilde{TM} \rightarrow M \times M$ ,  $(\pi, \exp)(\mathbf{v}) = (\pi(\mathbf{v}), \exp(\mathbf{v}))$ , where  $\pi$  is the bundle projection. We have already seen that any local chart  $(U, \mathbf{x})$  of  $M$  around a point  $\mathbf{p} \in M$  induces a chart  $(\pi^{-1}(U), \bar{\mathbf{x}})$  of  $TM$  around  $\pi^{-1}(\mathbf{p})$ , where

$$\bar{\mathbf{x}} = (\mathbf{x}, dx^1, \dots, dx^n) \circ \pi,$$

cf. (3.5.3). It also induces one of  $M \times M$  around  $(\mathbf{p}, \mathbf{p})$ , namely  $\mathbf{x} \times \mathbf{x}$ . Now, if  $\mathbf{v} = (\mathbf{p}, \mathbf{u}) \in TM \subset M \times \mathbb{R}^{n+k}$ , then  $\pi(\mathbf{p}, \mathbf{u}) = \mathbf{p}$ ; i.e.,  $\pi$  is projection onto the first factor. This means that the matrix of the derivative of  $(\pi, \exp)$  in the corresponding basis of coordinate



vector fields at any point of the zero section has the same top  $n$  rows as the  $2n \times 2n$  identity matrix  $I_{2n}$ . Furthermore, we claim that

$$\frac{\partial}{\partial \bar{x}^{n+j}}(\mathbf{0}_p) = (\iota_* \circ \mathcal{I}_{\mathbf{0}_p}) \frac{\partial}{\partial x^j}(\mathbf{p}),$$

where  $\iota : M_p \hookrightarrow TM$  is inclusion and  $\mathbf{0}_p$  denotes the zero vector in  $M_p$ . If we accept this claim for the moment, then Theorem 3.7.3(2) implies that

$$\exp_* \frac{\partial}{\partial \bar{x}^{n+j}}(\mathbf{0}_p) = (\exp_* \circ \iota_* \circ \mathcal{I}_{\mathbf{0}_p}) \frac{\partial}{\partial x^j}(\mathbf{p}) = \exp_{p_*} \mathcal{I}_{\mathbf{0}_p} \frac{\partial}{\partial x^j}(\mathbf{p}) = \frac{\partial}{\partial x^j}(\mathbf{p}),$$

so that the matrix of  $(\pi, \exp)_*$  has the form

$$\begin{bmatrix} I_n & 0 \\ * & I_n \end{bmatrix}.$$

Here all four entries are  $n \times n$  matrices, with  $I_n$  the identity and  $0$  the zero matrix.

To verify the claim, consider the ray  $\gamma_j : \mathbb{R} \rightarrow M_p$  with  $\gamma_j(t) = t \partial/\partial x^j(\mathbf{p})$ . Then

$$(\iota \circ \gamma_j)_* \mathbf{D}(0) = \sum_{i=1}^{2n} D(0)(\bar{x}^i \circ \iota \circ \gamma_j) \frac{\partial}{\partial \bar{x}^i}(\mathbf{0}_p).$$

For  $i \leq n$ ,  $\bar{x}^i \circ \iota \circ \gamma_j$  is the constant map  $x^i(\mathbf{p})$ , whereas

$$(\bar{x}^{n+i} \circ \iota \circ \gamma_j)(t) = dx^i(t \frac{\partial}{\partial x^j}(\mathbf{p})) = t$$

when  $i = j$  and  $0$  otherwise. Thus,  $(\iota \circ \gamma_j)_* \mathbf{D}(0) = (\partial/\partial \bar{x}^{n+j})(\mathbf{0}_p)$ . The claim now follows, since by definition,  $(\iota \circ \gamma_j)_* \mathbf{D}(0) = (\iota_* \circ \mathcal{I}_{\mathbf{0}_p}) \partial/\partial x^j(\mathbf{p})$ .

In conclusion, the derivative has maximal rank at each point of the zero section  $\mathbf{Z}$ , so that  $(\pi, \exp)$  is a diffeomorphism in a neighborhood of each  $\mathbf{Z}(\mathbf{p})$ ,  $\mathbf{p} \in M$ . Furthermore, the map itself is one-to-one when restricted to the zero section, since  $(\pi, \exp)\mathbf{Z}(\mathbf{p}) = (\mathbf{p}, \mathbf{p})$ . It can be shown that under these conditions,  $(\pi, \exp)$  is actually one-to-one on a neighborhood of the zero section. The proof of this last fact involves technical details outside the scope of this text. We nevertheless summarize this for future use:  $(\pi, \exp) : \widetilde{TM} \rightarrow M \times M$  maps a neighborhood of the zero section in  $TM$  diffeomorphically onto a neighborhood of the diagonal  $\Delta = \{(\mathbf{p}, \mathbf{p}) \mid \mathbf{p} \in M\}$  in  $M \times M$ .

### 3.8 The second fundamental tensor

Recall that given vector spaces  $V$  and  $W$ , a map  $M : V^k \rightarrow W$  is said to be multilinear if it is linear in each component; i.e. if for any  $(k-1)$ -tuple  $(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \in V^{k-1}$  and any  $i \in \{1, \dots, k\}$ , the map  $T \circ \iota_{i, \mathbf{v}_1, \dots, \mathbf{v}_{k-1}} : V \rightarrow W$  is linear, where

$$\begin{aligned} \iota_{i, \mathbf{v}_1, \dots, \mathbf{v}_{k-1}} : V &\rightarrow V^k, \\ \mathbf{u} &\mapsto (\mathbf{v}_1, \dots, \mathbf{v}_{i-1}, \mathbf{u}, \mathbf{v}_i, \dots, \mathbf{v}_{k-1}). \end{aligned}$$

We adopt the convention that  $V^0 = \mathbb{R}$ . For  $l = 0$  or  $1$ , a multilinear map  $T : V^k \rightarrow V^l$  is called a *tensor of type  $(k, l)$*  on  $V$ .

**Definition 3.8.1.** Let  $l = 0$  or  $1$ . A *tensor field of type  $(k, l)$*  on a manifold  $M$  is a map  $T$  such that for each  $\mathbf{p} \in M$ ,  $T(\mathbf{p}) : (M_{\mathbf{p}})^k \rightarrow (M_{\mathbf{p}})^l$  is a tensor of type  $(k, l)$  on the tangent space of  $M$  at  $\mathbf{p}$ .  $T$  is furthermore assumed to be smooth in the sense that for vector fields  $\mathbf{X}_1, \dots, \mathbf{X}_k$  on  $M$  the function (if  $l = 0$ ) or the vector field (if  $l = 1$ )  $T(\mathbf{X}_1, \dots, \mathbf{X}_k)$ , given by

$$T(\mathbf{X}_1, \dots, \mathbf{X}_k)(\mathbf{p}) = T(\mathbf{p})(\mathbf{X}_1(\mathbf{p}), \dots, \mathbf{X}_k(\mathbf{p})), \quad (3.8.1)$$

is smooth in the usual sense.

For example, on any manifold  $M$ , the map  $\mathbf{p} \mapsto g(\mathbf{p})$ , where  $g(\mathbf{p})$  denotes the inner product on the tangent space of  $M$  at  $\mathbf{p}$ , is a tensor field of type  $(2, 0)$ , called the *first fundamental tensor field of  $M$* .

We will denote by  $\mathfrak{X}M$  the vector space of vector fields on  $M$ , and by  $(\mathfrak{X}M)^0$  the vector space of all functions  $f : M \rightarrow \mathbb{R}$ . Thus, any tensor field  $T$  of type  $(k, l)$  defines a multilinear map  $T : (\mathfrak{X}M)^k \rightarrow (\mathfrak{X}M)^l$  via (3.8.1). This map is actually linear over functions: given  $f : M \rightarrow \mathbb{R}$ ,

$$T(\mathbf{X}_1, \dots, f\mathbf{X}_i, \dots, \mathbf{X}_k) = fT(\mathbf{X}_1, \dots, \mathbf{X}_i, \dots, \mathbf{X}_k). \quad (3.8.2)$$

It turns out that this property characterizes tensor fields:

**Theorem 3.8.1.** A multilinear map  $T : (\mathfrak{X}M)^k \rightarrow (\mathfrak{X}M)^l$  is a tensor field if and only if it is linear over functions; i.e., if and only if it satisfies (3.8.2).

*Proof.* The condition is necessary by the above remark. Conversely, suppose  $T$  is multilinear and linear over functions. We claim that if  $\mathbf{X}_i(\mathbf{p}) = \mathbf{Y}_i(\mathbf{p})$  for all  $i$ , then  $T(\mathbf{X}_1, \dots, \mathbf{X}_k)(\mathbf{p}) = T(\mathbf{Y}_1, \dots, \mathbf{Y}_k)(\mathbf{p})$ . To see this, assume for simplicity that  $k = 1$  and  $l = 0$ , the general case being analogous. It suffices to establish that if  $\mathbf{X}(\mathbf{p}) = \mathbf{0}$ , then  $T(\mathbf{X})(\mathbf{p}) = 0$ : indeed, if  $\mathbf{X}$  and  $\mathbf{Y}$  are as above, then  $\mathbf{X} - \mathbf{Y}$  is a vector field that vanishes at  $\mathbf{p}$ , so that  $T(\mathbf{X})(\mathbf{p}) - T(\mathbf{Y})(\mathbf{p}) = T(\mathbf{X} - \mathbf{Y})(\mathbf{p}) = 0$ .

So consider a chart  $(U, \mathbf{x})$  around  $\mathbf{p}$ , and write  $\mathbf{X}_{|U} = \sum f^i \partial/\partial x^i$  with  $f^i(\mathbf{p}) = 0$ . By Proposition 3.2.2, there exists a function  $\varphi$  on  $M$  with support in  $U$ , that equals one on a neighborhood of  $\mathbf{p}$  and takes values in  $[0, 1]$ . Define vector fields  $\mathbf{X}_i$  on  $M$  by setting them equal to  $\varphi \partial/\partial x^i$  on  $U$  and to  $0$  outside  $U$ . Each  $\mathbf{X}_i$  is differentiable on  $M$ : its restriction to  $U$  is given by a differentiable vector field, and if  $\mathbf{q}$  lies outside  $U$ , there is an open neighborhood of  $\mathbf{q}$  on which  $\mathbf{X}_i$  is identically zero, implying that  $\mathbf{X}_i$  is smooth at  $\mathbf{q}$ . Similarly, let  $g^i$  be the functions that equal  $\varphi f^i$  on  $U$  and  $0$  outside  $U$ . Notice that  $\sum g^i \mathbf{X}_i = \varphi^2 \mathbf{X}$ , because  $\sum_i g^i \mathbf{X}_i = \sum_i g^i \varphi \partial/\partial x^i = \varphi^2 \sum_i f^i \partial/\partial x^i = \varphi^2 \mathbf{X}$  on  $U$ , and both fields are zero outside  $U$ . Thus,  $\mathbf{X} = \varphi^2 \mathbf{X} + (1 - \varphi^2) \mathbf{X} = \sum_i g^i \mathbf{X}_i + (1 - \varphi^2) \mathbf{X}$ . Recalling that  $g^i(\mathbf{p}) = 0$  and  $\varphi(\mathbf{p}) = 1$ , we conclude that

$$(T\mathbf{X})(\mathbf{p}) = \left( \sum_i g^i(\mathbf{p})(T\mathbf{X}_i)(\mathbf{p}) \right) + (1 - \varphi^2(\mathbf{p}))(T\mathbf{X})(\mathbf{p}) = 0,$$

establishing the claim. We may therefore define for each  $\mathbf{p} \in M$  a multilinear map  $T(\mathbf{p})$  by

$$T(\mathbf{p})(\mathbf{u}_1, \dots, \mathbf{u}_k) := T(\mathbf{X}_1, \dots, \mathbf{X}_k)(\mathbf{p})$$

for any vector fields  $\mathbf{X}_i$  with  $\mathbf{X}_i(\mathbf{p}) = \mathbf{u}_i$ . The map  $\mathbf{p} \mapsto T(\mathbf{p})$  is clearly smooth.  $\square$

Let  $M^n \subset \mathbb{R}^{n+k}$  be a submanifold of Euclidean space, and consider the map  $S : \mathfrak{X}(M) \times \mathfrak{X}(M) \rightarrow \mathfrak{X}(\mathbb{R}^{n+k})|_M$  from the Cartesian product of the space of vector fields on  $M$  with itself to the space of vector fields on  $\mathbb{R}^{n+k}$  restricted to  $M$  given by

$$S(\mathbf{X}, \mathbf{Y}) = (D_{\mathbf{X}}\mathbf{Y})^\perp = D_{\mathbf{X}}\mathbf{Y} - \nabla_{\mathbf{X}}\mathbf{Y}. \quad (3.8.3)$$

$S$  is bilinear, and linear over functions in the first component by properties of the covariant derivative. Furthermore,  $S$  is symmetric, because

$$S(\mathbf{X}, \mathbf{Y}) - S(\mathbf{Y}, \mathbf{X}) = (D_{\mathbf{X}}\mathbf{Y} - D_{\mathbf{Y}}\mathbf{X})^\perp = [\mathbf{X}, \mathbf{Y}]^\perp = 0,$$

since the bracket of fields tangent to  $M$  is a field tangent to  $M$  by Theorem 2.9.3. Thus, it is also linear over functions in the second component, and in light of Theorem 3.8.1, the formula

$$S(\mathbf{x}, \mathbf{y}) = (D_{\mathbf{X}}\mathbf{Y})^\perp(\mathbf{p}), \quad \mathbf{x}, \mathbf{y} \in M_{\mathbf{p}}, \quad \mathbf{p} \in M, \quad (3.8.4)$$

where  $\mathbf{X}$  and  $\mathbf{Y}$  are any vector fields on  $M$  that equal  $\mathbf{x}$  and  $\mathbf{y}$  at  $\mathbf{p}$ , defines a tensor field on  $M$ , called the *second fundamental tensor of  $M$* . It is worth emphasizing that the fields  $\mathbf{X}$  and  $\mathbf{Y}$  only need to be defined in a neighborhood of the point being considered: just as in the proof of Theorem 3.8.1, if  $U$  is the domain of a chart around  $\mathbf{p}$ ,  $\mathbf{X}$  is a vector field on  $U$ , and  $\varphi$  a nonnegative function with support in  $U$  that equals 1 on a neighborhood of  $\mathbf{p}$ , then the field  $\tilde{\mathbf{X}}$  that equals  $\varphi\mathbf{X}$  on  $U$  and zero outside  $U$  is a vector field on  $M$  that equals  $\mathbf{X}$  on a neighborhood of  $\mathbf{p}$ .

**Definition 3.8.2.** Let  $\mathbf{u}$  denote a vector orthogonal to  $M_{\mathbf{p}}$ . The *second fundamental tensor with respect to  $\mathbf{u}$*  is the self-adjoint linear transformation  $S_{\mathbf{u}} : M_{\mathbf{p}} \rightarrow M_{\mathbf{p}}$  given by

$$S_{\mathbf{u}}\mathbf{x} = -(D_{\mathbf{x}}\mathbf{U})^\top, \quad \mathbf{x} \in M_{\mathbf{p}},$$

where  $\mathbf{U}$  is any locally defined normal vector field to  $M$  with  $\mathbf{U}(\mathbf{p}) = \mathbf{u}$  (here, normal to  $M$  means that  $\mathbf{U}(\mathbf{q}) \perp M_{\mathbf{q}}$  for all  $\mathbf{q}$  in the domain of  $\mathbf{U}$ ).

The above does not depend on the particular extension  $\mathbf{U}$  because it is tensorial in  $\mathbf{U}$ : If  $f : M \rightarrow \mathbb{R}$  is a function, then

$$(D_{\mathbf{v}}(f\mathbf{U}))^\top = (\mathbf{v}f)\mathbf{U}^\top + f(\mathbf{p})(D_{\mathbf{v}}\mathbf{U})^\top = f(\mathbf{p})(D_{\mathbf{v}}\mathbf{U})^\top,$$

since  $\mathbf{U}$  is normal to  $M$ , so that  $\mathbf{U}^\top$  is zero. Another way of seeing this is to consider the associated symmetric (as we will see in a moment) bilinear form on  $M_{\mathbf{p}}$  given by

$$(\mathbf{x}, \mathbf{y}) \mapsto \langle S_{\mathbf{u}}\mathbf{x}, \mathbf{y} \rangle, \quad \mathbf{x}, \mathbf{y} \in M_{\mathbf{p}}, \quad (3.8.5)$$

which is called the *second fundamental form with respect to  $\mathbf{u}$* . If  $\mathbf{Y}$  is a local extension of  $\mathbf{y}$ , then

$$\langle S_{\mathbf{u}}\mathbf{x}, \mathbf{y} \rangle = -\langle D_{\mathbf{x}}\mathbf{U}, \mathbf{y} \rangle = -D_{\mathbf{x}}\langle \mathbf{U}, \mathbf{Y} \rangle + \langle \mathbf{U}(\mathbf{p}), D_{\mathbf{x}}\mathbf{Y} \rangle = \langle \mathbf{u}, S(\mathbf{x}, \mathbf{y}) \rangle,$$

which, once again, only depends on the value of  $\mathbf{U}$  at the point  $\mathbf{p}$ . This also shows that the second fundamental form is symmetric as claimed.

**Example 3.8.1.** Consider the sphere  $M = S^n(r) \subset \mathbb{R}^{n+1}$  of radius  $r$ . One of the two unit normal fields is  $\mathbf{N} = (1/r)\mathbf{P}$ , where  $\mathbf{P}$  is the position vector field. The second fundamental tensor at  $\mathbf{p}$  is given by

$$\begin{aligned} S(\mathbf{x}, \mathbf{y}) &= (D_{\mathbf{x}}\mathbf{Y})^\perp = \langle D_{\mathbf{x}}\mathbf{Y}, \frac{1}{r}\mathbf{P}(\mathbf{p}) \rangle \frac{1}{r}\mathbf{P}(\mathbf{p}) \\ &= (\mathbf{x}\langle \mathbf{Y}, \mathbf{P} \rangle - \langle \mathbf{Y}(\mathbf{p}), D_{\mathbf{x}}\mathbf{P} \rangle) \frac{1}{r^2}\mathbf{P}(\mathbf{p}) \\ &= -\langle \mathbf{y}, D_{\mathbf{x}}\mathbf{P} \rangle \frac{1}{r^2}\mathbf{P}(\mathbf{p}), \end{aligned}$$

so that by Examples 2.8.2 (ii),

$$S(\mathbf{x}, \mathbf{y}) = -\frac{1}{r}\langle \mathbf{x}, \mathbf{y} \rangle \mathbf{n}, \quad (3.8.6)$$

where  $\mathbf{n}$  is the outward-pointing unit normal  $\mathbf{N}(\mathbf{p})$ . Thus, for any  $\mathbf{y} \in M_{\mathbf{p}}$ ,

$$\langle S_{\mathbf{n}}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{n}, S(\mathbf{x}, \mathbf{y}) \rangle = -\frac{1}{r}\langle \mathbf{x}, \mathbf{y} \rangle,$$

which implies

$$S_{\mathbf{n}}\mathbf{x} = -\frac{1}{r}\mathbf{x}, \quad \mathbf{x} \in M_{\mathbf{p}}. \quad (3.8.7)$$

### 3.9 Curvature

**Definition 3.9.1.** The *curvature tensor* of  $M$  is the tensor field  $R$  of type (3,1) given by

$$R(\mathbf{v}, \mathbf{w})\mathbf{z} = S_{S(\mathbf{w}, \mathbf{z})}\mathbf{v} - S_{S(\mathbf{v}, \mathbf{z})}\mathbf{w}, \quad \mathbf{v}, \mathbf{w}, \mathbf{z} \in M_{\mathbf{p}}, \quad \mathbf{p} \in M,$$

where  $S$  denotes the second fundamental tensor of  $M$ .

Since  $S$  is a tensor field, it is clear that  $R$  is one also.

**Example 3.9.1** (Curvature tensor of a sphere). Let  $M = S^n(r)$  denote the sphere of radius  $r$  centered at the origin. The second fundamental tensor  $S$  of  $M$  is given by (3.8.6), and for the outer unit normal vector  $\mathbf{n}$ ,  $S_{\mathbf{n}}\mathbf{v} = -\frac{1}{r}\mathbf{v}$ . Thus,

$$\begin{aligned} R(\mathbf{v}, \mathbf{w})\mathbf{z} &= S_{S(\mathbf{w}, \mathbf{z})}\mathbf{v} - S_{S(\mathbf{v}, \mathbf{z})}\mathbf{w} = S_{-\frac{1}{r}\langle \mathbf{w}, \mathbf{z} \rangle \mathbf{n}}\mathbf{v} - S_{-\frac{1}{r}\langle \mathbf{v}, \mathbf{z} \rangle \mathbf{n}}\mathbf{w} \\ &= \frac{1}{r^2} (\langle \mathbf{w}, \mathbf{z} \rangle \mathbf{v} - \langle \mathbf{v}, \mathbf{z} \rangle \mathbf{w}). \end{aligned}$$

There is another characterization of the curvature tensor that is often useful. First, though, observe that if  $\mathbf{X}$ ,  $\mathbf{Y}$ , and  $\mathbf{Z}$  are vector fields on an open set  $U \subset \mathbb{R}^n$ , then

$$(D_X D_Y - D_Y D_X - D_{[X, Y]})\mathbf{Z} = 0. \quad (3.9.1)$$

To see this, write  $\mathbf{Z} = \sum_i f^i \mathbf{D}_i$ , and recall that the coordinate vector fields  $\mathbf{D}_i$  are parallel. Thus,

$$\begin{aligned} (D_X D_Y - D_Y D_X - D_{[X, Y]})\mathbf{Z} &= \sum_i ((\mathbf{X}\mathbf{Y} - \mathbf{Y}\mathbf{X} - [\mathbf{X}, \mathbf{Y}])f^i) \mathbf{D}_i \\ &= \mathbf{0} \end{aligned}$$

by definition of the Lie bracket. In light of the following theorem, (3.9.1) just says that the curvature tensor of Euclidean space is identically zero; i.e., Euclidean space is *flat*.

**Theorem 3.9.1.** *For vector fields  $\mathbf{X}$ ,  $\mathbf{Y}$ ,  $\mathbf{Z}$  on a manifold  $M$ ,*

$$R(\mathbf{X}, \mathbf{Y})\mathbf{Z} = \nabla_X \nabla_Y \mathbf{Z} - \nabla_Y \nabla_X \mathbf{Z} - \nabla_{[X, Y]}\mathbf{Z}.$$

*Proof.* For the sake of clarity, we will sometimes write  $D_X^\perp \mathbf{Y}$  instead of  $(D_X \mathbf{Y})^\perp$ , and similarly for  $^\top$ . With this in mind, notice that

$$\nabla_X \nabla_Y \mathbf{Z} = \nabla_X (D_Y^\top \mathbf{Z}) = D_X^\top (D_Y \mathbf{Z} - S(\mathbf{Y}, \mathbf{Z})) = D_X^\top D_Y \mathbf{Z} + S_{S(\mathbf{Y}, \mathbf{Z})}\mathbf{X}.$$

Thus, the right side of the identity in the statement equals

$$\begin{aligned} &D_X^\top D_Y \mathbf{Z} + S_{S(\mathbf{Y}, \mathbf{Z})}\mathbf{X} - D_Y^\top D_X \mathbf{Z} - S_{S(\mathbf{X}, \mathbf{Z})}\mathbf{Y} - D_{[X, Y]}^\top \mathbf{Z} \\ &= (D_X D_Y - D_Y D_X - D_{[X, Y]})^\top \mathbf{Z} + S_{S(\mathbf{Y}, \mathbf{Z})}\mathbf{X} - S_{S(\mathbf{X}, \mathbf{Z})}\mathbf{Y} \\ &= S_{S(\mathbf{Y}, \mathbf{Z})}\mathbf{X} - S_{S(\mathbf{X}, \mathbf{Z})}\mathbf{Y} \\ &= R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \end{aligned}$$

as claimed. □

More generally, we have the following:

**Theorem 3.9.2.** *Let  $\mathbf{f} : N \rightarrow M$  be differentiable,  $\mathbf{U}$ ,  $\mathbf{V}$  vector fields on  $N$ , and  $\mathbf{X}$  a vector field along  $\mathbf{f}$ . Then*

$$R(\mathbf{f}_* \mathbf{U}, \mathbf{f}_* \mathbf{V})\mathbf{X} = \nabla_U \nabla_V \mathbf{X} - \nabla_V \nabla_U \mathbf{X} - \nabla_{[U, V]}\mathbf{X}.$$

*Proof.* The statement being a local one, we may work in the domain  $U$  of a chart on  $M$  with coordinate vector fields  $\mathbf{X}_i$ . The restriction of any vector field along  $\mathbf{f}$  to  $\mathbf{f}^{-1}(U)$  may be written as  $\sum_i h^i \mathbf{X}_i \circ \mathbf{f}$  for some functions  $h^i : \mathbf{f}^{-1}(U) \rightarrow \mathbb{R}$ . Since both sides of the identity we are proving are tensorial, we need only consider vector fields of the form  $\mathbf{X}_i \circ \mathbf{f}$ ; i.e., we may assume that there exist vector fields  $\tilde{\mathbf{U}}$ ,  $\tilde{\mathbf{V}}$ , and  $\tilde{\mathbf{X}}$  on  $U$  such that  $\mathbf{f}_* \mathbf{U} = \tilde{\mathbf{U}} \circ \mathbf{f}$ ,  $\mathbf{f}_* \mathbf{V} = \tilde{\mathbf{V}} \circ \mathbf{f}$ , and  $\mathbf{X} = \tilde{\mathbf{X}} \circ \mathbf{f}$  on  $\mathbf{f}^{-1}(U)$ . Thus, on  $\mathbf{f}^{-1}(U)$ ,

$$\begin{aligned} R(\mathbf{f}_* \mathbf{U}, \mathbf{f}_* \mathbf{V})\mathbf{X} &= R(\tilde{\mathbf{U}} \circ \mathbf{f}, \tilde{\mathbf{V}} \circ \mathbf{f})\tilde{\mathbf{X}} \circ \mathbf{f} = [R(\tilde{\mathbf{U}}, \tilde{\mathbf{V}})\tilde{\mathbf{X}}] \circ \mathbf{f} \\ &= [\nabla_{\tilde{\mathbf{U}}}\nabla_{\tilde{\mathbf{V}}}\tilde{\mathbf{X}} - \nabla_{\tilde{\mathbf{V}}}\nabla_{\tilde{\mathbf{U}}}\tilde{\mathbf{X}} - \nabla_{[\tilde{\mathbf{U}}, \tilde{\mathbf{V}}]}\tilde{\mathbf{X}}] \circ \mathbf{f}. \end{aligned}$$

By repeated use of Theorem 3.6.2 (6), the first term on the above line may be rewritten as:

$$\begin{aligned} [\nabla_{\tilde{U}}\nabla_{\tilde{V}}\tilde{\mathbf{X}}] \circ \mathbf{f} &= \nabla_{\tilde{U}\circ\mathbf{f}}\nabla_{\tilde{V}}\tilde{\mathbf{X}} = \nabla_{\mathbf{f}_*\tilde{U}}\nabla_{\tilde{V}}\tilde{\mathbf{X}} = \nabla_{\tilde{U}}(\nabla_{\tilde{V}}\tilde{\mathbf{X}} \circ \mathbf{f}) \\ &= \nabla_{\tilde{U}}\nabla_{\mathbf{f}_*\tilde{V}}\tilde{\mathbf{X}} = \nabla_{\tilde{U}}\nabla_{\tilde{V}}(\tilde{\mathbf{X}} \circ \mathbf{f}) \\ &= \nabla_{\tilde{U}}\nabla_{\tilde{V}}\mathbf{X}, \end{aligned}$$

and a similar expression holds for the second term. Finally, the last term may, in view of Theorem 2.9.2, be rewritten as follows:

$$(\nabla_{[\tilde{U},\tilde{V}]} \tilde{\mathbf{X}}) \circ \mathbf{f} = \nabla_{[\tilde{U},\tilde{V}]\circ\mathbf{f}} \tilde{\mathbf{X}} = \nabla_{\mathbf{f}_*[\tilde{U},\tilde{V}]} \tilde{\mathbf{X}} = \nabla_{[\tilde{U},\tilde{V}]}(\tilde{\mathbf{X}} \circ \mathbf{f}) = \nabla_{[\tilde{U},\tilde{V}]} \mathbf{X}.$$

This establishes the result.  $\square$

**Proposition 3.9.1.** *The following identities hold for vector fields  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \mathbf{U}$  on  $M$ :*

- (1)  $R(\mathbf{X}, \mathbf{Y})\mathbf{Z} = -R(\mathbf{Y}, \mathbf{X})\mathbf{Z}$ ;
- (2)  $R(\mathbf{X}, \mathbf{Y})\mathbf{Z} + R(\mathbf{Y}, \mathbf{Z})\mathbf{X} + R(\mathbf{Z}, \mathbf{X})\mathbf{Y} = \mathbf{0}$ ;
- (3)  $\langle R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{U} \rangle = -\langle R(\mathbf{X}, \mathbf{Y})\mathbf{U}, \mathbf{Z} \rangle$ , and
- (4)  $\langle R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{U} \rangle = \langle R(\mathbf{Z}, \mathbf{U})\mathbf{X}, \mathbf{Y} \rangle$ .

*Proof.* (1) is an immediate consequence of the definition. Since  $R$  is a tensor field, and is therefore linear over functions, it is enough to prove the identities for coordinate vector fields, or more generally, for vector fields with vanishing Lie brackets. Thus,

$$\nabla_{\mathbf{X}}\mathbf{Y} = \nabla_{\mathbf{Y}}\mathbf{X}, \quad R(\mathbf{X}, \mathbf{Y})\mathbf{Z} = \nabla_{\mathbf{X}}\nabla_{\mathbf{Y}}\mathbf{Z} - \nabla_{\mathbf{Y}}\nabla_{\mathbf{X}}\mathbf{Z},$$

with similar identities holding for the other terms. (2) now easily follows. For (3), it is enough to show that  $\langle R(\mathbf{X}, \mathbf{Y})\mathbf{V}, \mathbf{V} \rangle = 0$  for any vector field  $\mathbf{V}$ , since this will imply that

$$\begin{aligned} 0 &= \langle R(\mathbf{X}, \mathbf{Y})(\mathbf{Z} + \mathbf{U}), \mathbf{Z} + \mathbf{U} \rangle \\ &= \langle R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{Z} \rangle + \langle R(\mathbf{X}, \mathbf{Y})\mathbf{U}, \mathbf{U} \rangle + \langle R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{U} \rangle + \langle R(\mathbf{X}, \mathbf{Y})\mathbf{U}, \mathbf{Z} \rangle \\ &= \langle R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{U} \rangle + \langle R(\mathbf{X}, \mathbf{Y})\mathbf{U}, \mathbf{Z} \rangle. \end{aligned}$$

Now,

$$\begin{aligned} \langle \nabla_{\mathbf{X}}\nabla_{\mathbf{Y}}\mathbf{Z}, \mathbf{Z} \rangle &= \mathbf{X}\langle \nabla_{\mathbf{Y}}\mathbf{Z}, \mathbf{Z} \rangle - \langle \nabla_{\mathbf{Y}}\mathbf{Z}, \nabla_{\mathbf{X}}\mathbf{Z} \rangle \\ &= \frac{1}{2}\mathbf{XY}\langle \mathbf{Z}, \mathbf{Z} \rangle - \langle \nabla_{\mathbf{Y}}\mathbf{Z}, \nabla_{\mathbf{X}}\mathbf{Z} \rangle, \end{aligned}$$

so that

$$\langle R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{Z} \rangle = \frac{1}{2}(\mathbf{XY} - \mathbf{YX})\langle \mathbf{Z}, \mathbf{Z} \rangle = \frac{1}{2}[\mathbf{X}, \mathbf{Y}]\langle \mathbf{Z}, \mathbf{Z} \rangle = \mathbf{0}$$

because  $[\mathbf{X}, \mathbf{Y}] = \mathbf{0}$ .

The last identity can be derived from the other three:

$$\begin{aligned} \langle R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{U} \rangle &= -\langle R(\mathbf{Y}, \mathbf{X})\mathbf{Z}, \mathbf{U} \rangle \text{ by (1)} \\ &= \langle R(\mathbf{X}, \mathbf{Z})\mathbf{Y}, \mathbf{U} \rangle + \langle R(\mathbf{Z}, \mathbf{Y})\mathbf{X}, \mathbf{U} \rangle \text{ by (2)}. \end{aligned} \tag{3.9.2}$$

Similarly,

$$\begin{aligned}\langle R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{U} \rangle &= -\langle R(\mathbf{X}, \mathbf{Y})\mathbf{U}, \mathbf{Z} \rangle \text{ by (3)} \\ &= \langle R(\mathbf{Y}, \mathbf{U})\mathbf{X}, \mathbf{Z} \rangle + \langle R(\mathbf{U}, \mathbf{X})\mathbf{Y}, \mathbf{Z} \rangle \text{ by (2)}.\end{aligned}\tag{3.9.3}$$

Adding (3.9.2) and (3.9.3) implies

$$\begin{aligned}\langle 2R(\mathbf{X}, \mathbf{Y})\mathbf{Z}, \mathbf{U} \rangle &= \langle R(\mathbf{X}, \mathbf{Z})\mathbf{Y}, \mathbf{U} \rangle + \langle R(\mathbf{Z}, \mathbf{Y})\mathbf{X}, \mathbf{U} \rangle + \langle R(\mathbf{Y}, \mathbf{U})\mathbf{X}, \mathbf{Z} \rangle \\ &\quad + \langle R(\mathbf{U}, \mathbf{X})\mathbf{Y}, \mathbf{Z} \rangle.\end{aligned}\tag{3.9.4}$$

Swapping  $\mathbf{X}$  with  $\mathbf{Z}$  and  $\mathbf{Y}$  with  $\mathbf{U}$  in (3.9.4) yields

$$\begin{aligned}\langle 2R(\mathbf{Z}, \mathbf{U})\mathbf{X}, \mathbf{Y} \rangle &= \langle R(\mathbf{Z}, \mathbf{X})\mathbf{U}, \mathbf{Y} \rangle + \langle R(\mathbf{X}, \mathbf{U})\mathbf{Z}, \mathbf{Y} \rangle + \langle R(\mathbf{U}, \mathbf{Y})\mathbf{Z}, \mathbf{X} \rangle \\ &\quad + \langle R(\mathbf{Y}, \mathbf{Z})\mathbf{U}, \mathbf{X} \rangle.\end{aligned}\tag{3.9.5}$$

Finally, applying (1) and (3) to each term on the right side of (3.9.5) yields the right side of (3.9.4). This proves (4).  $\square$

Let  $\mathbf{p} \in M$ ,  $E$  a 2-dimensional subspace of  $M_{\mathbf{p}}$ . If  $\mathbf{v}$  and  $\mathbf{w}$  form a basis of  $E$ , then  $|\mathbf{v}|^2|\mathbf{w}|^2 - \langle \mathbf{v}, \mathbf{w} \rangle^2 > 0$  by the Cauchy-Schwartz inequality and Exercise 1.20. The *sectional curvature* of  $E$  is the number

$$K_E = \frac{\langle R(\mathbf{v}, \mathbf{w})\mathbf{w}, \mathbf{v} \rangle}{|\mathbf{v}|^2|\mathbf{w}|^2 - \langle \mathbf{v}, \mathbf{w} \rangle^2}\tag{3.9.6}$$

To see that this definition does not depend on the particular basis, denote by  $k$  the *curvature form* of  $R$ ; i.e., the map  $k : \mathfrak{X}M \times \mathfrak{X}M \rightarrow \mathbb{R}$ , where  $k(\mathbf{X}, \mathbf{Y}) = \langle R(\mathbf{X}, \mathbf{Y})\mathbf{Y}, \mathbf{X} \rangle$ . Notice that if  $R_S$  denotes the curvature tensor of the unit sphere, then the corresponding curvature form  $k_S$  is precisely the denominator of (3.9.6) by Example 3.9.1, and

$$K_E = \frac{k(\mathbf{v}, \mathbf{w})}{k_S(\mathbf{v}, \mathbf{w})}.$$

Suppose  $\mathbf{u} = a_{11}\mathbf{v} + a_{12}\mathbf{w}$  and  $\mathbf{v} = a_{21}\mathbf{v} + a_{22}\mathbf{w}$  form another basis of  $E$ . A straightforward computation shows that

$$k(\mathbf{u}, \mathbf{v}) = (\det A)^2 k(\mathbf{v}, \mathbf{w}),$$

where

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

is the transition matrix between the two bases. Since this identity does not depend on the particular curvature tensor, it also holds for  $k_S$ , and we conclude that the sectional curvature is indeed well-defined. Furthermore, if  $\mathbf{v}$  and  $\mathbf{w}$  are orthonormal, then

$$K_E = \langle R(\mathbf{v}, \mathbf{w})\mathbf{w}, \mathbf{v} \rangle.$$

Thus, for example, the sectional curvature of any plane tangent to a sphere of radius  $r$  is  $1/r^2$ . We therefore say that the sphere of radius  $r$  has *constant curvature*  $1/r^2$ . More generally, a *space of constant curvature*  $\kappa \in \mathbb{R}$  is a space where every plane has sectional curvature  $\kappa$ . As remarked earlier, Euclidean space has constant zero curvature.

### 3.10 Sectional curvature and the length of small circles

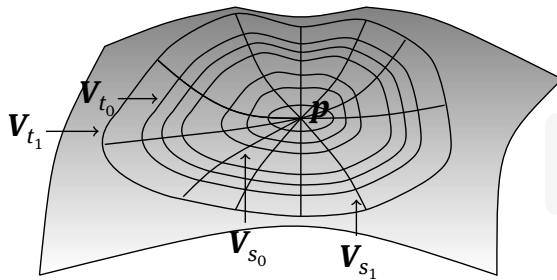
There is also a more geometric interpretation of the sectional curvature  $K_p$  of a plane  $P \subset M_p$ . Suppose  $\{\mathbf{v}, \mathbf{w}\}$  is an orthonormal basis of  $M_p$ , and consider the variation  $V : [0, a] \times [0, 2\pi] \rightarrow M$  by geodesics, where

$$V(t, s) = \exp_p[t(\cos s\mathbf{v} + \sin s\mathbf{w})].$$

When  $M$  is 2-dimensional, then for small  $t$ , the image of  $V_t$ , where  $V_t(s) = V(t, s)$ , is the set of points at distance  $t$  from  $p$ . By abuse of notation,  $V_t$  will be called the *circle of radius  $t$* , and  $V_s$ , where  $V_s(t) = V(t, s)$ , a “radial” geodesic. The length of the circle of radius  $t$  is

$$L(t) = \int_0^{2\pi} |\dot{V}_t(s)| ds.$$

$\dot{V}_t(s)$  is also  $Y_s(t)$ , where  $Y_s$  is the vector field  $Y_s(t) = V_*D_2(t, s)$  along the geodesic  $V_s$ .



Radial geodesics  $V_s$  and “circles”  $V_t$  of points at constant distance from  $p$

In order to estimate  $|Y_s|$ , we will need the following

**Lemma 3.10.1.** *Let  $Y$  be a vector field along a curve  $c$ . If  $E_i$  is the parallel vector field along  $c$  that equals  $Y^{(i)}(0)$  when  $t = 0$ , then*

$$Y(t) = \sum_{i=0}^k \frac{t^i}{i!} E_i(t) + Z(t),$$

where  $Z$  is a vector field along  $c$  satisfying  $|Z(t)|/t^k \rightarrow 0$  as  $t \rightarrow 0$ .

*Proof.* Choose a basis  $F_1, \dots, F_n$  of orthonormal parallel fields along  $c$ , so that  $Y = \sum f_l F_l$ , where  $f_l = \langle Y, F_l \rangle$ . By Taylor’s theorem,

$$f_l(t) = \sum_{i=1}^k \frac{t^i}{i!} f_l^{(i)}(0) + o(t^k),$$

with  $o(t^k)$  designating some expression that, when divided by  $t^k$ , goes to zero as  $t \rightarrow 0$ . Thus,

$$Y(t) = \sum_{i,l} \left( \frac{t^i}{i!} f_l^{(i)}(0) + o(t^k) \right) F_l(t).$$

The claim now follows, since  $E_i = \sum_l f_l^{(i)}(0) F_l$ . □



We will use the Lemma with  $k = 3$ : by definition,  $\mathbf{Y}_s(0) = \mathbf{0}$ , and

$$\mathbf{Y}'_s(0) = \nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_2(0, s) = \nabla_{\mathbf{D}_2} \mathbf{V}_* \mathbf{D}_1(0, s) = -\sin s \mathbf{v} + \cos s \mathbf{w},$$

since  $\mathbf{V}_* \mathbf{D}_1(0, s) = \cos s \mathbf{v} + \sin s \mathbf{w}$ . In order to compute  $\mathbf{Y}'_s(0)$ , notice that by Theorem 3.9.2,

$$\begin{aligned} R(\mathbf{V}_* \mathbf{D}_2, \mathbf{V}_* \mathbf{D}_1) \mathbf{V}_* \mathbf{D}_1 &= \nabla_{\mathbf{D}_2} \nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_1 - \nabla_{\mathbf{D}_1} \nabla_{\mathbf{D}_2} \mathbf{V}_* \mathbf{D}_1 \\ &\quad - \nabla_{[\mathbf{D}_2, \mathbf{D}_1]} \mathbf{V}_* \mathbf{D}_1 \\ &= -\nabla_{\mathbf{D}_1} \nabla_{\mathbf{D}_2} \mathbf{V}_* \mathbf{D}_1 = -\nabla_{\mathbf{D}_1} \nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_2. \end{aligned}$$

Thus,  $\mathbf{Y}''_s = -R(\mathbf{Y}_s, \dot{\mathbf{V}}_s) \dot{\mathbf{V}}_s$ , and since  $\mathbf{Y}_s(0) = \mathbf{0}$ ,  $\mathbf{Y}''_s(0) = \mathbf{0}$  also. It remains to evaluate  $\mathbf{Y}^{(3)} = -(R(\mathbf{Y}_s, \dot{\mathbf{V}}_s) \dot{\mathbf{V}}_s)'$  at 0. To do so, consider a basis  $\mathbf{E}_i$  of parallel fields along  $\mathbf{V}_s$  and write  $\mathbf{Y}_s = \sum f_i \mathbf{E}_i$  in terms of this basis. Then

$$\begin{aligned} \mathbf{Y}^{(3)}(0) &= -\sum (f_i R(\mathbf{E}_i, \dot{\mathbf{V}}_s) \dot{\mathbf{V}}_s)'(0) \\ &= -\sum_i \left( f_i(0) (R(\mathbf{E}_i, \dot{\mathbf{V}}_s) \dot{\mathbf{V}}_s)'(0) + f_i'(0) (R(\mathbf{E}_i, \dot{\mathbf{V}}_s) \dot{\mathbf{V}}_s)(0) \right) \\ &= -\sum f_i'(0) R(\mathbf{E}_i, \dot{\mathbf{V}}_s) \dot{\mathbf{V}}_s(0) = -R(\mathbf{Y}'_s(0), \dot{\mathbf{V}}_s(0)) \dot{\mathbf{V}}_s(0) \end{aligned}$$

since  $f_i(0) = 0$ . Lemma 3.10.1 now implies that  $\mathbf{Y}_s$  may be written as

$$\mathbf{Y}_s(t) = t \mathbf{E}_1(t) + \frac{t^3}{3!} \mathbf{E}_3(t) + \mathbf{Z}(t),$$

where  $\mathbf{E}_1(0) = \mathbf{w}_s := -\sin s \mathbf{v} + \cos s \mathbf{w}$ ,  $\mathbf{E}_3(0) = -R(\mathbf{w}_s, \mathbf{v}_s) \mathbf{v}_s$ , and  $\mathbf{v}_s := \cos s \mathbf{v} + \sin s \mathbf{w}$ . Furthermore,  $|\mathbf{Z}(t)|/t^3 \rightarrow 0$  as  $t \rightarrow 0$ . It follows that

$$\begin{aligned} |\mathbf{Y}_s|^2(t) &= t^2 - \frac{2t^4}{3!} K_P + \left( \frac{t^3}{3!} \right)^2 |R(\mathbf{w}_s, \mathbf{v}_s) \mathbf{v}_s|^2 + |\mathbf{Z}|^2(t) + 2t \langle \mathbf{E}_1, \mathbf{Z} \rangle(t) \\ &\quad + \frac{2t^3}{3!} \langle \mathbf{E}_3, \mathbf{Z} \rangle(t) \\ &= t^2 \left( 1 - \frac{2t^2}{3!} K_P - g(t) \right), \end{aligned}$$

where  $g$  is a function such that  $g(t)/t^2 \rightarrow 0$  as  $t \rightarrow 0$ . Since  $(1-x)^{1/2} = 1 - (x/2) + o(x)$ , we obtain for  $t > 0$

$$\begin{aligned} |\mathbf{Y}_s(t)| &= t \left( 1 - \frac{2t^2}{3!} K_P - g(t) \right)^{1/2} = t \left( 1 - \frac{t^2}{3!} K_P + o(t^2) \right) \\ &= t - \frac{t^3}{3!} K_P + o(t^3). \end{aligned} \tag{3.10.1}$$

We are now able to state our main result, which roughly says that  $M$  has positive (resp. negative) sectional curvature at  $\mathbf{p} \in M$  if and only if small circles centered at  $\mathbf{p}$  have smaller (resp. larger) circumference than circles of the same radius in Euclidean space:

**Theorem 3.10.1.** Let  $M$  be a manifold,  $\mathbf{p}$  a point in  $M$ , and  $P$  a plane in  $M_{\mathbf{p}}$ . For  $t > 0$ , denote by  $C_t$  the “circle” in  $M$  obtained by exponentiating the circle of radius  $t$  centered at the origin in  $P$ . Then the length of  $C_t$  is

$$L(C_t) = 2\pi t \left( 1 - \frac{K_P}{3!} t^2 + o(t^2) \right).$$

In particular, the sectional curvature of  $P$  is

$$K_P = \lim_{t \rightarrow 0} \frac{6}{t^2} \left( 1 - \frac{L(C_t)}{2\pi t} \right).$$

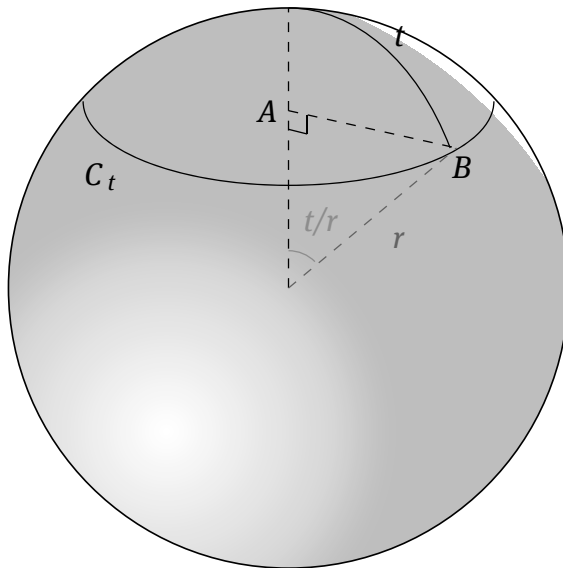
*Proof.* This is immediate from (3.10.1), since  $L(C_t) = \int_0^{2\pi} |\mathbf{Y}_s(t)| ds$  as remarked earlier.  $\square$

Even though the above formula for the sectional curvature is more qualitative than anything, there are some special cases where it can be used to actually evaluate the curvature:

**Example 3.10.1.** Consider a sphere of radius  $r > 0$ . A geodesic of length  $t < \pi r$  is an arc of a circle with radius  $r$  and central angle  $t/r$ . Elementary trigonometry implies that  $C_t$  has length  $2\pi r \sin(t/r)$ . Thus, any plane has curvature

$$K = \lim_{t \rightarrow 0} 6 \frac{t - \sin(t/r)}{t^3} = \frac{1}{r^2}$$

after applying l’Hospital’s rule three times.



the circle  $C_t$  has radius  
 $AB = r \sin \frac{t}{r}$

### 3.11 Isometries

Maps that preserve the inner product play a fundamental role in differential geometry:

**Definition 3.11.1.** A map  $f : M \rightarrow N$  between manifolds is said to be *isometric* if for all  $p \in M$  and  $v, w \in M_p$ ,  $\langle f_*v, f_*w \rangle = \langle v, w \rangle$ . An isometric diffeomorphism is called an *isometry*.

Equivalently,  $f : M \rightarrow N$  is isometric if  $|f_*v| = |v|$  for all  $v \in TM$ : indeed, if  $f$  preserves norms, then it must preserve inner products, for

$$\langle f_*u, f_*v \rangle = \frac{1}{2} (|f_*(u+v)|^2 - |f_*u|^2 - |f_*v|^2)$$

by linearity of  $f_*$ . The converse is also clear. In particular,  $f_{*p} : M_p \rightarrow N_{f(p)}$  is injective in this case, and  $\dim M \leq \dim N$ .

If  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is an isometry, and  $M$  is a submanifold of  $\mathbb{R}^n$ , then the image  $N := f(M)$  is also a manifold, and the restriction  $f : M \rightarrow N$  is an isometry: in fact, any local parametrization  $h$  of  $M$  induces a local parametrization  $f \circ h$  of  $N$ , so that  $N$  is a submanifold. Since the inner product on the tangent bundle of  $M$  is the restriction of that on  $\mathbb{R}^n$ , the restriction of  $f$  to  $M$  is clearly an isometry. The isometries of Euclidean space are called *Euclidean motions*, and are easy to describe:

**Proposition 3.11.1.**  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a Euclidean motion if and only if it is the composition of an orthogonal transformation  $A$  with a translation  $v \mapsto v + a$  by some vector  $a$ ; i.e.,  $f(v) = Av + a$ .

*Proof.* Clearly, compositions of orthogonal transformations with translations are isometries: If  $f(v) = Av + a$ , then for  $u = \mathcal{I}_p v \in \mathbb{R}_p^n$ ,

$$|f_*u| = |Df(p)v| = |Av| = |v| = |u|.$$

Conversely, suppose  $f$  is a Euclidean motion. Notice first of all that  $f$  preserves distances between points: by hypothesis,  $|f \circ c| = |f_* \circ \dot{c}| = |\dot{c}|$ , so that  $\ell(f \circ c) = \ell(c)$ , and  $f$  preserves the length of curves. Thus, if  $c$  is the line segment from  $p$  to  $q$ , then  $f \circ c$  is a curve from  $f(p)$  to  $f(q)$  of length  $|p - q|$ . This implies that  $|f(p) - f(q)| \leq |p - q|$  for any  $p, q$ . But  $f^{-1}$  is also an isometry, so by the same reasoning,

$$|f^{-1}(f(p)) - f^{-1}(f(q))| \leq |f(p) - f(q)| \leq |p - q|,$$

and all terms in the above inequality are the same; i.e.,  $f$  preserves distances as claimed.

If  $a = f(0)$ , and  $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$  denotes the map given by  $g(p) = f(p) - a$ , then  $Dg = Df$ , so that  $g$  is also an isometry. In particular,

$$|g(v)| = |g(v) - g(0)| = |v - 0| = |v|,$$

and it only remains to show that  $\mathbf{g}$  is linear. But if  $\mathbf{g}$  preserves norms, it also preserves inner products, because

$$\begin{aligned}\langle \mathbf{g}(\mathbf{v}), \mathbf{g}(\mathbf{w}) \rangle &= \frac{1}{2} (|\mathbf{g}(\mathbf{v})|^2 + |\mathbf{g}(\mathbf{w})|^2 - |\mathbf{g}(\mathbf{v}) - \mathbf{g}(\mathbf{w})|^2) \\ &= \frac{1}{2} (|\mathbf{v}|^2 + |\mathbf{w}|^2 - |\mathbf{v} - \mathbf{w}|^2) \\ &= \langle \mathbf{v}, \mathbf{w} \rangle.\end{aligned}$$

It follows that for any  $a \in \mathbb{R}$  and  $\mathbf{z} \in \mathbb{R}^n$ ,

$$\begin{aligned}\langle \mathbf{g}(a\mathbf{v} + \mathbf{w}) - a\mathbf{g}(\mathbf{v}) - \mathbf{g}(\mathbf{w}), \mathbf{z} \rangle &= \langle \mathbf{g}(a\mathbf{v} + \mathbf{w}) - a\mathbf{g}(\mathbf{v}) - \mathbf{g}(\mathbf{w}), \mathbf{g}(\mathbf{g}^{-1}(\mathbf{z})) \rangle \\ &= \langle \mathbf{g}(a\mathbf{v} + \mathbf{w}), \mathbf{g}(\mathbf{g}^{-1}(\mathbf{z})) \rangle \\ &\quad - a\langle \mathbf{g}(\mathbf{v}), \mathbf{g}(\mathbf{g}^{-1}(\mathbf{z})) \rangle - \langle \mathbf{g}(\mathbf{w}), \mathbf{g}(\mathbf{g}^{-1}(\mathbf{z})) \rangle \\ &= \langle a\mathbf{v} + \mathbf{w}, \mathbf{g}^{-1}(\mathbf{z}) \rangle - a\langle \mathbf{v}, \mathbf{g}^{-1}(\mathbf{z}) \rangle \\ &\quad - \langle \mathbf{w}, \mathbf{g}^{-1}(\mathbf{z}) \rangle \\ &= 0.\end{aligned}$$

Now take  $\mathbf{z} = \mathbf{g}(a\mathbf{v} + \mathbf{w}) - a\mathbf{g}(\mathbf{v}) - \mathbf{g}(\mathbf{w})$  in the above equation to deduce that  $\mathbf{g}(a\mathbf{v} + \mathbf{w}) - a\mathbf{g}(\mathbf{v}) - \mathbf{g}(\mathbf{w}) = \mathbf{0}$ . Thus,  $\mathbf{g}$  is linear, as claimed.  $\square$

We've already remarked that if  $M^n$  is a submanifold of  $\mathbb{R}^{n+k}$  and  $\mathbf{f}$  is a Euclidean motion, then the restriction  $\mathbf{f} : M \rightarrow \mathbf{f}(M)$  of  $\mathbf{f}$  to  $M$  is an isometry. It is by no means true that every isometry of  $M$  with another manifold is the restriction of a Euclidean motion: for example, let  $M = (0, 2\pi) \times \{0\} \subset \mathbb{R}^2$ . Then  $M$  is isometric to a circle of unit radius with one point removed via  $\mathbf{f}$ , where  $\mathbf{f}(t, 0) = (\cos t, \sin t)$ , but  $\mathbf{f}$  is not a Euclidean motion. We shall, however, see in a later chapter that if  $k = 1$  and  $\mathbf{f}$  also preserves the second fundamental form, then it is the restriction of a Euclidean motion (this is also true when  $k > 1$  under additional assumptions). In order to do so, we must first investigate how covariant derivatives behave under isometric maps. This will in fact lead us to an alternative expression for the covariant derivative, one which is of interest in its own right.

Let  $\mathbf{X}$ ,  $\mathbf{Y}$ , and  $\mathbf{Z}$  denote vector fields on  $M$ . Since  $\mathbf{Z}$  is tangent to  $M$ ,  $\langle \nabla_{\mathbf{X}}\mathbf{Y}, \mathbf{Z} \rangle = \langle D_{\mathbf{X}}\mathbf{Y}, \mathbf{Z} \rangle$ . Thus, by Theorem 3.6.2,

$$\mathbf{X}\langle \mathbf{Y}, \mathbf{Z} \rangle = \langle D_{\mathbf{X}}\mathbf{Y}, \mathbf{Z} \rangle + \langle \mathbf{Y}, D_{\mathbf{X}}\mathbf{Z} \rangle. \quad (3.11.1)$$

In the same way, but also using Proposition 2.9.1,

$$\begin{aligned}\mathbf{Y}\langle \mathbf{Z}, \mathbf{X} \rangle &= \langle D_{\mathbf{Y}}\mathbf{Z}, \mathbf{X} \rangle + \langle \mathbf{Z}, D_{\mathbf{Y}}\mathbf{X} \rangle \\ &= \langle D_{\mathbf{Y}}\mathbf{Z}, \mathbf{X} \rangle + \langle \mathbf{Z}, D_{\mathbf{X}}\mathbf{Y} \rangle - \langle \mathbf{Z}, [\mathbf{X}, \mathbf{Y}] \rangle,\end{aligned} \quad (3.11.2)$$

and

$$\mathbf{Z}\langle \mathbf{X}, \mathbf{Y} \rangle = \langle D_{\mathbf{X}}\mathbf{Z}, \mathbf{Y} \rangle + \langle [\mathbf{Z}, \mathbf{X}], \mathbf{Y} \rangle + \langle \mathbf{X}, D_{\mathbf{Y}}\mathbf{Z} \rangle - \langle \mathbf{X}, [\mathbf{Y}, \mathbf{Z}] \rangle. \quad (3.11.3)$$

Adding the first two equations (3.11.1), (3.11.2), and subtracting the third (3.11.3) yields

$$\begin{aligned} \mathbf{X}\langle \mathbf{Y}, \mathbf{Z} \rangle + \mathbf{Y}\langle \mathbf{Z}, \mathbf{X} \rangle - \mathbf{Z}\langle \mathbf{X}, \mathbf{Y} \rangle &= 2\langle D_{\mathbf{X}}\mathbf{Y}, \mathbf{Z} \rangle + \langle \mathbf{X}, [\mathbf{Y}, \mathbf{Z}] \rangle - \langle \mathbf{Y}, [\mathbf{Z}, \mathbf{X}] \rangle \\ &\quad - \langle \mathbf{Z}, [\mathbf{X}, \mathbf{Y}] \rangle, \end{aligned}$$

so that

$$\begin{aligned} \langle D_{\mathbf{X}}\mathbf{Y}, \mathbf{Z} \rangle = \langle \nabla_{\mathbf{X}}\mathbf{Y}, \mathbf{Z} \rangle &= \frac{1}{2} \left\{ \mathbf{X}\langle \mathbf{Y}, \mathbf{Z} \rangle + \mathbf{Y}\langle \mathbf{Z}, \mathbf{X} \rangle - \mathbf{Z}\langle \mathbf{X}, \mathbf{Y} \rangle \right. \\ &\quad \left. + \langle \mathbf{Z}, [\mathbf{X}, \mathbf{Y}] \rangle + \langle \mathbf{Y}, [\mathbf{Z}, \mathbf{X}] \rangle - \langle \mathbf{X}, [\mathbf{Y}, \mathbf{Z}] \rangle \right\}. \end{aligned} \quad (3.11.4)$$

Notice that (3.11.4) actually characterizes the covariant derivative  $\nabla_{\mathbf{X}}\mathbf{Y}$ , since for any  $\mathbf{p} \in M$  one can find a local orthonormal basis  $\mathbf{Z}_i$  of vector fields in a neighborhood of  $\mathbf{p}$  in  $M$ , and

$$\nabla_{\mathbf{X}(\mathbf{p})}\mathbf{Y} = \sum_i \langle \nabla_{\mathbf{X}(\mathbf{p})}\mathbf{Y}, \mathbf{Z}_i(\mathbf{p}) \rangle \mathbf{Z}_i(\mathbf{p}).$$

Now, (3.11.1) actually holds for vector fields along maps. In particular, given a map  $\mathbf{f} : N \rightarrow M$  and vector fields  $\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}$  on  $N$ ,

$$\tilde{\mathbf{X}}\langle \mathbf{f}_*\tilde{\mathbf{Y}}, \mathbf{f}_*\tilde{\mathbf{Z}} \rangle = \langle \nabla_{\tilde{\mathbf{X}}}\mathbf{f}_*\tilde{\mathbf{Y}}, \mathbf{f}_*\tilde{\mathbf{Z}} \rangle + \langle \mathbf{f}_*\tilde{\mathbf{Y}}, \nabla_{\tilde{\mathbf{X}}}\mathbf{f}_*\tilde{\mathbf{Z}} \rangle.$$

Furthermore, by Theorem 2.9.1,

$$\mathbf{f}_*[\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}] = \nabla_{\tilde{\mathbf{X}}}\mathbf{f}_*\tilde{\mathbf{Y}} - \nabla_{\tilde{\mathbf{Y}}}\mathbf{f}_*\tilde{\mathbf{X}}.$$

Thus, the same argument that lead to (3.11.4) yields the following generalization of it:

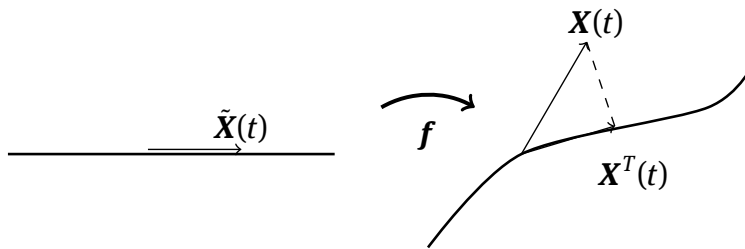
$$\begin{aligned} \langle \nabla_{\tilde{\mathbf{X}}}\mathbf{f}_*\tilde{\mathbf{Y}}, \mathbf{f}_*\tilde{\mathbf{Z}} \rangle &= \frac{1}{2} \left\{ \tilde{\mathbf{X}}\langle \mathbf{f}_*\tilde{\mathbf{Y}}, \mathbf{f}_*\tilde{\mathbf{Z}} \rangle + \tilde{\mathbf{Y}}\langle \mathbf{f}_*\tilde{\mathbf{Z}}, \mathbf{f}_*\tilde{\mathbf{X}} \rangle - \tilde{\mathbf{Z}}\langle \mathbf{f}_*\tilde{\mathbf{X}}, \mathbf{f}_*\tilde{\mathbf{Y}} \rangle \right. \\ &\quad \left. + \langle \mathbf{f}_*\tilde{\mathbf{Z}}, \mathbf{f}_*[\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}] \rangle + \langle \mathbf{f}_*\tilde{\mathbf{Y}}, \mathbf{f}_*[\tilde{\mathbf{Z}}, \tilde{\mathbf{X}}] \rangle - \langle \mathbf{f}_*\tilde{\mathbf{X}}, \mathbf{f}_*[\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}] \rangle \right\}. \end{aligned} \quad (3.11.5)$$

The next lemma essentially identifies the tangential component of a vector field along an isometric map:

**Lemma 3.11.1.** *Given an isometric map  $\mathbf{f} : N \rightarrow M$ , there exists, for any vector field  $\mathbf{X}$  along  $\mathbf{f}$ , a unique vector field  $\mathbf{X}^T$  along  $\mathbf{f}$  such that  $\mathbf{X}^T = \mathbf{f}_*\tilde{\mathbf{X}}$  for some vector field  $\tilde{\mathbf{X}}$  on  $N$ , and*

$$\langle \mathbf{X}, \mathbf{f}_*\tilde{\mathbf{Y}} \rangle = \langle \mathbf{X}^T, \mathbf{f}_*\tilde{\mathbf{Y}} \rangle$$

for any vector field  $\tilde{\mathbf{Y}}$  on  $N$ .



*Proof.* We shall construct  $\mathbf{X}^T$  locally. Given  $\mathbf{p} \in N$ , consider an orthonormal basis  $\tilde{\mathbf{X}}_i$  of vector fields in a neighborhood  $U$  of  $\mathbf{p}$  in  $N$ , and set

$$\varphi_i = \langle \mathbf{X}, \mathbf{f}_* \tilde{\mathbf{X}}_i \rangle : U \rightarrow \mathbb{R}.$$

The restriction  $\tilde{\mathbf{Y}}|_U$  to  $U$  of a vector field  $\tilde{\mathbf{Y}}$  on  $N$  then equals  $\sum_i \langle \tilde{\mathbf{Y}}, \tilde{\mathbf{X}}_i \rangle \tilde{\mathbf{X}}_i$ , and

$$\mathbf{f}_* \tilde{\mathbf{Y}}|_U = \sum_i \langle \tilde{\mathbf{Y}}, \tilde{\mathbf{X}}_i \rangle \mathbf{f}_* \tilde{\mathbf{X}}_i.$$

Thus,

$$\langle \mathbf{X}, \mathbf{f}_* \tilde{\mathbf{Y}} \rangle = \sum_i \langle \tilde{\mathbf{Y}}, \tilde{\mathbf{X}}_i \rangle \langle \mathbf{X}, \mathbf{f}_* \tilde{\mathbf{X}}_i \rangle = \langle \sum_i \varphi_i \tilde{\mathbf{X}}_i, \tilde{\mathbf{Y}} \rangle,$$

and we may take  $\mathbf{X}|_U^T = \mathbf{f}_* \tilde{\mathbf{X}}$ , where  $\tilde{\mathbf{X}} = \sum_i \varphi_i \tilde{\mathbf{X}}_i$ .

For uniqueness, if  $\mathbf{Z} = \mathbf{f}_* \tilde{\mathbf{Z}}$  satisfies  $\langle \mathbf{Z}, \mathbf{f}_* \tilde{\mathbf{Y}} \rangle = \langle \mathbf{X}^T, \mathbf{f}_* \tilde{\mathbf{Y}} \rangle$  for every vector field  $\tilde{\mathbf{Y}}$  on  $U$ , then

$$\langle \tilde{\mathbf{Z}}, \tilde{\mathbf{Y}} \rangle = \langle \mathbf{Z}, \mathbf{f}_* \tilde{\mathbf{Y}} \rangle = \langle \mathbf{X}^T, \mathbf{f}_* \tilde{\mathbf{Y}} \rangle = \langle \tilde{\mathbf{X}}, \tilde{\mathbf{Y}} \rangle,$$

so that  $\langle \tilde{\mathbf{Z}} - \tilde{\mathbf{X}}, \tilde{\mathbf{Y}} \rangle = 0$  for any such  $\tilde{\mathbf{Y}}$ . Now take  $\tilde{\mathbf{Y}} = \tilde{\mathbf{Z}} - \tilde{\mathbf{X}}$  to conclude that  $\tilde{\mathbf{Z}} = \tilde{\mathbf{X}}$ , and therefore also  $\mathbf{Z} = \mathbf{X}^T$ .  $\square$

The next theorem roughly says that covariant derivatives are preserved under isometric maps:

**Theorem 3.11.1.** *Suppose  $\mathbf{f} : N \rightarrow M$  is isometric, and  $\mathbf{X}, \mathbf{Y}$  are vector fields on  $N$ . Then*

$$(\nabla_{\mathbf{X}} \mathbf{f}_* \mathbf{Y})^T = \mathbf{f}_* \nabla_{\mathbf{X}} \mathbf{Y}.$$

*Proof.* Notice that although the same notation is used in the equation above, the covariant derivative on the left is the one on  $M$ , and the one on the right is that on  $N$ . It must be shown that  $\langle \nabla_{\mathbf{X}} \mathbf{f}_* \mathbf{Y}, \mathbf{f}_* \mathbf{Z} \rangle = \langle \nabla_{\mathbf{X}} \mathbf{Y}, \mathbf{Z} \rangle$  for any vector field  $\mathbf{Z}$  on  $N$ . The first inner product is given by (3.11.5), if one deletes all tildes in that equation. Since  $\mathbf{f}$  is isometric, all  $\mathbf{f}_*$  on the right side of that identity may be dropped. What remains is the right side of (3.11.4), thereby establishing the claim.  $\square$

Suppose furthermore that  $N$  and  $M$  have the same dimension. The above theorem then says that

$$\nabla_{\mathbf{X}} \mathbf{f}_* \mathbf{Y} = \mathbf{f}_* \nabla_{\mathbf{X}} \mathbf{Y} \tag{3.11.6}$$

for any vector fields  $\mathbf{X}$  and  $\mathbf{Y}$  on  $N$ . This implies that if  $\mathbf{X}$  is a parallel vector field along a curve  $\mathbf{c}$  in  $N$ , then  $\mathbf{f}_* \mathbf{X}$  is a parallel vector field along  $\mathbf{f} \circ \mathbf{c}$ . In particular,  $\mathbf{f}$  maps a geodesic  $\mathbf{c}$  in  $N$  to a geodesic  $\mathbf{f} \circ \mathbf{c}$  in  $M$ ; equivalently,  $\exp_M \circ \mathbf{f}_* = \mathbf{f} \circ \exp_N$ . Moreover, by Theorem 3.9.2,  $\mathbf{f}$  preserves the curvature tensors in the sense that  $\mathbf{f}_* R_N(\mathbf{X}, \mathbf{Y})\mathbf{Z} = R_M(\mathbf{f}_* \mathbf{X}, \mathbf{f}_* \mathbf{Y})\mathbf{f}_* \mathbf{Z}$  for any vector fields  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$  on  $N$ . This, of course, also implies that the sectional curvature of a plane  $P \subset TN$  equals that of  $\mathbf{f}_*(P) \subset TM$ . Summarizing, we have proved:

**Theorem 3.11.2.** Suppose  $f : N \rightarrow M$  is an isometric map between manifolds of the same dimension. Then

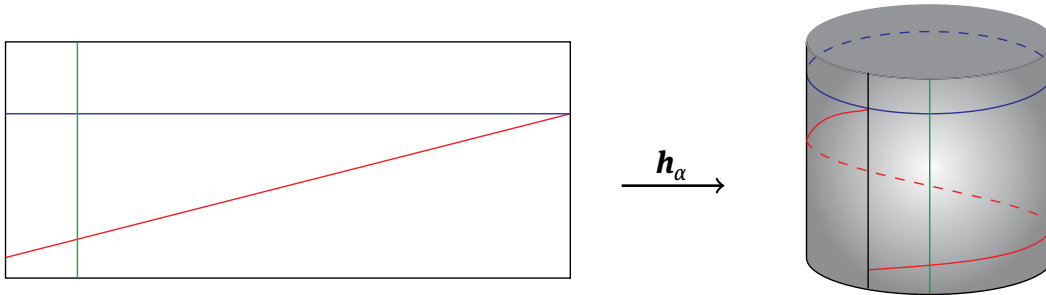
- (1)  $\exp_M \circ f_* = f \circ \exp_N$ .
- (2)  $f_* R_N(X, Y)Z = R_M(f_* X, f_* Y)f_* Z$  for any vector fields  $X, Y, Z$  on  $N$ .
- (3) The sectional curvature of a plane  $P \subset TN$  equals that of  $f_*(P) \subset TM$ .

**Example 3.11.1.** Consider the cylinder  $M = \{(x, y, z) \mid x^2 + y^2 = r^2\}$  of radius  $r > 0$  in  $\mathbb{R}^3$ . Given any  $\alpha \in \mathbb{R}$ , the map  $h_\alpha : (r\alpha, r(\alpha + 2\pi)) \times \mathbb{R} \rightarrow M$ , where  $h_\alpha(u, v) = (r \cos(u/r), r \sin(u/r), v)$  is a local parametrization of  $M$ . Its image is all of  $M$  except for a line parallel to the  $z$ -axis.  $M$  can therefore be covered by two parametrizations corresponding to judiciously chosen values of  $\alpha$ . The domain of  $h_\alpha$  is an open subset of  $\mathbb{R}^2$  and has constant sectional curvature zero. Furthermore, the two vector fields

$$h_{\alpha*} D_1(u, v) = -\sin \frac{u}{r} D_1 \circ h_\alpha(u, v) + \cos \frac{u}{r} D_2 \circ h_\alpha(u, v),$$

$$h_{\alpha*} D_2(u, v) = D_3 \circ h_\alpha(u, v),$$

are everywhere orthonormal, so that  $h_\alpha$  is isometric. It follows that  $M$  is flat. The geodesics of  $M$  are the images of the geodesics in the plane and are therefore helices in  $M$ .



Isometries can be useful in identifying geodesics of  $M$ : Suppose  $N$  is a manifold contained in  $M$ .  $N$  is said to be *totally geodesic in  $M$*  if every geodesic of  $N$  is also a geodesic of  $M$ . For example,  $S^n \times \{0\}$  is totally geodesic in  $S^{n+1} \subset \mathbb{R}^{n+2}$  (but not in  $\mathbb{R}^{n+2}$ ).

**Proposition 3.11.2.** Let  $M$  be a manifold and  $f : M \rightarrow M$  an isometry. If the fixed point set  $N = \{p \in M \mid f(p) = p\}$  of  $f$  is a manifold, then it is totally geodesic in  $M$ .

*Proof.* Consider any  $p \in N$  and  $u \in N_p$ . Denote by  $c_u$  the geodesic of  $M$  with  $\dot{c}_u(0) = u$ . Since  $f$  is an isometry,  $f \circ c_u$  is also a geodesic, and  $f \circ \dot{c}_u(0) = f_* u = u$  because  $u \in N_p$  and  $f_*$  is the identity on  $N_p$ . By uniqueness of geodesics,  $f \circ c_u = c_u$ , and  $c_u$  has its image contained in  $N$ . By Theorem 6.1.1, it is also a geodesic in  $N$ . Since geodesics are uniquely determined by their tangent vector at the origin and  $u$  was arbitrary, the result follows. □

The above proposition provides for example another way of describing geodesics on spheres: We may first of all assume the sphere is centered at the origin, since translations are isometries, and thus map geodesics to geodesics. So let  $p \in M = S^n(r)$ ,

$\mathbf{u} \in \mathbb{R}^{n+1}$  a unit vector orthogonal to  $\mathbf{p}$ , so that  $\mathcal{I}_{\mathbf{p}}\mathbf{u} \in M_{\mathbf{p}}$ . We wish to describe the geodesic starting at  $\mathbf{p}$  with initial tangent vector  $\mathcal{I}_{\mathbf{p}}\mathbf{u}$ .

If  $E$  is a subspace of  $\mathbb{R}^{n+1}$ , recall that *reflection* in  $E$  is the linear transformation of Euclidean space that maps  $\mathbf{v} + \mathbf{w} \in E \oplus E^\perp = \mathbb{R}^{n+1}$  to  $\mathbf{v} - \mathbf{w}$ . It is an isometry because  $\mathbf{v} \perp \mathbf{w}$ , so that  $|\mathbf{v} + \mathbf{w}| = |\mathbf{v} - \mathbf{w}|$ . Reflection in the plane  $P$  spanned by  $\mathbf{p}$  and  $\mathbf{u}$  is then an isometry of Euclidean space which maps  $M$  to itself, so that its restriction to  $M$  is an isometry of  $M$ . Furthermore, the fixed point set of this restriction is the great circle  $P \cap M$  on the sphere, and this circle is therefore totally geodesic. Since this is true for any point  $\mathbf{p} \in M$  and unit vector in the tangent space at  $\mathbf{p}$ , all great circles are the images of geodesics. As seen earlier, the geodesic above may be parametrized by  $t \mapsto (\cos(t/r))\mathbf{p} + (\sin(t/r))(r\mathbf{u})$ .

### 3.12 Exercises

**3.1.** Let  $R, h > 0$ . An *open cone* with base  $R$  and height  $h$  is the subset of  $\mathbb{R}^3$  obtained by rotating the half-open line segment

$$\{(t, 0, (-h/R)t + h) \mid 0 \leq t < R\}$$

joining  $(0, 0, h)$  and  $(R, 0, 0)$  about the  $z$ -axis. Notice that the line segment contains the first point but not the second.

- (a) Show that even though a cone is commonly called a surface, it is not a submanifold of  $\mathbb{R}^3$ .
- (b) If the tip  $(0, 0, h)$  of the cone is removed, is the resulting set a manifold?

**3.2.** Recall that a *hypersurface* is a manifold with dimension one less than that of the ambient Euclidean space.

- (a) Let  $\mathbf{f} : \mathbb{R}^{n+k} \supset U \rightarrow \mathbb{R}^k$ . Show that if  $\mathbf{0}$  is a regular value of  $\mathbf{f}$ , then  $\mathbf{0}$  is a regular value of each  $u^i \circ \mathbf{f}$ ,  $i = 1, \dots, k$ . Conclude that  $\mathbf{f}^{-1}(\mathbf{0})$  is the intersection of  $k$  hypersurfaces.
- (b) Suppose  $f_1, \dots, f_k : U \rightarrow \mathbb{R}$  all have  $0$  as regular value, and set  $M_i = f_i^{-1}(0)$ . Show that  $M = \bigcap_{i=1}^k M_i$  is an  $n$ -dimensional manifold of  $\mathbb{R}^{n+k}$  if and only if the vectors  $\nabla f_1(\mathbf{p}), \dots, \nabla f_k(\mathbf{p})$  are linearly independent for every  $\mathbf{p} \in M$ .

**3.3.** Prove the “if” part of Corollary 3.2.1:  $M \subset \mathbb{R}^{n+k}$  is an  $n$ -dimensional manifold if for every  $\mathbf{p} \in M$ , there exists a neighborhood  $U$  of  $\mathbf{p}$  in  $\mathbb{R}^{n+k}$ , and a map  $\mathbf{f} : U \rightarrow \mathbb{R}^k$  having  $\mathbf{0}$  as a regular value, such that  $U \cap M = \mathbf{f}^{-1}(\mathbf{0})$ .

**3.4.** Let  $M^n$  be a manifold,  $(U, \mathbf{x})$  a chart of  $M$  around some  $\mathbf{p} \in M$ . Show that  $\{dx^i(\mathbf{p}) \mid i = 1, \dots, n\}$  is the basis of  $M_{\mathbf{p}}^*$  dual to  $\{\partial/\partial x^i(\mathbf{p})\}$ .

**3.5.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be given by  $f(\mathbf{a}) = |\mathbf{a}|^2$ .

- (a) Find  $df$ .
- (b) Find  $dg$ , where  $g := f|_{S^{n-1}}$ .



**3.6.** Let  $f, g : M \rightarrow \mathbb{R}$ ,  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ . Prove that

- (a)  $d(f + g) = df + dg$ .
- (b)  $d(fg) = f dg + g df$ .
- (c)  $d(\varphi \circ f) = (\varphi' \circ f) df$ .

**3.7.** If  $r, \theta$  denote polar coordinates on the plane, express  $dr$  and  $d\theta$  in terms of  $du^1$  and  $du^2$ .

**3.8.** Determine  $df$ , if  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  is given by:

- (a)  $f(x, y, z) = \cos e^{x^2 y + z}$ .
- (b)  $f(\mathbf{a}) = \langle \mathbf{p}, \mathbf{a} \rangle$  for some  $\mathbf{p} \in \mathbb{R}^3$ .

**3.9.** A *differential 1-form* on a manifold  $M$  is just a tensor field of type (1,0).

- (a) Let  $\alpha$  be a differential 1-form on  $M$ , and  $(U, \mathbf{x})$  a chart of  $M$ . Prove that the restriction of  $\alpha$  to  $U$  can be written

$$\alpha|_U = \sum_i \alpha_i dx^i$$

for some functions  $\alpha_i : U \rightarrow \mathbb{R}$ . How is  $\alpha_i$  defined?

- (b) Use the first fundamental tensor field of  $M$  to show that there is a one-to-one correspondence between differential 1-forms and vector fields on  $M$ , see also Corollary 1.4.2.

**3.10.** Let  $M$  be a manifold,  $f : M \rightarrow \mathbb{R}$ , and  $(U, \mathbf{x})$  a chart of  $M$ , so that the restriction of  $df$  to  $U$  can be written

$$df|_U = \sum_i df \left( \frac{\partial}{\partial x^i} \right) \frac{\partial}{\partial x^i} = \sum_i \frac{\partial}{\partial x^i} (f) \frac{\partial}{\partial x^i}.$$

Show that  $\partial/\partial x^i (f) = D_i(f \circ \mathbf{x}^{-1}) \circ \mathbf{x}$ .

**3.11.** (a) Use the method of Lagrange multipliers to determine the maximum of  $g = \sum_{i=1}^n u^i u^{n+i} : \mathbb{R}^{2n} \rightarrow \mathbb{R}$  on the manifold  $\mathbf{f}^{-1}(\mathbf{0})$ , where

$$\mathbf{f} = \left( \sum_{i=1}^n (u^i)^2 - 1, \sum_{i=1}^n (u^{n+i})^2 - 1 \right) : \mathbb{R}^{2n} \rightarrow \mathbb{R}^2.$$

- (b) Given any  $(x_1, \dots, x_n, y_1, \dots, y_n) \in \mathbb{R}^n \times \mathbb{R}^n$  with  $\sum x_i^2, \sum y_i^2 \neq 0$ , the point

$$\mathbf{p} = \left( \frac{x_1}{(\sum x_i^2)^{1/2}}, \dots, \frac{x_n}{(\sum x_i^2)^{1/2}}, y_1/(\sum y_i^2)^{1/2}, \dots, y_n/(\sum y_i^2)^{1/2} \right)$$

belongs to  $\mathbf{f}^{-1}(\mathbf{0})$ . Use this to give another proof of the Cauchy-Schwarz inequality from Theorem 1.4.1.

**3.12.** Given a manifold  $M$ , the *bundle projection* is the map  $\pi : TM \rightarrow M$  that maps  $\mathbf{u} \in TM$  to  $\mathbf{p}$ , if  $\mathbf{u} \in M_{\mathbf{p}}$ . Show that the bundle projection is differentiable.

**3.13.** An  $n$ -dimensional manifold  $M$  is said to be *parallelizable* if its tangent bundle is diffeomorphic to  $M \times \mathbb{R}^n$  by means of a diffeomorphism  $f$  which makes the diagram

$$\begin{array}{ccc} TM & \xrightarrow{f} & M \times \mathbb{R}^n \\ & \searrow \pi & \downarrow \pi_1 \\ & & M \end{array}$$

commute. Here,  $\pi$  is the bundle projection, and  $\pi_1$  is projection onto the first factor. Show that  $M$  is parallelizable if and only if it admits  $n$  vector fields which are linearly independent at every point of  $M$ . Conclude that every Lie group is parallelizable. In contrast, it can be shown that a 2-dimensional sphere is not parallelizable; in fact, every vector field on it must vanish somewhere.

**3.14.** Let  $M, N$  be manifolds. Given  $p \in M$ , define  $J_p : N \rightarrow M \times N$  by  $J_p(q) = (p, q)$ . Similarly, for  $q \in N$ , define  $\iota_q : M \rightarrow M \times N$  by  $\iota_q(p) = (p, q)$ .

- (a) Prove that the map  $\iota_{q*} + J_{p*} : M_p \times N_q \xrightarrow{\cong} (M \times N)_{(p,q)}$  given by  $(\iota_{q*} + J_{p*})(u, v) = \iota_{q*}u + J_{p*}v$ , is an isomorphism.
- (b) Show that the isomorphism from (a) has as inverse  $(\pi_{M*}(p,q), \pi_{N*}(p,q))$ , where  $\pi_M : M \times N \rightarrow M$  and  $\pi_N : M \times N \rightarrow N$  denote the projections.

**3.15.** Let  $(U, x)$  denote a chart on  $M$ . Show that if  $X$  is a vector field on  $U$ , then

$$X = \sum_i dx^i(X) \frac{\partial}{\partial x^i}.$$

Conclude that  $dx^i$  is a tensor field on  $U$  (i.e., that it is differentiable).

**3.16.** Show that if  $(U, x)$  and  $(V, y)$  are two charts of  $M$ , then on  $U \cap V$ ,

$$\frac{\partial}{\partial x^i} = \sum_j (D_i(y^j \circ x^{-1}) \circ x) \frac{\partial}{\partial y^j}.$$

Notice that according to Exercise 3.10, this is equivalent to

$$\frac{\partial}{\partial x^i} = \sum_j \frac{\partial y^j}{\partial x^i} \frac{\partial}{\partial y^j}.$$

**3.17.** Let  $X$  be a vector field on  $\mathbb{R}^n$  with flow  $\Phi_t$ . Define a map  $F : U \rightarrow \mathbb{R}^n$  on a neighborhood  $U$  of the origin by

$$F(a_1, \dots, a_n) = \Phi_{a_1}(0, a_2, \dots, a_n).$$

- (a) Show that  $F_*D_1 = X \circ F.t'$
- (b) Show that  $F_*D_i(\mathbf{0}) = D_i(\mathbf{0})$  for  $i > 1$ .
- (c) Let  $M$  be a manifold,  $p \in M$ , and  $X$  a vector field on  $M$  with  $X(p) \neq \mathbf{0}$ . Prove that there exists a chart  $(U, x)$  around  $p$  such that

$$X|_U = \frac{\partial}{\partial x^1}.$$

**3.18.** Let  $X$  be a vector field on a manifold  $M$  that vanishes outside a compact set. Show that  $X$  is complete.

**3.19.** Let  $a < b$  denote two regular values of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . Suppose that  $f^{-1}(a)$  and  $f^{-1}(b)$  are nonempty, so that they are  $(n - 1)$ -dimensional manifolds. Prove that if  $f^{-1}[a, b]$  is compact and contains no critical points of  $f$ , then  $f^{-1}(a)$  is diffeomorphic to  $f^{-1}(b)$ . *Hint:* Consider a vector field  $X$  that equals  $\nabla f / |\nabla f|$  on  $f^{-1}[a, b]$  and vanishes outside a compact set.  $X$  is then complete by Exercise 3.18, so that  $\Phi_{b-a}$  is defined, where  $\Phi_t$  is the flow of  $X$ .

**3.20.** Recall that  $M_{m,n} \cong \mathbb{R}^{mn}$  denotes the space of all  $m \times n$  matrices. The goal of this exercise is to show that the subset  $M_{m,n}(k)$  consisting of all matrices of rank  $k$  is a submanifold of  $\mathbb{R}^{mn}$  with dimension  $k(m + n - k)$ .

(a) Show that if  $M \in M_{m,n}$  has rank at least  $k$ , then there exist elementary matrices  $P \in M_{m,m}$  and  $Q \in M_{n,n}$  such that

$$PMQ = \begin{bmatrix} A & B \\ C & D \end{bmatrix}$$

where  $A$  is an invertible  $k \times k$  matrix,  $B \in M_{k,n-k}$ ,  $C \in M_{m-k,k}$ , and  $D \in M_{m-k,n-k}$ .

(b) Let  $A$  be a nonsingular  $k \times k$  matrix. Prove that

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \in M_{m,n}$$

has rank  $k$  if and only if  $D = CA^{-1}B$ . *Hint:* the matrix

$$\begin{bmatrix} A & B \\ 0 & -CA^{-1}B + D \end{bmatrix} = \begin{bmatrix} I_k & 0 \\ -CA^{-1} & I_{m-k} \end{bmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix}$$

has the same rank as

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}.$$

(c) Given  $M_0 \in M_{m,n}(k)$ , choose elementary matrices  $P$  and  $Q$  such that

$$PM_0Q = \begin{bmatrix} A_0 & B_0 \\ C_0 & D_0 \end{bmatrix}$$

with  $A_0$  invertible. There exists an open neighborhood  $V$  of  $A_0$  in  $M_{k,k}$  such that every  $A$  in  $V$  is invertible. Denote by  $U$  the open neighborhood of  $M_0$  in  $M_{m,n}(k)$  consisting of all matrices of the form

$$M = P^{-1} \begin{bmatrix} A & B \\ C & D \end{bmatrix} Q^{-1},$$

with  $A \in V$ ,  $B \in M_{k,n-k}$ ,  $C \in M_{m-k,k}$ , and  $D \in M_{m-k,n-k}$ .

Identify  $\mathbb{R}^{k(m+n-k)} = \mathbb{R}^{k^2+k(n-k)+k(m-k)}$  with the subspace of  $M_{m,n}$  that consists of all matrices of the form

$$\begin{bmatrix} A & B \\ C & 0 \end{bmatrix}, \quad A \in M_{k,k}, \quad B \in M_{k,n-k}, \quad C \in M_{m-k,k}.$$

If

$$\begin{aligned} \pi : \mathbb{R}^{mn} &\longrightarrow \mathbb{R}^{k(m+n-k)} \\ \begin{bmatrix} A & B \\ C & D \end{bmatrix} &\longmapsto \begin{bmatrix} A & B \\ C & 0 \end{bmatrix} \end{aligned}$$

denotes projection, define

$$\begin{aligned} \mathbf{x} : U &\longrightarrow \mathbb{R}^{k(m+n-k)} \\ M &\longmapsto \pi(PMQ). \end{aligned} \tag{3.12.1}$$

Show that the collection of all  $(U, \mathbf{x})$  in (3.12.1) defines an atlas on  $M_{m,n}(k)$ .

**3.21.** Let  $m : G \times G \rightarrow G$  denote the multiplication map on a Lie group  $G$ .

- (a) Show that for  $(\mathbf{v}, \mathbf{0}) \in G_e \times G_e$ ,  $m_*(\mathbf{v}, \mathbf{0}) = \mathbf{v}$ .  
 (b) Suppose  $\mathbf{c}_1$  and  $\mathbf{c}_2$  are two curves in  $G$  that pass through  $e$  at time zero. Show that if  $\mathbf{c}(t) = \mathbf{c}_1(t) \cdot \mathbf{c}_2(t)$ , then

$$\dot{\mathbf{c}}(0) = \dot{\mathbf{c}}_1(0) + \dot{\mathbf{c}}_2(0).$$

**3.22.** Suppose  $G$  is a Lie group,  $\mathbf{X}, \mathbf{Y} \in \mathfrak{g}$  with flows  $\Phi_t$  and  $\Psi_t$  respectively.

- (a) Prove that  $(\Phi_t \circ \Psi_t)(e) = \Psi_t(e) \cdot \Phi_t(e)$ .  
 (b) Let  $\mathbf{c}_X$  denote the integral curve of  $\mathbf{X} \in \mathfrak{g}$  that passes through  $e$  at time 0. Show that if  $[\mathbf{X}, \mathbf{Y}] = \mathbf{0}$  for  $\mathbf{X}, \mathbf{Y} \in \mathfrak{g}$ , then  $\mathbf{c}_X(t) \cdot \mathbf{c}_Y(t) = \mathbf{c}_Y(t) \cdot \mathbf{c}_X(t)$  for all  $t$ .  
 (c) Suppose  $\mathbf{X}, \mathbf{Y} \in \mathfrak{g}$  have vanishing bracket. Use Exercise 3.21 to prove that the curve  $\mathbf{c}_X \cdot \mathbf{c}_Y : \mathbb{R} \rightarrow G$  is a homomorphism with velocity vector  $\mathbf{X}(e) + \mathbf{Y}(e)$  at zero. Conclude that  $\exp(\mathbf{X} + \mathbf{Y}) = \exp(\mathbf{X}) \exp(\mathbf{Y})$ .

**3.23.** If  $g$  is an element of a Lie group  $G$ , *conjugation by  $g$*  is the diffeomorphism  $\tau_g : G \rightarrow G$  given by  $\tau_g a = gag^{-1}$ ,  $a \in G$ . Its derivative at the identity is therefore an isomorphism  $\tau_{g*e} : G_e \rightarrow G_e$ . In analogy with  $GL(n)$ , the space of all isomorphisms of  $G_e \cong \mathfrak{g}$  with itself is denoted  $GL(\mathfrak{g})$ , and is a Lie group under composition. The map

$$\begin{aligned} \text{Ad} : G &\longrightarrow GL(\mathfrak{g}), \\ g &\longmapsto \text{Ad}_g = \tau_{g*e} \end{aligned}$$

is called the *adjoint representation* of  $G$ .

- (a) Prove that  $\text{Ad}$  is a Lie group homomorphism.  
 (b) Suppose  $G = GL(n)$ , and identify  $\mathfrak{g}$  with  $M_n$ . Show that for  $A \in G$  and  $M \in \mathfrak{g}$ ,  $\text{Ad}_A M = AMA^{-1}$ . In particular,  $\text{Ad}$  is differentiable. It can be shown that the adjoint representation in any Lie group is differentiable.

(c) The derivative of  $\text{Ad}$  at the identity is denoted  $\text{ad} : \mathfrak{g} \rightarrow M_n(\mathfrak{g})$ . Prove that for  $\mathbf{X}, \mathbf{Y} \in \mathfrak{g}$ ,  $\text{ad}_{\mathbf{X}}(\mathbf{Y}) = [\mathbf{X}, \mathbf{Y}]$  when  $G = GL(n)$ . It can be shown that this identity is valid in any Lie group.

**3.24.** Prove that on a compact manifold, the exponential map at any point is defined on the whole tangent space; equivalently, any geodesic has all of  $\mathbb{R}$  as its domain.

**3.25.** A *geodesic vector field* on a manifold  $M$  is a vector field  $\mathbf{X}$  that satisfies  $\nabla_{\mathbf{X}}\mathbf{X} \equiv \mathbf{0}$ .

- (a) Prove that  $\mathbf{X}$  is geodesic if and only if its integral curves are geodesics in  $M$ .
- (b) Let  $\mathbf{X}$  be a vector field on an open set  $U \subset \mathbb{R}^n$  represented by  $\mathbf{f} : U \rightarrow \mathbb{R}^n$ , so that  $\mathbf{X}(\mathbf{p}) = \mathcal{I}_{\mathbf{p}}\mathbf{f}(\mathbf{p})$ . Show that  $\mathbf{X}$  is geodesic if and only if

$$\langle \nabla f^i, \mathbf{f} \rangle \equiv \mathbf{0}, \quad 1 \leq i \leq n.$$

In particular, any parallel vector field on  $U$  is geodesic. Are these the only ones?  
*Hint:* Consider  $\mathbf{P}/|\mathbf{P}|$  on  $\mathbb{R}^n \setminus \{\mathbf{0}\}$ , where  $\mathbf{P}$  is the position vector field.

**3.26.** A vector field  $\mathbf{X}$  on a manifold  $M$  is called a *Killing field* if its flow  $\{\Phi_t\}$  consists of local isometries of  $M$ .

- (a) Show that  $\mathbf{X}$  is Killing if and only if for any vector field  $\mathbf{Y}$  that is  $\Phi_t$ -related to itself for all  $t$ ,  $\langle \mathbf{Y}, \mathbf{Y} \rangle$  is constant.
- (b) Show that in (a) the condition of being  $\Phi_t$ -related to itself is equivalent to  $[\mathbf{X}, \mathbf{Y}] \equiv \mathbf{0}$ .
- (c) Show that  $\mathbf{X}$  is Killing if and only  $\langle D_u\mathbf{X}, \mathbf{u} \rangle = 0$  for all  $\mathbf{u} \in TM$ ; equivalently, the operator  $\mathbf{u} \mapsto D_u\mathbf{X}$  on  $M_{\mathbf{p}}$  is skew-adjoint for all  $\mathbf{p} \in M$ .

**3.27.** Let  $M^2$  be a 2-dimensional submanifold of  $\mathbb{R}^3$ ,  $\mathbf{n}$  a unit normal vector field to  $M$  in a neighborhood of  $\mathbf{p} \in M$ , and  $S$  the second fundamental tensor field of  $M$  with respect to  $\mathbf{n}$ .

- (a) Show that  $R(\mathbf{x}, \mathbf{y})\mathbf{z} = \langle \mathbf{S}\mathbf{y}, \mathbf{z} \rangle \mathbf{S}\mathbf{x} - \langle \mathbf{S}\mathbf{x}, \mathbf{z} \rangle \mathbf{S}\mathbf{y}$  for  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in M_{\mathbf{p}}$ .
- (b) Prove that if  $\mathbf{x}$  and  $\mathbf{y}$  form an orthonormal basis of  $M_{\mathbf{p}}$ , then the sectional curvature  $K(\mathbf{p})$  of  $M$  at  $\mathbf{p}$  is given by

$$K(\mathbf{p}) = \langle \mathbf{S}\mathbf{x}, \mathbf{x} \rangle \langle \mathbf{S}\mathbf{y}, \mathbf{y} \rangle - \langle \mathbf{S}\mathbf{x}, \mathbf{y} \rangle^2.$$

(c) Let  $(U, \mathbf{x})$  be a chart of  $M$  around  $\mathbf{p}$ , and consider the six functions on  $U$  given by

$$\begin{aligned} E &= \left\langle \frac{\partial}{\partial x^1}, \frac{\partial}{\partial x^1} \right\rangle & F &= \left\langle \frac{\partial}{\partial x^1}, \frac{\partial}{\partial x^2} \right\rangle & G &= \left\langle \frac{\partial}{\partial x^2}, \frac{\partial}{\partial x^2} \right\rangle \\ l &= \left\langle S \frac{\partial}{\partial x^1}, \frac{\partial}{\partial x^1} \right\rangle & n &= \left\langle S \frac{\partial}{\partial x^1}, \frac{\partial}{\partial x^2} \right\rangle & m &= \left\langle S \frac{\partial}{\partial x^2}, \frac{\partial}{\partial x^2} \right\rangle \end{aligned}$$

Prove that the sectional curvature on  $U$  is given by the function  $K : U \rightarrow \mathbb{R}$ , where

$$K = \frac{lm - n^2}{EG - F^2}.$$

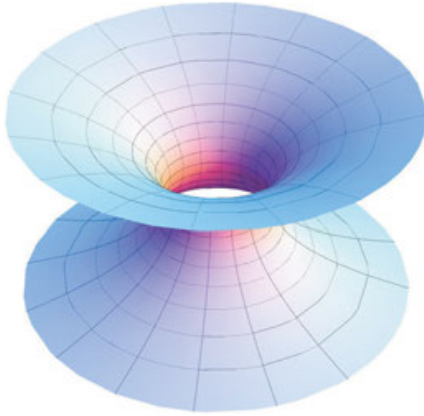


Fig. 3.5: A catenoid

**3.28.** Let  $M^2$  be a 2-dimensional submanifold of  $\mathbb{R}^3$ ,  $\mathbf{n}$  and  $S$  as in Exercise 3.27. The *mean curvature* of  $M$  at a point  $\mathbf{p} \in M$  is the trace of  $S(\mathbf{p})$ . A surface is said to be *minimal* if its mean curvature is zero everywhere.

Let  $a > 0$ . The image of the map  $\mathbf{h} : [0, 2\pi] \times \mathbb{R} \rightarrow \mathbb{R}^3$ , where

$$\mathbf{h}(u, v) = \left( a \cosh \frac{v}{a} \cos u, a \cosh \frac{v}{a} \sin u, v \right),$$

is called a *catenoid*. A catenoid can be simulated by dipping two circles in a soapy solution and pulling them slowly apart. Compute the curvature of a catenoid and show that it is a minimal surface.

**3.29.** Since  $S^1$  is a submanifold of  $\mathbb{R}^2$ ,  $S^1 \times S^1$  is in a natural way a submanifold of  $\mathbb{R}^4 = \mathbb{R}^2 \times \mathbb{R}^2$ , called a *flat torus*, and denoted  $T^2$ .

(a) Show that  $T^2$  is indeed flat; i.e., its curvature is identically zero.

(b) Let  $0 < r < R$ . Prove that the map  $\mathbf{h} : [0, 2\pi] \times [0, 2\pi] \rightarrow \mathbb{R}^3$ , where

$$\mathbf{h}(u, v) = ((R + r \cos u) \cos v, (R + r \cos u) \sin v, r \sin u),$$

is an immersion, and that its image  $M$  is an (imbedded) 2-dimensional submanifold of  $\mathbb{R}^3$ .  $M$  is also called a torus.

(c) Show that  $M$  and  $T^2$  are diffeomorphic but not isometric.

**3.30.** Given  $a, b, c > 0$ , consider the ellipsoid

$$M^2 = \left\{ (x, y, z) \in \mathbb{R}^3 \mid \frac{x^2}{a^2} + \frac{y^2}{b^2} + \frac{z^2}{c^2} = 1 \right\}$$

in  $\mathbb{R}^3$ .

(a) Show that

$$\mathbf{N} = \frac{u^1}{a^2} \mathbf{D}_1 + \frac{u^2}{b^2} \mathbf{D}_2 + \frac{u^3}{c^2} \mathbf{D}_3$$

is a vector field normal to  $M$ .

- (b) Let  $\mathbf{p} = (x, y, z) \in M$ , and denote by  $s$  the second fundamental form of  $M$  at  $\mathbf{p}$  with respect to  $\mathbf{N}(\mathbf{p})/|\mathbf{N}(\mathbf{p})|$ . Prove that

$$s(\mathbf{x}, \mathbf{y}) = \left( \frac{x^2}{a^4} + \frac{y^2}{b^4} + \frac{z^2}{c^4} \right)^{-\frac{1}{2}} \left( \frac{x_1 y_1}{a^2} + \frac{x_2 y_2}{b^2} + \frac{x_3 y_3}{c^2} \right),$$

for  $\mathbf{x} = \sum_{i=1}^3 x_i \mathbf{D}_i(\mathbf{p})$ ,  $\mathbf{y} = \sum_{i=1}^3 y_i \mathbf{D}_i(\mathbf{p})$ .

- (c) Show that the sectional curvature of  $M$  at  $\mathbf{p} = (x, y, z)$  equals

$$K = \left[ abc \left( \frac{x^2}{a^4} + \frac{y^2}{b^4} + \frac{z^2}{c^4} \right) \right]^{-2}.$$

**3.31.** Given  $a, b, c > 0$ , consider the hyperboloid

$$M^2 = \left\{ (x, y, z) \in \mathbb{R}^3 \mid \frac{x^2}{a^2} + \frac{y^2}{b^2} - \frac{z^2}{c^2} = 1 \right\}$$

in  $\mathbb{R}^3$ .

- (a) Show that

$$\mathbf{N} = \frac{u^1}{a^2} \mathbf{D}_1 + \frac{u^2}{b^2} \mathbf{D}_2 - \frac{u^3}{c^2} \mathbf{D}_2$$

is a vector field normal to  $M$ .

- (b) Let  $\mathbf{p} = (x, y, z) \in M$ , and denote by  $s$  the second fundamental form of  $M$  at  $\mathbf{p}$  with respect to  $\mathbf{N}(\mathbf{p})/|\mathbf{N}(\mathbf{p})|$ . Prove that

$$s(\mathbf{x}, \mathbf{y}) = \left( \frac{x^2}{a^4} + \frac{y^2}{b^4} + \frac{z^2}{c^4} \right)^{-\frac{1}{2}} \left( \frac{x_1 y_1}{a^2} + \frac{x_2 y_2}{b^2} - \frac{x_3 y_3}{c^2} \right),$$

for  $\mathbf{x} = \sum_{i=1}^3 x_i \mathbf{D}_i(\mathbf{p})$ ,  $\mathbf{y} = \sum_{i=1}^3 y_i \mathbf{D}_i(\mathbf{p})$ .

- (c) Show that the sectional curvature of  $M$  at  $\mathbf{p} = (x, y, z)$  equals

$$K = - \left[ abc \left( \frac{x^2}{a^4} + \frac{y^2}{b^4} + \frac{z^2}{c^4} \right) \right]^{-2}.$$

**3.32.** Describe three mutually orthogonal geodesics (i.e., their tangent vectors are orthogonal at the points of intersection) on the ellipsoid from exercise 3.30.

**3.33.** Let  $0 \leq a < b$ , and  $f : (a, b) \rightarrow \mathbb{R}$  be a smooth function. If  $a = 0$ , suppose that  $f$  is extendable to  $a$  with  $f^{(n)}(a) = 0$  for all  $n$ .

- (a) Show that  $M = \{(x_1, \dots, x_n) \in \mathbb{R}^n \mid x_n = f((x_1^2 + \dots + x_n^2)^{1/2})\}$  is an  $(n - 1)$ -dimensional submanifold of  $\mathbb{R}^n$ .

- (b) Let  $A$  denote the annulus  $\{\mathbf{u} \in \mathbb{R}^{n-1} \mid a < |\mathbf{u}| < b\}$ . Show that for  $\mathbf{u} \in A$ , the curve  $t \mapsto (t\mathbf{u}, f(|t\mathbf{u}|))$  is, after reparametrization, a geodesic of  $M$ .

**3.34.** A more general surface of revolution than the one introduced in Section 3.1 is that generated by a curve rather than by the graph of a function: let  $\mathbf{c} = (c^1, 0, c^2) : I \rightarrow$

$\mathbb{R}^3$  denote a curve in the  $x$ - $z$  plane. Its image, when rotated about the  $z$ -axis, generates a surface  $M$  which can be parametrized by  $\mathbf{h} : I \times [0, 2\pi] \rightarrow \mathbb{R}^3$ , with

$$\mathbf{h}(u, v) = (c^1(u) \cos v, c^1(u) \sin v, c^2(u)).$$

- (a) Show that any meridian  $t \mapsto \mathbf{h}(t, v_0)$ ,  $v_0 \in [0, 2\pi]$ , is a geodesic.
- (b) Let  $\gamma : [0, 2\pi] \rightarrow M$  denote a parallel; i.e.,  $\gamma(t) = \mathbf{h}(u_0, t)$  for some  $u_0 \in I$ . Prove that  $\gamma$  is a geodesic if and only if  $u_0$  is a critical point of  $c^1$ . Geometrically, this means that the tangent line to  $\mathbf{c}$  is vertical at the point  $\mathbf{c}(u_0)$  where  $\mathbf{c}$  intersects the parallel circle.
- (c) If  $\mathbf{c}$  is given by

$$\mathbf{c}(t) = (t, 0, -\frac{h}{R}t + h), \quad t \in (0, R), \quad h > 0,$$

so that the image of  $\mathbf{c}$  is a line segment joining  $(0, 0, h)$  to  $(R, 0, 0)$ , then the resulting surface is a cone, and by (b), parallels are not geodesics. Determine explicitly the geodesic of  $M$  that points in the parallel direction (i.e., the geodesic whose initial tangent vector equals that of some parallel  $\gamma$  in (b)). *Hint:* Show first that the cone is flat.



## 4 Integration on Euclidean space

We have postponed the subject of integration until now because our aim is to not only integrate real-valued functions on Euclidean space, but also differential forms on manifolds. The latter are particularly suited to formulate Stokes' theorem. We begin with the former, which closely parallels integration of single-variable functions.

### 4.1 The integral of a function over a box

Recall that a partition  $P$  of an interval  $[a, b] \subset \mathbb{R}$  is just a finite subset of  $[a, b]$  that contains the endpoints of the interval. This subset is then ordered:  $a = t_0 < t_1 < \dots < t_k = b$ , and each  $[t_{i-1}, t_i]$ ,  $i = 1, \dots, k$ , is referred to as a *subinterval* of the partition. Given partitions  $P_i$  of  $[a_i, b_i]$ ,  $i = 1, \dots, n$ , the product  $P = P_1 \times \dots \times P_n$  is called a *partition* of the box  $A = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_n, b_n] \subset \mathbb{R}^n$ . If  $P_i$  partitions  $[a_i, b_i]$  into  $k_i$  subintervals, then  $P$  partitions  $A$  into  $k_1 k_2 \dots k_n$  boxes, which will be called the *subboxes* of  $P$ . Each subbox is then a Cartesian product  $J_1 \times \dots \times J_n$ , where  $J_i$  is a subinterval of the partition  $P_i$ . By abuse of notation, we write  $B \in P$  if  $B$  is a subbox of  $P$ . Definition 1.6.1 implies that the volume of  $A$  equals  $\prod_{i=1}^n (b_i - a_i)$ .

Just as the integral  $\int_a^b f$  of a nonnegative function  $f$  is meant to represent the area under the graph of  $f$  between  $a$  and  $b$ , the integral  $\int_A f$  of a nonnegative function  $f$  of  $n$  variables will be interpreted as the volume of the  $(n + 1)$ -dimensional solid that lies under the graph of  $f$  and above  $A$ . In the same vein, we begin by approximating this volume by sums of volumes of boxes.

Let  $A$  be a box as above,  $P$  a partition of  $A$ , and  $f$  a bounded real-valued function whose domain contains  $A$ . For each  $B \in P$ , set

$$m_B(f) = \inf\{f(\mathbf{b}) \mid \mathbf{b} \in B\}, \quad M_B(f) = \sup\{f(\mathbf{b}) \mid \mathbf{b} \in B\}.$$

The *lower* and *upper sums* of  $f$  for  $P$  are respectively given by

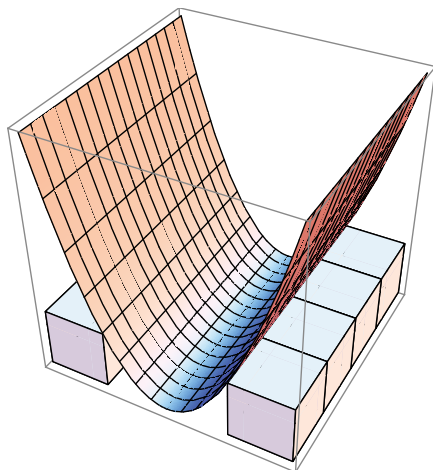
$$L(f, P) = \sum_{B \in P} m_B(f) \text{vol}(B), \quad U(f, P) = \sum_{B \in P} M_B(f) \text{vol}(B).$$

Of course  $L(f, P) \leq U(f, P)$ , but a stronger property holds. A partition  $\tilde{P}$  is said to be a *refinement* of  $P$  if  $P \subset \tilde{P}$ . Two easy but important observations are in order: firstly, if  $\tilde{P}$  refines  $P$ , then  $L(f, P) \leq L(f, \tilde{P})$  and  $U(f, \tilde{P}) \leq U(f, P)$ . To establish the first inequality, notice that any  $B \in P$  equal a union of subboxes  $\tilde{B}_1, \dots, \tilde{B}_k$  of  $\tilde{P}$ , so that

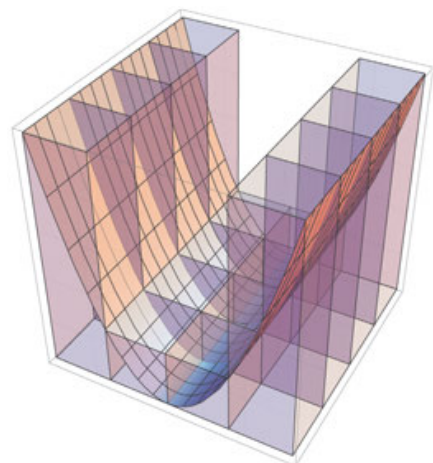
$$m_B(f) \text{vol}(B) = m_B(f) \sum_{i=1}^k \text{vol}(\tilde{B}_i) \leq \sum_{i=1}^k m_{\tilde{B}_i}(f) \text{vol}(\tilde{B}_i).$$

(The last inequality follows from the fact  $m_B(f) \leq m_{\tilde{B}_i}(f)$  since  $\tilde{B}_i \subset B$ ). The sum over all  $B$  in  $P$  of the left side is  $L(f, P)$ , and the corresponding sum of the right side equals  $L(f, \tilde{P})$ . The argument for upper sums is similar.

The second observation is that given any two partitions of  $A$ , there exists one that refines both: indeed, if  $P = P_1 \times \cdots \times P_n$  and  $\tilde{P} = \tilde{P}_1 \times \cdots \times \tilde{P}_n$  are the two partitions, then  $(P_1 \cup \tilde{P}_1) \times \cdots \times (P_n \cup \tilde{P}_n)$  is one such refinement.



**Fig. 4.1:** A lower sum for  $(x, y) \mapsto x^2$  on  $[-2, 2] \times [-2, 2]$ , using a partition with squares of side length 1.



**Fig. 4.2:** The corresponding upper sum.

**Proposition 4.1.1.** For any two partitions  $P$  and  $\tilde{P}$  of a box  $A$ ,  $L(f, P) \leq U(f, \tilde{P})$ .

*Proof.* If  $\tilde{P}$  is a refinement of both  $P$  and  $\tilde{P}$ , then

$$L(f, P) \leq L(f, \tilde{P}) \leq U(f, \tilde{P}) \leq U(f, \tilde{P}). \quad \square$$

By Proposition 4.1.1, the supremum  $L(f, A)$  of the collection of lower sums of  $f$  for all possible partitions of  $A$  exists, and is called the *lower integral of  $f$  over  $A$* . Similarly, the infimum  $U(f, A)$  of all upper sums exists, and is called the *upper integral of  $f$  over  $A$* . Again by the proposition, we always have that  $L(f, A) \leq U(f, A)$ . Furthermore, by a fundamental property of suprema – see Appendix A, given any  $\varepsilon > 0$ , there exists a partition  $P$  for which  $L(f, P) > L(f, A) - \varepsilon$ . A similar property holds for infima.

**Definition 4.1.1.** Let  $f$  be a bounded function on a box  $A$ .  $f$  is said to be *integrable over*  $A$  if the lower and upper integrals of  $f$  over  $A$  are equal. In this case, their common value is defined to be the *integral of  $f$  over  $A$* , denoted  $\int_A f$ .

The following is often a useful tool in deciding whether or not a function is integrable:

**Theorem 4.1.1.** Let  $f$  denote a bounded function over a box  $A$ .  $f$  is integrable over  $A$  if and only if for any  $\varepsilon > 0$ , there exists a partition  $P$  of  $A$  such that  $U(f, P) - L(f, P) < \varepsilon$ .

*Proof.* Suppose  $f$  is integrable over  $A$ , and  $\varepsilon > 0$  is given. There exists a partition  $P$  of  $A$  such that  $U(f, P) < U(f, A) + \varepsilon/2$ . Similarly, there exist a partition  $P'$  of  $A$  with  $L(f, P') > L(f, A) - \varepsilon/2$ . Since  $U(f, A) = L(f, A)$ ,  $U(f, P) - L(f, P') < \varepsilon$ . Thus, if  $\tilde{P}$  refines both  $P$  and  $P'$ , then

$$U(f, \tilde{P}) - L(f, \tilde{P}) \leq U(f, P) - L(f, P') < \varepsilon.$$

Conversely, suppose that for any  $\varepsilon > 0$ , there exists a partition  $P$  of  $A$  such that  $U(f, P) - L(f, P) < \varepsilon$ . Then the nonnegative number  $U(f, A) - L(f, A)$  is less than  $\varepsilon$  for any  $\varepsilon > 0$ , and must therefore equal 0; i.e.,  $f$  is integrable.  $\square$

**Examples and Remarks 4.1.1.** (i) Let  $c \in \mathbb{R}$ , and  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  denote the constant function  $f(\mathbf{a}) = c$ ,  $\mathbf{a} \in \mathbb{R}^n$ . Then for any box  $B$ ,  $m_B(f) = M_B(f) = c$ , so that for any partition  $P$  of  $A$ ,  $L(f, P) = U(f, P) = c \cdot \text{vol}(A)$ .  $f$  is therefore integrable, and  $\int_A f = c \cdot \text{vol}(A)$ .

(ii) Consider the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by  $f(x, y) = xy$ , and the square  $A = [0, 1] \times [0, 1]$ . Denote by  $P_n$  the partition of  $A$  into  $n^2$  squares of equal sides; i.e.,  $P_n = \{(i/n, j/n) \mid 0 \leq i, j \leq n\}$ . If  $B_{ij}$  denotes the subsquare  $[(i-1)/n, i/n] \times [(j-1)/n, j/n]$ , then the minimum of  $f$  on  $B_{ij}$  occurs at the lower left corner  $((i-1)/n, (j-1)/n)$ , and the maximum at the upper right corner  $(i/n, j/n)$ , so that  $m_{B_{ij}}(f) = (i-1)(j-1)/n^2$  and  $M_{B_{ij}}(f) = ij/n^2$ . Thus,

$$\begin{aligned} L(f, P_n) &= \sum_{i,j=1}^n \frac{(i-1)(j-1)}{n^4} = \frac{1}{n^4} \sum_{j=1}^n (j-1) \frac{n(n-1)}{2} = \frac{n^2(n-1)^2}{4n^4} \\ &= \frac{1}{4} \left(1 - \frac{1}{n}\right)^2, \end{aligned}$$

and similarly,

$$U(f, P_n) = \sum_{i,j=1}^n \frac{ij}{n^4} = \frac{1}{4} \left(1 + \frac{1}{n}\right)^2.$$

It follows that  $U(f, P_n) - L(f, P_n) = 1/n$ , and  $f$  is integrable by Theorem 4.1.1. Furthermore,  $L(f, P_n) \leq 1/4 \leq U(f, P_n)$  for all  $n$ , so that  $\int_A f = 1/4$ .

(iii) Let  $A = [0, 1] \times [0, 1]$  as in (ii), and  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by  $f(x, y) = 1$  if  $x$  and  $y$  are both rational, and 0 otherwise. Given any partition  $P$  of  $A$ , an arbitrary box  $B$  of  $P$  will contain points whose coordinates are rational, and points with at least one

irrational coordinate. Thus,  $m_B(f) = 0$  and  $M_B(f) = 1$ . This implies that  $L(f, A) = L(f, P) = 0$  and  $U(f, A) = U(f, P) = 1$ , so that  $f$  is not integrable.

(iv) A word about notation: It is common practice, when given a specific formula for a function, to include that formula in the integral. For example, the integral of the function  $f$  in (ii) is also denoted  $\int_A xy \, dx \, dy$ , or  $\int_A u^1 u^2 \, du^1 \, du^2$ . This is rather unfortunate since  $du^1$  is really the differential of the function  $u^1$  and has a completely different meaning here. Nevertheless, this practice is so widespread that we will conform to it.

**Theorem 4.1.2.** *Suppose  $f$  and  $g$  are integrable over  $A$ .*

- (1) *If  $f \leq g$ , then  $\int_A f \leq \int_A g$ ;*
- (2) *If  $m$  and  $M$  are lower and upper bounds respectively of  $f$  on  $A$ , then  $m \operatorname{vol}(A) \leq \int_B f \leq M \operatorname{vol}(A)$ ;*
- (3)  *$|f|$  is integrable, and  $|\int_A f| \leq \int_A |f|$ ;*
- (4)  *$f + g$  is integrable, and  $\int_A (f + g) = \int_A f + \int_A g$ ;*
- (5) *Given  $c \in \mathbb{R}$ ,  $cf$  is integrable, and  $\int_A cf = c \int_A f$ .*

*Proof.* If  $f \leq g$ , then  $L(f, P) \leq L(g, P)$  for any partition  $P$  of  $A$ . The first statement follows by taking suprema over all partitions of  $A$ . The second one follows from the first, since  $f$  is squeezed in between the constant functions  $m$  and  $M$ . For the third one, observe first that for any box  $B$ ,

$$|f(\mathbf{a})| - |f(\mathbf{b})| \leq M_B(f) - m_B(f), \quad \mathbf{a}, \mathbf{b} \in B. \tag{4.1.1}$$

This inequality is (tediously) verified on a case by case basis: for example, suppose  $f(\mathbf{a}) \leq 0 \leq f(\mathbf{b})$ . Then  $M_B(f) \geq 0$ , and

$$|f(\mathbf{a})| - |f(\mathbf{b})| = -f(\mathbf{a}) - f(\mathbf{b}) \leq -f(\mathbf{a}) \leq -m_B(f) \leq M_B(f) - m_B(f).$$

Fixing an arbitrary  $\mathbf{b}$  in (4.1.1) and taking the supremum over all  $\mathbf{a} \in B$  implies that  $M_B(|f|) - |f(\mathbf{b})| \leq M_B(f) - m_B(f)$ , or equivalently, that

$$|f(\mathbf{b})| \geq M_B(|f|) - (M_B(f) - m_B(f)), \quad \mathbf{b} \in B.$$

Finally, taking the infimum over all  $\mathbf{b} \in B$  and rearranging terms yields

$$M_B(|f|) - m_B(|f|) \leq M_B(f) - m_B(f).$$

Thus, given any partition  $P$  of  $A$ , we have

$$U(|f|, P) - L(|f|, P) \leq U(f, P) - L(f, P).$$

It now follows from Theorem 4.1.1 that  $|f|$  is integrable whenever  $f$  is. The inequality  $|\int_A f| \leq \int_A |f|$  is obtained by applying the first assertion of the theorem to  $-|f| \leq f \leq |f|$ . The last two statements – like the others – have proofs that are similar to those for functions of one variable, and are left as an exercise.  $\square$

In the context of Theorem 4.1.2, it is also true that a product of integrable functions is integrable, as well as a quotient (provided the bottom function is nonzero). This turns out to be an immediate consequence of results from the next section.

## 4.2 Integrability and discontinuities

Our next goal is to characterize those functions that are integrable, and define integration over regions that are more complicated than boxes. This, in turn, will lead us to seek a more general concept of integral.

**Definition 4.2.1.** A set  $A \subset \mathbb{R}^n$  is said to have *measure zero* if for any  $\varepsilon > 0$ ,  $A$  can be covered by a countable collection  $B_1, B_2, \dots$  of open boxes with  $\sum_{i=1}^{\infty} \text{vol}(B_i) < \varepsilon$ .

Since a closed box has the same volume as its interior, a set of measure zero may also be covered by a countable collection of closed boxes satisfying the above volume condition. The converse is easily verified, so that one may replace open by closed in the definition. Given a collection  $\{B_i\}$  of boxes, we will refer to the number  $\sum_i \text{vol}(B_i)$  (if it exists) as the *total volume* of the collection. This is just for convenience's sake, since already for a finite collection, the total volume will in general be larger than the volume of the union if the boxes intersect.

**Examples 4.2.1.** (i) Any countable set of points has measure zero: if  $\mathbf{a}_1, \mathbf{a}_2, \dots$  is a sequence and  $\varepsilon$  is a positive number, let  $B_i$  be a box of volume smaller than  $\varepsilon/2^i$  containing  $\mathbf{a}_i$ . Then this collection has total volume less than  $\varepsilon$ . For example, the set  $\mathbb{Q}$  of rationals has measure zero.

(ii) A countable union of sets  $A_i$  of measure zero has measure zero: given  $\varepsilon > 0$ , cover each  $A_i$  by boxes  $B_{i1}, B_{i2}, \dots$  such that  $\sum_j \text{vol}(B_{ij}) < \varepsilon/2^i$ . The collection  $\{B_{ij} \mid i, j = 1, 2, \dots\}$  is then a countable cover of  $\cup A_i$  (see Appendix A), and may therefore be arranged into a sequence  $C_1, C_2, \dots$ . There are, of course, many ways to do this, but if  $\sum \text{vol}(C_i)$  converges, then any rearrangement also converges to the same limit, since the convergence is absolute. To see that it does converge, let  $k \in \mathbb{N}$ . Then there exists an integer  $l$  such that all the *terms*  $C_1, \dots, C_k$  appear in the *sequences*  $B_{1i}$  through  $B_{li}$ ,  $i = 1, \dots, \infty$ . Thus,

$$\begin{aligned} \text{vol}(C_1) + \dots + \text{vol}(C_k) &\leq \sum_{i=1}^{\infty} \text{vol}(B_{1i}) + \dots + \sum_{i=1}^{\infty} \text{vol}(B_{li}) \\ &< \frac{\varepsilon}{2} + \dots + \frac{\varepsilon}{2^l} \\ &< \varepsilon. \end{aligned}$$

This means that the sequence of partial sums is bounded above by  $\varepsilon$ , and therefore  $\sum_{i=1}^{\infty} \text{vol}(C_i) \leq \varepsilon$ .

(iii) A (non-degenerate) box does not have measure zero. This may seem an obvious statement, but with only the definition to work with, it does require an argument. To see this, it suffices to consider a closed box  $A$ . Suppose  $B_1, B_2, \dots$  is a collection of boxes covering  $A$ . We claim that  $\sum \text{vol}(B_i) \geq \text{vol}(A)$ . Indeed,  $A$  is compact, so a finite subcollection, which we may rename as  $B_1, \dots, B_k$ , covers  $A$ . This subcollection defines a partition  $P$  of  $A$ ; namely,  $P$  consists of all points whose  $i$ -th

coordinate is either the beginning or the end point of the projection of some  $B_j \cap A$  onto the  $i$ -th coordinate axis. Specifically, if  $B_j \cap A = \Pi_{i=1}^n [a_i^j, b_i^j]$ , then

$$P = \{(x_1, \dots, x_n) \mid \text{for each } i = 1, \dots, n, x_i = a_i^j \text{ or } b_i^j \text{ for some } j\}.$$

By construction, any box in the partition  $P$  is contained in each of the  $B_j$  it intersects, and the union of these boxes equals  $A$ , so that

$$\sum_i \text{vol}(B_i) \geq \sum_{i=1}^k \text{vol}(B_i) \geq \sum_{B \in P} \text{vol}(B) = \text{vol}(A).$$

This establishes the assertion, for if  $A$  had measure zero, then by the above inequality,  $A$  would have volume less than any positive number  $\epsilon$ , and would then be a degenerate box.

- (iv) The collection of all irrational numbers between 0 and 1 does not have measure zero: if it did, then by (i) and (ii), the interval  $[0, 1]$  would have measure zero, contradicting (iii).
- (v) We will later see that a ball of radius  $r$  in  $\mathbb{R}^n$  has volume  $k_n r^n$ , where  $k_n$  is a constant depending on  $n$ , see also Exercise 4.25. Assuming this, we could use open or closed balls instead of boxes in the definition of measure zero, since

$$[-r, r]^n \subset B_{r\sqrt{n}}(\mathbf{0}) \subset [-r\sqrt{n}, r\sqrt{n}]^n.$$

**Remark 4.2.1.** The type of construction referred to in part (iii) of the above examples will be used many times. One variant of it is worth emphasizing: If  $B_1, \dots, B_k$  are boxes contained in a larger box  $B$ , then there exists a partition  $P$  of  $B$  such that each  $B_i$  is a union of subboxes of  $P$ . The proof is the same as the one used in (iii).

Now that the preliminaries have been dealt with, we are in a position to characterize those functions that are integrable. In light of our work with upper and lower sums, it shouldn't come as a surprise to learn that they are the ones with not too many jumps. Specifically:

**Theorem 4.2.1.** *Let  $f : A \rightarrow \mathbb{R}$  denote a function that is bounded on the closed box  $A$ . Then  $f$  is integrable over  $A$  if and only if the set of discontinuities of  $f$  has measure zero.*

*Proof.* Denote by  $D$  the set of points where  $f$  is discontinuous, and for  $\delta > 0$ , by  $D_\delta$  the set of all  $\mathbf{a} \in A$  such that  $M_U(f) - m_U(f) \geq \delta$  for any neighborhood  $U$  of  $\mathbf{a}$ . Then  $D_\delta \subset D$  for any  $\delta$ , and  $D = \bigcup_{i=1}^\infty D_{1/i}$  by Exercise 1.53.

Suppose first that  $f$  is integrable. To show that  $D$  has measure zero, it suffices to show that each  $D_{1/k}$  has measure zero. Now, given any partition  $P$  of  $A$ ,  $D_{1/k}$  may be written as a union

$$D_{1/k} = \left( D_{1/k} \cap \bigcup_{B \in P} B^0 \right) \cup \left( D_{1/k} \cap \bigcup_{B \in P} \partial B \right)$$

of two sets, the second of which has measure zero. Thus, given  $\epsilon > 0$ , it suffices to produce a partition  $P$  of  $A$  for which those boxes that intersect  $D_{1/k}$  in their interior have

total volume less than  $\varepsilon$ . We claim that any partition  $P$  of  $A$  for which  $U(f, P) - L(f, P) < \varepsilon/k$  will do. Indeed, denote by  $\mathcal{C}$  the collection of boxes in  $P$  that intersect  $D_{1/k}$  in their interior. Then  $M_B(f) - m_B(f) \geq 1/k$  for each  $B \in \mathcal{C}$ , so that

$$\begin{aligned} \sum_{B \in \mathcal{C}} \text{vol}(B) &\leq k \cdot \sum_{B \in \mathcal{C}} (M_B(f) - m_B(f)) \text{vol}(B) \\ &\leq k \cdot \sum_{B \in P} (M_B(f) - m_B(f)) \text{vol}(B) = k(U(f, P) - L(f, P)) \\ &< \varepsilon, \end{aligned}$$

and  $D_{1/k}$  has measure zero, as claimed.

Conversely, suppose  $D$  has measure zero. Given  $\varepsilon > 0$ , we must exhibit a partition  $P$  of  $A$  satisfying  $U(f, P) - L(f, P) < \varepsilon$ . Let  $M$  be an upper bound of  $|f|$  on  $A$ , and set  $\delta := \varepsilon/(2M + \text{vol}(A))$ . Cover  $D$  by open boxes  $B_i$  with  $\sum_i \text{vol}(B_i) < \delta$ . Similarly, for each  $\mathbf{a} \in A \setminus D$ , choose some box  $C_{\mathbf{a}}$  that contains  $\mathbf{a}$  in its interior and satisfies  $M_{C_{\mathbf{a}}}(f) - m_{C_{\mathbf{a}}}(f) < \delta$ . The two collections yield a cover of  $A$ . Choose a finite subcover, and a partition  $P$  of  $A$  such that each subbox in  $P$  is contained in one of the boxes from the subcover, see Remark 4.2.1.

If  $B$  is a subbox contained in some  $B_i$ , then  $(M_B(f) - m_B(f)) \text{vol}(B) < 2M \text{vol}(B)$ , and if  $B$  is contained in some  $C_{\mathbf{a}}$ , then  $(M_B(f) - m_B(f)) \text{vol}(B)$  is less than  $\delta \text{vol}(B)$ . Thus,

$$U(f, P) - L(f, P) < 2M\delta + \delta \text{vol}(A) = \varepsilon. \quad \square$$

We are now in a position to integrate over regions that are more general than boxes:

**Definition 4.2.2.** The *characteristic function*  $\chi_A$  of a bounded set  $A \subset \mathbb{R}^n$  is defined by  $\chi_A(\mathbf{a}) = 1$  if  $\mathbf{a} \in A$  and  $\chi_A(\mathbf{a}) = 0$  otherwise.

$A$  is said to be *Jordan-measurable* if  $\int_B \chi_A$  exists, with  $B$  denoting some box that contains  $A$ . In this case, the value of this integral is called the  *$n$ -dimensional volume* of  $A$ .

**Remark 4.2.2.** If  $A$  is a bounded Jordan-measurable set of measure zero, then it has zero volume: consider any partition  $P$  of some box  $B$  containing  $A$ . The infimum of  $\chi_A$  on any subbox of  $P$  must be zero, since the only other possibility is 1. The latter is ruled out, for it would imply that the subbox is contained in  $A$ , contradicting the fact that  $A$  has measure zero. Thus  $L(\chi_A, P) = 0$  for any  $P$ , and the claim follows.

Notice, however, that there exist bounded sets of measure zero which are not Jordan-measurable. One such is  $A = \mathbb{Q} \cap [0, 1] \subset \mathbb{R}$ .

**Definition 4.2.3.** A bounded  $f : A \rightarrow \mathbb{R}$  is said to be *integrable over  $A$*  if  $\int_B f \cdot \chi_A$  exists, with  $B$  denoting any box containing  $A$ . The value of this integral is denoted  $\int_A f$ .

Clearly, the above definition does not depend on the particular box  $B$ , and  $f$  is integrable if  $\int_B f \cdot \chi_A$  exists for *some* box  $B$  containing  $A$ .

**Theorem 4.2.2.** A bounded set  $A$  is Jordan-measurable if and only if the boundary of  $A$  has measure zero.

*Proof.* If  $B$  is a closed box whose interior contains  $A$ , then the restriction of  $\chi_A$  to  $B$  has exactly the boundary of  $A$  as its set of discontinuities, since any neighborhood of a point on the boundary contains points of  $A$  (where  $\chi_A$  equals 1) and points outside  $A$  (where it equals zero). The claim now follows from Theorem 4.2.1.  $\square$

**Examples 4.2.2.** (i) Any region in  $\mathbb{R}^2$  whose boundary is (the image of) a bounded curve is Jordan-measurable: let  $\mathbf{c} : [0, b] \rightarrow \mathbb{R}^2$  be a piecewise-smooth curve parametrizing the boundary, and let  $L = \int_0^b |\mathbf{c}'|$  denote its length. Given  $\varepsilon > 0$ , we must find a cover of the boundary by rectangles with total volume (which we call area in dimension 2) less than  $\varepsilon$ . So let  $n$  be an integer large enough so that  $(n + 1)/n^2 < \varepsilon/9L^2$ . Since the function  $f : [0, b] \rightarrow \mathbb{R}$  given by  $f(t) = \int_0^t |\mathbf{c}'|$  is continuous, there exists a partition  $t_0 = 0 < t_1 < \dots < t_n = b$  of  $[0, b]$  such that  $f(t_i) = iL/n$ ,  $0 \leq i \leq n$ ; i.e., the portion of the curve joining the points  $\mathbf{p}_{i-1} = \mathbf{c}(t_{i-1})$  and  $\mathbf{p}_i = \mathbf{c}(t_i)$  has length  $L/n$ . In particular, the distance between these two is no larger than  $L/n$ , so that any point on the boundary is at distance less than or equal to  $L/n$  from some  $\mathbf{p}_i$ . It follows that the collection  $\mathcal{C}$  of squares centered at  $\mathbf{p}_i$  with sides of length  $3L/n$  covers  $\mathcal{C}$ , and

$$\sum_{S \in \mathcal{C}} \text{vol}(S) = (n + 1) \frac{9L^2}{n^2} < \varepsilon,$$

which establishes the claim.

(ii) One would hope that any bounded open set in  $\mathbb{R}^n$  is Jordan-measurable. This is not the case, however, even in dimension 1, a fact that will lead us to generalize the concept of integral in order to avoid these exceptions. One such is obtained by arranging all rational numbers in  $(0, 1)$  in a sequence  $x_1, x_2, \dots$ , and setting

$$U_i = \left( x_i - \frac{1}{2^{i+2}}, x_i + \frac{1}{2^{i+2}} \right) \cap (0, 1), \quad U = \cup_{i=1}^{\infty} U_i.$$

Notice that  $U$  is open, so that its boundary  $\partial U$  is disjoint from  $U$ . Furthermore,  $\partial U \subset [0, 1]$ , since any  $x \notin [0, 1]$  has a neighborhood (namely  $(-\infty, 0)$  or  $(1, \infty)$ ) that does not intersect  $[0, 1]$ , and therefore does not intersect  $U$  either. Thus,  $\partial U \subset [0, 1] \setminus U$ . Conversely, if  $x \in [0, 1] \setminus U$ , then any neighborhood of  $x$  intersects both  $U$  (since it must contain some rational) and the complement of  $U$  (since it contains  $x$ ). Thus,  $\partial U = [0, 1] \setminus U$ , and  $[0, 1]$  equals the disjoint union of  $U$  and  $\partial U$ . This means that  $\partial U$  does not have measure zero: indeed,

$$\text{vol}(U) \leq \sum_i \text{vol}(U_i) \leq \sum_{i=1}^{\infty} \frac{1}{2^{i+1}} = \frac{1}{2},$$

so that if  $\partial U$  could be covered by intervals of total length less than  $1/2$ , then  $[0, 1]$  could be covered by intervals of total length smaller than 1, contradicting Examples 4.2.1 (iii).

(iii) Let  $U, V$  denote open sets in  $\mathbb{R}^n$ . A map  $h : U \rightarrow V$  is said to be a *homeomorphism* if it is bijective, continuous, and has continuous inverse. If  $h$  is a homeomorphism



and  $A \subset U$  is a compact Jordan-measurable set, then so is  $h(A)$ . The compactness property was proved in Chapter 1. The Jordan-measurability can be seen as follows: the fact that  $h$  is a homeomorphism easily implies that  $h(\partial A) = \partial h(A)$ . The restriction of  $h$  to  $A$  is uniformly continuous. Thus, for any  $\varepsilon > 0$ , there exists some  $\delta > 0$  such that any ball of radius  $\delta$  centered at a point of  $\partial A$  is mapped by  $h$  into a ball of radius  $\varepsilon$  centered at a point of  $\partial h(A)$ . The claim now follows from Examples 4.2.1 (v); alternatively, one could avoid any mention of volume of balls by rephrasing continuity in terms of boxes.

In order to avoid the awkward situation from part (ii) in the above example, we will use partitions of unity when the region of integration is an open set. Let  $f$  be a function with open domain  $U$  that is bounded in some neighborhood of any point in  $U$ , and suppose the set of discontinuities of  $f$  has measure zero. Let  $\{U_i\}$  denote an *admissible* open cover of  $U$  (meaning that each  $U_i$  is contained in  $U$ ), and  $\Phi = \{\varphi_i\}$  a countable partition of unity subordinate to the cover. An arbitrary  $\varphi \in \Phi$  is zero outside some compact set in  $U$ , so that  $\varphi|f|$  is integrable on any closed box that contains  $U$ . Suppose that the series  $\sum_i \int_U \varphi_i|f|$  converges. Since  $|\int_U \varphi_i f| \leq \int_U \varphi_i|f|$ , this in turn implies absolute convergence of  $\sum_i \int_U \varphi_i f$ . Now,  $\sum \varphi_i \equiv 1$ , so it is tempting to define  $\int_U f$  as the sum of this series. It must first be checked, though, that a different partition of unity  $\Psi = \{\psi_j\}$  produces the same sum. This is not difficult to show, for if  $\varphi \in \Phi$ , then  $\int_U \varphi f = \int_U \sum_j \psi_j \varphi f$ . Furthermore, each point in the support of  $\varphi$  has a neighborhood on which only finitely many  $\psi_j$  are nonzero. Since the support is compact, only finitely many  $\psi_j$  are nonzero on it. Thus,

$$\int_U \varphi f = \sum_{j=1}^{\infty} \int_U \psi_j \varphi f,$$

and consequently,

$$\sum_{i=1}^{\infty} \int_U \varphi_i f = \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \int_U \psi_j \varphi_i f.$$

Applying this identity to  $|f|$  shows that the convergence is absolute. Similarly,

$$\sum_{j=1}^{\infty} \int_U \psi_j f = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} \int_U \psi_j \varphi_i f.$$

Absolute convergence implies that the order of summation may be interchanged (see Examples 4.2.1 (ii)), so that

$$\sum_{\varphi \in \Phi} \int_U \varphi f = \sum_{\psi \in \Psi} \int_U \psi f.$$

Notice that the sum is not only independent of the partition of unity, it is also independent of the open cover to which it is subordinate, since this cover was not used in the argument.

**Lemma 4.2.1.** *Let  $U \subset \mathbb{R}^n$  be an open bounded set,  $f : U \rightarrow \mathbb{R}$  a function that is bounded and has a set of discontinuities of measure zero. Given any admissible open cover of  $U$  and partition of unity  $\{\varphi_i\}_{i \in \mathbb{N}}$  subordinate to this cover, the series*

$$\sum_{i=1}^{\infty} \int_U \varphi_i |f|$$

*converges, and the sum is independent of the particular cover and partition of unity chosen.*

*Proof.* Let  $B$  denote a closed box that contains  $U$ . Each  $\varphi_i |f|$  extends to a function (denoted in the same way) on  $B$  which is smooth outside  $U$  by setting it equal to zero there. If  $M$  denotes an upper bound of  $|f|$  on  $U$ , the partial sums of the series are bounded above, because

$$\sum_{i=1}^n \int_U \varphi_i |f| = \int_U \sum_{i=1}^n \varphi_i |f| = \int_B \sum_{i=1}^n \varphi_i |f| \leq \int_B |f| \leq M \operatorname{vol}(B).$$

Thus, the series converges. Independence of the cover and partition of unity was established earlier.  $\square$

**Definition 4.2.4.** Let  $f : U \rightarrow \mathbb{R}$ , where  $U$  and  $f$  satisfy the hypotheses of Lemma 4.2.1. The *generalized integral of  $f$  over  $U$*  is defined to be the number

$$\int_U f := \sum_{i=1}^{\infty} \int_U \varphi_i f,$$

where  $\{\varphi_i\}$  is any partition of unity subordinate to some admissible open cover of  $U$ .

Of course, if this definition is to be of any value, it should coincide with the old one when the region of integration is Jordan-measurable.

**Theorem 4.2.3.** *If  $f$  is a bounded function which is integrable over some bounded Jordan-measurable set  $A$ , then the original  $\int_A f$  coincides with the generalized one from Definition 4.2.4.*

*Proof.* Let  $\Phi = \{\varphi_1, \varphi_2, \dots\}$  be a partition of unity for  $A$ ,  $\varepsilon > 0$ , and  $\int_A f$  denote the original integral. The claim follows once we show that there exists a natural number  $N$  such that

$$\left| \int_A f - \sum_{j=1}^k \int_A \varphi_j f \right| < \varepsilon \text{ for all } k \geq N.$$

So consider a box  $B$  that contains  $A$ , and let  $M$  denote an upper bound of  $|f|$  on  $B$ . Since  $\partial A$  is compact and has measure zero, there exist boxes  $B_1, \dots, B_k \subset B$  covering  $\partial A$  with total volume less than  $\varepsilon/M$ . Take a partition  $P$  of  $B$  where each  $B_i$  is a union of subboxes of  $P$  (see Remark 4.2.1), and denote by  $C$  the union of the subboxes of  $P$  that

lie inside  $A$ . Then

$$\text{vol}(A \setminus C) \leq \sum_{i=1}^k \text{vol}(B_i) < \frac{\varepsilon}{M}.$$

Furthermore, by compactness of  $C$ , the collection of all  $\varphi \in \Phi$  that are not identically zero on  $C$  is finite. This means that there exists an integer  $N$  such that for  $k \geq N$ , the restriction of  $\varphi_k$  to  $C$  is zero, and therefore,  $\int_A \varphi_k f = \int_{A \setminus C} \varphi_k f$ . But then for any  $k \geq N$ ,

$$\begin{aligned} \left| \int_A f - \sum_{j=1}^k \int_A \varphi_j f \right| &= \left| \int_A \left( f - \sum_{j=1}^k \varphi_j f \right) \right| \leq \int_A \left| f - \sum_{j=1}^k \varphi_j f \right| \\ &= \int_A \left( 1 - \sum_{j=1}^k \varphi_j \right) |f| \leq M \int_A \left( 1 - \sum_{j=1}^k \varphi_j \right) \\ &= M \int_A \sum_{j>k} \varphi_j \leq M \int_{A \setminus C} 1 = M \text{vol}(A \setminus C) \\ &< \varepsilon. \end{aligned}$$

□

As before, we define the *volume*  $\text{vol}(A)$  of  $A \subset \mathbb{R}^n$  to be  $\int_A 1$ . Thus, the volume of a bounded Jordan-measurable set coincides with the original notion of volume, but Lemma 4.2.1 now guarantees that *any bounded open set*, even a non Jordan-measurable one (see Examples 4.2.2), has a well-defined volume.

**Remarks 4.2.3.** (i) The proof of Theorem 4.2.3 shows that if  $A$  is a Jordan-measurable bounded set, then for any  $\varepsilon > 0$  there exists a compact Jordan-measurable subset  $C$  of  $A$  such that  $\text{vol}(A \setminus C) < \varepsilon$ .

(ii) If  $U$  is open in  $\mathbb{R}^n$ , and  $f, g : U \rightarrow \mathbb{R}$  are equal everywhere except on a set of measure zero, then  $f$  is integrable on  $U$  if and only if  $g$  is, and in that case, both integrals are equal. This is because  $h = f - g$  vanishes except on a set of measure zero and is therefore integrable. Thus, if  $f$  is integrable, then so is  $g = f - h$ , and vice-versa.

### 4.3 Fubini's theorem

Any box  $B$  in  $\mathbb{R}^{m+n}$  decomposes as a product  $B_1 \times B_2 = \pi_1(B) \times \pi_2(B) \subset \mathbb{R}^m \times \mathbb{R}^n$  of boxes, with  $\pi_i$  denoting the orthogonal projection of  $\mathbb{R}^m \times \mathbb{R}^n$  onto each factor,  $i = 1, 2$ . Given  $f : B \rightarrow \mathbb{R}$ , each  $\mathbf{x} \in B_1$  induces a function  $f_{\mathbf{x}} : B_2 \rightarrow \mathbb{R}$  by setting  $f_{\mathbf{x}}(\mathbf{y}) = f(\mathbf{x}, \mathbf{y})$ . As an elementary example, let  $B = [0, 2] \times [0, 2]$ , and  $f(x, y) = (x + 1)y^2$ . Then  $f_0(x) = y^2$  and  $f_1(x) = 2y^2$ .

Suppose next that for each  $\mathbf{x} \in B_1$ ,  $f_{\mathbf{x}}$  is integrable on  $B_2$ . This yields a new function  $g$  on  $B_1$  by setting  $g(\mathbf{x}) = \int_{B_2} f_{\mathbf{x}}$ . This expression is commonly denoted  $\int_{B_2} f(\mathbf{x}, \mathbf{y}) \, d\mathbf{y}$ . In the above example,

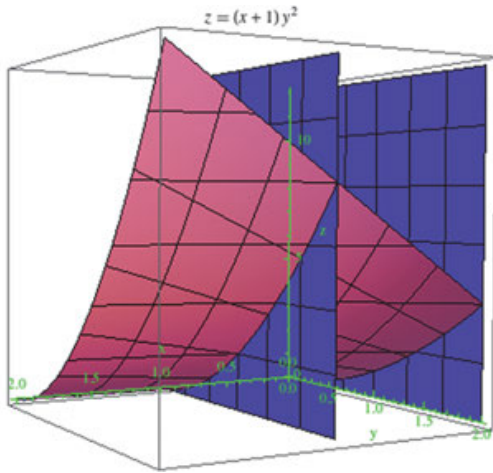


Fig. 4.3: Cross-sectional areas  $\int_{B_2} f(x, y) dy$  for  $x = 1$  and  $x = 0$ .

$$\int_{B_2} f(x, y) dy = \int_0^2 (x + 1)y^2 dy = (x + 1) \frac{y^3}{3} \Big|_{y=0}^{y=2} = \frac{8(x + 1)}{3}.$$

If this function  $g$  is integrable over  $B_1$ , its integral is called an *iterated integral*, and is denoted  $\int_{B_1} \int_{B_2} f(x, y) dy dx$ . In our example, the reader can verify that this integral equals  $32/3$ , and that reversing the order of integration (that is, evaluating  $\int_{B_2} \int_{B_1} f(x, y) dx dy$  instead) yields the same result. This is no coincidence. Fubini's theorem implies that for continuous  $f$ , both iterated integrals are equal, and furthermore equal the integral  $\int_B f$  of the original function over  $B$ , thereby substantially simplifying the evaluation of these integrals. The theorem actually holds for arbitrary integrable functions, although the statement becomes more complicated because iterated integrals do not always exist.

Recall that a bounded function  $f : B \rightarrow \mathbb{R}$  always admits lower and upper integrals  $L(f, B)$ ,  $U(f, B)$ , regardless of whether it is integrable. Also recall that a partition  $P$  of  $B_1 \times B_2$  induces partitions  $P_i = \pi_i(P)$  of  $B_i$ ,  $i = 1, 2$ , and  $P = P_1 \times P_2$ .

**Lemma 4.3.1.** *Let  $B_1$  and  $B_2$  denote rectangles in  $\mathbb{R}^m$  and  $\mathbb{R}^n$  respectively, and  $f : B_1 \times B_2 \rightarrow \mathbb{R}$  be a bounded function. For each  $\mathbf{x} \in B_1$ , define  $f_{\mathbf{x}} : B_2 \rightarrow \mathbb{R}$  by  $f_{\mathbf{x}}(\mathbf{y}) = f(\mathbf{x}, \mathbf{y})$ , and denote by  $g_1 : B_1 \rightarrow \mathbb{R}$  the function given by  $g_1(\mathbf{x}) = L(f_{\mathbf{x}}, B_2)$ . Then for any partition  $P$  of  $B$ ,*

$$L(f, P) \leq L(g_1, P_1) \leq U(g_1, P_1) \leq U(f, P),$$

where  $P_1 = \pi_1(P)$  is the induced partition of  $B_1$ .

*Proof.* Observe that if  $S_i$  are subboxes of  $P_i$ ,  $i = 1, 2$ , and  $\mathbf{x} \in S_1$ , then

$$m_{S_1 \times S_2}(f) \leq m_{\{\mathbf{x}\} \times S_2}(f) = m_{S_2}(f_{\mathbf{x}}).$$

Consequently,

$$\sum_{S_2 \in P_2} m_{S_1 \times S_2}(f) \text{ vol}(S_2) \leq \sum_{S_2 \in P_2} m_{S_2}(f_{\mathbf{x}}) \text{ vol}(S_2) = L(f_{\mathbf{x}}, P_2) \leq g_1(\mathbf{x}).$$

Since this holds for every  $\mathbf{x} \in S_1$ ,

$$\sum_{S_2 \in \mathcal{P}_2} m_{S_1 \times S_2}(f) \operatorname{vol}(S_2) \leq m_{S_1}(g_1),$$

so that

$$\begin{aligned} L(f, P) &= \sum_{S_1 \times S_2} m_{S_1 \times S_2}(f) \operatorname{vol}(S_1 \times S_2) \\ &= \sum_{S_1 \in \mathcal{P}_1} \left( \sum_{S_2 \in \mathcal{P}_2} m_{S_1 \times S_2}(f) \operatorname{vol}(S_2) \right) \operatorname{vol}(S_1) \leq \sum_{S_1 \in \mathcal{P}_1} m_{S_1}(g_1) \operatorname{vol}(S_1) \\ &= L(g_1, P_1). \end{aligned} \quad (4.3.1)$$

Next, if  $h_1 : B_1 \rightarrow \mathbb{R}$  is given by  $h_1(\mathbf{x}) = U(f_{\mathbf{x}}, P_2)$ , then a similar argument, using  $h_1$  instead of  $g_1$ , shows that  $U(h_1, P_1) \leq U(f, P)$ . But  $g_1 \leq h_1$ , so that

$$U(g_1, P_1) \leq U(f, P).$$

This, together with (4.3.1), yields the claim.  $\square$

**Remarks 4.3.1.** (i) With minor modifications, the proof of the lemma shows that

$$L(f, P) \leq L(h_1, P_1) \leq U(h_1, P_1) \leq U(f, P),$$

where  $h_1$  is the function that was used in the proof,  $h_1(\mathbf{x}) = U(f_{\mathbf{x}}, P_2)$ .

(ii) Define, for each  $\mathbf{y} \in B_2$ , a function  $f_{\mathbf{y}} : B_1 \rightarrow \mathbb{R}$  by  $f_{\mathbf{y}}(\mathbf{x}) = f(\mathbf{x}, \mathbf{y})$ , and denote by  $g_2, h_2 : B_2 \rightarrow \mathbb{R}$  the functions given by

$$g_2(\mathbf{y}) = L(f_{\mathbf{y}}, B_1), \quad h_2(\mathbf{y}) = U(f_{\mathbf{y}}, B_1).$$

Arguments similar to those used in the proof of the lemma imply that

$$L(f, P) \leq L(g_2, P_2) \leq U(g_2, P_2) \leq U(f, P),$$

and

$$L(f, P) \leq L(h_2, P_2) \leq U(h_2, P_2) \leq U(f, P).$$

**Theorem 4.3.1 (Fubini's Theorem).** *Suppose  $f : B_1 \times B_2 \subset \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$  is integrable. Then, with notation as in Lemma 4.3.1 and the remark following it, the functions  $g_i, h_i$  are integrable on  $B_i$ , and*

$$\int_{B_1 \times B_2} f = \int_{B_i} g_i = \int_{B_i} h_i, \quad i = 1, 2.$$

*In particular if  $f$  is continuous, then*

$$\int_{B_1 \times B_2} f = \int_{B_1} \int_{B_2} f(\mathbf{x}, \mathbf{y}) \, d\mathbf{y} \, d\mathbf{x} = \int_{B_2} \int_{B_1} f(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} \, d\mathbf{y}.$$

*Proof.* We limit ourselves to establishing  $\int_{A \times B} f = \int_{B_1} g_1$  in the first identity, since they are all proved in essentially the same way. By Lemma 4.3.1,

$$L(f, B) \leq L(g_1, B_1) \leq U(g_1, B_1) \leq U(f, B).$$

But  $f$  is integrable, so  $L(f, B) = U(f, B)$ , and the claim follows.

For the second identity, observe that if  $f$  is continuous, then so is  $f_x$  for any  $x \in B_1$ , and

$$g_1(x) = L(f_x, B_2) = \int_{B_2} f_x = \int_{B_2} f(x, y) dy.$$

Therefore,

$$\int_{B_1 \times B_2} f = \int_{B_1} g_1 = \int_{B_1} \int_{B_2} f(x, y) dy dx.$$

Using  $g_2$  instead of  $g_1$  shows that the order of integration may be reversed.  $\square$

**Examples 4.3.1.** (i) The function  $f$ , where  $f(x, y) = 2x^3 e^{x^2 y}$ , is integrable over  $R = [0, 1] \times [0, 1]$  since it is continuous. Trying to evaluate the iterated integral  $\int_0^1 \int_0^1 2x^3 e^{x^2 y} dx dy$  doesn't look auspicious. If we reverse the order of integration, however, then

$$\int_0^1 2x^3 e^{x^2 y} dy = 2x e^{x^2 y} \Big|_{y=0}^{y=1} = 2x(e^{x^2} - 1),$$

so that  $\int_R f = \int_0^1 2x(e^{x^2} - 1) dx = e^{x^2} - x^2 \Big|_0^1 = e - 2$ .

(ii) It may well happen that a function is integrable, but one of the iterated integrals does not exist: consider

$$f : [0, 1] \times [0, 1] \rightarrow \mathbb{R},$$

$$(x, y) \mapsto \begin{cases} 1 & \text{if } x = y = 0, \\ 0 & \text{if } x \text{ or } y \text{ is irrational,} \\ \frac{1}{n} & \text{if } y \text{ is rational and } x = \frac{m}{n} \neq 0, \end{cases}$$

where  $m, n$  are integers with no common factor and  $n > 0$ . We claim that the set of discontinuities of  $f$  is  $(\mathbb{Q} \cap [0, 1]) \times [0, 1]$ , so that  $f$  is integrable. To see this, consider a point  $(x_0, y_0)$  in this set, with  $x_0 = p/q$ . If  $y_0$  is irrational, then  $f(x_0, y_0) = 0$ , but  $f(x_0, y_k) = 1/q$  for any sequence  $\{y_k\}$  of rationals that converges to  $y_0$ . This shows that  $f$  is discontinuous at that point. If  $y_0$  is rational, then  $f(x_0, y_0) \neq 0$  but any neighborhood of  $(x_0, y_0)$  contains points where  $f$  vanishes, so again  $f$  is discontinuous there. Next, we check that  $f$  is continuous at any other point; i.e., at any  $(x_0, y_0)$  with  $x_0$  irrational. So consider a sequence  $(x_k, y_k) \rightarrow (x_0, y_0)$ . If  $x_k$  is rational for only finitely many  $k$ , then for large enough  $k$ ,  $f(x_k, y_k) = 0 = f(x_0, y_0)$ . Otherwise, there exists a subsequence  $x_{k_n} = p_n/q_n$  of  $\{x_k\}$  consisting of rationals. It is not difficult to prove that in this case,  $q_n \rightarrow \infty$  (see Examples and Remarks (ii))

in Appendix A), so that  $f(x_{k_n}, y_{k_n})$  is either zero or  $1/q_n$ , and in any case converges to zero. Since those terms that are not in the subsequence are mapped to zero,  $f$  is indeed continuous at  $(x_0, y_0)$ .

Thus,  $f$  is integrable, and its integral is zero, since any lower sum is. However, if  $x = m/n$  is a nonzero rational, then  $f_x(y) = 1/n$  when  $y$  is rational and 0 otherwise, so that  $f_x$  is not integrable and  $\int_0^1 f(x, y) dy$  does not exist. The other integral  $\int_0^1 f(x, y) dx$  does, however, exist, and equals zero: this is clear if  $y$  is irrational, since  $f_y$  is then identically zero, and follows from the example in Appendix A mentioned above when  $y$  is rational.

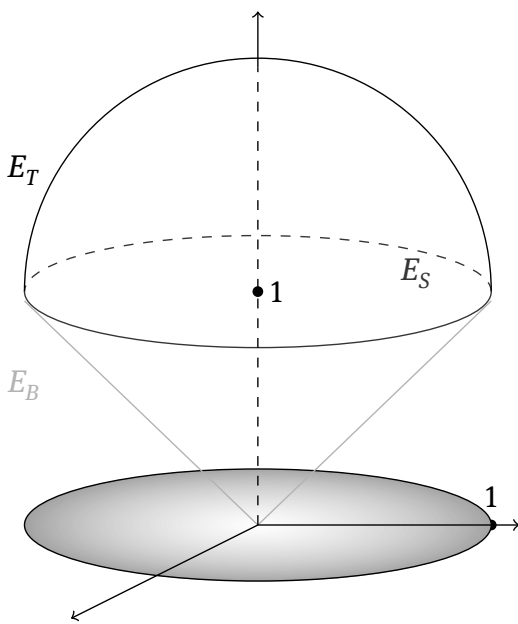
Fubini's theorem is applicable to regions that are more general than boxes: The  $j$ -th coordinate plane in  $\mathbb{R}^n$ ,  $j = 1, \dots, n$ , is the set  $\Pi_j$  of all points  $\mathbf{a} \in \mathbb{R}^n$  with  $u^j(\mathbf{a}) = 0$ . It is canonically isomorphic with  $\mathbb{R}^{n-1}$  via

$$\begin{aligned} \iota_j : \mathbb{R}^{n-1} &\rightarrow \Pi_j, \\ (a_1, \dots, a_{n-1}) &\mapsto (a_1, \dots, a_{j-1}, 0, a_j, \dots, a_{n-1}), \end{aligned}$$

and we routinely identify the two. Define the projection  $\pi_j$  of  $\mathbb{R}^n$  onto the  $j$ -th coordinate plane by  $\pi_j(\mathbf{a}) = (u^1(\mathbf{a}), \dots, u^{j-1}(\mathbf{a}), u^{j+1}(\mathbf{a}), \dots, u^n(\mathbf{a}))$ . A nonempty set  $E \subset \mathbb{R}^n$  is said to be *projectable* if there exists a compact Jordan-measurable set  $\tilde{E} \subset \mathbb{R}^{n-1}$ , some  $j \in \{1, \dots, n\}$  and continuous functions  $g_i : \tilde{E} \rightarrow \mathbb{R}$ ,  $i = 1, 2$ , such that

$$E = \{\mathbf{a} \in \mathbb{R}^n \mid \pi_j(\mathbf{a}) \in \tilde{E} \text{ and } (g_1 \circ \pi_j)(\mathbf{a}) \leq u^j(\mathbf{a}) \leq (g_2 \circ \pi_j)(\mathbf{a})\}.$$

Geometrically speaking, a projectable set  $E$  consists of the region lying between the graphs of two functions defined on the projection of  $E$  onto some coordinate plane. Notice that such a set is entirely determined (and we say it is generated) by  $j$ ,  $\tilde{E}$ ,  $g_1$ , and  $g_2$ .



Projectable set with  
 $g_1(x, y) = \sqrt{x^2 + y^2},$   
 $g_2(x, y) = 1 + \sqrt{1 - x^2 - y^2},$   
 $\tilde{E} = \bar{B}_1(\mathbf{0}) \subset \mathbb{R}^2, j = 3.$

**Lemma 4.3.2.** *A projectable set is Jordan-measurable.*

*Proof.* The claim follows from Theorem 4.2.2 once we establish that the boundary of  $E$  has measure zero. So suppose  $E$  is generated by  $j$ ,  $\tilde{E}$ ,  $g_1$ , and  $g_2$ . We may assume without loss of generality that  $j = n$ . Now, the boundary of  $E$  decomposes as a union of three compact overlapping sets, a “top”  $E_T = \{(\mathbf{x}, g_2(\mathbf{x})) \mid \mathbf{x} \in \tilde{E}\}$ , a “bottom”  $E_B = \{(\mathbf{x}, g_1(\mathbf{x})) \mid \mathbf{x} \in \tilde{E}\}$ , and a “side”  $E_S = \{(\mathbf{x}, x_n) \mid \mathbf{x} \in \partial E, g_1(\mathbf{x}) \leq x_n \leq g_2(\mathbf{x})\}$ .

Consider first the top. Let  $\varepsilon > 0$ , and choose some box  $B \subset \mathbb{R}^{n-1}$  that contains  $\tilde{E}$ . Since  $g_2$  is uniformly continuous on the compact set  $\tilde{E}$ , there exists some  $\delta > 0$  such that  $|g_2(\mathbf{a}) - g_2(\mathbf{b})| < \varepsilon/(4 \operatorname{vol}(B))$  whenever  $\mathbf{a}, \mathbf{b} \in \tilde{E}$  are at a distance less than  $\delta$  from each other. Next, divide  $B$  into subboxes  $B_1, \dots, B_k$  of diameter less than  $\delta$ , and choose some  $\mathbf{a}_i \in B_i, i = 1, \dots, k$ . Then the sets

$$B_i \times \left[ g_2(\mathbf{a}_i) - \frac{\varepsilon}{3 \operatorname{vol}(B)}, g_2(\mathbf{a}_i) + \frac{\varepsilon}{3 \operatorname{vol}(B)} \right], \quad i = 1, \dots, k$$

form a collection of  $n$ -dimensional boxes covering  $E_T$  with total volume  $\sum_i 2\varepsilon \operatorname{vol}(B_i)/(3 \operatorname{vol}(B)) = 2\varepsilon/3 < \varepsilon$ . A similar argument shows that the bottom  $E_B$  can also be covered by boxes with arbitrarily small total volume.

It remains to consider the side  $E_S$ . Set

$$m = \min\{g_1(\mathbf{x}) \mid \mathbf{x} \in \partial\tilde{E}\}, \quad M = \max\{g_2(\mathbf{x}) \mid \mathbf{x} \in \partial\tilde{E}\}.$$

Let  $\varepsilon > 0$ . Since  $\tilde{E}$  is Jordan-measurable, there exist  $(n-1)$ -dimensional boxes  $\tilde{B}_1, \dots, \tilde{B}_l$  covering the boundary of  $\tilde{E}$  with total volume less than  $\varepsilon/(M - m)$  (at least if  $m < M$ . If  $m = M$ , then  $g_1 \equiv g_2$  equal a constant  $\alpha$ , and  $E_S = \partial\tilde{E} \times \{\alpha\}$  certainly has measure zero). Then

$$E_S \subset \bigcup_{i=1}^l \tilde{B}_i \times [m, M],$$

and  $\sum_i \operatorname{vol}(\tilde{B}_i \times [m, M]) < \varepsilon$ . □

**Proposition 4.3.1.** *Suppose  $E \subset \mathbb{R}^n$  is projectable, generated by  $j, \tilde{E}, g_1$ , and  $g_2$ . For each  $\mathbf{x} = (x_1, \dots, x_{n-1}) \in \tilde{E}$ , define  $f_{\mathbf{x}} : [g_1(\mathbf{x}), g_2(\mathbf{x})] \rightarrow \mathbb{R}$  by*

$$f_{\mathbf{x}}(t) = f(x_1, \dots, x_{j-1}, t, x_j, \dots, x_{n-1}).$$

*If  $f : E \rightarrow \mathbb{R}$  is integrable, then*

$$\int_E f = \int_{\tilde{E}} \tilde{f}, \quad \text{where } \tilde{f}(\mathbf{x}) = \int_{g_1(\mathbf{x})}^{g_2(\mathbf{x})} f_{\mathbf{x}}, \quad \mathbf{x} \in \tilde{E}.$$

*Proof.* For simplicity of notation, assume that  $j = n$ . Consider any box  $B = [a_1, b_1] \times \dots \times [a_n, b_n]$  that contains  $E$ , and define  $g : B \rightarrow \mathbb{R}$  by setting it equal to  $f$  inside  $E$ , and



zero otherwise. If  $\tilde{B} = \pi_n(B) \subset \mathbb{R}^{n-1}$ , then by Fubini's theorem,

$$\begin{aligned} \int_E f &= \int_B g = \int_{\tilde{B}} \left( \int_{a_n}^{b_n} g(\mathbf{x}, t) dt \right) d\mathbf{x} = \int_{\tilde{E}} \left( \int_{a_n}^{b_n} g(\mathbf{x}, t) dt \right) d\mathbf{x} \\ &= \int_{\tilde{E}} \left( \int_{g_1(\mathbf{x})}^{g_2(\mathbf{x})} f(\mathbf{x}, t) dt \right) d\mathbf{x} \\ &= \int_{\tilde{E}} \tilde{f}. \end{aligned} \quad \square$$

**Examples 4.3.2.** (i) Suppose we wish to find the volume of the 3-dimensional region  $R$  that lies inside the cylinder  $x^2 + y^2 = 9$  and between the planes  $z = 1$  and  $y + z = 5$ . Let  $D$  denote the disk  $\{(x, y) \mid x^2 + y^2 \leq 9\}$  in  $\mathbb{R}^2$ . Then

$$R = \{(x, y, z) \mid (x, y) \in D \text{ and } 1 \leq z \leq 5 - y\},$$

so that

$$\text{vol}(R) = \int_D \left( \int_1^{5-y} dz \right) dy dx = \int_D (4 - y) dy dx.$$

Now,  $D$  itself is projectable; in fact,

$$D = \{(x, y) \mid -3 \leq x \leq 3, -\sqrt{9-x^2} \leq y \leq \sqrt{9-x^2}\},$$

and

$$\begin{aligned} \text{vol}(R) &= \int_{-3}^3 \int_{-\sqrt{9-x^2}}^{\sqrt{9-x^2}} (4 - y) dy dx = \int_{-3}^3 8\sqrt{9-x^2} dx \\ &= 8 \left( \frac{x}{2} \sqrt{9-x^2} + \frac{9}{2} \arcsin \frac{x}{3} \right) \Big|_{-3}^3 = 36\pi. \end{aligned}$$

(ii) We have so far dealt with integration over bounded regions. More generally, if  $A \subset \mathbb{R}^n$  and  $f : A \rightarrow \mathbb{R}$ , we define

$$\int_A f = \lim_{k \rightarrow \infty} \int_{A \cap [-k, k]^n} f,$$

provided the sequence on the right converges.

As a simple application that will be used in the proof of Sard's theorem, suppose  $A \subset \mathbb{R}^n = \mathbb{R} \times \mathbb{R}^{n-1}$  is a set such that its intersection  $A_t = A \cap (\{t\} \times \mathbb{R}^{n-1})$  with each hyperplane  $u^1 = t$ ,  $t \in \mathbb{R}$ , has  $((n-1)$ -dimensional) volume zero. Then  $A$  has volume zero. To see this, notice first of all that by the above definition,  $A$  may be assumed to be bounded: indeed, if the claim holds for bounded sets, and  $A$  is

an unbounded set satisfying the above hypothesis, then  $\text{vol}(A \cap [-k, k]^n) = 0$  for every  $k$ , and so is the limit as  $k \rightarrow \infty$ . So consider a box  $B = [a_1, b_1] \times B_1$  in  $\mathbb{R}^n$  that contains  $A$ . Then

$$\text{vol}(A) = \int_A 1 = \int_B \chi_A = \int_{a_1}^{b_1} \left( \int_{B_1} \chi_{A_t} \right) dt$$

is zero since by assumption each  $\int_{B_1} \chi_{A_t} = 0$ .

- (ii) In Section 1.6, we defined the volume of the parallelepiped  $P$  spanned by  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^n$  to be the absolute value of the determinant of the matrix that has  $\mathbf{x}_1, \dots, \mathbf{x}_n$  as columns. This may be rephrased as follows: notice that  $P = L(B)$ , where  $B = [0, 1]^n$ , and  $L : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is the linear transformation determined by  $L\mathbf{e}_i = \mathbf{x}_i$ ,  $i = 1, \dots, n$ . The above formula for the volume of  $P$  may therefore be written as

$$\text{vol} L(B) = |\det L| \cdot \text{vol}(B). \quad (4.3.2)$$

For the sake of consistency, it must be checked that this definition coincides with the definition of volume as an integral. We will establish, somewhat more generally, that (4.3.2) holds for any box  $B = [a_1, b_1] \times \dots \times [a_n, b_n]$  in  $\mathbb{R}^n$ . First of all, notice that  $L$  may be assumed to be an isomorphism, since otherwise both sides of the above identity vanish. Next, since any isomorphism is a composition of elementary transformations (see Appendix B) and the determinant of a composition is the product of the determinants, it is enough to prove (4.3.2) for elementary transformations. Recall that these transformations come in three flavors:

- (1)  $L\mathbf{e}_i = \mathbf{e}_j$ ,  $L\mathbf{e}_j = \mathbf{e}_i$  for some  $1 \leq i < j \leq n$ , and  $L\mathbf{e}_k = \mathbf{e}_k$  if  $k \neq i, j$ ;
- (2) There exists  $1 \leq i \leq n$  and  $\alpha \neq 0$  such that  $L\mathbf{e}_i = \alpha\mathbf{e}_i$ , and  $L\mathbf{e}_j = \mathbf{e}_j$  for  $j \neq i$ ;
- (3) There exist distinct  $i, j \in \{1, \dots, n\}$ ,  $\alpha \in \mathbb{R}$ , such that  $L\mathbf{e}_i = \mathbf{e}_i + \alpha\mathbf{e}_j$ , and  $L\mathbf{e}_k = \mathbf{e}_k$  when  $k \neq i$ .

If  $L$  is of type 1, then  $L(B)$  is obtained by interchanging edges  $i$  and  $j$  in  $B$ , and both have the same volume because the latter equals the product of the lengths of its edges. Since  $|\det L| = 1$ , (4.3.2) holds in this case. If  $L$  is of type 2, then  $L(B)$  is a box that has the same edges as  $B$  except for the  $i$ -th one, which has length  $|\alpha|$  times that of  $B$ ; but  $\alpha = \det L$ , so again (4.3.2) is true in this case. Finally, if  $L$  is of type 3, then by what was established for type 1, we may assume that  $L\mathbf{e}_i = \mathbf{e}_i$  if  $i \neq n-1$  and  $L\mathbf{e}_{n-1} = \mathbf{e}_{n-1} + \alpha\mathbf{e}_n$ . Set

$$B_1 = [a_1, b_1] \times \dots \times [a_{n-2}, b_{n-2}], \quad B_2 = [a_{n-1}, b_{n-1}] \times [a_n, b_n],$$

so that  $B = B_1 \times B_2$ . By Fubini's theorem,

$$\text{vol} L(B) = \int_{L(B_1 \times B_2)} 1 = \int_{B_1 \times L(B_2)} 1 = \text{vol} B_1 \cdot \text{vol} L(B_2),$$

and since  $\det L = 1$ , it remains to show that  $B_2$  and  $L(B_2)$  have the same volume; i.e., that (4.3.2) holds for a type 3 elementary transformation  $L : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ . But if  $B = [a, b] \times [c, d]$ , then  $L(B)$  is the projectable region

$$L(B) = \{(x, y) \mid a \leq x \leq b, ax + c \leq y \leq ax + d\},$$

and

$$\text{vol } L(B) = \int_a^b \left( \int_{ax+c}^{ax+d} 1 \right) dx = \int_a^b (d - c) = \text{vol } B.$$

## 4.4 Sard's theorem

Recall from the section on Taylor polynomials that  $\mathbf{a} \in U \subset \mathbb{R}^n$  is said to be a *critical point* of a map  $\mathbf{f} : U \rightarrow \mathbb{R}^m$  if  $D\mathbf{f}(\mathbf{a}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is either not onto or does not exist, and that in this case,  $\mathbf{f}(\mathbf{a})$  is called a *critical value* of  $\mathbf{f}$ . As an application of Fubini's theorem, we discuss a remarkable result of Sard, which asserts that the set of critical values of a  $C^k$  map  $\mathbf{f} : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  has measure zero, if  $k$  is large enough (depending on  $n$  and  $m$ ). This is of course not surprising (and not difficult to prove) when  $n < m$ : after all, even in the worst case scenario when  $\mathbf{f}$  has maximal rank everywhere, the set of critical values is all of  $\mathbf{f}(U)$ , which is locally an immersed  $n$ -dimensional submanifold of  $\mathbb{R}^m$ . Our approach of the argument follows that of [11]. It has the advantage of being relatively short, but comes at additional cost, namely  $\mathbf{f}$  will be assumed to have continuous partial derivatives of any order.

**Theorem 4.4.1 (Sard).** *Let  $U \subset \mathbb{R}^n$ , and  $\mathbf{f} : U \rightarrow \mathbb{R}^m$  a  $C^\infty$  map; i.e., the component functions  $u^i \circ \mathbf{f}$  have continuous partial derivatives of any order. Then the set of critical values of  $\mathbf{f}$  has measure zero.*

*Proof.* The argument will be by induction on the dimension  $n$  of the domain, beginning with  $n = 0$ . Recall that the zero-dimensional space  $\mathbb{R}^0 = \{0\}$ , so the statement is true in this case. Assume then that the statement holds in dimensions less than  $n$ .

Denote by  $C$  the set of critical points of  $\mathbf{f}$ , and by  $C_k$  the subset consisting of those points where all partial derivatives of (the components of)  $\mathbf{f}$  of order  $\leq k$  vanish,  $k \geq 1$ . Thus,  $C \supset C_k \supset C_{k+1}$ ,  $k \in \mathbb{N}$ . We will repeatedly make use of the following observation: in order to show that a given set  $\mathbf{f}(A)$  has measure zero, it is enough to show that every point in  $A$  admits a neighborhood  $V$  such that  $\mathbf{f}(V \cap A)$  has measure zero. This is because the resulting cover of  $A$  contains a countable subcover by Theorem 1.7.5, and a countable union of measure zero sets has measure zero.

*Claim 1:  $\mathbf{f}(C \setminus C_1)$  has measure zero.* To see this, consider  $\mathbf{p} \in C \setminus C_1$ . By assumption,  $D_i f^j(\mathbf{p}) \neq 0$  for some  $1 \leq i, j \leq n$ . It may be assumed that  $i = j = 1$ , since interchanging two coordinates in the domain or in the range of  $\mathbf{f}$  amounts to composing  $\mathbf{f}$  on the right

or on the left with a diffeomorphism  $F$ ; in the former case,  $f \circ F$  has the same critical values as  $f$ , and in the latter case,  $(F \circ f)(C)$  has measure zero if and only if  $f(C)$  has measure zero. Now, the map  $h : U \rightarrow \mathbb{R}^n$  given by

$$h(\mathbf{a}) = (f^1(\mathbf{a}), a_2, \dots, a_n), \quad a_i = u^i(\mathbf{a}),$$

has rank  $n$  at  $\mathbf{p}$ : indeed, its Jacobian at that point is a matrix whose entries below the diagonal all vanish, so that its determinant equals the product of the diagonal elements, namely  $D_1 f^1(\mathbf{p}) \neq 0$ . By the inverse function theorem, there exists a neighborhood  $V$  of  $\mathbf{p}$  in  $U$  such that the restriction  $h : V \rightarrow h(V)$  is a diffeomorphism. Set  $\mathbf{g} := f \circ h^{-1} : h(V) \rightarrow \mathbb{R}^m$ . Since  $h$  is a diffeomorphism, the set  $A$  of critical values of  $\mathbf{g}$  coincides with the set of critical values of the restriction of  $f$  to  $V$ ; i.e.,  $A = f(V \cap C)$ , and it remains to show that  $A$  has measure zero. Notice that  $u^1 \circ h = f^1$ , so that

$$u^1 \circ \mathbf{g} = u^1 \circ f \circ h^{-1} = f^1 \circ h^{-1} = u^1 \circ h \circ h^{-1} = u^1.$$

In other words,  $\mathbf{g}$  maps each hyperplane  $u^1 = \text{constant}$  into itself. Denote by  $\mathbf{g}_\alpha$  the restriction of  $\mathbf{g}$  to the (portion of the) hyperplane

$$H_\alpha = \{\mathbf{p} \in U \mid u^1(\mathbf{p}) = \alpha\}, \quad \alpha \in \mathbb{R}.$$

Observe that for each  $\mathbf{p} \in H_\alpha$ ,

- (1) the Jacobian matrix of  $\mathbf{g}$  at  $\mathbf{p}$  has  $\mathbf{e}_1^T$  as its first row, and
- (2)  $D_i g^j(\mathbf{p}) = D_i g_\alpha^j(\mathbf{p})$  for all  $i, j \neq 1$ .

The latter statement says that if  $\mathbf{g}_\alpha$  is viewed as a map from  $\mathbb{R}^{n-1} \cong H_\alpha$  to itself, then the Jacobian of  $\mathbf{g}_\alpha$  at  $\mathbf{p}$  is obtained by deleting the first row and column from that of  $\mathbf{g}$  at  $\mathbf{p}$ :

$$D\mathbf{g}(\mathbf{p}) = \begin{bmatrix} 1 & 0 & \dots & 0 \\ * & & & \\ \vdots & & D\mathbf{g}_\alpha(\mathbf{p}) & \\ * & & & \end{bmatrix}$$

Together with the former, this means that the columns of  $[D\mathbf{g}(\mathbf{p})]$  are linearly dependent if and only if those of  $[D\mathbf{g}_\alpha(\mathbf{p})]$  are: for if  $M^i$  denotes the  $i$ -th column of the first matrix, and  $M_\alpha^i$  that of the second, then the first entry of the column vector  $\sum a_i M^i$  is  $a_1$ . Thus, if  $\sum a_i M^i = \mathbf{0}$ , then  $a_1 = 0$  and  $\sum a_i M^i = \sum_{i>1} a_i M^i \in \mathbb{R}^m$ . This vector has zero as its first entry, and deleting this entry yields  $\sum_i a_i M_\alpha^{i-1} = \mathbf{0} \in \mathbb{R}^{m-1}$ . Summarizing,  $\mathbf{p} \in H_\alpha$  is a critical point of  $\mathbf{g}_\alpha$  iff it is a critical point of  $\mathbf{g}$ . According to our induction hypothesis, the set of critical values of  $\mathbf{g}_\alpha$  has measure zero, so that the intersection of the set of critical values of  $\mathbf{g}$  with each hyperplane  $u^j = \alpha$  has measure zero. By Examples 4.3.2 (ii),  $A$  itself then has measure zero.

*Claim 2:*  $f(C_k \setminus C_{k+1})$  has measure zero for  $k \geq 1$ . Given any  $\mathbf{p} \in C_k \setminus C_{k+1}$ , we shall once again exhibit an open neighborhood  $V$  of  $\mathbf{p}$  such that  $f(V \cap C_k)$  has measure zero. By

assumption, there exists a partial derivative of (a component function of)  $\mathbf{f}$  of order  $k$ , which we denote by  $\varphi$ , such that  $\varphi(\mathbf{p}) = 0$  but  $D_i\varphi(\mathbf{p}) \neq 0$  for some  $i$  between 1 and  $n$ . As before, the map  $\mathbf{h} : U \rightarrow \mathbb{R}^n$ , given by

$$\mathbf{h}(\mathbf{a}) = (a_1, \dots, a_{i-1}, \varphi(\mathbf{a}), a_{i+1}, \dots, a_n),$$

has rank  $n$  at  $\mathbf{p}$  and there exists a neighborhood  $V$  of  $\mathbf{p}$  on which the restriction  $\mathbf{h} : V \rightarrow \mathbf{h}(V)$  is a diffeomorphism. Set once again  $\mathbf{g} = \mathbf{f} \circ \mathbf{h}^{-1} : \mathbf{h}(V) \rightarrow \mathbb{R}^m$ . Since  $\varphi$  vanishes when restricted to  $V \cap C_k$ ,  $\mathbf{h}(V \cap C_k)$  lies in the coordinate hyperplane  $u^i = 0$ . If  $\mathbf{g}_0$  denotes the restriction of  $\mathbf{g}$  to this hyperplane, then every point of  $\mathbf{h}(V \cap C_k)$  is a critical point of  $\mathbf{g}_0$  because  $V \cap C_k$  is contained in the set of critical points of  $\mathbf{f}$ . By the induction hypothesis, the critical values of  $\mathbf{g}_0$  form a set of measure zero. Thus,  $\mathbf{f}(V \cap C_k) = \mathbf{g}_0(\mathbf{h}(V \cap C_k))$  has measure zero, as claimed.

*Claim 3:  $\mathbf{f}(C_k)$  has measure zero if  $k$  is large enough.* To see this, consider a box  $B = \prod_{i=1}^n [a_i, b_i]$  with edges of common length  $R = b_i - a_i$  small enough that  $B$  is contained in  $U$ . It suffices to show that if  $k > (n/m) - 1$ , then  $\mathbf{f}(C_k \cap B)$  has measure zero. Applying Remark 2.7.1 to each component function of  $\mathbf{f}$ , we see that there exists  $\beta > 0$  such that

$$|R(\mathbf{x}, \mathbf{h})| := |\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x})| \leq \beta |\mathbf{h}|^{k+1}, \quad \mathbf{x} \in C_k \cap B, \quad \mathbf{x} + \mathbf{h} \in B. \quad (4.4.1)$$

Given any natural number  $l > 1$ , partition each  $[a_i, b_i]$  into  $l$  intervals of equal length  $R/l$ . The corresponding partition of  $B$  consists of  $l^n$  subboxes with diameter  $\sqrt{n}R/l$ . Let  $\mathbf{x} \in C_k \cap B$ , and  $\tilde{B}$  be a subbox containing  $\mathbf{x}$ . Then any other point in  $\tilde{B} \cap C_k$  is of the form  $\mathbf{x} + \mathbf{h}$ , where  $|\mathbf{h}| \leq \sqrt{n}R/l$ . It follows from (4.4.1) that  $\mathbf{f}(\tilde{B} \cap C_k)$  is contained in a box centered at  $\mathbf{f}(\mathbf{x})$  with all sides of common length  $2\beta(\sqrt{n}R/l)^{k+1}$ . Since this box lives in  $\mathbb{R}^m$ , it has volume  $\alpha/l^{m(k+1)}$ , where  $\alpha$  is a constant that does not depend on  $l$ . Thus,  $\mathbf{f}(B \cap C_k)$  is contained in  $l^n$  boxes with total volume

$$V \leq \alpha l^{n-m(k+1)}.$$

The exponent of  $l$  is a negative integer, so that  $V$  can be made arbitrarily small by choosing  $l$  large enough. As observed earlier, this implies that  $\mathbf{f}(C_k)$  has measure zero. To conclude, set  $C_0 := C$ . If  $k > (n/m) - 1$ , then

$$\mathbf{f}(C) = \left( \bigcup_{i=1}^k \mathbf{f}(C_{i-1} \setminus C_i) \right) \cup \mathbf{f}(C_k)$$

is a finite union of sets of measure zero. This completes the proof.  $\square$

**Remark 4.4.1.** Sard's theorem holds more generally for maps between manifolds: A subset  $A$  of a manifold  $M^k$  is said to have measure zero if  $A$  can be written as a countable union  $\cup A_n$  where each  $A_n$  lies in the domain of some chart  $(U_n, \mathbf{x}_n)$ , and  $\mathbf{x}_n(A_n)$  has measure zero in  $\mathbb{R}^k$ . Now, if  $\mathbf{f} : M \rightarrow N$  is a  $C^\infty$  map between manifolds  $M$  and  $N$ , then for any charts  $(U, \mathbf{x})$  of  $M$  and  $(V, \mathbf{y})$  of  $N$ , the set of critical values of  $\mathbf{y} \circ \mathbf{f} \circ \mathbf{x}^{-1}$  has measure zero. It follows that the set of critical values of  $\mathbf{f}$  has measure zero.

## 4.5 The change of variables theorem

For functions of a single variable, the chain rule implies the well-known change of variables theorem:

If  $g : [a, b] \rightarrow \mathbb{R}$  is continuously differentiable, and  $f$  is a continuous function whose domain contains  $g([a, b])$ , then

$$\int_{g(a)}^{g(b)} f = \int_a^b (f \circ g)g'. \quad (4.5.1)$$

The proof is easy: if  $F$  is an antiderivative of  $f$ , then  $F \circ g$  is an antiderivative of  $(f \circ g)g'$ , so that both sides equal  $F(g(b)) - F(g(a))$ . Before stating the generalization of this theorem to higher dimensions – let alone proving it, which turns out to be much more involved – we observe that it may be reformulated as follows:

$$\int_{g(a,b)} f = \int_{(a,b)} (f \circ g)|g'|, \quad (4.5.2)$$

at least if  $g$  is one-to-one. Indeed,  $g$  is then either increasing, in which case  $g' \geq 0$  and both identities are the same, or decreasing, so that  $g(a, b) = (g(b), g(a))$ ,  $g' \leq 0$ , and both sides of (4.5.2) are the negative of those in (4.5.1).

The higher-dimensional analogue of (4.5.2) is given by the following:

**Theorem 4.5.1.** *Let  $U \subset \mathbb{R}^n$  be bounded, and suppose  $\mathbf{g} : U \rightarrow \mathbb{R}^n$  is a continuously differentiable map, which is injective on an open subset whose complement in  $U$  has measure zero. If  $f : \mathbf{g}(U) \rightarrow \mathbb{R}$  is integrable, then*

$$\int_{\mathbf{g}(U)} f = \int_U (f \circ \mathbf{g})|\det D\mathbf{g}|.$$

The proof will be handled in a series of steps. First of all,  $U$  may be assumed to be open and  $\mathbf{g}$  injective on  $U$  by Remarks 4.2.3 (ii) together with the fact that the image  $\mathbf{g}(A)$  of a set  $A$  of measure zero has measure zero if  $\mathbf{g}$  is continuous and  $A$  lies in some compact set, see Exercise 4.3. Similarly, by Sard's theorem, we may suppose that the Jacobian determinant of  $\mathbf{g}$  is nowhere zero. Next, we point out two observations that will be used repeatedly in the proof:

**Observation 4.5.1.** The class of maps  $\mathbf{g}$  for which the theorem holds is closed under composition.

To see this, let  $\mathbf{g}_1 : U \rightarrow \mathbb{R}^n$ ,  $\mathbf{g}_2 : V \rightarrow \mathbb{R}^n$ , where  $\mathbf{g}_1(U) \subset V$ . If the claim is true for each of these maps, then

$$\begin{aligned} \int_{(\mathbf{g}_2 \circ \mathbf{g}_1)(U)} f &= \int_{\mathbf{g}_2(\mathbf{g}_1(U))} f = \int_{\mathbf{g}_1(U)} (f \circ \mathbf{g}_2) |\det D\mathbf{g}_2| \\ &= \int_U (f \circ \mathbf{g}_2 \circ \mathbf{g}_1) |(\det D\mathbf{g}_2) \circ \mathbf{g}_1| \det D\mathbf{g}_1| \\ &= \int_U (f \circ \mathbf{g}_2 \circ \mathbf{g}_1) |\det D(\mathbf{g}_2 \circ \mathbf{g}_1)|. \end{aligned}$$

**Observation 4.5.2.** It suffices to prove that for any  $\mathbf{a} \in U$ , there exists an open neighborhood  $W \subset U$  of  $\mathbf{a}$  for which the theorem holds.

Indeed, if the statement is true in this case, then there exists an admissible open cover  $\{W_k \mid k \in A \subset \mathbb{N}\}$  of  $U$  such that the theorem holds on each  $W_k$ . Then  $\{\mathbf{g}(W_k) \mid k \in A\}$  is an admissible open cover of  $\mathbf{g}(U)$ . If  $\{\varphi_k \mid k \in A\}$  is a partition of unity subordinate to  $\{\mathbf{g}(W_k) \mid k \in A\}$ , then again by hypothesis,

$$\int_{\mathbf{g}(W_k)} \varphi_k f = \int_{W_k} (\varphi_k f) \circ \mathbf{g} |\det D\mathbf{g}|. \quad (4.5.3)$$

Now,  $\varphi_k$  vanishes outside  $\mathbf{g}(W_k)$ , so that, since  $\mathbf{g}$  is injective,  $\varphi_k \circ \mathbf{g}$  is zero outside  $W_k$  (and in particular,  $\{\varphi_k \circ \mathbf{g} \mid k \in A\}$  is a partition of unity subordinate to  $U$ ). Thus, (4.5.3) becomes

$$\int_{\mathbf{g}(U)} \varphi_k f = \int_U (\varphi_k f) \circ \mathbf{g} |\det D\mathbf{g}|. \quad (4.5.4)$$

As noted above, the collection  $\{\varphi_k \circ \mathbf{g} \mid k \in A\}$  is a partition of unity subordinate to  $U$ . The definition of the generalized integral then implies that

$$\int_{\mathbf{g}(U)} f = \sum_{k \in A} \int_{\mathbf{g}(U)} \varphi_k f = \sum_{k \in A} \int_U (\varphi_k \circ \mathbf{g})(f \circ \mathbf{g}) |\det D\mathbf{g}| = \int_U (f \circ \mathbf{g}) |\det D\mathbf{g}|.$$

Now that these observations are out of the way, we proceed with the proof of the change of variables formula by first considering constant functions. Since the integral of a constant function  $c$  over a set is  $c$  times the volume of the set, it suffices to consider the constant function 1:

**Lemma 4.5.1.** *Let  $U$  be open in  $\mathbb{R}^n$ ,  $\mathbf{g} : U \rightarrow \mathbb{R}^n$  a one-to-one continuously differentiable map with nowhere zero Jacobian determinant. Then*

$$\text{vol } \mathbf{g}(U) = \int_{\mathbf{g}(U)} 1 = \int_U |\det D\mathbf{g}|.$$

*Proof.* We proceed by induction on  $n$ . The case  $n = 1$  follows from (4.5.2), so assume the claim holds in dimension  $n - 1$ . By Observation 4.5.2, it suffices to establish the

result for some neighborhood of an arbitrary point  $\mathbf{a} \in U$ . Furthermore, we may suppose that  $[D\mathbf{g}(\mathbf{a})] = I_n$ ; if the lemma is valid for maps with Jacobian matrix equal to the identity at  $\mathbf{a}$ , then it must hold for the map  $D\mathbf{g}(\mathbf{a})^{-1} \circ \mathbf{g}$ . On the other hand, it is also true for the linear transformation  $D\mathbf{g}(\mathbf{a})$  by Examples and Remarks 4.3.2 (iv). By Observation 4.5.1, it then holds for  $\mathbf{g} = D\mathbf{g}(\mathbf{a}) \circ (D\mathbf{g}(\mathbf{a})^{-1} \circ \mathbf{g})$ .

In order to use the induction hypothesis, we rewrite  $\mathbf{g}$  as a composition  $\mathbf{f} \circ \mathbf{h}$  as follows: let  $\pi_1 : \mathbb{R}^n = \mathbb{R}^{n-1} \times \mathbb{R} \rightarrow \mathbb{R}^{n-1}$  denote projection, and define  $\mathbf{h} : U \rightarrow \mathbb{R}^n$  by  $\mathbf{h} = (\pi_1 \circ \mathbf{g}, u^n)$ ; i.e., for  $\mathbf{x} = (x_1, \dots, x_n)$ ,  $\mathbf{h}(\mathbf{x}) = (g^1(\mathbf{x}), \dots, g^{n-1}(\mathbf{x}), x_n)$ . Since  $D\mathbf{g}(\mathbf{a}) = 1_{\mathbb{R}^n}$ ,  $D\mathbf{h}(\mathbf{a}) = 1_{\mathbb{R}^n}$  and  $\mathbf{h}$  is invertible in some neighborhood  $W$  of  $\mathbf{a}$ . Set  $\tilde{W} = \mathbf{h}(W)$ , and define  $\mathbf{f} : \tilde{W} \rightarrow \mathbb{R}^n$  by  $\mathbf{f} = (\pi_1, g^n \circ (\mathbf{h}|_W)^{-1})$ . Then on  $W$

$$\mathbf{f} \circ \mathbf{h} = (\pi_1 \circ \mathbf{h}, g^n) = (\pi_1 \circ \mathbf{g}, g^n) = \mathbf{g}.$$

Furthermore,

$$D\mathbf{f}(\mathbf{h}(\mathbf{a})) = D\mathbf{f}(\mathbf{h}(\mathbf{a})) \circ D\mathbf{h}(\mathbf{a}) = D\mathbf{g}(\mathbf{a}) = 1_{\mathbb{R}^n},$$

so that  $\mathbf{f}$  is injective with Jacobian matrix of rank  $n$  on some neighborhood  $V$  of  $\mathbf{h}(\mathbf{a})$ . Set  $U = \mathbf{h}^{-1}(V)$ , restrict  $\mathbf{h}$  to  $U$  and  $\mathbf{f}$  to  $V$ . Then  $\mathbf{h} : U \rightarrow V$ ,  $\mathbf{f} : V \rightarrow \mathbf{f}(V)$  are diffeomorphisms,  $\mathbf{g}|_U = \mathbf{f} \circ \mathbf{h}$ , and it suffices, by the two observations above, to establish the claim for  $\mathbf{f}$  and  $\mathbf{h}$  on open boxes containing  $\mathbf{h}(\mathbf{a})$  and  $\mathbf{a}$  respectively. We begin with  $\mathbf{h}$ : consider an open box  $B = B_1 \times (a_n, b_n) \subset U$  containing  $\mathbf{a}$ , where  $B_1 = \pi_1(B)$ . By Examples 4.2.2 (iii) and Exercise 1.29,  $\mathbf{h}(B \times [a_n, b_n])$  and  $\mathbf{h}(B \times \{t\})$  are Jordan-measurable, and Fubini's theorem implies

$$\int_{\mathbf{h}(B)} 1 = \int_{a_n}^{b_n} \left( \int_{\mathbf{h}(B_1 \times \{t\})} 1 \, d\mathbf{x} \right) dt. \tag{4.5.5}$$

Now, for each  $t \in (a_n, b_n)$ , the map  $\mathbf{h}_t : B_1 \rightarrow \mathbb{R}^{n-1}$ , where  $\mathbf{h}_t(\mathbf{x}) = (\pi_1 \circ \mathbf{g})(\mathbf{x}, t)$ , is injective, and

$$[D\mathbf{h}(\mathbf{x}, t)] = \begin{bmatrix} & & & * \\ & & & \vdots \\ & D\mathbf{h}_t(\mathbf{x}) & & * \\ 0 & \dots & 0 & 1 \end{bmatrix},$$

so that  $\det D\mathbf{h}_t(\mathbf{x}) \neq 0$ . By the induction hypothesis,

$$\int_{\mathbf{h}_t(B_1)} 1 \, d\mathbf{x} = \int_{B_1} |D\mathbf{h}_t| \, d\mathbf{x} = \int_{B_1} |D\mathbf{h}(\mathbf{x}, t)| \, d\mathbf{x}.$$

But  $\mathbf{h}_t(B_1) = \mathbf{h}(B_1 \times \{t\})$ , so that (4.5.5) now yields

$$\int_{\mathbf{h}(B)} 1 = \int_{a_n}^{b_n} \left( \int_{B_1} |D\mathbf{h}(\mathbf{x}, t)| \, d\mathbf{x} \right) dt = \int_B |D\mathbf{h}|.$$



This proves the claim for  $\mathbf{h}$ , and we now turn our attention to  $\mathbf{f}$ . Let  $B = B_1 \times (a_n, b_n)$  be a box in  $\mathbb{R}^{n-1} \times \mathbb{R}$ , and for each  $\mathbf{x} \in B$ , define  $f_{\mathbf{x}} : (a_n, b_n) \rightarrow \mathbb{R}$  by  $f_{\mathbf{x}}(t) = f^n(\mathbf{x}, t) = (\mathbf{g}^n \circ \mathbf{h}^{-1})(\mathbf{x}, t)$ . Since  $\mathbf{f}(\mathbf{x}, t) = (\mathbf{x}, f_{\mathbf{x}}(t))$ ,

$$[D\mathbf{f}(\mathbf{x}, t)] = \begin{bmatrix} & & 0 \\ & I_{n-1} & \vdots \\ * & \dots & * & f'_{\mathbf{x}}(t) \end{bmatrix},$$

and in particular,  $|\det D\mathbf{f}(\mathbf{x}, t)| = |f'_{\mathbf{x}}(t)|$ . Now,

$$\mathbf{f}(B) = \bigcup_{\mathbf{x} \in B_1} \{\mathbf{x}\} \times f_{\mathbf{x}}(a_n, b_n),$$

so that

$$\begin{aligned} \int_{\mathbf{f}(B)} 1 &= \int_{B_1} \left( \int_{f_{\mathbf{x}}(a_n, b_n)} 1 \, dt \right) d\mathbf{x} = \int_{B_1} \left( \int_{(a_n, b_n)} |f'_{\mathbf{x}}(t)| \, dt \right) d\mathbf{x} \\ &= \int_{B_1} \left( \int_{a_n}^{b_n} |\det D\mathbf{f}(\mathbf{x}, t)| \, dt \right) d\mathbf{x} \\ &= \int_B |\det D\mathbf{f}|. \end{aligned}$$

This completes the proof of the lemma.  $\square$

*Proof of Theorem 4.5.1.* By Observation 4.5.2, it suffices to show that for any  $\mathbf{a} \in U$ , there exists a neighborhood  $W \subset U$  of  $\mathbf{a}$  such that the theorem holds on  $W$ . So consider an open box  $B$  that contains  $\mathbf{g}(\mathbf{a})$ , and let  $W = \mathbf{g}^{-1}(B)$ . If  $P$  is a partition of  $B$ , then

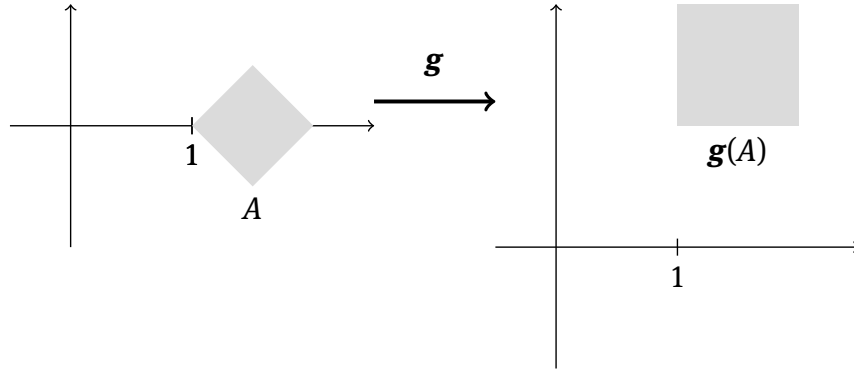
$$\begin{aligned} L(f, P) &= \sum_{\tilde{B} \in P} m_{\tilde{B}}(f) \operatorname{vol}(\tilde{B}) = \sum_{\tilde{B} \in P} m_{\tilde{B}}(f) \int_{\tilde{B}^0} 1 \\ &= \sum_{\tilde{B} \in P} m_{\tilde{B}}(f) \int_{\mathbf{g}^{-1}(\tilde{B}^0)} (1 \circ \mathbf{g}) |\det D\mathbf{g}| \\ &= \sum_{\tilde{B} \in P} \int_{\mathbf{g}^{-1}(\tilde{B}^0)} (m_{\tilde{B}}(f)) \circ \mathbf{g} |\det D\mathbf{g}| \\ &\leq \int_W (f \circ \mathbf{g}) |\det D\mathbf{g}|, \end{aligned}$$

and therefore  $\int_{\mathbf{g}(W)} f \leq \int_W (f \circ \mathbf{g}) |\det D\mathbf{g}|$ . A similar argument using  $U(f, P)$  instead of  $L(f, P)$  and  $M_{\tilde{B}}(f)$  instead of  $m_{\tilde{B}}(f)$  implies that the inequality holds in the other direction. Thus,

$$\int_{\mathbf{g}(W)} f = \int_W (f \circ \mathbf{g}) |\det D\mathbf{g}|.$$

This shows that the theorem holds for  $W$  and establishes the result.  $\square$

**Example 4.5.1.** Consider the integral  $\int_A (x - y)/(x + y) dx dy$ , where  $A$  denotes the region bounded by the square with vertices  $(1, 0)$ ,  $(3/2, 1/2)$ ,  $(2, 0)$  and  $(3/2, -1/2)$ . Although it can be evaluated with Fubini's theorem (notice that the region is projectable), the computation is lengthy.



Alternatively, observe that  $\mathbf{g} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , where  $\mathbf{g}(x, y) = (x - y, x + y)$ , is an isomorphism, so we may apply the change of variables theorem to the function  $f(u, v) = u/v$  to obtain

$$\begin{aligned} \int_A \frac{x-y}{x+y} dx dy &= \int_A (f \circ \mathbf{g}) = \frac{1}{|\det D\mathbf{g}|} \int_A (f \circ \mathbf{g}) |\det D\mathbf{g}| \\ &= \frac{1}{|\det D\mathbf{g}|} \int_{\mathbf{g}(A)} f. \end{aligned}$$

Now, the Jacobian determinant of  $\mathbf{g}$  equals 2, and  $A$  is bounded by the lines  $x - y = 1$ ,  $x - y = 2$ ,  $x + y = 1$ , and  $x + y = 2$ . Thus,  $\mathbf{g}(A) = [1, 2] \times [1, 2]$ , and

$$\int_A \frac{x-y}{x+y} dx dy = \frac{1}{2} \int_1^2 \left( \int_1^2 \frac{u}{v} du \right) dv = \frac{1}{2} \int_1^2 u \ln 2 du = \frac{3}{4} \ln 2.$$

## 4.6 Cylindrical and spherical coordinates

As an application of the results from the previous section, we discuss some changes of variables that are useful when working with regions that exhibit radial or spherical symmetry.

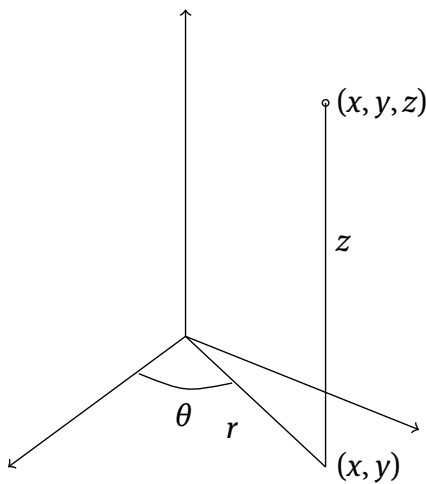
### 4.6.1 Cylindrical coordinates

*Polar coordinates* in  $\mathbb{R}^2$  are given by the map  $(r, \theta) : \mathbb{R}^2 \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}^2$ , where  $r$  assigns to a point  $\mathbf{p}$  its distance to the origin  $\mathbf{0}$ , and  $\theta$  is the angle in  $[0, 2\pi)$  between the positive  $x$ -axis and the line segment  $\mathbf{0p}$ . The polar angle  $\theta$  requires some care in writing down,

because  $\tan$  is only invertible when restricted to an interval of length  $\pi$ . Formally,  $r = \sqrt{(u^1)^2 + (u^2)^2}$ , whereas

$$\theta = \begin{cases} \arctan(u^2/u^1) & \text{if } u^1 > 0, u^2 \geq 0, \\ \pi/2 & \text{if } u^1 = 0 \text{ and } u^2 > 0, \\ \pi + \arctan(u^2/u^1) & \text{if } u^1 < 0, \\ 3\pi/2 & \text{if } u^1 = 0 \text{ and } u^2 < 0, \\ 2\pi + \arctan(u^2/u^1) & \text{if } u^1 > 0 \text{ and } u^2 < 0. \end{cases}$$

Polar coordinates extend to *cylindrical coordinates* in  $\mathbb{R}^3$  via  $(r, \theta, u^3) : \{(x, y, z) \in \mathbb{R}^3 \mid (x, y) \neq (0, 0)\} \rightarrow \mathbb{R}^3$ .



They are particularly useful when working with surfaces that are invariant under rotation about the  $z$ -axis. If  $V \subset \mathbb{R}^3$  is a region that can be expressed as  $U$  in cylindrical coordinates, then  $V = \mathbf{g}(U)$ , where  $\mathbf{g}$  is the inverse of cylindrical coordinates,

$$\mathbf{g} : [0, \infty) \times [0, 2\pi) \times \mathbb{R} \rightarrow \mathbb{R}^3, \\ (r, \theta, z) \mapsto (r \cos \theta, r \sin \theta, z).$$

$\mathbf{g}$  is injective on  $(0, \infty) \times [0, 2\pi) \times \mathbb{R}$  and differentiable on  $(0, \infty) \times (0, 2\pi) \times \mathbb{R}$  with Jacobian

$$[D\mathbf{g}(r, \theta, z)] = \begin{bmatrix} \cos \theta & -r \sin \theta & 0 \\ \sin \theta & r \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Thus, if  $V \subset \mathbb{R}^3$  is expressed as  $U$  in cylindrical coordinates, i.e., if  $V = \mathbf{g}(U)$ , and if  $f$  is integrable on  $V$ , then

$$\int_V f = \int_U (f \circ \mathbf{g})(r, \theta, z) r \, dr \, d\theta \, dz. \quad (4.6.1)$$

Integration in polar coordinates is a special case of the above: suppose  $f$  is a function of 2 variables defined on a region  $V$  in the plane. Since the inverse  $\mathbf{h}$  of polar coordinates  $(r, \theta)$  equals  $\pi_1 \circ \mathbf{g}|_{\mathbb{R}^2 \times \{0\}}$ , where  $\pi_1 = (u^1, u^2) : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  is projection, we have

$$[D\mathbf{h}(r, \theta)] = \begin{bmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{bmatrix}.$$

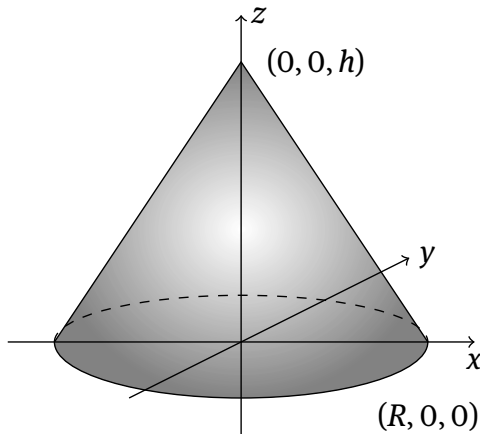
Thus, if  $V \subset \mathbb{R}^2$  is expressed as  $U$  in polar coordinates, i.e., if  $V = \mathbf{g}(U)$ , and if  $f$  is integrable on  $V$ , then

$$\int_V f = \int_U (f \circ \mathbf{g})(r, \theta) r \, dr \, d\theta. \quad (4.6.2)$$

**Examples 4.6.1.** (i) A *right circular cone* of height  $h$  and radius  $R$  is the surface (the term is used loosely here) in  $\mathbb{R}^3$  obtained by rotating the line segment

$$\{(t, 0, (-h/R)t + h) \mid 0 \leq t \leq R\}$$

joining  $(R, 0, 0)$  and  $(0, 0, h)$  about the  $z$ -axis.



A cone with base radius  $R$  and height  $h$

The volume of this cone (or rather the volume of the 3-dimensional region bounded by the cone and the plane  $z = 0$ ) is therefore equal to  $\text{vol } A$ , where

$$A = \left\{ (x, y, z) \mid 0 \leq \sqrt{x^2 + y^2} \leq R, \quad 0 \leq z \leq -\frac{h}{R} \sqrt{x^2 + y^2} + h \right\}.$$

Thus, the interior  $A^0$  of  $A$  can be expressed as  $\mathbf{g}(U)$ , where  $U$  is the projectable region

$$U = \left\{ (r, \theta, z) \mid 0 \leq r < R, \quad 0 \leq \theta < 2\pi, \quad 0 < z < h \left( 1 - \frac{r}{R} \right) \right\}.$$

By Fubini's and the change of variables theorems,

$$\begin{aligned} \text{vol}(A) &= \text{vol}(A^0) = \text{vol}(\mathbf{g}(U)) = \int_{\mathbf{g}(U)} 1 = \int_U r \, dz \, dr \, d\theta \\ &= \int_0^{2\pi} \int_0^R \int_0^{h(1-\frac{r}{R})} r \, dz \, dr \, d\theta = 2\pi \int_0^R \left( rh - r^2 \frac{h}{R} \right) dr \\ &= \frac{\pi R^2 h}{3}. \end{aligned}$$

- (ii) Suppose we are asked to determine the area of the planar region  $R$  bounded by the curve  $(x^2 + y^2)^3 = (x^2 - y^2)^2$ . This region is better visualized in polar coordinates, where the equation of the curve becomes

$$r^6 = (r^2 \cos^2 \theta - r^2 \sin^2 \theta)^2 = r^4 \cos^2(2\theta), \quad \text{or } r = |\cos 2\theta|$$

since  $r \geq 0$ .

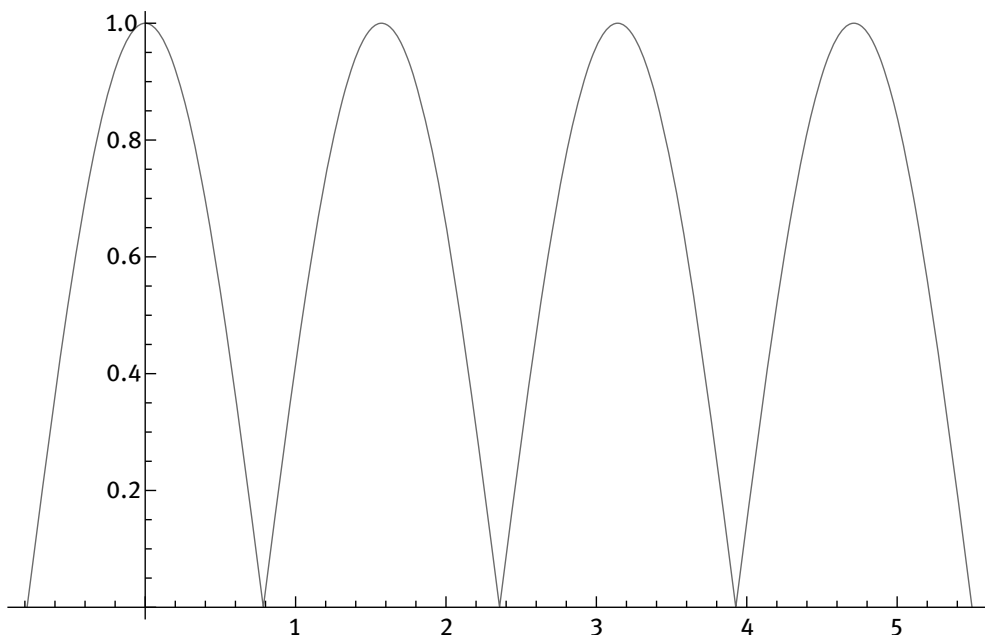


Fig. 4.4:  $r = |\cos 2\theta|$

The curve is invariant under reflection in both coordinate axes (that is, the equation is unchanged when  $\theta$  is replaced by  $-\theta$  or by  $\pi - \theta$ ), and represents the boundary of a 4-leaved rose. For example, the leaf

$$A = \{(r, \theta) \mid -\frac{\pi}{4} \leq \theta \leq \frac{\pi}{4}, \quad 0 \leq r \leq \cos 2\theta\}$$

in the region  $-\pi/4 \leq \theta \leq \pi/4$  starts out at the origin when  $\theta = -\pi/4$ , and  $r$  increases until it reaches a maximum of 1 at  $\theta = 0$ , which corresponds to the point (1,0) in Cartesian coordinates. This is the bottom half of the leaf, and the top

is obtained by reflecting in the  $x$ -axis. Reflection in the  $y$ -axis yields a second leaf. Finally, notice that the curve is also invariant under reflection in the line  $y = x$  because the equation is unchanged when  $\theta$  is replaced by  $\pi/2 - \theta$ . Said reflection then reveals the two remaining leaves.

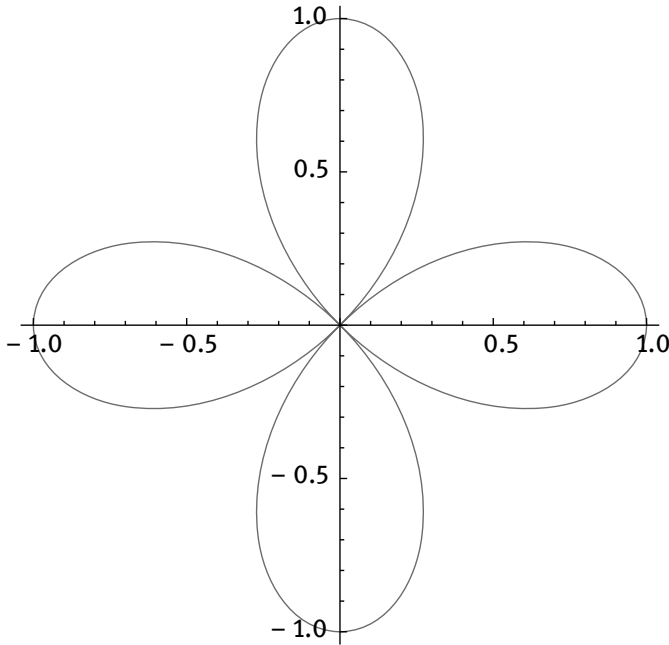


Fig. 4.5:  $r = |\cos 2\theta|$  in polar coordinates

Thus,

$$\begin{aligned} \text{area}(R) &= \int_R 1 \cdot r \, dr \, d\theta = 4 \int_A r \, dr \, d\theta = 4 \int_{-\pi/4}^{\pi/4} \int_0^{\cos 2\theta} r \, dr \, d\theta \\ &= 4 \int_{-\pi/4}^{\pi/4} \frac{1}{2} \cos^2(2\theta) \, d\theta = \int_{-\pi/4}^{\pi/4} (1 + \cos 4\theta) \, d\theta = \left[ \theta + \frac{1}{4} \sin 4\theta \right]_{-\pi/4}^{\pi/4} \\ &= \frac{\pi}{2}. \end{aligned}$$

#### 4.6.2 Spherical coordinates

Like cylindrical coordinates, spherical coordinates generalize polar ones, but in a different direction. Denote by  $\arccos : [-1, 1] \rightarrow [0, \pi]$  the inverse of the restriction of  $\cos$  to the interval  $[0, \pi]$ . Let  $\rho : \mathbb{R}^3 \rightarrow \mathbb{R}$  denote the distance function to the origin,  $\rho(\mathbf{x}) = |\mathbf{x}|$ , and  $\varphi : \mathbb{R}^3 \setminus \{\mathbf{0}\} \rightarrow [0, \pi]$  the angle with the positive  $z$ -axis,

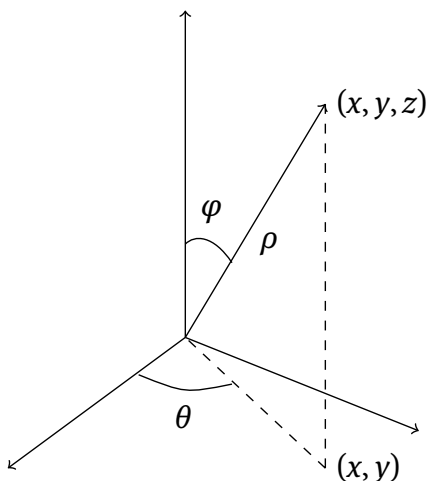
$$\varphi(\mathbf{x}) = \arccos(\langle \mathbf{x}, \mathbf{e}_3 \rangle / |\mathbf{x}|), \quad \mathbf{x} \in \mathbb{R}^3 \setminus \{\mathbf{0}\}.$$

Thus,

$$\mathbf{x} = (\pi(\mathbf{x}), \rho(\mathbf{x}) \cos \varphi(\mathbf{x})),$$

where  $\pi : \mathbb{R}^3 = \mathbb{R}^2 \times \mathbb{R} \rightarrow \mathbb{R}^2$  is projection onto the  $x$ - $y$  plane. Since  $\pi(\mathbf{x}) \in \mathbb{R}^2$ , it may be expressed in polar coordinates; in fact,  $r(\pi(\mathbf{x})) = \rho(\mathbf{x}) \sin \varphi(\mathbf{x})$ , so that

$$\mathbf{x} = (\rho(\mathbf{x}) \cos \theta(\pi(\mathbf{x})) \sin \varphi(\mathbf{x}), \rho(\mathbf{x}) \sin \theta(\pi(\mathbf{x})) \sin \varphi(\mathbf{x}), \rho(\mathbf{x}) \cos \varphi(\mathbf{x})). \quad (4.6.3)$$



Spherical coordinates

*Spherical coordinates* are given by the map  $(\rho, \theta, \varphi) : \mathbb{R}^3 \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}^3$ , with  $\rho, \varphi$  as above, and  $\theta$  the same function from cylindrical coordinates. If

$$\begin{aligned} \mathbf{g} : [0, \infty) \times [0, 2\pi) \times [0, \pi] &\rightarrow \mathbb{R}^3, \\ (a_1, a_2, a_3) &\mapsto a_1(\sin a_3 \cos a_2, \sin a_3 \sin a_2, \cos a_3), \end{aligned}$$

then (4.6.3) says that  $(\rho, \theta, \varphi)$  is invertible with inverse  $\mathbf{g}$ .  $\mathbf{g}$  has Jacobian matrix

$$[D\mathbf{g}(\rho, \theta, \varphi)] = \begin{bmatrix} \sin \varphi \cos \theta & -\rho \sin \varphi \sin \theta & \rho \cos \varphi \cos \theta \\ \sin \varphi \sin \theta & \rho \sin \varphi \cos \theta & \rho \cos \varphi \sin \theta \\ \cos \varphi & 0 & -\rho \sin \varphi \end{bmatrix},$$

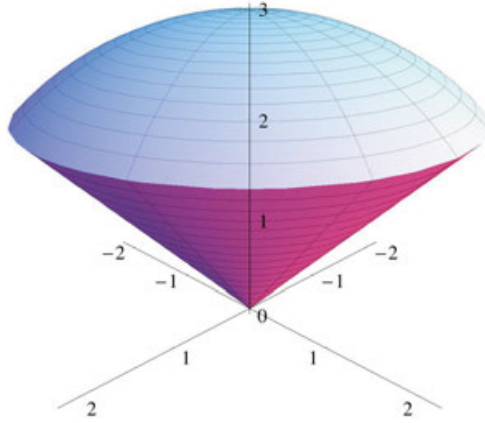
which has determinant  $-\rho^2 \sin \varphi$ . Thus, if  $V = \mathbf{g}(U)$  and  $f$  is integrable over  $V$ , then

$$\int_V f = \int_U (f \circ \mathbf{g})(\rho, \theta, \varphi) \rho^2 \sin \varphi \, d\rho \, d\theta \, d\varphi. \quad (4.6.4)$$

**Examples and Remarks 4.6.1.** (i) Suppose we are asked to determine the volume of the ‘ice cream cone’ consisting of the region  $V$  bounded by the sphere  $x^2 + y^2 + z^2 = R^2$  and the cone  $z = \sqrt{x^2 + y^2}$ . First of all, notice that in spherical coordinates, the sphere and cone have equations  $\rho = R$  and  $\varphi = \pi/4$  respectively (more generally the equation  $\varphi = \alpha$ ,  $\alpha \in (0, \pi)$  describes a cone with tip at the origin and axis the  $z$ -axis). This means that  $V^0 = \mathbf{g}(U)$ , where  $U = (0, R) \times [0, 2\pi) \times [0, \pi/4)$ ,

so that

$$\begin{aligned}
 \text{vol } V &= \text{vol } V^0 = \int_{\mathbf{g}(U)} 1 = \int_U \rho^2 \sin \varphi \, d\rho \, d\varphi \, d\theta \\
 &= \int_0^{2\pi} \int_0^{\pi/4} \int_0^R \rho^2 \sin \varphi \, d\rho \, d\varphi \, d\theta = 2\pi \frac{R^3}{3} \int_0^{\pi/4} \sin \varphi \, d\varphi \\
 &= 2\pi \frac{R^3}{3} \left(1 - \frac{\sqrt{2}}{2}\right).
 \end{aligned}$$



Notice that the volume of the whole ball is obtained by replacing the upper bound of  $\pi/4$  in the last integral with  $\pi$ , which results in  $4\pi R^3/3$ .

- (ii) More generally, one can compute the volume of the ball  $B^n(R)$  of radius  $R > 0$  around  $\mathbf{0}$  in  $\mathbb{R}^n$  as follows: first of all, notice that the diffeomorphism  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $\mathbf{g}(\mathbf{x}) = R\mathbf{x}$ , maps  $B^n(1)$  onto  $B^n(R)$ . Thus, by the change of variables theorem,

$$\text{vol}(B^n(R)) = \int_{B^n(1)} |\det D\mathbf{g}| = \int_{B^n(1)} R^n = R^n \text{vol}(B^n(1)). \quad (4.6.5)$$

Next, write  $B^n(1) = \{(\mathbf{x}, \mathbf{a}) \mid \mathbf{x} \in B^2(1), \mathbf{a} \in \mathbb{R}^{n-2}, |\mathbf{a}|^2 + |\mathbf{x}|^2 \leq 1\}$ , and express  $B^2(1)$  in polar coordinates to obtain with (4.6.5)

$$\begin{aligned}
 \text{vol}(B^n(1)) &= \int_0^{2\pi} \int_0^1 \text{vol}(B^{n-2}(\sqrt{1-r^2})) r \, dr \, d\theta \\
 &= 2\pi \int_0^1 (1-r^2)^{\frac{n-2}{2}} \text{vol}(B^{n-2}(1)) r \, dr \\
 &= 2\pi \text{vol}(B^{n-2}(1)) \int_0^1 (1-r^2)^{\frac{n-2}{2}} r \, dr \\
 &= \frac{2\pi}{n} \text{vol}(B^{n-2}(1)).
 \end{aligned}$$



An easy induction argument together with (4.6.5) then yields

$$\text{vol } B^{2n+1}(R) = \frac{\pi^n 2^{n+1} R^{2n+1}}{(2n+1)(2n-1)\cdots 5 \cdot 3}, \quad \text{vol } B^{2n}(R) = \frac{\pi^n R^{2n}}{n!}.$$

- (iii) Spherical coordinates in  $\mathbb{R}^3$  can be extended in exactly the same way to  $\mathbb{R}^4$ : change the notation for  $\varphi$  to  $\varphi_1$ , and define  $\varphi_2 : \mathbb{R}^4 \setminus \{\mathbf{0}\} \rightarrow [0, \pi]$  by

$$\varphi_2(\mathbf{x}) = \arccos(\langle \mathbf{x}, \mathbf{e}_4 \rangle / |\mathbf{x}|), \quad \mathbf{x} \in \mathbb{R}^4 \setminus \{\mathbf{0}\}.$$

Then  $\mathbf{x} = (\pi(\mathbf{x}), \rho(\mathbf{x}) \cos \varphi_2(\mathbf{x}))$ , where  $\rho$  denotes the distance function to the origin in  $\mathbb{R}^4$  and  $\pi : \mathbb{R}^3 \times \mathbb{R} \rightarrow \mathbb{R}^3$  the projection. Using 3-dimensional spherical coordinates for  $\pi(\mathbf{x})$ , we obtain

$$\mathbf{x} = (\rho \cos \theta \sin \varphi_1 \sin \varphi_2, \rho \sin \theta \sin \varphi_1 \sin \varphi_2, \rho \cos \varphi_1 \sin \varphi_2, \rho \cos \varphi_2)(\mathbf{x}).$$

A straightforward if tedious induction argument now yields spherical coordinates  $(\rho, \theta, \varphi_1, \dots, \varphi_{n-2}) : \mathbb{R}^n \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}^n$  on  $\mathbb{R}^n$  with inverse  $\mathbf{g}$ , where

$$\mathbf{g}(\rho, \theta, \varphi_1, \dots, \varphi_{n-2}) = \rho(\cos \theta \sin \varphi_1 \cdots \sin \varphi_{n-2}, \sin \theta \sin \varphi_1 \cdots \sin \varphi_{n-2}, \cos \varphi_1 \sin \varphi_2 \cdots \sin \varphi_{n-2}, \dots, \cos \varphi_{n-2}), \quad (4.6.6)$$

and

$$|\det D\mathbf{g}|(\rho, \theta, \varphi_1, \dots, \varphi_{n-2}) = \rho^{n-1} \sin^{n-2} \varphi_{n-2} \sin^{n-3} \varphi_{n-3} \cdots \sin \varphi_1. \quad (4.6.7)$$

This provides an alternative method for deriving the volume of a ball in  $\mathbb{R}^n$ , see Exercise 4.25.

- (iv) Recall that for a continuous function  $f$  of one variable, the *improper integral*  $\int_a^\infty f$  is defined to be  $\lim_{x \rightarrow \infty} \int_a^x f$ , provided the limit exists, in which case the improper integral is said to *converge*. Similarly,  $\int_{-\infty}^a f = \lim_{x \rightarrow -\infty} \int_x^a f$ , and if both exist, we define  $\int_{-\infty}^\infty f = \int_{-\infty}^a f + \int_a^\infty f$  (notice that the particular value of  $a$  is irrelevant). Consider the function  $f$  given by  $f(x) = e^{-x^2}$ . Even though  $f$  has an antiderivative  $F$ , it is well known that there is no formula for  $F$  other than  $F(x) = \int_a^x e^{-t^2}$  for some  $a \in \mathbb{R}$ . Nevertheless, we will compute explicitly  $\int_{-\infty}^\infty f$ . First of all, observe that  $\int_0^\infty f$  exists, because  $e^{-x^2} \leq 1/x^2$  for large  $x$ , and  $\int_a^\infty 1/x^2 dx$  converges. Furthermore,  $\int_{-\infty}^\infty f = \lim_{R \rightarrow \infty} \int_{-R}^R f$  since  $f$  is even. If  $S_R = [-R, R] \times [-R, R]$ , then

$$\left( \int_{-R}^R f \right)^2 = \left( \int_{-R}^R e^{-x^2} dx \right) \left( \int_{-R}^R e^{-y^2} dy \right) = \int_{S_R} e^{-(x^2+y^2)} dx dy \quad (4.6.8)$$

by Fubini's theorem. Let  $B_R$  denote the disk of radius  $R$  centered at the origin. Observe that if  $g(x, y) = e^{-(x^2+y^2)}$ , then  $\lim_{R \rightarrow \infty} \int_{B_R} g = \lim_{R \rightarrow \infty} \int_{S_R} g$  in the sense

that if one exists, then so does the other and the two are equal: indeed,  $S_{R/\sqrt{2}} \subset B_R \subset S_{2R}$ , and  $|g| \leq e^{-R^2}$  on  $S_{2R} \setminus B_R$ . Similarly,  $|g| \leq e^{-R^2/2}$  on  $B_R \setminus S_{R/\sqrt{2}}$ , so that

$$\int_{B_R} g - \int_{S_{R/\sqrt{2}}} g \leq e^{-R^2/2} R^2 (\pi - 2) \xrightarrow{R \rightarrow \infty} 0,$$

and

$$\int_{S_{2R}} g - \int_{B_R} g \leq e^{-R^2} R^2 (4 - \pi) \xrightarrow{R \rightarrow \infty} 0,$$

thereby proving the claim. Using polar coordinates,

$$\int_{B_R} g = \int_0^{2\pi} \int_0^R e^{-r^2} r \, dr \, d\theta = \pi(1 - e^{-R^2}),$$

so that by (4.6.8),

$$\int_{-\infty}^{\infty} e^{-x^2} \, dx = \left( \lim_{R \rightarrow \infty} \int_{B_R} g \right)^{\frac{1}{2}} = \sqrt{\pi}.$$

## 4.7 Some applications

Before discussing applications of integration to concepts from physics, we observe that if  $f : [a, b] \rightarrow \mathbb{R}$  is integrable, then its integral may be evaluated by considering only partitions  $P_n$  of  $[a, b]$  into  $n$  subintervals

$$I_j^n = \left[ a + (j-1) \frac{b-a}{n}, a + j \frac{b-a}{n} \right]$$

of equal length  $l(I_j^n) = (b-a)/n$ . If  $x_j^n$  is any point in  $I_j^n$ , the expression

$$\sum_{j=1}^n f(x_j^n) l(I_j^n)$$

is called a *Riemann sum* of  $f$  for  $P_n$ . Since this sum is sandwiched between  $L(f, P_n)$  and  $U(f, P_n)$ , and the latter two sequences converge to the integral of  $f$ ,

$$\int_a^b f = \lim_{n \rightarrow \infty} \sum_{j=1}^n f(x_j^n) l(I_j^n).$$

The same is of course true for integrals of functions of more than one variable, if one considers partitions by boxes of the same size. This turns out to be useful in defining physical concepts as limits of “approximations”.

### 4.7.1 Mass

Consider a body occupying a region  $E$  in 3-space, made of a not necessarily homogeneous material. Its *density*  $\rho(\mathbf{r})$  at an interior point  $\mathbf{r} \in E$  is

$$\rho(\mathbf{r}) = \lim_{\varepsilon \rightarrow 0} \frac{\text{mass}(B_\varepsilon(\mathbf{r}))}{\text{vol}(B_\varepsilon(\mathbf{r}))},$$

provided the limit exists. We assume that the density is defined everywhere and is, in fact, continuous. In order to evaluate the mass of the whole body, we begin by approximating it with Riemann sums. Suppose first that  $E$  is a box, and consider a partition  $P_n$  of  $E$  by  $n^3$  boxes of equal size  $B_1, \dots, B_{n^3}$ , obtained by partitioning each side of the box into  $n$  subintervals. If  $n$  is large enough that each  $B_i$  is very small, the continuity of  $\rho$  ensures that for any  $\mathbf{r}_i \in B_i$ ,  $\rho(\mathbf{r}_i) \text{vol}(B_i)$  is a fair approximation of the mass of  $B_i$ . Thus,  $\sum_{i=1}^{n^3} \rho(\mathbf{r}_i) \text{vol}(B_i)$  is an approximation to the mass  $m$  of the body, and it is expected that the approximation gets better as  $n$  is larger. Since this is a Riemann sum of the continuous function  $\rho$  for  $P_n$ , its limit as  $n \rightarrow \infty$  equals the integral of  $\rho$  over  $E$ , and we define the mass  $m$  of the body to equal

$$m = \int_E \rho.$$

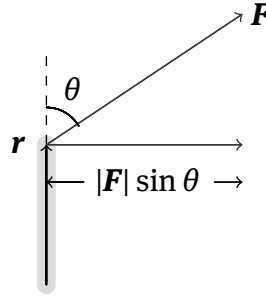
This is easily extended to a more general body, by enclosing  $E$  inside some box and integrating  $\rho \chi_E$  over the box. The above formula is therefore still valid.

### 4.7.2 Center of mass

When trying to place a two-dimensional rigid object such as a tray on the tip of a thin vertical rod, a little experimentation shows that there is one and only one point on the object which, when in contact with the rod, leaves the object balanced. This point is called the center of mass of the object, and can be defined more generally for a three-dimensional body. In order to determine this point, the concept of *torque* is useful.

In physics, the torque  $\boldsymbol{\tau}$  of a force  $\mathbf{F}$  is a vector that measures the tendency of that force to rotate an object about an axis or a point. If the point where the torque is measured is located at the origin, and the object is at the end of a lever arm at position  $\mathbf{r}$ , then the magnitude of the torque is experimentally determined to be proportional to both the length of the lever arm and to the magnitude of the component of the force orthogonal to the arm. In other words, if  $\theta$  is the angle between  $\mathbf{F}$  and  $\mathbf{r}$ , the magnitude of the torque is proportional to  $|\mathbf{r}||\mathbf{F}| \sin \theta$ . This is also the magnitude of  $\mathbf{r} \times \mathbf{F}$ , and the torque is therefore defined to be

$$\boldsymbol{\tau} = \mathbf{r} \times \mathbf{F}.$$



If we now have a discrete system of point masses  $m_i$  at positions  $\mathbf{r}_i$ ,  $i = 1, \dots, n$ , then the center of mass of the system is that point  $\mathbf{R}$  with respect to which the total torque of gravity on the point masses is zero, so that rotationally speaking, the system behaves as though all mass were concentrated at the center of mass. Since the torque due to the  $i$ -th object is  $(\mathbf{r}_i - \mathbf{R}) \times -m_i g \mathbf{k}$ , where  $g$  is the gravitational constant, the total torque is  $\sum_i (\mathbf{r}_i - \mathbf{R}) \times -m_i g \mathbf{k}$ . Setting this equal to zero implies that the vector  $\sum_i m_i (\mathbf{r}_i - \mathbf{R})$  is parallel to  $\mathbf{k}$ . If this is to hold for any rotation of the system about its center of mass, then the vector itself must be zero, since this amounts to replacing  $\mathbf{k}$  by an arbitrary vector. Thus, the center of mass is located at

$$\mathbf{R} = \frac{1}{M} \sum_i m_i \mathbf{r}_i,$$

where  $M = \sum_i m_i$  is the total mass.

The above discussion for discrete mass distributions generalizes to continuous mass distributions: consider a solid occupying a Jordan-measurable region  $E$  in space, with possibly variable density  $\rho$ . Assume first that  $E$  is a box, and partition it into subboxes  $B_i$  of equal size,  $i = 1, \dots, n^3$ . If  $\mathbf{r}_i$  denotes the center of gravity of  $B_i$  and  $B_i$  is small enough, we approximate the mass of  $B_i$  by  $\rho(\mathbf{r}_i) \Delta V$ , where  $\Delta V$  is the common volume of the subboxes. Denoting by  $\mathbf{R}$  the position vector of the center of mass, the torque due to  $B_i$  experienced at the center of mass is approximately  $(\mathbf{r}_i - \mathbf{R}) \times \rho(\mathbf{r}_i) \Delta V g \mathbf{k}$ . Eliminating  $g \mathbf{k}$  as above, we see that the vector  $\sum_i (\mathbf{r}_i - \mathbf{R}) \rho(\mathbf{r}_i) \Delta V$  must go to zero as  $n \rightarrow \infty$ . Its three components are limits of Riemann sums, and if  $\mathbf{R} = [\bar{x} \ \bar{y} \ \bar{z}]^T$ , then the first component equals  $\int_E (x - \bar{x}) \rho(x, y, z) dx dy dz$ . Setting this equal to zero and doing the same with the other components, we obtain for the coordinates of the center of mass

$$\begin{aligned} \bar{x} &= \frac{1}{M} \int_E x \rho(x, y, z) dx dy dz, \\ \bar{y} &= \frac{1}{M} \int_E y \rho(x, y, z) dx dy dz, \\ \bar{z} &= \frac{1}{M} \int_E z \rho(x, y, z) dx dy dz. \end{aligned}$$

The same argument used in discussing the mass of an object shows that the above identities hold for any Jordan-measurable body  $E$ , not just box-like ones.

When the object is homogeneous, that is, when the density is constant, the center of mass is called the *centroid*.

**Example 4.7.1.** Let us find the centroid of a solid hemisphere  $E$  of radius  $a$ . Since the density is constant, the mass  $M$  equals density times the volume of  $E$ , and the first coordinate of the centroid is

$$\bar{x} = \frac{1}{\text{vol}(E)} \int_E x \, dx \, dy \, dz = \frac{3}{2\pi a^3} \int_E x \, dx \, dy \, dz,$$

with similar formulae for the other coordinates. Assuming the solid is the Northern hemisphere of a sphere centered at the origin, we would expect, by symmetry, that  $\bar{x} = \bar{y} = 0$ . This is indeed the case: using spherical coordinates,

$$\begin{aligned} \int_E x \, dx \, dy \, dz &= \int_0^{2\pi} \int_0^{\pi/2} \int_0^a (\rho \sin \varphi \cos \theta) \cdot \rho^2 \sin \varphi \, d\theta \, d\varphi \, d\rho \\ &= \int_0^{2\pi} \cos \theta \, d\theta \int_0^{\pi/2} \sin^2 \varphi \, d\varphi \int_0^a \rho^3 \, d\rho \\ &= 0 \end{aligned}$$

because the first integral on the second to last line vanishes. A similar calculation yields  $\bar{y} = 0$ . Finally,

$$\begin{aligned} \int_E z \, dx \, dy \, dz &= \int_0^{2\pi} \int_0^{\pi/2} \int_0^a (\rho \cos \varphi) \cdot \rho^2 \sin \varphi \, d\theta \, d\varphi \, d\rho \\ &= 2\pi \int_0^{\pi/2} \sin \varphi \cos \varphi \, d\varphi \int_0^a \rho^3 \, d\rho = 2\pi \left. \frac{\sin^2 \varphi}{2} \right|_0^{\pi/2} \left. \frac{\rho^4}{4} \right|_0^a \\ &= \frac{\pi a^4}{4}, \end{aligned}$$

so that  $\bar{z} = (3a)/8$ .

### 4.7.3 Moment of inertia

According to Newton's second law, a force must be applied to an object moving along a straight line in order to change its velocity. The magnitude  $F$  of the force is proportional to the object's acceleration  $a$ , and the constant of proportionality is the mass  $m$  of the object:  $F = ma$ .

A similar law holds for objects that are rotating about an axis. In this case, torque has to be applied in order to change the object's angular velocity: recall that if the

object is at distance  $r$  from the axis, and  $F$  is the (norm of the) component of the force tangent to the object's circular path, then the amount of torque measured at the axis is  $\tau = rF$ . Now, if  $m$  denotes the object's mass, then its acceleration is  $a = F/m$ , and its angular acceleration is  $\alpha = a/r = F/(mr)$ . Therefore, the applied torque

$$\tau = (mr^2)\alpha$$

is proportional to the angular acceleration. The constant of proportionality,  $I = mr^2$ , is called the *moment of inertia* of the object about the axis.

The discussion generalizes to discrete mass distributions and continuous ones in exactly the same way as was done for the center of mass. Thus, the moment of inertia of a solid body occupying a region  $E$  with density  $\rho$  about an axis is given by

$$I = \int_E \rho r^2, \quad (4.7.1)$$

where  $r : E \rightarrow [0, \infty)$  is the distance function to the axis.

**Example 4.7.2.** Let us find the moment of inertia of a solid ball with constant density and radius  $a$  about any axis through its center. By symmetry, the ball may be assumed to be centered at the origin and the axis is the  $z$ -axis. Using spherical coordinates, and denoting by  $\bar{\rho}$  the density of the ball to avoid confusion with the spherical coordinate  $\rho$ ,

$$\begin{aligned} I &= \int_{\{(x,y,z)|x^2+y^2+z^2 \leq a^2\}} \bar{\rho}(x^2 + y^2) dx dy dz \\ &= \bar{\rho} \int_0^{2\pi} \int_0^a \int_0^\pi (\rho^2 \sin^2 \varphi) \cdot \rho^2 \sin \varphi d\varphi d\rho d\theta = 2\pi \bar{\rho} \int_0^a \rho^4 d\rho \int_0^\pi \sin^3 \varphi d\varphi \\ &= 2\pi \bar{\rho} \cdot \frac{a^5}{5} \cdot \frac{4}{3}. \end{aligned}$$

Since  $\bar{\rho}$  is constant, the mass  $m$  of the ball equals  $(4/3)\pi a^3 \bar{\rho}$ , so the moment of inertia may also be written as  $I = (2/5)ma^2$ .

## 4.8 Exercises

**4.1.** Prove the last two assertions of Theorem 4.1.2.

**4.2.** Prove or disprove:

- If  $f$  is a bounded function on  $A \subset \mathbb{R}^n$  and  $|f|$  is integrable on  $A$ , then so is  $f$ .
- If  $U$  is open in  $\mathbb{R}^n$  and  $f : U \rightarrow \mathbb{R}$  is continuous at some  $\mathbf{p} \in U$ , then  $f$  is integrable on  $B_r(\mathbf{p})$  for small enough  $r > 0$ .
- If  $A \subset \mathbb{R}^n$  has measure zero, then  $A$  is bounded.
- If  $A \subset \mathbb{R}^n$  has measure zero, then the boundary of  $A$  has measure zero.

**4.3.** Let  $A$  be a compact set of measure zero in  $\mathbb{R}^n$ . Show that if  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^k$  is continuous, then  $\mathbf{g}(A)$  has measure zero in  $\mathbb{R}^k$ .

**4.4.** It was shown in Remark 4.2.2 that if  $A \subset \mathbb{R}^n$  is Jordan-measurable and has measure zero, then it has volume zero. Show that the converse is also true, so that for Jordan-measurable sets, the two concepts are equivalent.

**4.5.** Show that any compact submanifold of  $\mathbb{R}^n$  has measure zero in  $\mathbb{R}^n$ .

**4.6.**  $A \subset \mathbb{R}^n$  is said to have *content zero* if for any  $\varepsilon > 0$ , there exists a finite cover  $\{B_1, \dots, B_k\}$  of  $A$  by boxes with  $\sum_i \text{vol}(B_i) < \varepsilon$ .

(a) Show that if  $A$  is compact and has measure zero, then it has content zero.

(b) Give examples of sets of measure zero that do not have content zero.

**4.7.** Let  $f : U \rightarrow \mathbb{R}^n$  be a function which is integrable over some box  $A = \prod_{i=1}^n [a_i, b_i]$  in  $U$ . For each positive integer  $k$ , consider the partition  $P_k$  of  $A$  obtained by partitioning each  $[a_i, b_i]$  into  $k$  intervals of equal length, and choose some point  $\mathbf{a}_B$  in each subbox  $B \in P_k$ . Prove that

$$\lim_{k \rightarrow \infty} \sum_{B \in P_k} f(\mathbf{a}_B) \text{vol}(B) = \int_A f.$$

**4.8.** Prove that if  $U$  is a bounded open set in  $\mathbb{R}^n$ , then there exists a sequence of smooth functions  $f_k : \mathbb{R}^n \rightarrow \mathbb{R}$  such that

$$\lim_{k \rightarrow \infty} \int_{\mathbb{R}^n} f_k = \text{vol}(U).$$

**4.9.** Given  $f : [a, b] \rightarrow \mathbb{R}$ , define  $g : [a, b] \times [c, d] \rightarrow \mathbb{R}$  by  $g(x, y) = f(x)$ . Using only the definition of integral, show that  $g$  is integrable if and only if  $f$  is, and if it is, then  $\int_{[a,b] \times [c,d]} g = (d - c) \int_{[a,b]} f$ .

**4.10.** Let  $U$  be an open set in  $\mathbb{R}^n$ , and suppose  $f : U \rightarrow \mathbb{R}$  is integrable on  $U$  and continuous at some  $\mathbf{p} \in U$ .

(a) If  $C_r(\mathbf{p})$  denotes the cube

$$[u^1(\mathbf{p}) - r/2, u^1(\mathbf{p}) + r/2] \times \cdots \times [u^n(\mathbf{p}) - r/2, u^n(\mathbf{p}) + r/2]$$

with sides of length  $r$  centered at  $\mathbf{p}$ , show that

$$\lim_{r \rightarrow 0^+} \frac{1}{\text{vol}(C_r(\mathbf{p}))} \int_{C_r(\mathbf{p})} f = f(\mathbf{p}).$$

(b) Prove that

$$\lim_{r \rightarrow 0^+} \frac{1}{\text{vol}(B_r(\mathbf{p}))} \int_{B_r(\mathbf{p})} f = f(\mathbf{p}).$$

**4.11.** Let  $A \subset \mathbb{R}^n$  be compact, connected, and Jordan-measurable. If  $f : A \rightarrow \mathbb{R}$  is continuous, show that there exists some  $\mathbf{a} \in A$  such that

$$\int_A f = f(\mathbf{a}) \cdot \text{vol}(A).$$

**4.12.** Let  $f : [a, b] \times [c, d] \rightarrow \mathbb{R}$  be continuous.

(a) If  $g(x, t) = \int_c^t f(x, y) dy$ , find  $D_1 G$  and  $D_2 G$ .

(b) If  $h(s, t) = \int_a^s g(x, t) dx = \int_a^s (\int_c^t f(x, y) dy) dx$ , find  $D_{12} h$ .

**4.13.** Let  $f : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$  be continuous. Prove that

$$\int_0^1 (\int_0^x f(x, y) dy) dx = \int_0^1 (\int_y^1 f(x, y) dx) dy.$$

**4.14.** Let  $A$  denote a Jordan-measurable set in  $\mathbb{R}^n$ , and  $B = A \times [0, 1] \subset \mathbb{R}^{n+1}$ . Given  $\mathbf{u} \in \mathbb{R}^n$ , define

$$B_{\mathbf{u}} = \{\mathbf{a} + t(\mathbf{u} + \mathbf{e}_{n+1}) \mid \mathbf{a} \in A, t \in [0, 1]\} \subset \mathbb{R}^{n+1}.$$

Thus,  $B = B_{\mathbf{0}}$  is a right solid cylinder over  $A$  of height 1, and  $B_{\mathbf{u}}$  is a slanted cylinder over  $A$  of the same height. Prove that  $\text{vol } B_{\mathbf{u}} = \text{vol } B$  for any  $\mathbf{u} \in \mathbb{R}^n$ .

**4.15.** (a) Define  $f : (0, 1) \rightarrow \mathbb{R}$  by

$$f(x) = -n, \quad x \in (1 - \frac{1}{n}, 1 - \frac{1}{n+1}).$$

Show that  $\lim_{x \rightarrow 1} \int_{(0,x)} f$  exists, but  $\int_{(0,1)} f$  does not.

(b) Prove that if  $f : (0, 1) \rightarrow \mathbb{R}$  is continuous, then  $\lim_{x \rightarrow 1} \int_{(0,x)} f$  exists if and only if  $\int_{(0,1)} f$  does, and if they do exist, then they coincide. Give an example of a continuous  $f : (0, 1) \rightarrow \mathbb{R}$  that is not integrable over  $(0, 1)$ .

**4.16.** Define  $f : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$  by

$$f(x, y) = \begin{cases} 1 & \text{if } x \text{ is rational,} \\ y & \text{otherwise.} \end{cases}$$

(a) Show that  $\int_0^1 \int_0^1 f(x, y) dy dx$  exists, and find it.

(b) Prove that  $\int_0^1 \int_0^1 f(x, y) dx dy$  does not exist.

(c) Show that  $f$  is not integrable over its domain.

**4.17.** Explain why one may assume that the Jacobian determinant of  $\mathbf{g}$  is nowhere zero in the proof of the change of variables theorem.

**4.18.** Prove that Euclidean motions preserve volume; i.e., if  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a Euclidean motion, and  $A \subset \mathbb{R}^n$  is Jordan-measurable, then  $\text{vol}(\mathbf{f}(A)) = \text{vol}(A)$ .



**4.19.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be continuously differentiable.

(a) Suppose that  $|\det Df(\mathbf{p})| > 1$  for some  $\mathbf{p}$ . Show that  $\text{vol}(f(B_r(\mathbf{p}))) > \text{vol}(B_r(\mathbf{p}))$  for sufficiently small  $r > 0$ .

(b) More generally, prove that if  $|\det Df(\mathbf{p})| \neq 0$ , then

$$|\det Df(\mathbf{p})| = \lim_{r \rightarrow 0^+} \frac{\text{vol}(f(B_r(\mathbf{p})))}{\text{vol}(B_r(\mathbf{p}))}.$$

**4.20.** Evaluate  $\int_A f$ , if  $f(x, y) = \cos \frac{x}{x+y}$  and  $A \subset \mathbb{R}^2$  is the region inside the triangle with vertices  $(0, 0)$ ,  $(0, 1)$ , and  $(1, 0)$ . *Hint:* Use the change of variables  $u = x$ ,  $v = x + y$ .

**4.21.** The set of points  $(x, y)$  in the plane satisfying  $(x^2 + y^2 - y)^2 - x^2 - y^2 = 0$  is called a *cardioid*. Sketch this curve and determine the area of the region it encloses.

**4.22.** Use an appropriate change of variables to evaluate  $\int_A f$ , if  $f(x, y) = \sin(2x^2 + y^2)$ , and  $A = \{(x, y) \in \mathbb{R}^2 \mid y \geq 0, 2x^2 + y^2 \leq 2\}$ .

**4.23.** Evaluate once again the volume of the cone from Examples 4.6.1 (i), but using spherical coordinates rather than cylindrical ones.

**4.24.** Determine the volume of the region in  $\mathbb{R}^3$  that is bounded by the paraboloid  $z = x^2 + y^2$  and the sphere  $x^2 + y^2 + z^2 = 2$ .

**4.25.** (a) Use induction and integration by parts to show that

$$\int_0^{\pi} \sin^{2n+1} x \, dx = \frac{2^{n+1} n!}{(2n+1)(2n-1)\cdots 5 \cdot 3},$$

$$\int_0^{\pi} \sin^{2n} x \, dx = \pi \frac{(2n-1)(2n-3)\cdots 5 \cdot 3}{2^n n!}.$$

(b) Use (4.6.7) to prove that the volume of a ball  $B^n(R)$  of radius  $R$  in  $\mathbb{R}^n$  is given by

$$\text{vol } B^{2n+1}(R) = \frac{\pi^n 2^{n+1} R^{2n+1}}{(2n+1)(2n-1)\cdots 5 \cdot 3}, \quad \text{vol } B^{2n}(R) = \frac{\pi^n R^{2n}}{n!}.$$

**4.26.** Let  $a_i > 0$ ,  $1 \leq i \leq n$ . Use an appropriate change of variables to evaluate the volume of the region in  $\mathbb{R}^n$  bounded by the ellipsoid

$$\{(x_1, \dots, x_n) \in \mathbb{R}^n \mid \sum_{i=1}^n \frac{x_i^2}{a_i^2} = 1\}.$$

**4.27.** The *gamma function* is defined by

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} \, dt, \quad 0 < x < \infty.$$

- (a) Show that this improper integral converges for all  $x > 0$ . (Notice that if  $x < 1$ , the integral is also improper at 0; i.e., it must be shown that

$$\int_0^a t^{x-1} e^{-t} dt := \lim_{\varepsilon \rightarrow 0^+} \int_{\varepsilon}^a t^{x-1} e^{-t} dt$$

converges for some and hence all  $a > 0$ ).

- (b) Show that  $\Gamma(1) = 1$ .  
 (c) Use the change of variables  $t = u^2$  and Examples and Remarks 4.6.1 (iv) to prove that  $\Gamma(1/2) = \sqrt{\pi}$ .  
 (d) Use integration by parts to show that  $\Gamma(x + 1) = x\Gamma(x)$ , and deduce that  $\Gamma(n) = (n - 1)!$  for all  $n \in \mathbb{N}$ .  
 (e) In Examples and Remarks 4.6.1 (ii), a formula was derived for the volume  $\text{vol}(B^n(R))$  of a ball of radius  $R$  in  $\mathbb{R}^n$ . There were actually two formulas, one for even  $n$  and the other for odd  $n$ . Prove that they can be unified into one by means of the gamma function:

$$\text{vol}(B^n(R)) = \frac{2R^n \pi^{n/2}}{n\Gamma(n/2)}.$$

- 4.28.** (a) Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be given by

$$f(x) = \begin{cases} 1 & \text{if } |x| \leq 1, \\ \frac{1}{x^2} & \text{if } |x| \geq 1. \end{cases}$$

Show that  $\lim_{r \rightarrow \infty} \int_{-r}^r f$  exists.

- (b) Let  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$  be given by

$$g(\mathbf{x}) = \begin{cases} 1 & \text{if } |\mathbf{x}| \leq 1, \\ \frac{1}{|\mathbf{x}|^2} & \text{if } |\mathbf{x}| \geq 1. \end{cases}$$

Show that  $\lim_{r \rightarrow \infty} \int_{[-r,r] \times [-r,r]} g$  does not exist.

- 4.29.** Reprove Theorem 2.3.1 by using the methods from this chapter; i.e., show that if  $f, D_2f : [a, b] \times [c, d] \rightarrow \mathbb{R}$  are continuous, then the function

$$\begin{aligned} \varphi : [c, d] &\longrightarrow \mathbb{R}, \\ y &\longmapsto \int_a^b f(x, y) dx \end{aligned}$$

is differentiable on  $(c, d)$  and

$$\varphi'(y_0) = \int_a^b D_2f(x, y_0) dx, \quad y_0 \in (c, d).$$

*Hint:*  $\varphi(y) = \int_a^b \left( \int_c^y D_2f(x, t) dt + f(x, c) \right) dx$ . Notice also that the hypotheses may be somewhat weakened.

**4.30.** This exercise uses Sard's theorem to show that any smooth map  $\mathbf{f} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^m$  can be “approximated” by an immersion if  $m \geq 2n$ ; more specifically, given any  $\varepsilon > 0$ , there exists an  $m \times n$  matrix  $A = (a_{ij})$  with  $|a_{ij}| < \varepsilon$  for all  $i$  and  $j$ , such that the map

$$\begin{aligned} U &\longrightarrow \mathbb{R}^m \\ \mathbf{x} &\longmapsto \mathbf{f}(\mathbf{x}) + A\mathbf{x} \end{aligned}$$

is an immersion. Recall from Exercise 3.20 that the collection  $M_{m,n}(k)$  of all  $m \times n$  matrices of rank  $k$  is a manifold of dimension  $k(m + n - k)$ . Define

$$\begin{aligned} \mathbf{g}_k : U \times M_{m,n}(k) &\longrightarrow M_{m,n} \\ (\mathbf{u}, B) &\longmapsto B - [D\mathbf{f}(\mathbf{u})]. \end{aligned}$$

- Show that if  $m \geq 2n$  and  $k < n$ , then the dimension of the domain of  $\mathbf{g}_k$  is less than that of its image. *Hint:* the function  $t \mapsto t(m + n - t)$  is increasing if  $t < n < m$ .
- Show that the image of  $\mathbf{g}_k$  has measure zero under the assumptions from part (a), and prove the claim made at the beginning of the exercise.
- Show that if  $\mathbf{f} : \mathbb{R}^n \supset U \rightarrow \mathbb{R}^m$  is smooth, with  $U$  bounded and  $m \geq 2n$ , then for any  $\varepsilon > 0$  there exists an immersion  $\mathbf{g} : U \rightarrow \mathbb{R}^m$  such that  $|\mathbf{g}(\mathbf{x}) - \mathbf{f}(\mathbf{x})| < \varepsilon$  for all  $\mathbf{x} \in U$ .

This result can be used in proving the so-called “Whitney imbedding theorem” which states that any  $n$ -dimensional manifold  $M$  can be imbedded in  $\mathbb{R}^{2n+1}$  (see for example [14]); i.e., there exists an injective map  $\mathbf{f} : M \rightarrow \mathbb{R}^{2n+1}$  of maximal rank everywhere with continuous inverse. In particular  $\mathbf{f}(M)$  is a submanifold of  $\mathbb{R}^{2n+1}$  (any parametrization  $\mathbf{h}$  of  $M$  generates a parametrization of  $\mathbf{f} \circ \mathbf{h}$  of  $\mathbf{f}(M)$ ) which is diffeomorphic to  $M$  (via  $\mathbf{f}$ ), so that  $M$  itself may be considered to be sitting inside  $\mathbb{R}^{2n+1}$  regardless of the dimension of the original Euclidean space it is a submanifold of.

It should be noted that manifolds can be defined as abstract sets that are not contained in Euclidean space. Whitney's theorem asserts that our definition is equivalent to that one.



## 5 Differential Forms

Now that we are familiar with integration on Euclidean space, we would like to translate this to manifolds. Functions will be replaced by more exotic objects called differential forms. Their main advantage lies in that they have the change of variables formula built-in, once the concept of orientation is introduced. We begin by recalling the more general concept of tensor field.

### 5.1 Tensors and tensor fields

**Definition 5.1.1.** Let  $V$  be a vector space. Given  $k \in \mathbb{N}$ , a  $k$ -tensor on  $V$  is a multilinear map  $T : V^k \rightarrow \mathbb{R}$ , where  $V^k$  denotes the  $k$ -fold Cartesian product of  $V$  with itself.

The collection  $\mathcal{T}_k(V)$  of all  $k$ -tensors on  $V$  is clearly a vector space under the usual addition of maps and scalar multiplication. For example,  $\mathcal{T}_1(V)$  is just the dual space  $V^*$ . By convention,  $\mathcal{T}_0(V)$  is defined to be  $\mathbb{R}$ . An inner product on  $V$  is a 2-tensor (with some additional properties).

**Remark 5.1.1.**  $k$ -tensors are actually a special case of what we defined to be tensors in Chapter 3; i.e., what we call a  $k$ -tensor here is just a tensor of order  $(0, k)$  in the previous terminology. We will only deal with this particular subset in this chapter.

**Definition 5.1.2.** The *tensor product* of  $T \in \mathcal{T}_k(V)$  with  $\tilde{T} \in \mathcal{T}_l(V)$  is the  $(k + l)$ -tensor  $T \otimes \tilde{T}$  given by

$$(T \otimes \tilde{T})(\mathbf{v}_1, \dots, \mathbf{v}_{k+l}) = T(\mathbf{v}_1, \dots, \mathbf{v}_k) \cdot \tilde{T}(\mathbf{v}_{k+1}, \dots, \mathbf{v}_{k+l}), \quad \mathbf{v}_i \in V.$$

The following properties are easy consequences of the definition, and their verification is left to the reader: for tensors  $T_i$ ,  $a \in \mathbb{R}$ ,

$$\begin{aligned} (T_1 \otimes T_2) \otimes T_3 &= T_1 \otimes (T_2 \otimes T_3); \\ T_1 \otimes (T_2 + T_3) &= T_1 \otimes T_2 + T_1 \otimes T_3; \\ (T_1 + T_2) \otimes T_3 &= T_1 \otimes T_3 + T_2 \otimes T_3; \\ a(T_1 \otimes T_2) &= (aT_1) \otimes T_2 = T_1 \otimes (aT_2), \end{aligned}$$

where of course the tensors being added in the second and third identities are assumed to have the same order. In view of the first property, either side will be denoted  $T_1 \otimes T_2 \otimes T_3$ . Observe, though, that the tensor product is, in general, not commutative; i.e.,  $T_1 \otimes T_2$  need not equal  $T_2 \otimes T_1$ .

Tensor products yield explicit bases for the spaces  $\mathcal{T}_k(V)$ :

**Theorem 5.1.1.** Let  $\mathbf{v}_1, \dots, \mathbf{v}_n$  denote a basis of  $V$ . If  $\alpha^1, \dots, \alpha^n \in V^* = \mathcal{T}_1(V)$  is the dual basis (i.e.,  $\alpha^i(\mathbf{v}_j) = \delta_{ij}$ ), then

$$\{\alpha^{i_1} \otimes \dots \otimes \alpha^{i_k} \mid 1 \leq i_1, \dots, i_k \leq n\}$$

is a basis of  $\mathcal{T}_k(V)$ , which therefore has dimension  $n^k$ .

*Proof.* Set  $T_{i_1 \dots i_k} = T(\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k}) \in \mathbb{R}$ . We claim that

$$T = \sum_{i_1, \dots, i_k=1}^n T_{i_1 \dots i_k} \alpha^{i_1} \otimes \dots \otimes \alpha^{i_k},$$

which will show that the set in the statement indeed spans  $\mathcal{T}_k(V)$ . In order to establish the claim, it suffices to check that both sides agree when evaluated on basis elements, since both are multilinear: for if  $\mathbf{w}_i = \sum_j a_{ij} \mathbf{v}_j$  and  $M$  is any  $k$ -tensor, then

$$\begin{aligned} M(\mathbf{w}_1, \dots, \mathbf{w}_k) &= M\left(\sum_{j_1=1}^n a_{1j_1} \mathbf{v}_{j_1}, \dots, \sum_{j_k=1}^n a_{kj_k} \mathbf{v}_{j_k}\right) \\ &= \sum_{j_1, \dots, j_k=1}^n a_{1j_1} \cdots a_{kj_k} M(\mathbf{v}_{j_1}, \dots, \mathbf{v}_{j_k}), \end{aligned}$$

so that  $M$  is entirely determined by what it does to basis vectors. Now,  $\alpha^i(\mathbf{v}_j) = \delta_{ij}$ , so that

$$\begin{aligned} &\sum_{i_1, \dots, i_k=1}^n T_{i_1 \dots i_k} \alpha^{i_1} \otimes \dots \otimes \alpha^{i_k}(\mathbf{v}_{j_1}, \dots, \mathbf{v}_{j_k}) \\ &= \sum_{i_1, \dots, i_k=1}^n T_{i_1 \dots i_k} \alpha^{i_1}(\mathbf{v}_{j_1}) \cdots \alpha^{i_k}(\mathbf{v}_{j_k}) = T_{j_1 \dots j_k} \\ &= T(\mathbf{v}_{j_1}, \dots, \mathbf{v}_{j_k}), \end{aligned}$$

thereby establishing the claim. Linear independence is similar: suppose that  $\sum a_{i_1 \dots i_k} \alpha^{i_1} \otimes \dots \otimes \alpha^{i_k} = 0$ . As in the last calculation, given any  $j_1, \dots, j_k \in \{1, \dots, n\}$ , applying both sides to  $\mathbf{v}_{j_1}, \dots, \mathbf{v}_{j_k}$  yields  $a_{j_1 \dots j_k} = 0$ . This concludes the proof of the theorem.  $\square$

As in Chapter 3, we define a  $k$ -tensor field on a manifold  $M^n$  to be a map  $T$  that assigns to each  $\mathbf{p} \in M$  a  $k$ -tensor  $T(\mathbf{p}) \in \mathcal{T}_k(M_{\mathbf{p}})$  which is smooth in the sense that for any vector fields  $\mathbf{X}_1, \dots, \mathbf{X}_k$  on  $M$ , the function  $T(\mathbf{X}_1, \dots, \mathbf{X}_k)$  which assigns to  $\mathbf{p} \in M$  the number  $T(\mathbf{p})(\mathbf{X}_1(\mathbf{p}), \dots, \mathbf{X}_k(\mathbf{p}))$  is smooth in the usual sense. If  $(U, \mathbf{x})$  is a local chart of  $M$ , then the restriction of  $T$  to  $U$  is smooth if and only if the functions

$$T_{i_1 \dots i_k} := T\left(\frac{\partial}{\partial x^{i_1}}, \dots, \frac{\partial}{\partial x^{i_k}}\right) : U \rightarrow \mathbb{R}, \quad 1 \leq i_1, \dots, i_k \leq n$$

are differentiable. In fact, they must by definition be differentiable if  $T$  is smooth. Conversely, if these functions are differentiable, and  $x^i$  denotes as usual  $u^i \circ \mathbf{x}$ , then

$$T|_U = \sum_{i_1, \dots, i_k=1}^n T_{i_1 \dots i_k} dx^{i_1} \otimes \dots \otimes dx^{i_k}$$

is smooth because each  $dx^i$  is: recall that, given any vector field  $\mathbf{X}$  on  $U$ ,  $\mathbf{X} = \sum_i dx^i(\mathbf{X}) \partial/\partial x^i$  (see Exercise 3.15), so that  $dx^i(\mathbf{X})$  is indeed a differentiable function on  $U$ .

The collection of all  $k$ -tensor fields on  $M$  is clearly a vector space, denoted  $\mathcal{T}_k(M)$ .

**Definition 5.1.3.** Let  $f : M \rightarrow N$  be a differentiable map,  $T$  a  $k$ -tensor field on  $N$ ,  $k > 0$ . The *pullback* of  $T$  by  $f$  is the  $k$ -tensor field  $f^*T$  on  $M$  given by

$$(f^*T)(\mathbf{p})(\mathbf{v}_1, \dots, \mathbf{v}_k) = T(f(\mathbf{p}))(f_*\mathbf{v}_1, \dots, f_*\mathbf{v}_k), \quad \mathbf{p} \in M, \quad \mathbf{v}_i \in M_{\mathbf{p}}.$$

A zero-tensor field on  $N$  is just a function  $h : N \rightarrow \mathbb{R}$ . In this case, we define  $f^*h$  to equal  $h \circ f$ .

The following properties of the pullback are easily verified and left as an exercise:

**Proposition 5.1.1.** Let  $f : M \rightarrow N$  be a differentiable map between manifolds  $M$  and  $N$ ,  $S, T$   $k$ -tensor fields on  $N$ ,  $R$  an  $l$ -tensor field on  $N$ , and  $h : N \rightarrow \mathbb{R}$  a function. Then

- (1)  $f^*(aS + bT) = af^*S + bf^*T$ ,  $a, b \in \mathbb{R}$ ,
- (2)  $f^*(R \otimes S) = f^*R \otimes f^*S$ ,
- (3)  $f^*(hT) = (h \circ f)f^*T$ .

**Example 5.1.1.** A *Riemannian metric* on  $M$  is a 2-tensor field  $\mathbf{g}$  on  $M$  such that  $\mathbf{g}(\mathbf{p})$  is an inner product on  $M_{\mathbf{p}}$  for all  $\mathbf{p} \in M$ . The *standard Riemannian metric* on  $\mathbb{R}^n$  is  $\mathbf{g} = \sum_i du^i \otimes du^i$ . This is the inner product we've been using all along on each tangent space, since by definition  $\mathbf{g}(\mathbf{D}_i, \mathbf{D}_j) = \delta_{ij}$ , so that  $\{\mathbf{D}_i \mid i = 1, \dots, n\}$  is an orthonormal basis when evaluated at any point. If  $M \subset \mathbb{R}^n$ , the *standard Riemannian metric* on  $M$  is  $\iota^*\mathbf{g}$ , where  $\iota : M \hookrightarrow \mathbb{R}^n$  is the inclusion map and  $\mathbf{g}$  is the standard Riemannian metric on Euclidean space. Again, it is by definition the restriction to each  $M_{\mathbf{p}}$  of the inner product on  $\mathbb{R}_{\mathbf{p}}^n$ , and coincides with the first fundamental tensor field on  $M$  defined in Chapter 3.

**Remark 5.1.2.** In physics, tensors are usually defined in a more convoluted way. For (relative) simplicity, we only discuss Cartesian tensors; i.e., tensors defined on Euclidean space  $\mathbb{R}^n$ . Consider two ordered bases  $\mathbf{v}_i$  and  $\mathbf{w}_j$  of  $\mathbb{R}^n$ , and let  $T$  be a tensor of order  $k$  on  $\mathbb{R}^n$  in our sense of the word. Denote by  $T_{i_1 \dots i_k} = T(\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k})$  the components of  $T$  in the first basis, and by  $\tilde{T}_{j_1 \dots j_k} = T(\mathbf{w}_{j_1}, \dots, \mathbf{w}_{j_k})$  the components of the same tensor in the second one. If  $L = [L_{ij}]$  is the change of basis matrix, then

$$\begin{aligned} \tilde{T}_{j_1 \dots j_k} &= T(\mathbf{w}_{j_1}, \dots, \mathbf{w}_{j_k}) = T\left(\sum_{i_1} L_{j_1 i_1} \mathbf{v}_{i_1}, \dots, \sum_{i_k} L_{j_k i_k} \mathbf{v}_{i_k}\right) \\ &= \sum_{i_1, i_2, \dots, i_k} L_{j_1 i_1} L_{j_2 i_2} \cdots L_{j_k i_k} T(\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k}), \end{aligned}$$

so that

$$\tilde{T}_{j_1 \dots j_k} = \sum_{i_1, i_2, \dots, i_k} L_{j_1 i_1} L_{j_2 i_2} \cdots L_{j_k i_k} T_{i_1 \dots i_k}. \quad (5.1.1)$$

A Cartesian tensor of order  $k$  is then traditionally defined as an “object” determined by  $n^k$  components in any given “coordinate system” (meaning basis), such that components in different coordinate systems are related by (5.1.1). Needless to say, such a condition can be quite difficult to check in a given instance.

## 5.2 Alternating tensors and forms

Recall from Chapter 1 that a  $k$ -tensor  $T$  on a vector space  $V$  is said to be *alternating* or *skew-symmetric* if

$$T(\mathbf{v}_1, \dots, \mathbf{v}_i, \dots, \mathbf{v}_j, \dots, \mathbf{v}_k) = -T(\mathbf{v}_1, \dots, \mathbf{v}_j, \dots, \mathbf{v}_i, \dots, \mathbf{v}_k),$$

for all  $1 \leq i \neq j \leq k$ ,  $\mathbf{v}_l \in V$ ,  $l = 1, \dots, k$ , and *symmetric* if the above equation holds with the minus sign removed. It is easily seen that any 2-tensor  $T$  can be written as a sum  $T_s + T_a$  of a symmetric tensor  $T_s$  and an alternating one  $T_a$ : set

$$T_s(\mathbf{x}, \mathbf{y}) = \frac{1}{2}((T(\mathbf{x}, \mathbf{y}) + T(\mathbf{y}, \mathbf{x}))), \quad T_a(\mathbf{x}, \mathbf{y}) = \frac{1}{2}(T(\mathbf{x}, \mathbf{y}) - T(\mathbf{y}, \mathbf{x})),$$

for  $\mathbf{x}, \mathbf{y} \in V$ . Both symmetric and alternating parts can be generalized to  $k$ -tensors for arbitrary  $k$ . We only outline the latter, since we are particularly interested in alternating tensors for now: given a  $k$ -tensor  $T$ , define a new tensor  $T_a$  by

$$T_a(\mathbf{v}_1, \dots, \mathbf{v}_k) = \frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) T(\mathbf{v}_{\sigma(1)}, \dots, \mathbf{v}_{\sigma(k)}), \quad \mathbf{v}_1, \dots, \mathbf{v}_k \in V,$$

with the terminology from Chapter 1:  $S_k$  denotes the set of permutations on  $k$  elements, and  $\varepsilon(\sigma)$  is the sign of the permutation  $\sigma$ . Notice that for  $k = 2$ , the formula coincides with the 2-tensor  $T_a$  defined earlier. We claim that  $T_a$  is alternating. To see this, let  $\tau$  denote the transposition  $(i, j)$ . Since  $S_k = \{\sigma \circ \tau \mid \sigma \in S_k\}$ ,

$$\begin{aligned} T_a(\mathbf{v}_1, \dots, \mathbf{v}_j, \dots, \mathbf{v}_i, \dots, \mathbf{v}_k) &= T_a(\mathbf{v}_{\tau(1)}, \dots, \mathbf{v}_{\tau(i)}, \dots, \mathbf{v}_{\tau(j)}, \dots, \mathbf{v}_{\tau(k)}) \\ &= \frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) T(\mathbf{v}_{(\sigma \circ \tau)(1)}, \dots, \mathbf{v}_{(\sigma \circ \tau)(k)}) \\ &= -\frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) \varepsilon(\tau) T(\mathbf{v}_{(\sigma \circ \tau)(1)}, \dots, \mathbf{v}_{(\sigma \circ \tau)(k)}) \\ &= -\frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma \circ \tau) T(\mathbf{v}_{(\sigma \circ \tau)(1)}, \dots, \mathbf{v}_{(\sigma \circ \tau)(k)}) \\ &= -\frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) T(\mathbf{v}_{\sigma(1)}, \dots, \mathbf{v}_{\sigma(k)}) \\ &= -T_a(\mathbf{v}_1, \dots, \mathbf{v}_i, \dots, \mathbf{v}_j, \dots, \mathbf{v}_k). \end{aligned}$$

The collection of all alternating  $k$ -tensors on  $V$  is a vector space denoted  $\Lambda_k(V)$ . Its elements are called *k-forms*. The following properties are easily proven, and their verification is left to the reader:

- Proposition 5.2.1.** (1) If  $\alpha \in \Lambda_k(V)$ , then  $\alpha_a = \alpha$ . In particular, for any tensor  $T$  on  $V$ ,  
 $(T_a)_a = T_a$ ;  
 (2) The map  $\mathcal{T}_k(V) \rightarrow \Lambda_k(V)$  which sends  $T$  to  $T_a$  is linear; i.e.,  $(S + T)_a = S_a + T_a$ ,  
 $(cT)_a = cT_a$  for all  $S, T \in \mathcal{T}_k(V)$ ,  $c \in \mathbb{R}$ ;



(3)

$$\begin{aligned}(S \otimes (T_1 + T_2))_a &= (S \otimes T_1)_a + (S \otimes T_2)_a, \\ ((T_1 + T_2) \otimes S)_a &= (T_1 \otimes S)_a + (T_2 \otimes S)_a, \\ c(S \otimes T)_a &= (cS \otimes T)_a = (S \otimes cT)_a\end{aligned}$$

for all  $T, T_1, T_2 \in \mathcal{T}_k(V)$ ,  $S \in \mathcal{T}_l(V)$ ,  $c \in \mathbb{R}$ ;

Since the tensor product of two forms is, in general, no longer alternating, we modify it as follows:

**Definition 5.2.1.** Given  $\alpha \in \Lambda_k(V)$ ,  $\beta \in \Lambda_l(V)$ , their *wedge product* or *exterior product* is the  $(k+l)$ -form  $\alpha \wedge \beta \in \Lambda_{k+l}(V)$  given by

$$\alpha \wedge \beta = \frac{(k+l)!}{k!l!}(\alpha \otimes \beta)_a.$$

The reason for including the factorial term in the definition of the wedge product will be revealed shortly. Notice that if  $\alpha, \alpha^1, \alpha^2 \in \Lambda_k(V)$ ,  $\beta \in \Lambda_l(V)$ , and  $c \in \mathbb{R}$ , then

$$\begin{aligned}(\alpha^1 + \alpha^2) \wedge \beta &= \alpha^1 \wedge \beta + \alpha^2 \wedge \beta, \\ \beta \wedge (\alpha^1 + \alpha^2) &= \beta \wedge \alpha^1 + \beta \wedge \alpha^2, \\ c(\alpha \wedge \beta) &= (c\alpha) \wedge \beta = \alpha \wedge (c\beta).\end{aligned}$$

by the third statement in Proposition 5.2.1. It is also true that  $(\alpha \wedge \beta) \wedge \gamma = \alpha \wedge (\beta \wedge \gamma)$ , but a few more properties are needed in order to prove this:

**Lemma 5.2.1.** For  $\alpha^1, \dots, \alpha^k \in \Lambda_1(V) = \mathcal{T}_1(V)$ ,  $\sigma \in S_k$ ,

$$(\alpha^1 \otimes \dots \otimes \alpha^k)_a = \frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) \alpha^{\sigma(1)} \otimes \dots \otimes \alpha^{\sigma(k)},$$

$$(\alpha^{\sigma(1)} \otimes \dots \otimes \alpha^{\sigma(k)})_a = \varepsilon(\sigma) (\alpha^1 \otimes \dots \otimes \alpha^k)_a,$$

$$\begin{aligned}(\alpha^1 \otimes \dots \otimes \alpha^k)_a &= ((\alpha^1 \otimes \dots \otimes \alpha^{k-l})_a \otimes \alpha^{k-l+1} \otimes \dots \otimes \alpha^k)_a \\ &= (\alpha^1 \otimes \dots \otimes \alpha^{k-l} \otimes (\alpha^{k-l+1} \otimes \dots \otimes \alpha^k))_a.\end{aligned}$$

In particular,  $(S_a \otimes T)_a = (S \otimes T_a)_a = (S \otimes T)_a$  for all  $S \in \mathcal{T}_k(V)$ ,  $T \in \mathcal{T}_l(V)$ .

*Proof.* For the first identity, observe that

$$\begin{aligned}(\alpha^1 \otimes \dots \otimes \alpha^k)_a(\mathbf{v}_1, \dots, \mathbf{v}_k) &= \frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) \alpha^1(\mathbf{v}_{\sigma(1)}) \dots \alpha^k(\mathbf{v}_{\sigma(k)}) \\ &= \frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) (\alpha^{\sigma^{-1}(1)} \otimes \dots \otimes \alpha^{\sigma^{-1}(k)})(\mathbf{v}_1, \dots, \mathbf{v}_k).\end{aligned}$$

Furthermore, a permutation has the same sign as its inverse, so that

$$\begin{aligned} (\alpha^1 \otimes \cdots \otimes \alpha^k)_a &= \frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) \alpha^{\sigma^{-1}(1)} \otimes \cdots \otimes \alpha^{\sigma^{-1}(k)} \\ &= \frac{1}{k!} \sum_{\sigma^{-1} \in S_k} \varepsilon(\sigma^{-1}) \alpha^{\sigma^{-1}(1)} \otimes \cdots \otimes \alpha^{\sigma^{-1}(k)} \\ &= \frac{1}{k!} \sum_{\sigma \in S_k} \varepsilon(\sigma) \alpha^{\sigma(1)} \otimes \cdots \otimes \alpha^{\sigma(k)}. \end{aligned}$$

For the second identity, we have by the one just proved,

$$\begin{aligned} (\alpha^{\sigma(1)} \otimes \cdots \otimes \alpha^{\sigma(k)})_a &= \frac{1}{k!} \sum_{\tau \in S_k} \varepsilon(\tau) \alpha^{(\tau \circ \sigma)(1)} \otimes \cdots \otimes \alpha^{(\tau \circ \sigma)(k)} \\ &= \varepsilon(\sigma) \frac{1}{k!} \sum_{\tau \in S_k} \varepsilon(\tau \circ \sigma) \alpha^{(\tau \circ \sigma)(1)} \otimes \cdots \otimes \alpha^{(\tau \circ \sigma)(k)} \\ &= \varepsilon(\sigma) \frac{1}{k!} \sum_{\tau \circ \sigma \in S_k} \varepsilon(\tau \circ \sigma) \alpha^{(\tau \circ \sigma)(1)} \otimes \cdots \otimes \alpha^{(\tau \circ \sigma)(k)} \\ &= \varepsilon(\sigma) (\alpha^1 \otimes \cdots \otimes \alpha^k)_a. \end{aligned}$$

The two equalities in the third identity are proved in the same way, so we only argue the first one:

$$\begin{aligned} &[(\alpha^1 \otimes \cdots \otimes \alpha^{k-l})_a \otimes \alpha^{k-l+1} \otimes \cdots \otimes \alpha^l]_a \\ &= \frac{1}{(k-l)!} \sum_{\sigma \in S_{k-l}} \varepsilon(\sigma) (\alpha^{\sigma(1)} \otimes \cdots \otimes \alpha^{\sigma(k-l)} \otimes \alpha^{k-l+1} \otimes \cdots \otimes \alpha^l)_a \\ &= \frac{1}{(k-l)!} \sum_{\sigma \in S_{k-l}} \varepsilon^2(\sigma) (\alpha^1 \otimes \cdots \otimes \alpha^{k+l})_a \\ &= (\alpha^1 \otimes \cdots \otimes \alpha^{k+l})_a, \end{aligned}$$

using the second identity on the third line. Finally, the last statement follows from the third identity, since the map  $T \mapsto T_a$  is linear and each  $k$ -tensor  $T$  is a linear combination of basis elements of the form  $\alpha^{i_1} \otimes \cdots \otimes \alpha^{i_k}$  by Theorem 5.1.1.  $\square$

The final statement in the above Lemma also implies that the wedge product is associative, meaning:

**Theorem 5.2.1.** *If  $\alpha \in \Lambda_k(V)$ ,  $\beta \in \Lambda_l(V)$ , and  $\gamma \in \Lambda_m(V)$ , then*

$$(\alpha \wedge \beta) \wedge \gamma = \alpha \wedge (\beta \wedge \gamma) = \frac{(k+l+m)!}{k! l! m!} (\alpha \otimes \beta \otimes \gamma)_a.$$

*Proof.*

$$\begin{aligned}
 (\alpha \wedge \beta) \wedge \gamma &= \frac{(k+l+m)!}{(k+l)!m!} ((\alpha \wedge \beta) \otimes \gamma)_a \\
 &= \frac{(k+l+m)!}{(k+l)!m!} \frac{(k+l)!}{k!l!} ((\alpha \otimes \beta)_a \otimes \gamma)_a \\
 &= \frac{(k+l+m)!}{k!l!m!} (\alpha \otimes \beta \otimes \gamma)_a.
 \end{aligned}$$

A similar argument shows that  $\alpha \wedge (\beta \wedge \gamma)$  also equals the last term in the above identity.  $\square$

In view of the above theorem, we write  $\alpha \wedge \beta \wedge \gamma$  for  $(\alpha \wedge \beta) \wedge \gamma$  or  $\alpha \wedge (\beta \wedge \gamma)$ , and similarly for products of higher order. Notice that the first two identities from Lemma 5.2.1, together with associativity of the wedge product, immediately imply

$$\begin{aligned}
 \alpha^1 \wedge \cdots \wedge \alpha^k &= \sum_{\sigma \in S_k} \varepsilon(\sigma) \alpha^{\sigma(1)} \otimes \cdots \otimes \alpha^{\sigma(k)} \\
 \alpha^{\sigma(1)} \wedge \cdots \wedge \alpha^{\sigma(k)} &= \varepsilon(\sigma) \alpha^1 \wedge \cdots \wedge \alpha^k
 \end{aligned} \tag{5.2.1}$$

for  $\alpha^1, \dots, \alpha^k \in \Lambda_1(V)$ ,  $\sigma \in S_k$ . The second identity, in particular, lets us identify a basis of  $\Lambda_k(V)$ :

**Theorem 5.2.2.** *Let  $\mathbf{v}_1, \dots, \mathbf{v}_n$  denote a basis of  $V$ , and  $\alpha^1, \dots, \alpha^n \in \Lambda_1(V)$  the dual basis. Then the set*

$$\{\alpha^{i_1} \wedge \cdots \wedge \alpha^{i_k} \mid 1 \leq i_1 < i_2 < \cdots < i_k \leq n\}$$

*is a basis of  $\Lambda_k(V)$ , which therefore has dimension  $\binom{n}{k}$ , where*

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} = \frac{n(n-1)\cdots(n-k+1)}{k!}.$$

*Proof.* If  $\alpha \in \Lambda_k(V)$ , then  $\alpha$  is a  $k$ -tensor, and by Theorem 5.1.1,

$$\alpha = \sum_{i_1, \dots, i_k=1}^n \alpha_{i_1 \dots i_k} \alpha^{i_1} \otimes \cdots \otimes \alpha^{i_k}.$$

Thus,

$$\begin{aligned}
 \alpha &= \alpha_a = \sum_{j_1, \dots, j_k=1}^n \alpha_{j_1 \dots j_k} (\alpha^{j_1} \otimes \cdots \otimes \alpha^{j_k})_a \\
 &= \frac{1}{k!} \sum_{j_1, \dots, j_k=1}^n \alpha_{j_1 \dots j_k} \alpha^{j_1} \wedge \cdots \wedge \alpha^{j_k}.
 \end{aligned}$$

Let  $i_1 < i_2 < \cdots < i_k$  denote  $j_1, \dots, j_k$  written in increasing order. Then  $\alpha^{i_1} \wedge \cdots \wedge \alpha^{i_k} = \pm \alpha^{j_1} \wedge \cdots \wedge \alpha^{j_k}$  by (5.2.1), and the collection in the statement of the theorem spans  $\Lambda_k(V)$ .

To see that it is linearly independent, suppose that

$$\sum_{1 \leq j_1 < \cdots < j_k \leq n} \alpha_{j_1 \dots j_k} \alpha^{j_1} \wedge \cdots \wedge \alpha^{j_k} = 0. \tag{5.2.2}$$

Observe that if  $i_1 < \dots < i_k$ , then

$$\alpha^{j_1} \wedge \dots \wedge \alpha^{j_k}(\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k}) = \sum_{\sigma \in S_k} \varepsilon(\sigma) \alpha^{j_1}(\mathbf{v}_{\sigma(i_1)}) \dots \alpha^{j_k}(\mathbf{v}_{\sigma(i_k)}), \quad (5.2.3)$$

where  $S_k$  denotes all the permutations of  $\{i_1, \dots, i_k\}$ . Since  $\alpha^j(\mathbf{v}_i) = \delta_{ij}$ , the only nonzero terms in this sum are those for which  $\sigma(i_l) = j_l$ ,  $l = 1, \dots, k$ , in which case they equal  $\varepsilon(\sigma)$ . In particular, as sets,  $\{i_1, \dots, i_k\} = \{\sigma(i_1), \dots, \sigma(i_k)\} = \{j_1, \dots, j_k\}$ . But  $i_1 < \dots < i_k$  and  $j_1 < \dots < j_k$ , so all terms vanish except when  $i_1 = j_1, \dots, i_k = j_k$ , and the sum then equals 1. Thus, for any  $i_1 < \dots < i_k$ , applying both sides of (5.2.2) to  $\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k}$  yields

$$0 = \sum_{1 \leq j_1 < \dots < j_k \leq n} \alpha_{j_1 \dots j_k} \alpha^{j_1} \wedge \dots \wedge \alpha^{j_k}(\mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_k}) = \alpha_{i_1 \dots i_k},$$

which establishes linear independence.

To conclude the proof, it remains to show that the dimension is  $\binom{n}{k}$  as claimed. The size of the basis is the number of distinct  $k$  elements chosen from a set of  $n$  elements written in increasing order. Now the number of distinct ordered sets of  $k$  elements is  $n(n-1)\dots(n-k+1)$ , since there are  $n$  choices for the first one,  $n-1$  for the second, and so on. Since there are  $k!$  ways of ordering a given set of  $k$  elements, the claim follows.  $\square$

The theorem says in particular that  $\Lambda_n(V)$ , where  $n$  is the dimension of  $V$ , is one-dimensional. This we already knew of course from Theorem 1.3.1: on  $\mathbb{R}^n$ , any  $n$ -form is a multiple of the determinant, so  $\Lambda_n(\mathbb{R}^n)$  is one-dimensional. A choice of basis for  $V$  induces an isomorphism of  $V$  with  $\mathbb{R}^n$ , which in turn induces one of  $\Lambda_n(V)$  with  $\Lambda_n(\mathbb{R}^n)$ .

Observe also that if  $\{\varepsilon^i\}$  is the basis of  $\mathbb{R}^{n*}$  dual to  $\{\mathbf{e}_i\}$ , then

$$\varepsilon^1 \wedge \dots \wedge \varepsilon^n = \det. \quad (5.2.4)$$

Indeed, the form on the left equals  $c$  times  $\det$  for some  $c \in \mathbb{R}$ . Thus,

$$\varepsilon^1 \wedge \dots \wedge \varepsilon^n(\mathbf{e}_1, \dots, \mathbf{e}_n) = c \cdot \det(\mathbf{e}_1, \dots, \mathbf{e}_n) = c.$$

But by (5.2.3), the term on the left equals 1. The factorials in the definition of the wedge product were chosen in order for (5.2.4) to hold.

**Definition 5.2.2.** An *orientation* on an  $n$ -dimensional vector space  $V$  is a choice of a nonzero  $\omega \in \Lambda_n(V)$ . Given any two nonzero  $\omega_1, \omega_2 \in \Lambda_n(V)$ ,  $\omega_1 = c \cdot \omega_2$  for some  $c \neq 0$ . If  $c > 0$ , we say  $\omega_1$  and  $\omega_2$  determine the *same orientation*. Otherwise, they determine opposite orientations.

Clearly, there are only two possible orientations on a given vector space. The *standard orientation* on  $\mathbb{R}^n$  is the one induced by the determinant. In an oriented vector space  $V$ , if  $\omega \in \Lambda_n(V)$  induces the given orientation, an ordered basis  $(\mathbf{v}_1, \dots, \mathbf{v}_n)$  of  $V$  (i.e., an  $n$ -tuple consisting of basis elements) is said to be *positively oriented* if  $\omega(\mathbf{v}_1, \dots, \mathbf{v}_n) > 0$ .

The following result deals with *decomposable* elements of  $\Lambda_k(V)$ ; i.e., elements of the form  $\alpha^1 \wedge \cdots \wedge \alpha^k$ , where  $\alpha^i \in V^*$  (as opposed to a sum of such terms):

**Corollary 5.2.1.** *Let  $V$  be an  $n$ -dimensional vector space,  $\mathbf{v}_i \in V$ ,  $\alpha^i \in V^*$ ,  $1 \leq i, j \leq k \leq n$ . Then*

- (1)  $(\alpha^1 \wedge \cdots \wedge \alpha^k)(\mathbf{v}_1, \dots, \mathbf{v}_k) = \det(\alpha_i(\mathbf{v}_j))$ ;
- (2)  $\alpha^1 \wedge \cdots \wedge \alpha^k \neq 0$  if and only if  $\alpha^1, \dots, \alpha^k$  are linearly independent;
- (3) *There is a bijective correspondence between one-dimensional decomposable subspaces of  $\Lambda_k(V)$  and  $k$ -dimensional subspaces of  $V$ .*

*Proof.* For (1), (5.2.1) implies that

$$(\alpha^1 \wedge \cdots \wedge \alpha^k)(\mathbf{v}_1, \dots, \mathbf{v}_k) = \sum_{\sigma \in \mathcal{S}_k} \varepsilon(\sigma) \alpha^{\sigma(1)}(\mathbf{v}_1) \cdots \alpha^{\sigma(k)}(\mathbf{v}_k) = \det(\alpha^i(\mathbf{v}_j))$$

by definition of the determinant.

For (2), if  $\alpha^1, \dots, \alpha^k$  are linearly dependent, then one of them, say  $\alpha^1$ , is a linear combination  $\sum c_j \alpha^j$  of the others. But then

$$\alpha^1 \wedge \cdots \wedge \alpha^k = \sum_{j=2}^k c_j \alpha^j \wedge \alpha_2 \wedge \cdots \wedge \alpha^k = 0$$

since  $\alpha^j$  appears twice in each term inside the summation. On the other hand, if they are linearly independent, then they can be extended to a basis of  $V^*$ , and by Theorem 5.2.2,  $\alpha^1 \wedge \cdots \wedge \alpha^k$  is one of the corresponding basis elements of  $\Lambda_k(V)$ . Being a basis element, it cannot be zero.

For the third statement, the correspondence

$$\text{span}\{\alpha^1 \wedge \cdots \wedge \alpha^k\} \longleftrightarrow \text{span}\{\alpha^1, \dots, \alpha^k\}$$

where  $\alpha^1, \dots, \alpha^k$  are linearly independent in  $V^*$  is a bijective one between one-dimensional decomposable subspaces of  $\Lambda_k(V)$  and  $k$ -dimensional subspaces of  $V^*$  by (2). The claim follows in view of the isomorphism between a space and its dual.  $\square$

Let us return for a moment to the topic of orientation. Suppose that in addition to an orientation,  $V$  is endowed with an inner product. If  $(\mathbf{v}_1, \dots, \mathbf{v}_n)$  is a positively oriented orthonormal basis of  $V$ , then there exists a unique  $n$ -form  $\alpha$  on  $V$  such that  $\alpha(\mathbf{v}_1, \dots, \mathbf{v}_n) = 1$ ; in fact,  $\alpha = \alpha^1 \wedge \cdots \wedge \alpha^n$ , where  $\alpha^1, \dots, \alpha^n$  is the basis of  $V^*$  dual to  $\mathbf{v}_1, \dots, \mathbf{v}_n$ .  $\alpha$  is called the *volume form* of  $V$ . The reason we used a definite rather than indefinite article for volume form is because  $\alpha(\mathbf{w}_1, \dots, \mathbf{w}_n) = 1$  for *any* positively oriented orthonormal basis  $(\mathbf{w}_1, \dots, \mathbf{w}_n)$  of  $V$ , which means that such an  $\alpha$  is unique:

**Theorem 5.2.3.** *Let  $V$  be an  $n$ -dimensional vector space.*

- (1) *If  $\mathbf{v}_1, \dots, \mathbf{v}_n$  is a basis of  $V$ , and  $\alpha \in \Lambda_n(V)$ , then*

$$\alpha(\mathbf{w}_1, \dots, \mathbf{w}_n) = \det(a_{ij}) \alpha(\mathbf{v}_1, \dots, \mathbf{v}_n),$$

for any  $\mathbf{w}_i = \sum_j a_{ij} \mathbf{v}_j \in V$ .

- (2) Suppose that in addition,  $V$  is oriented and endowed with an inner product. If  $\omega$  is the corresponding volume form, then for any positively oriented orthonormal basis  $\mathbf{w}_1, \dots, \mathbf{w}_n$  of  $V$ ,  $\omega(\mathbf{w}_1, \dots, \mathbf{w}_n) = 1$ .

*Proof.* Define an  $n$ -form  $\beta$  on  $\mathbb{R}^n$  by

$$\beta \left( \begin{bmatrix} b_{11} \\ \vdots \\ b_{n1} \end{bmatrix}, \dots, \begin{bmatrix} b_{1n} \\ \vdots \\ b_{nn} \end{bmatrix} \right) = \alpha \left( \sum_i b_{i1} \mathbf{v}_i, \dots, \sum_i b_{in} \mathbf{v}_i \right),$$

for  $b_{ij} \in \mathbb{R}$ ,  $1 \leq i, j \leq n$ . By Theorem 1.3.1,  $\beta = c \cdot \det$  for some  $c \in \mathbb{R}$ , and

$$c = c \cdot \det(\mathbf{e}_1, \dots, \mathbf{e}_n) = \beta(\mathbf{e}_1, \dots, \mathbf{e}_n) = \alpha(\mathbf{v}_1, \dots, \mathbf{v}_n).$$

Thus,

$$\alpha(\mathbf{w}_1, \dots, \mathbf{w}_n) = \beta \left( \begin{bmatrix} a_{11} \\ \vdots \\ a_{n1} \end{bmatrix}, \dots, \begin{bmatrix} a_{1n} \\ \vdots \\ a_{nn} \end{bmatrix} \right) = \det(a_{ij}) \alpha(\mathbf{v}_1, \dots, \mathbf{v}_n),$$

which establishes the first claim. For the second one, let  $\omega = \omega^1 \wedge \dots \wedge \omega^n$ , and  $\mathcal{B} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  the positively oriented orthonormal basis of  $V$  dual to  $(\omega^1, \dots, \omega^n)$ . If  $\mathcal{C}$  is the basis  $(\mathbf{w}_1, \dots, \mathbf{w}_n)$ , then by what was just proved,

$$\omega(\mathbf{w}_1, \dots, \mathbf{w}_n) = \det[1_V]_{\mathcal{C}, \mathcal{B}} \cdot \omega(\mathbf{v}_1, \dots, \mathbf{v}_n) = \det[1_V]_{\mathcal{C}, \mathcal{B}},$$

and it remains to show this determinant equals 1. Now,  $[1_V]_{\mathcal{C}, \mathcal{B}}$  is also the matrix with respect to  $\mathcal{B}$  of the operator  $L$  which maps  $\mathbf{v}_i$  to  $\mathbf{w}_i$ ,  $i = 1, \dots, n$ , because

$$[\mathbf{w}_i]_{\mathcal{B}} = [1_V]_{\mathcal{C}, \mathcal{B}} [\mathbf{w}_i]_{\mathcal{C}} = [1_V]_{\mathcal{C}, \mathcal{B}} \mathbf{e}_i, \quad i = 1, \dots, n,$$

on the one hand, and  $[\mathbf{w}_i]_{\mathcal{B}} = [L]_{\mathcal{B}} \mathbf{e}_i$  by definition of  $L$  on the other. But  $L$  maps an orthonormal basis to another one, so it is a linear isometry. In particular,  $[L]_{\mathcal{B}} [L]_{\mathcal{B}}^T = I_n$ , which implies that  $(\det L)^2 = 1$ , so that  $\det L = \pm 1$ . Summarizing,  $\omega(\mathbf{w}_1, \dots, \mathbf{w}_n) = \pm 1$ . Finally,  $\omega(\mathbf{w}_1, \dots, \mathbf{w}_n) > 0$  since  $\mathcal{C}$  is positively oriented, so it must equal one. This completes the proof.  $\square$

2-forms on  $V$  possess an additional property, which in turn has several consequences explored in the exercises:

**Proposition 5.2.2.** *Given any nonzero form  $\alpha \in \Lambda_2(V)$ , there exists a basis  $\alpha^1, \dots, \alpha^n$  of  $\Lambda_1(V)$  such that*

$$\alpha = \alpha^1 \wedge \alpha^2 + \alpha^3 \wedge \alpha^4 + \dots + \alpha^{2k-1} \wedge \alpha^{2k}.$$

*Proof.* It suffices to construct a basis  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of  $V$  such that

$$\alpha(\mathbf{v}_{2i-1}, \mathbf{v}_{2i}) = 1 \text{ if } i \leq k, \quad \alpha(\mathbf{v}_i, \mathbf{v}_j) = 0 \text{ if } i \text{ or } j > k, \text{ or if } |j - i| > 1,$$

since the dual basis will then satisfy the claim. We argue by induction on the dimension of  $V$ . If the dimension is 1, there is nothing to prove, so assume the claim is true in all dimensions less than  $n$ . Since  $\alpha \neq 0$ , there exist  $\mathbf{v}_1, \mathbf{w}_2 \in V$  such that  $\alpha(\mathbf{v}_1, \mathbf{w}_2) \neq 0$ . If  $\mathbf{v}_2 = \mathbf{w}_2/\alpha(\mathbf{v}_1, \mathbf{w}_2)$ , then  $\alpha(\mathbf{v}_1, \mathbf{v}_2) = 1$ . Let  $W$  be the subspace of  $V$  spanned by  $\mathbf{v}_1$  and  $\mathbf{v}_2$ , and define

$$Z = \{\mathbf{v} \in V \mid \alpha(\mathbf{v}, \mathbf{v}_1) = \alpha(\mathbf{v}, \mathbf{v}_2) = 0\}.$$

By construction,  $Z$  is a subspace of  $V$  that intersects  $W$  in the zero vector only. Furthermore,  $Z$  has dimension at least  $n - 2$ : indeed,  $Z = \ker \beta^1 \cap \ker \beta^2$ , where  $\beta^i \in V^*$  is given by  $\beta^i(\mathbf{v}) = \alpha(\mathbf{v}, \mathbf{v}_i)$ ,  $i = 1, 2$ ,  $\mathbf{v} \in V$ , and each kernel has dimension  $n - 1$ , see also Exercise 1.21. By that same exercise,

$$\begin{aligned} \dim Z &= \dim \ker \beta^1 + \dim \ker \beta^2 - \dim(\ker \beta^1 + \ker \beta^2) \\ &\geq \dim \ker \beta^1 + \dim \ker \beta^2 - \dim V \\ &= n - 2. \end{aligned}$$

Since  $W$  is two-dimensional and  $Z \cap W = \{\mathbf{0}\}$ ,  $V = W \oplus Z$ . The result now follows by the induction hypothesis applied to the restriction of  $\alpha$  to  $Z$ .  $\square$

Even though the tensor product is not commutative, i.e.,  $S \otimes T$  need not equal  $T \otimes S$ , the wedge product is “almost” commutative:

$$\alpha \wedge \beta = (-1)^{kl} \beta \wedge \alpha, \quad \alpha \in \Lambda_k(V), \quad \beta \in \Lambda_l(V). \quad (5.2.5)$$

This is easily seen by writing both forms in terms of a basis and using (5.2.1), see Exercise 5.5.

**Examples 5.2.1.** (i) Let  $\alpha \in \Lambda_1(V)$ ,  $\beta \in \Lambda_2(V)$ . Then  $\alpha \wedge \beta = \beta \wedge \alpha$ , and for any  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$ ,

$$\begin{aligned} (\alpha \wedge \beta)(\mathbf{x}, \mathbf{y}, \mathbf{z}) &= \frac{1}{2} [\alpha(\mathbf{x})\beta(\mathbf{y}, \mathbf{z}) - \alpha(\mathbf{x})\beta(\mathbf{z}, \mathbf{y}) - \alpha(\mathbf{y})\beta(\mathbf{x}, \mathbf{z}) + \alpha(\mathbf{y})\beta(\mathbf{z}, \mathbf{x}) \\ &\quad + \alpha(\mathbf{z})\beta(\mathbf{x}, \mathbf{y}) - \alpha(\mathbf{z})\beta(\mathbf{y}, \mathbf{x})] \\ &= \alpha(\mathbf{x})\beta(\mathbf{y}, \mathbf{z}) + \alpha(\mathbf{y})\beta(\mathbf{z}, \mathbf{x}) + \alpha(\mathbf{z})\beta(\mathbf{x}, \mathbf{y}) \\ &= \cup (\alpha \otimes \beta)(\mathbf{x}, \mathbf{y}, \mathbf{z}), \end{aligned}$$

where  $\cup$  denotes cyclic summation: given a 3-tensor  $T$ ,

$$\cup T(\mathbf{x}, \mathbf{y}, \mathbf{z}) = T(\mathbf{x}, \mathbf{y}, \mathbf{z}) + T(\mathbf{y}, \mathbf{z}, \mathbf{x}) + T(\mathbf{z}, \mathbf{x}, \mathbf{y}).$$

(ii) Let  $\alpha, \beta \in \Lambda^1(V) = V^*$ . Recall from Corollary 5.2.1(2) that if  $\alpha \neq 0$  and  $\alpha \wedge \beta = 0$ , then  $\beta \in \text{span}\{\alpha\}$ . The following generalization is known as Cartan’s lemma: Suppose  $\alpha^1, \dots, \alpha^k \in V^*$  are linearly independent. If  $\beta^1, \dots, \beta^k \in V^*$  satisfy

$$\sum_{i=1}^k \alpha^i \wedge \beta^i = 0,$$

then each  $\beta^i \in \text{span}\{\alpha^1, \dots, \alpha^k\}$ . Furthermore, if  $\beta^i = \sum_j a_{ij} \alpha^j$ , then  $a_{ij} = a_{ji}$ .

To see this, extend the  $\alpha^i$ 's to a basis  $\alpha^1, \dots, \alpha^n$  of  $V^*$ , and write

$$\beta^i = \sum_{j=1}^k a_{ij} \alpha^j + \sum_{j=k+1}^n b_{ij} \alpha^j.$$

Then

$$0 = \sum_{i=1}^k \alpha^i \wedge \beta^i = \sum_{1 \leq i < j \leq k} (a_{ij} - a_{ji}) \alpha^i \wedge \alpha^j + \sum_{i \leq k < j} b_{ij} \alpha^i \wedge \alpha^j,$$

and the claim follows from linear independence of the  $\alpha^i \wedge \alpha^j, i < j$ .

### 5.3 Differential forms

The same procedure used in going from tensors to tensor fields (or from vectors to vector fields) can be applied to forms. The result is not, however, traditionally called a “form field”:

**Definition 5.3.1.** A *differential  $k$ -form*, or simply a differential form, on a manifold  $M$  is a map  $\omega$  that assigns to each  $\mathbf{p} \in M$  an element  $\omega(\mathbf{p})$  of  $\Lambda_k(M_{\mathbf{p}})$ .  $\omega$  is assumed to be smooth in the sense that for any vector fields  $\mathbf{X}_1, \dots, \mathbf{X}_k$  on  $M$ , the function

$$\begin{aligned} \omega(\mathbf{X}_1, \dots, \mathbf{X}_k) : M &\longrightarrow \mathbb{R}, \\ \mathbf{p} &\longmapsto \omega(\mathbf{p})(\mathbf{X}_1(\mathbf{p}), \dots, \mathbf{X}_k(\mathbf{p})) \end{aligned}$$

is differentiable.

Since a differential form is a tensor field, the results established for tensor fields hold in this context as well. Thus, for example, smoothness may be rephrased by requiring that for any  $\mathbf{p} \in M$ , there be a chart  $(U, \mathbf{x})$  around  $\mathbf{p}$  such that the functions  $\omega(\partial/\partial x^{i_1}, \dots, \partial/\partial x^{i_k})$  are differentiable on  $U$ . Furthermore, if we denote these functions by  $\omega_{i_1 \dots i_k}$ , then the restriction of  $\omega$  to  $U$  is given by

$$\omega|_U = \sum_{1 \leq i_1 < \dots < i_k \leq n} \omega_{i_1 \dots i_k} dx^{i_1} \wedge \dots \wedge dx^{i_k}.$$

Similarly, the collection of all differential  $k$ -forms on  $M$  is a vector space, denoted  $\Lambda_k(M)$ . Finally, Proposition 5.1.1 implies that given  $\mathbf{f} : M \rightarrow N$ ,  $\alpha, \beta$   $k$ -forms on  $N$ ,  $\gamma$  an  $l$ -form on  $N$ ,  $h : \mathbb{N} \rightarrow \mathbb{R}$ , and real numbers  $a, b$ ,

$$\begin{aligned} \mathbf{f}^*(a\alpha + b\beta) &= a\mathbf{f}^*\alpha + b\mathbf{f}^*\beta, \\ \mathbf{f}^*(\alpha \wedge \gamma) &= \mathbf{f}^*\alpha \wedge \mathbf{f}^*\gamma, \\ \mathbf{f}^*(h\alpha) &= (h \circ \mathbf{f})\mathbf{f}^*\alpha. \end{aligned} \tag{5.3.1}$$

There is one construction that is specific to forms as opposed to general tensor fields: let us agree to call a function on  $M$  a differential 0-form. Taking differentials of func-



tions is then a linear map  $d : \Lambda_0(M) \rightarrow \Lambda_1(M)$ . We now extend it to an operator  $d : \Lambda_k(M) \rightarrow \Lambda_{k+1}(M)$  for any nonnegative integer  $k$  as follows: let  $\alpha$  be a differential  $k$ -form,  $\mathbf{p} \in M$ . In order to define  $d\alpha(\mathbf{p})$ , consider a chart  $(U, \mathbf{x})$  of  $M$  around  $\mathbf{p}$ . Then as noted earlier, the restriction of  $\alpha$  to  $U$  may be written as  $\alpha|_U = \sum_{1 \leq i_1 < \dots < i_k \leq n} \alpha_{i_1 \dots i_k} dx^{i_1} \wedge \dots \wedge dx^{i_k}$ . Define

$$(d\alpha)|_U := \sum_{1 \leq i_1 < \dots < i_k \leq n} d\alpha_{i_1 \dots i_k} \wedge dx^{i_1} \wedge \dots \wedge dx^{i_k}. \quad (5.3.2)$$

It must be checked that this formula is independent of the chosen chart. To see this, we first claim that  $d$  satisfies the following properties at  $\mathbf{p}$ :

- (1)  $d\alpha(\mathbf{p}) \in \Lambda_{k+1}(M_{\mathbf{p}})$ ;
- (2) if  $\alpha = \beta$  on some neighborhood of  $\mathbf{p}$ , then  $d\alpha(\mathbf{p}) = d\beta(\mathbf{p})$ ;
- (3)  $d(a\alpha + b\beta)(\mathbf{p}) = a d\alpha(\mathbf{p}) + b d\beta(\mathbf{p})$ ,  $a, b \in \mathbb{R}$ ,  $\beta \in \Lambda_k(M)$ ;
- (4)  $d(\alpha \wedge \beta)(\mathbf{p}) = d\alpha(\mathbf{p}) \wedge \beta(\mathbf{p}) + (-1)^k \alpha(\mathbf{p}) \wedge d\beta(\mathbf{p})$ ,  $\beta \in \Lambda_l(M)$ ;
- (5)  $d(df)(\mathbf{p}) = 0$  for  $f \in \Lambda_0(M)$ .

The first three items follow directly from (5.3.2). To prove the fourth one, it suffices to check it in the case that  $\alpha = f dx^{i_1} \wedge \dots \wedge dx^{i_k}$  and  $\beta = g dx^{j_1} \wedge \dots \wedge dx^{j_l}$  by the linearity of  $d$  established in (3). Furthermore, the identity certainly holds if  $\alpha$  and/or  $\beta$  are functions (i.e., zero forms). Define  $\gamma = dx^{i_1} \wedge \dots \wedge dx^{i_k}$ , so that  $\beta = g\gamma$ . Then

$$\begin{aligned} d(\alpha \wedge \beta)(\mathbf{p}) &= d(fg dx^{i_1} \wedge \dots \wedge dx^{i_k} \wedge \gamma)(\mathbf{p}) \\ &= (df(\mathbf{p})g(\mathbf{p}) + f(\mathbf{p})dg(\mathbf{p})) \wedge (dx^{i_1} \wedge \dots \wedge dx^{i_k} \wedge \gamma)(\mathbf{p}) \\ &= (df \wedge dx^{i_1} \wedge \dots \wedge dx^{i_k})(\mathbf{p}) \wedge (g\gamma)(\mathbf{p}) \\ &\quad + (-1)^k (f dx^{i_1} \wedge \dots \wedge dx^{i_k})(\mathbf{p}) \wedge (dg \wedge \gamma)(\mathbf{p}) \\ &= d\alpha(\mathbf{p}) \wedge \beta(\mathbf{p}) + (-1)^k \alpha(\mathbf{p}) \wedge d\beta(\mathbf{p}). \end{aligned}$$

The last identity can be seen by writing  $(df)|_U = \sum_i \partial/\partial x^i(f) dx^i$ , so that

$$\begin{aligned} d(df)(\mathbf{p}) &= \left( \sum_i d\left(\frac{\partial}{\partial x^i} f\right) \wedge dx^i \right)(\mathbf{p}) = \sum_{ij} \left( \left( \frac{\partial}{\partial x^j} \frac{\partial}{\partial x^i} f \right) dx^j \wedge dx^i \right)(\mathbf{p}) \\ &= \sum_{j < i} \left( \left( \frac{\partial}{\partial x^i} \frac{\partial}{\partial x^j} - \frac{\partial}{\partial x^j} \frac{\partial}{\partial x^i} \right) f \right) dx^j \wedge dx^i(\mathbf{p}) \\ &= 0. \end{aligned}$$

We can now justify the claim made earlier that the definition of  $d$  does not depend on the chosen chart: indeed, suppose  $\tilde{d}$  is the operator obtained by using a different chart. It suffices to show that

$$\tilde{d}(f dx^{i_1} \wedge \dots \wedge dx^{i_k}) = df \wedge dx^{i_1} \wedge \dots \wedge dx^{i_k}$$

for any function  $f$  by linearity of  $\tilde{d}$ . But  $\tilde{d}f = df$  and  $\tilde{d}$  satisfies properties (1) through (5), so that

$$\tilde{d}(f dx^{i_1} \wedge \dots \wedge dx^{i_k}) = df \wedge dx^{i_1} \wedge \dots \wedge dx^{i_k} + f \tilde{d}(dx^{i_1} \wedge \dots \wedge dx^{i_k})$$

by (4). It remains to show that the last term in the above identity is zero. But again by (4), this term is a sum over  $j$  of expressions of the form

$$\pm dx^{i_1} \wedge \cdots \wedge \tilde{d}(dx^{j_i}) \wedge \cdots \wedge dx^{i_k},$$

and  $\tilde{d}(dx^{j_i}) = 0$  by (5). This shows that  $d$  is well defined.

**Theorem 5.3.1.** *There exists a unique linear map  $d : \Lambda_k(M) \rightarrow \Lambda_{k+1}(M)$ ,  $k = 0, 1, 2, \dots$ , called the exterior derivative operator, that satisfies*

- (1)  $d(\alpha \wedge \beta) = (d\alpha) \wedge \beta + (-1)^k \alpha \wedge (d\beta)$ ,  $\alpha \in \Lambda_k(M)$ ,
- (2)  $d \circ d = 0$ , and
- (3) for  $f \in \Lambda_0(M)$ ,  $df$  is the differential of  $f$ .

*Proof.* Existence has already been established. For uniqueness, it suffices to show that a linear map as in the statement necessarily satisfies the five properties listed earlier in the definition of  $d$ . All but the second property are clear. So let  $\tilde{d}$  be a map satisfying the conditions of the statement, and suppose that  $\alpha = \beta$  on a neighborhood  $U$  of  $\mathbf{p}$ . Set  $\gamma = \alpha - \beta$ . Then  $\gamma|_U \equiv 0$ , and we must show that  $\tilde{d}\gamma(\mathbf{p}) = 0$ . To see this is so, consider a neighborhood  $V$  of  $\mathbf{p}$  whose closure lies in  $U$ , and a function  $\varphi$  with values between zero and one, which equals 1 on  $V$  and has support in  $U$ . Then  $\gamma = (1 - \varphi)\gamma$ , so that

$$\tilde{d}\gamma(\mathbf{p}) = d(1 - \varphi)(\mathbf{p}) \wedge \gamma(\mathbf{p}) + (1 - \varphi)(\mathbf{p}) \tilde{d}\gamma(\mathbf{p}) = 0. \quad \square$$

In the case of  $\mathbb{R}^n$ , the definition we gave of exterior derivative is intrinsic, but on a generic manifold it is somewhat unsatisfactory to rely on a definition in terms of charts. The following property, which could have been used as a definition, seeks to remedy this:

**Theorem 5.3.2.** *For  $\alpha \in \Lambda_k(M)$ ,  $\mathbf{X}_i \in \mathfrak{X}M$ ,  $i = 0, \dots, k$ ,*

$$d\alpha(\mathbf{X}_0, \dots, \mathbf{X}_k) = \sum_{i=0}^k (-1)^i \mathbf{X}_i(\alpha(\mathbf{X}_0, \dots, \hat{\mathbf{X}}_i, \dots, \mathbf{X}_k)) \\ + \sum_{i < j} (-1)^{i+j} \alpha([\mathbf{X}_i, \mathbf{X}_j], \mathbf{X}_0, \dots, \hat{\mathbf{X}}_i, \dots, \hat{\mathbf{X}}_j, \dots, \mathbf{X}_k).$$

(We use  $\hat{\mathbf{X}}_i$  to indicate that  $\mathbf{X}_i$  does not appear in the list).

*Proof.* We prove the identity for  $k = 1$ , the general case being similar; i.e., we will show that

$$d\alpha(\mathbf{X}, \mathbf{Y}) = \mathbf{X}(\alpha(\mathbf{Y})) - \mathbf{Y}(\alpha(\mathbf{X})) - \alpha([\mathbf{X}, \mathbf{Y}]).$$

We first establish that the right side of the above equation is tensorial in  $\mathbf{X}$  and  $\mathbf{Y}$ . It is certainly linear over vector fields, so that by Theorem 3.8.1 we only need to check linearity over functions. Now, for  $f : M \rightarrow \mathbb{R}$ , replacing  $\mathbf{X}$  by  $f\mathbf{X}$  in the last two terms

(the first one is trivially linear) on the right yields

$$\begin{aligned} -\mathbf{Y}(\alpha(f\mathbf{X})) - \alpha[f\mathbf{X}, \mathbf{Y}] &= -\mathbf{Y}(f(\alpha\mathbf{X})) - \alpha[f\mathbf{X}, \mathbf{Y}] \\ &= -(\mathbf{Y}f)\alpha(\mathbf{X}) - f\mathbf{Y}(\alpha\mathbf{X}) - \alpha(f[\mathbf{X}, \mathbf{Y}] - (\mathbf{Y}f)\mathbf{X}) \\ &= -f(\mathbf{Y}(\alpha(\mathbf{X})) - \alpha[\mathbf{X}, \mathbf{Y}]), \end{aligned}$$

which shows that it is indeed tensorial in the first argument. Since  $\alpha$  is skew-symmetric, it is also tensorial in the second argument. In light of this, the identity need only be established for coordinate vector fields. So let  $(U, \mathbf{x})$  be a chart. Then the restriction of  $\alpha$  to  $U$  equals  $\sum_l \alpha(\partial/\partial x^l) dx^l$ , so that

$$\begin{aligned} d\alpha|_U &= \sum_l d\left(\alpha \frac{\partial}{\partial x^l}\right) \wedge dx^l = \sum_{k,l} \frac{\partial}{\partial x^k} \left(\alpha \frac{\partial}{\partial x^l}\right) dx^k \wedge dx^l \\ &= \sum_{k<l} \left(\frac{\partial}{\partial x^k} \left(\alpha \frac{\partial}{\partial x^l}\right) - \frac{\partial}{\partial x^l} \left(\alpha \frac{\partial}{\partial x^k}\right)\right) dx^k \wedge dx^l. \end{aligned}$$

This means that

$$d\alpha \left(\frac{\partial}{\partial x^i}, \frac{\partial}{\partial x^j}\right) = \frac{\partial}{\partial x^i} \left(\alpha \frac{\partial}{\partial x^j}\right) - \frac{\partial}{\partial x^j} \left(\alpha \frac{\partial}{\partial x^i}\right),$$

which is the desired identity for coordinate vector fields.  $\square$

Another important property of the exterior derivative is that it commutes with pull-backs:

**Theorem 5.3.3.** *Let  $M, N$  be manifolds with exterior derivative operators  $d^M$  and  $d^N$  respectively. Given a map  $f : M \rightarrow N$ ,  $f^* \circ d^N = d^M \circ f^*$ .*

*Proof.* Having emphasized the fact that the two exterior derivative operators live in different spaces, we omit the superscript for brevity. Let us first consider the case of a zero-form  $h : N \rightarrow \mathbb{R}$ ; it must be established that  $f^* dh = d(h \circ f)$ . This is nothing more than the chain rule together with an exercise in notation, once we recall from Chapter 3 that by definition, given  $\mathbf{q} \in N$ ,  $\mathbf{u} \in N_{\mathbf{q}}$ ,

$$dh(\mathbf{q})\mathbf{u} = \mathcal{I}_{h(\mathbf{q})}^{-1}(h_{*\mathbf{q}}\mathbf{u}).$$

Indeed, for  $\mathbf{p} \in M$  and  $\mathbf{v} \in M_{\mathbf{p}}$ ,

$$\begin{aligned} f^* dh(\mathbf{p})(\mathbf{v}) &= dh(f(\mathbf{p}))(\mathbf{f}_{*\mathbf{p}}\mathbf{v}) = \mathcal{I}_{(h \circ f)(\mathbf{p})}^{-1} h_{*}(\mathbf{f}_{*\mathbf{p}}\mathbf{v}) = \mathcal{I}_{(h \circ f)(\mathbf{p})}^{-1} ((h \circ f)_{*\mathbf{p}}\mathbf{v}) \\ &= d(h \circ f)(\mathbf{p})(\mathbf{v}), \end{aligned}$$

as claimed. In the general case, linearity of both sides in the identity allows us to only consider a  $k$ -form on  $N$  of type  $h dx^{i_1} \wedge \cdots \wedge dx^{i_k}$ . Then

$$\begin{aligned}
 d(\mathbf{f}^* h dx^{i_1} \wedge \cdots \wedge dx^{i_k}) &= d((h \circ \mathbf{f}) \mathbf{f}^* dx^{i_1} \wedge \cdots \wedge \mathbf{f}^* dx^{i_k}) \\
 &= d((h \circ \mathbf{f}) d\mathbf{f}^* x^{i_1} \wedge \cdots \wedge d\mathbf{f}^* x^{i_k}) \\
 &= d(h \circ \mathbf{f}) \wedge d\mathbf{f}^* x^{i_1} \wedge \cdots \wedge d\mathbf{f}^* x^{i_k} \text{ since } d \circ d = 0 \\
 &= \mathbf{f}^* dh \wedge \mathbf{f}^* dx^{i_1} \wedge \cdots \wedge \mathbf{f}^* dx^{i_k} \\
 &= \mathbf{f}^* (dh \wedge dx^{i_1} \wedge \cdots \wedge dx^{i_k}) \\
 &= \mathbf{f}^* d(h dx^{i_1} \wedge \cdots \wedge dx^{i_k}).
 \end{aligned}$$

This completes the proof of the Theorem.  $\square$

## 5.4 Integration on manifolds

Our goal in this section is to define the integral of a form of degree  $k$  on a differentiable manifold of dimension  $k$ . The form will always be assumed to have compact support, and the manifold will need to be endowed with an orientation, in the following sense:

**Definition 5.4.1.** An  $n$ -dimensional manifold  $M$  is said to be *orientable* if it admits a nowhere-zero  $n$ -form.

Given any two nowhere-zero  $n$ -forms  $\omega_1$  and  $\omega_2$  on  $M^n$ , there exists a function  $f : M \rightarrow \mathbb{R}$ , which is either positive everywhere or negative everywhere, such that  $\omega_1 = f\omega_2$ . The two forms are said to be equivalent if  $f > 0$ . Thus, the collection of all non-zero  $n$ -forms splits into two disjoint subsets, called equivalence classes. An *orientation* of  $M$  is a choice of one of these two classes. Clearly, an orientable manifold has exactly two orientations, and the property of being orientable may be thought of as being able to orient all tangent spaces in a continuous manner. The *standard orientation* of  $\mathbb{R}^n$  is the equivalence class containing  $du^1 \wedge \cdots \wedge du^n$ . This is often described as the orientation induced by  $du^1 \wedge \cdots \wedge du^n$ .

**Proposition 5.4.1.**  $M$  is orientable if and only if it admits an atlas consisting of charts whose transition functions have positive Jacobian determinant; i.e, for  $(U, \mathbf{x})$  and  $(V, \mathbf{y})$  in the atlas,  $\det D(\mathbf{y} \circ \mathbf{x}^{-1}) > 0$  if  $U \cap V \neq \emptyset$ .

*Proof.* Suppose  $M$  is orientable, and consider an  $n$ -form  $\omega$  on  $M$  which vanishes nowhere. Choose any atlas of  $M$ . By reordering the components of the coordinate maps, if necessary, it may be assumed that for any chart  $(U, \mathbf{x})$ ,  $f_{\mathbf{x}} := \omega(\partial/\partial x^1, \dots, \partial/\partial x^n) > 0$ . Since  $\omega$ , when restricted to  $U$ , equals  $f_{\mathbf{x}} dx^1 \wedge \cdots \wedge dx^n$ , we have for overlapping charts  $(U, \mathbf{x})$  and  $(V, \mathbf{y})$ ,

$$\frac{1}{f_{\mathbf{y}}}\omega = dy^1 \wedge \cdots \wedge dy^n = \frac{f_{\mathbf{x}}}{f_{\mathbf{y}}} dx^1 \wedge \cdots \wedge dx^n$$

on  $U \cap V$ . By Theorem 5.2.3,  $f_x/f_y = \det D(\mathbf{y} \circ \mathbf{x}^{-1})$ , and the latter is therefore positive. Conversely, suppose there exists an atlas  $\{(U_\alpha, \mathbf{x}_\alpha)\}$  whose transition functions have positive Jacobian. Consider a partition of unity  $\{\varphi_i\}$  subordinate to the atlas, and for each  $i$  choose a chart  $(U_i, \mathbf{x}_i)$  in the atlas that contains the support of  $\varphi_i$ . Define forms  $\omega_i$  on  $M$  by

$$\omega_i(\mathbf{p}) = \begin{cases} \varphi_i(\mathbf{p}) (dx_i^1 \wedge \cdots \wedge dx_i^n)(\mathbf{p}), & \text{if } \mathbf{p} \in U_i \\ 0 & \text{otherwise.} \end{cases}$$

Then the form  $\omega = \sum_i \omega_i$  cannot vanish anywhere, because if  $\mathbf{q} \in U_j$ , then  $\omega(\mathbf{q})$  applied to the  $n$ -tuple  $(\partial/\partial x_j^1, \dots, \partial/\partial x_j^n)$  of vector fields at  $\mathbf{q}$  is a sum of nonnegative terms, at least one of which is positive (namely  $\omega_j(\mathbf{q})$  for any  $i$  such that  $\varphi_i(\mathbf{q}) > 0$ ).  $\square$

One example of nonorientable manifold is the Möbius strip described in Chapter 7.

Given manifolds  $M_i$  with orientations respectively induced by  $n$ -forms  $\omega_i$ ,  $i = 1, 2$ , a diffeomorphism  $\mathbf{f} : M_1 \rightarrow M_2$  is said to be *orientation-preserving* if  $\mathbf{f}^* \omega_2$  induces the same orientation of  $M_1$  as  $\omega_1$ . An orientation of  $M$  also induces an orientation of any open subset  $U$  of  $M$ , namely the one defined by  $\iota^* \omega$ , where  $\omega$  represents the orientation of  $M$  and  $\iota : U \hookrightarrow M$  denotes the inclusion map. A chart  $(U, \mathbf{x})$  of an oriented manifold  $M$  is said to be *positive* if  $\mathbf{x} : U \rightarrow \mathbf{x}(U)$  is orientation-preserving for the standard orientation of  $\mathbb{R}^n$ .

In order to define integration of an  $n$ -form on a manifold  $M^n$ , we begin with the special case  $M = \mathbb{R}^n$ :

**Definition 5.4.2.** Let  $\omega$  denote an  $n$ -form with compact support on  $\mathbb{R}^n$ ,  $U \subset \mathbb{R}^n$ . The *integral of  $\omega$  over  $U$*  is defined to be the ordinary integral over  $U$  of the function  $\omega(\mathbf{D}_1, \dots, \mathbf{D}_n)$ .

In other words, if we write  $\omega = f du^1 \wedge \cdots \wedge du^n$ , then  $\int_U \omega = \int_U f$ .

Next, suppose  $M^n$  is an oriented manifold,  $(U, \mathbf{x})$  a positive chart of  $M$ , and  $\omega$  an  $n$ -form on  $M$  with support in  $U$ . Define

$$\int_M \omega = \int_U \omega := \int_{\mathbf{x}(U)} (\mathbf{x}^{-1})^* \omega, \quad (5.4.1)$$

with the last integral as in Definition 5.4.2. It must be checked that this definition makes sense; i.e., that it does not depend on the particular chart. This is essentially due to the change of variables theorem:

**Theorem 5.4.1.** If  $(U, \mathbf{x})$  and  $(V, \mathbf{y})$  are positive charts on  $M^n$ , and  $\omega$  is an  $n$ -form with compact support in  $U \cap V$ , then

$$\int_{\mathbf{x}(U \cap V)} (\mathbf{x}^{-1})^* \omega = \int_{\mathbf{y}(U \cap V)} (\mathbf{y}^{-1})^* \omega.$$

*Proof.* First of all, observe that

$$\begin{aligned} \int_{\mathbf{x}(U \cap V)} (\mathbf{x}^{-1})^* \omega &= \int_{\mathbf{x}(U \cap V)} \omega \circ \mathbf{x}^{-1} (\mathbf{x}_*^{-1} \mathbf{D}_1, \dots, \mathbf{x}_*^{-1} \mathbf{D}_n) \\ &= \int_{\mathbf{x}(U \cap V)} \omega \circ \mathbf{x}^{-1} \left( \frac{\partial}{\partial x^1} \circ \mathbf{x}^{-1}, \dots, \frac{\partial}{\partial x^n} \circ \mathbf{x}^{-1} \right) \\ &= \int_{\mathbf{x}(U \cap V)} \left[ \omega \left( \frac{\partial}{\partial x^1}, \dots, \frac{\partial}{\partial x^n} \right) \right] \circ \mathbf{x}^{-1}, \end{aligned}$$

and similarly for the term involving  $\mathbf{y}$ , so that the identity to be established becomes

$$\int_{\mathbf{x}(U \cap V)} \omega \left( \frac{\partial}{\partial x^1}, \dots, \frac{\partial}{\partial x^n} \right) \circ \mathbf{x}^{-1} = \int_{\mathbf{y}(U \cap V)} \omega \left( \frac{\partial}{\partial y^1}, \dots, \frac{\partial}{\partial y^n} \right) \circ \mathbf{y}^{-1}.$$

Now, by Exercise 3.16,

$$\frac{\partial}{\partial y^i} = \sum_j (D_i(u^j \circ \mathbf{x} \circ \mathbf{y}^{-1}) \circ \mathbf{y}) \frac{\partial}{\partial x^j}.$$

Together with Theorem 5.2.3, this yields

$$\omega \left( \frac{\partial}{\partial y^1}, \dots, \frac{\partial}{\partial y^n} \right) = (\det D(\mathbf{x} \circ \mathbf{y}^{-1}) \circ \mathbf{y}) \omega \left( \frac{\partial}{\partial x^1}, \dots, \frac{\partial}{\partial x^n} \right).$$

But then by the change of variables theorem,

$$\begin{aligned} \int_{\mathbf{y}(U \cap V)} \omega \left( \frac{\partial}{\partial y^1}, \dots, \frac{\partial}{\partial y^n} \right) \circ \mathbf{y}^{-1} \\ &= \int_{\mathbf{y}(U \cap V)} \det D(\mathbf{x} \circ \mathbf{y}^{-1}) \omega \left( \frac{\partial}{\partial x^1}, \dots, \frac{\partial}{\partial x^n} \right) \circ \mathbf{y}^{-1} \\ &= \int_{\mathbf{y}(U \cap V)} |\det D(\mathbf{x} \circ \mathbf{y}^{-1})| \left[ \omega \left( \frac{\partial}{\partial x^1}, \dots, \frac{\partial}{\partial x^n} \right) \circ \mathbf{x}^{-1} \right] \circ \mathbf{x} \circ \mathbf{y}^{-1} \\ &= \int_{\mathbf{x}(U \cap V)} \omega \left( \frac{\partial}{\partial x^1}, \dots, \frac{\partial}{\partial x^n} \right) \circ \mathbf{x}^{-1}. \end{aligned}$$

This establishes the claim. □

Notice how the above proof used the fact that both charts were positive.

The general case (when the form does not have support inside the domain of a single chart) uses a partition of unity:

**Definition 5.4.3.** Let  $\omega$  denote an  $n$ -form with compact support on  $M^n$ , and  $\{(U_i, \mathbf{x}_i)\}$  a countable atlas with subordinate partition of unity  $\{\varphi_i\}$ . For each  $i$ , define an  $n$ -form  $\omega_i$  by  $\omega_i = \varphi_i \omega$  on  $U_i$  and  $\omega_i = 0$  outside  $U_i$ . The *integral of  $\omega$  over  $M$*  is

$$\int_M \omega = \sum_i \int_M \omega_i.$$

The proof that this definition does not depend on the chosen partition of unity is similar to that in Chapter 4, but more straightforward since the sum is a finite one.

**Examples and Remarks 5.4.1.** (i) Let  $M^2$  be a surface – i.e., a two-dimensional manifold – in  $\mathbb{R}^3$ . Any chart  $(U, \mathbf{x})$  of  $M$  induces a unit vector field  $\mathbf{N}$  on  $U$  orthogonal to  $M$ , namely

$$\mathbf{N} = \frac{1}{\left| \frac{\partial}{\partial x^1} \times \frac{\partial}{\partial x^2} \right|} \frac{\partial}{\partial x^1} \times \frac{\partial}{\partial x^2}.$$

If  $(V, \mathbf{y})$  is another chart, then by Exercise 3.16,

$$\frac{\partial}{\partial x^i} = \sum_j (D_i(\mathbf{u}^j \circ \mathbf{y} \circ \mathbf{x}^{-1}) \circ \mathbf{x}) \frac{\partial}{\partial x^j},$$

so that

$$\frac{\partial}{\partial x^1} \times \frac{\partial}{\partial x^2} = \det(D(\mathbf{y} \circ \mathbf{x}^{-1}) \circ \mathbf{x}) \frac{\partial}{\partial y^1} \times \frac{\partial}{\partial y^2}.$$

It follows that positive charts determine the same unit normal field on the intersection of their domains. Thus, an orientable surface admits a global unit normal field. Conversely, if  $M$  admits a unit normal field  $\mathbf{N}$ , then it is orientable: the 2-form  $\omega$  on  $M$ , given by  $\omega(\mathbf{X}, \mathbf{Y}) = \det(\mathbf{N}, \mathbf{X}, \mathbf{Y})$  is nowhere zero. We will later generalize this to hypersurfaces.

- (ii) The *volume form* of an oriented  $n$ -manifold  $M$  is the  $n$ -form  $\eta$  on  $M$  such that  $\eta(\mathbf{p})(\mathbf{u}_1, \dots, \mathbf{u}_n) = 1$  for any positively oriented orthonormal basis  $(\mathbf{u}_1, \dots, \mathbf{u}_n)$  of  $M_{\mathbf{p}}$ ,  $\mathbf{p} \in M$ . It is well-defined by Theorem 5.2.3. In fact, if  $\alpha$  is any  $n$ -form inducing the orientation, then  $f\alpha$ , where  $f(\mathbf{p}) = 1/\alpha(\mathbf{p})(\mathbf{u}_1, \dots, \mathbf{u}_n)$  for any positive orthonormal basis  $(\mathbf{u}_1, \dots, \mathbf{u}_n)$  of  $M_{\mathbf{p}}$ , is the volume form of  $M$ . When  $M$  is compact, its *volume* is  $\int_M \eta$ .
- (iii) By Theorem 1.4.4, there is a bijection  $\flat : \mathfrak{X}M \rightarrow \Lambda_1(M)$  given by  $\mathbf{X}^\flat(\mathbf{Y}) = \langle \mathbf{X}, \mathbf{Y} \rangle$  for vector fields  $\mathbf{X}, \mathbf{Y}$  on  $M$ , see also Exercise 3.9. It is customary in physics to consider so-called ‘line integrals of vector fields in  $\mathbb{R}^n$  along oriented curves’. This can be interpreted in the present context as follows: let  $M$  be a one-dimensional oriented manifold in  $\mathbb{R}^n$ ,  $\mathbf{X}$  a vector field on some open set  $U$  containing  $M$ , so that  $\mathbf{X}^\flat$  is a one-form on  $U$ . If  $\iota : M \hookrightarrow \mathbb{R}^n$  denotes inclusion, then  $\iota^* \mathbf{X}^\flat$  is a one-form on  $M$ , and the line integral of  $\mathbf{X}$  along the curve is defined to be  $\int_M \iota^* \mathbf{X}^\flat$ . This, in turn, is easily extendible to the case when the image of the curve  $\mathbf{c} : I \rightarrow \mathbb{R}^n$  is not a manifold: define the integral of  $\mathbf{X}$  along  $\mathbf{c}$  to be  $\int_I \mathbf{c}^* \mathbf{X}^\flat$ .

In physics, a common kind of vector field is a force field  $\mathbf{F}$ , which, when acting on a body of mass  $m$ , produces an acceleration  $\dot{\mathbf{c}}'$  satisfying  $\mathbf{F}(\mathbf{c}(t)) = m \dot{\mathbf{c}}'(t)$  according to Newton’s second law of motion. The *work*  $W$  done by  $\mathbf{F}$  in moving the object from  $\mathbf{c}(a)$  to  $\mathbf{c}(b)$  is defined to be  $\int_{[a,b]} \mathbf{c}^* \mathbf{F}^\flat$ . The *kinetic energy* of the object at time  $t$  is  $K(\mathbf{c}(t)) = (1/2)m|\dot{\mathbf{c}}|^2$ . It follows that the work done by the force equals the

change in kinetic energy:

$$\begin{aligned} W &= \int_{[a,b]} \mathbf{c}^* \mathbf{F}^b = \int_a^b \mathbf{F}^b(\dot{\mathbf{c}}) = \int_a^b \langle m\dot{\mathbf{c}}', \dot{\mathbf{c}} \rangle = \frac{m}{2} \int_a^b |\dot{\mathbf{c}}|^2 \\ &= K(\mathbf{c}(b)) - K(\mathbf{c}(a)). \end{aligned}$$

(iv) Let  $\omega$  denote the one-form on  $\mathbb{R}^2 \setminus \{0\}$  given by

$$\omega = \frac{-u^2}{((u^1)^2 + (u^2)^2)^{1/2}} du^1 + \frac{u^1}{((u^1)^2 + (u^2)^2)^{1/2}} du^2.$$

If  $\iota_r : S^1(r) \hookrightarrow \mathbb{R}^2$  is the inclusion map from the circle of radius  $r$  and center the origin into the plane, then  $\iota_r^* \omega$  is the volume form of  $S^1(r)$  for the induced (counterclockwise orientation): indeed the metric dual

$$\omega^\sharp = \frac{-u^2}{((u^1)^2 + (u^2)^2)^{1/2}} \mathbf{D}_1 + \frac{u^1}{((u^1)^2 + (u^2)^2)^{1/2}} \mathbf{D}_2$$

of  $\omega$  is tangent to all circles (being orthogonal to the position vector field), has length 1, and  $\omega(\omega^\sharp)$  is easily computed to equal one.

The curve  $\mathbf{c} : (0, 2\pi) \rightarrow \mathbb{R}^2$ ,  $\mathbf{c}(t) = (r \cos t, r \sin t)$ , is a parametrization of  $S^1(r) \setminus \{r\mathbf{e}_1\}$ , and

$$\dot{\mathbf{c}} = -r \sin \mathbf{D}_1 \circ \mathbf{c} + r \cos \mathbf{D}_2 \circ \mathbf{c} = r\omega^\sharp \circ \mathbf{c}.$$

Thus, the volume of  $S^1(r)$  equals

$$\int_{S^1(r)} \iota_r^* \omega = \int_0^{2\pi} \omega(\dot{\mathbf{c}}) = \int_0^{2\pi} r\omega(\omega^\sharp) = \int_0^{2\pi} r = 2\pi r,$$

as expected.

## 5.5 Manifolds with boundary

The *upper half-space*  $H^n$  is the set of all points  $\mathbf{p} \in \mathbb{R}^n$  such that  $u^n(\mathbf{p}) \geq 0$ . Even though it is closed in  $\mathbb{R}^n$ , we extended the notion of differentiability to maps defined on sets such as these; namely, a map is differentiable on  $H^n$  if it is extendable to a map that is smooth on some open set containing  $H^n$ .

In order to discuss Stokes' theorem, we need the following concept, which generalizes that of manifold:

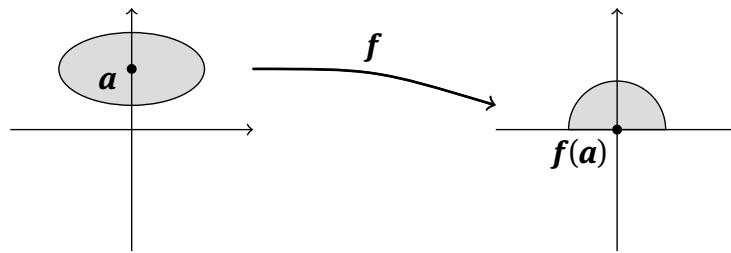
**Definition 5.5.1.** A subset  $M \subset \mathbb{R}^{n+k}$  of Euclidean space is said to be an *n-dimensional manifold with boundary* if every point  $\mathbf{p} \in M$  admits a neighborhood  $U$  in  $\mathbb{R}^{n+p}$  with either one of the following properties:



- (1) There exists an open set  $V$  in  $\mathbb{R}^n$ , a one-to-one differentiable map  $\mathbf{h} : V \rightarrow \mathbb{R}^{n+k}$  of maximal rank everywhere such that
- $\mathbf{h}(V) = U \cap M$ , and
  - $\mathbf{h}^{-1}$  is continuous;
- (2) There exists an open set  $V$  in  $H^n$ , a one-to-one differentiable map  $\mathbf{h} : V \rightarrow \mathbb{R}^{n+k}$  of maximal rank everywhere such that
- $\mathbf{h}(V) = U \cap M$ ,
  - $\mathbf{h}^{-1}$  is continuous, and
  - $(u^n \circ \mathbf{h}^{-1})(\mathbf{p}) = 0$ .

$\mathbf{h}$  is called a *local parametrization* of  $M$ , and its inverse a *chart*.

When all points of  $M$  satisfy the first condition, one recovers the usual definition of a manifold. It is important to realize that if a point  $\mathbf{p}$  satisfies the second condition for some parametrization  $\mathbf{h}$ , then it satisfies it for every parametrization: for if there were some other parametrization  $\tilde{\mathbf{h}}$  of the first type, with, say,  $\tilde{\mathbf{h}}(\mathbf{a}) = \mathbf{p}$  and  $u^n(\mathbf{a}) \neq 0$ , then  $\mathbf{f} := \mathbf{h}^{-1} \circ \tilde{\mathbf{h}}$  would be a map of maximal rank on a neighborhood of  $\mathbf{a}$  in  $\mathbb{R}^n$ . The image of arbitrarily small open neighborhoods of  $\mathbf{a}$  would not be open in  $\mathbb{R}^n$ , contradicting the inverse function theorem.



In view of this, we may define the *boundary*  $\partial M$  of  $M$  to be the set of all points satisfying the second condition. A word of caution is in order here: even though the notation is the same, the boundary of a manifold  $M$  need not coincide with the topological boundary of the set  $M$  as defined in Chapter 1, see Exercise 5.24.

One of the simplest examples of manifold with boundary is a closed metric ball in  $\mathbb{R}^n$ . The boundary sphere is then also a manifold, with dimension lower by 1. This is always the case:

**Proposition 5.5.1.** *Let  $M$  be an  $n$ -dimensional manifold with boundary. If  $\partial M$  is non-empty, then it is an  $(n - 1)$ -dimensional manifold.*

*Proof.* Consider  $\mathbf{p} \in \partial M$ . Any parametrization  $\mathbf{h} : U \rightarrow \mathbf{h}(U)$  of a neighborhood  $\mathbf{h}(U)$  of  $\mathbf{p}$  in  $M$  induces a parametrization  $\tilde{\mathbf{h}}$  of a neighborhood of  $\mathbf{p}$  in  $\partial M$ : Let  $\tilde{U} = \{\mathbf{u} \in \mathbb{R}^n \mid u^n(\mathbf{u}) = 0\}$ .  $\tilde{U}$  is naturally identified with a subset of  $\mathbb{R}^{n-1} = \mathbb{R}^{n-1} \times \{0\} \subset \mathbb{R}^n$ , and the restriction  $\tilde{\mathbf{h}}$  of  $\mathbf{h}$  to  $\tilde{U}$  is differentiable, has maximal rank, and has continuous inverse because  $\mathbf{h}$  enjoys these properties.  $\square$

In light of the above proposition, if  $\mathbf{p} \in \partial M$ , then the tangent space of  $\partial M$  at  $\mathbf{p}$  is well defined. It is also convenient to have a notion of tangent space of  $M$  itself at  $\mathbf{p}$ . In order to do so, consider a parametrization  $\mathbf{h}$  with  $\mathbf{h}(\mathbf{a}) = \mathbf{p}$ . Let  $\tilde{\mathbf{h}}$  be an extension of  $\mathbf{h}$  to an open set in  $\mathbb{R}^n$ , and define  $M_{\mathbf{p}} = \tilde{\mathbf{h}}_*(\mathbb{R}^n)$ . It must of course be checked that this is independent of the extension. If  $\bar{\mathbf{h}}$  is another extension of  $\mathbf{h}$ , choose a neighborhood  $U$  of  $\mathbf{a}$  in  $\mathbb{R}^n$  small enough so that both extensions are defined on  $U$ . Then  $\bar{\mathbf{h}}^{-1} \circ \tilde{\mathbf{h}}$  is the identity on  $U \cap H^n$ , and therefore the same is true for the derivative  $D(\bar{\mathbf{h}}^{-1} \circ \tilde{\mathbf{h}})$  on the intersection of  $U$  with the interior of  $H^n$ . Choosing any sequence  $\mathbf{a}_k \in U \cap (H^n)^0$  converging to  $\mathbf{a}$ , we have that  $D(\bar{\mathbf{h}}^{-1} \circ \tilde{\mathbf{h}})(\mathbf{a})$  is the identity by continuity, and thus,  $D(\bar{\mathbf{h}})(\mathbf{a}) = D(\tilde{\mathbf{h}})(\mathbf{a})$  as claimed. We will for convenience's sake denote either by  $D\mathbf{h}(\mathbf{a})$ .

Consider any  $\mathbf{p} \in \partial M$  and parametrization  $\mathbf{h}$  of a neighborhood of  $\mathbf{p}$ . Since  $M_{\mathbf{p}}$  has dimension  $n$ , and  $(\partial M)_{\mathbf{p}}$  is a subspace of  $M_{\mathbf{p}}$  of dimension  $n - 1$ , there are exactly two unit vectors in  $M_{\mathbf{p}}$  orthogonal to  $(\partial M)_{\mathbf{p}}$ . If  $\mathbf{x} = \mathbf{h}^{-1}$ , then they cannot lie in the kernel of  $d\mathbf{x}^n(\mathbf{p})$  because this kernel is precisely  $(\partial M)_{\mathbf{p}}$ . Thus, there is one, and only one unit vector  $\mathbf{n} \in M_{\mathbf{p}}$  orthogonal to  $(\partial M)_{\mathbf{p}}$  with  $\mathbf{n}(x^n) < 0$ .  $\mathbf{n}$  is called the *outward unit normal at  $\mathbf{p}$* . This procedure can be done at any point of  $\partial M$ , and the resulting map  $\mathbf{N} : \partial M \rightarrow TM$  which assigns to each element in the boundary the outward unit normal vector is called the *outward unit normal field*. This field can be used to induce an orientation of the boundary of  $M$ . In order to do so, we introduce the following terminology:

**Definition 5.5.2.** Let  $\mathbf{X}$  be a vector field on a manifold  $M$ . *Interior multiplication by  $\mathbf{X}$*  is the map  $i(\mathbf{X}) : \Lambda_k(M) \rightarrow \Lambda_{k-1}(M)$  given by

$$(i(\mathbf{X})\alpha)(\mathbf{X}_1, \dots, \mathbf{X}_{k-1}) = \alpha(\mathbf{X}, \mathbf{X}_1, \dots, \mathbf{X}_{k-1}), \quad \mathbf{X}_1, \dots, \mathbf{X}_{k-1} \in \mathfrak{X}M.$$

Suppose now that  $M$  is an oriented manifold with boundary, with orientation induced by some nowhere-zero  $n$ -form  $\omega$ . If  $\mathbf{N}$  is the outward unit normal field, then  $i(\mathbf{N})\omega$  is a nowhere-zero form on  $\partial M$ , and thus induces an orientation on  $\partial M$ . This orientation is said to be the one *induced by  $M$* . Notice that if  $\eta$  is the volume form of  $M$ , then  $i(\mathbf{N})\eta$  is the volume form of  $\partial M$ .

As a simple example, consider the closed disk  $D$  of radius 1 about the origin in  $\mathbb{R}^2$ . It is a 2-dimensional manifold with boundary. The open ball has a natural orientation as an open set in  $\mathbb{R}^2$ , namely that induced by the determinant. The induced orientation on the boundary  $S^1$  is the one induced by the (restriction of the) one-form  $\omega$  from Examples and Remarks 5.4.1 (iii), and is usually referred to as the counterclockwise orientation.

Another important special case is that of a one-dimensional manifold with boundary. For simplicity, let us consider the case when  $M^1$  is parametrized by a curve  $\mathbf{c} : [0, a] \rightarrow M$  with  $|\dot{\mathbf{c}}| \equiv 1$ .  $\partial M$  is a zero-dimensional manifold consisting of  $\mathbf{c}(0)$  and  $\mathbf{c}(1)$ . Since a zero-form is a function, two nowhere zero forms induce the same orientation on  $\partial M$  if they both have the same sign when evaluated at each point. Thus, an orientation is an assignment of a sign to each point in the boundary. If  $\dot{\mathbf{c}}$  represents

the orientation of  $M$ , then the induced orientation of  $\partial M$  is  $(\mathbf{c}(1), +)$ ,  $(\mathbf{c}(0), -)$ , because the outward unit normal is  $\dot{\mathbf{c}}(1)$  at  $\mathbf{c}(1)$  and  $-\dot{\mathbf{c}}(0)$  at  $\mathbf{c}(0)$ .

If  $\alpha$  is zero-form on an oriented compact zero-dimensional manifold  $M = \{\mathbf{a}_1, \dots, \mathbf{a}_k\}$ , define

$$\int_M \alpha := \sum_{i=1}^k \text{sign}(\mathbf{a}_i) \alpha(\mathbf{a}_i),$$

where  $\text{sign}$  denotes the sign induced by the orientation. The reader is invited to verify that with this notation, if  $M = [a, b] \subset \mathbb{R}$  with the standard orientation, and  $f : M \rightarrow \mathbb{R}$  is a zero-form, then the Fundamental Theorem of Calculus reads

$$\int_M df = \int_{\partial M} \iota^* f,$$

where  $\iota : \partial M \hookrightarrow M$  denotes inclusion (recall that for a zero-form  $f$ ,  $\iota^* f = f \circ \iota$  is the restriction of  $f$  to  $\partial M$ ).

## 5.6 Stokes' theorem

In the previous section, we rewrote the Fundamental Theorem of Calculus as

$$\int_M df = \int_{\partial M} \iota^* f.$$

The generalization to higher dimensions is known as Stokes' Theorem. Perhaps unsurprisingly, the key ingredients in its proof consist of the above theorem and Fubini's.

**Theorem 5.6.1** (Stokes' Theorem). *Let  $M$  be an oriented  $n$ -dimensional manifold with boundary,  $\omega$  an  $(n-1)$ -form with compact support in  $M$ . Then*

$$\int_M d\omega = \int_{\partial M} \iota^* \omega,$$

with  $\iota : \partial M \hookrightarrow M$  denoting inclusion.

*Proof.* By definition of the integral and the fact that  $\omega$  can be written as a finite sum of forms each of which has support inside some chart, it may be assumed that there is a local positive parametrization  $\mathbf{h} : H^n \supset U \rightarrow M$  whose image contains the support of  $\omega$ .

By Theorem 5.3.3,  $\int_M d\omega = \int_U \mathbf{h}^* d\omega = \int_U d\mathbf{h}^* \omega$ . Now,

$$\mathbf{h}^* \omega = \sum_{i=1}^n f^i du^1 \wedge \dots \wedge du^{i-1} \wedge du^{i+1} \wedge \dots \wedge du^n,$$

where  $f^i = \mathbf{h}^* \omega(\mathbf{D}_1, \dots, \mathbf{D}_{i-1}, \mathbf{D}_{i+1}, \dots, \mathbf{D}_n)$ . By Theorem 5.3.3,

$$\int_M d\omega = \int_U d\mathbf{h}^* \omega = \int_U \sum_i (-1)^i D_i f^i,$$

since  $d\mathbf{h}^*\omega = (\sum_i(-1)^i D_i f^i) du^1 \wedge \dots \wedge du^n$ . Extend  $f^i$  to smooth functions on  $H^n$  by setting them equal to zero outside  $U$ , and denote them both by the same symbol. We consider two possibilities:

Suppose first that  $\mathbf{h}(U)$  does not intersect the boundary of  $M$ . Then  $\iota^*\omega$ , and its integral over the boundary of  $M$  vanish. Consider a box  $B = [a_1, b_1] \times \dots \times [a_n, b_n]$  that contains  $U$  in its interior. Since the support of  $\omega$  is compact, it may be assumed that  $a_n > 0$ ; i.e., the box does not intersect the boundary of  $H^n$ . If  $B_i$  denotes the lower-dimensional box  $[a_1, b_1] \times \dots \times [a_{i-1}, b_{i-1}] \times [a_{i+1}, b_{i+1}] \times \dots \times [a_n, b_n]$ , then Fubini's theorem and the fundamental theorem of Calculus imply that  $\int_B D_i f^i = \int_{B_i} g^i$ , where

$$g^i(u^1, \dots, u^{n-1}) = f^i(u^1, \dots, u^{i-1}, b_i, u^i, \dots, u^{n-1}) - f^i(u^1, \dots, u^{i-1}, a_i, u^i, \dots, u^{n-1}).$$

But then  $g^i$  is identically zero, since  $f^i$  vanishes on the boundary of the box. This proves the theorem in this case.

Suppose next that  $\mathbf{h}(U)$  does intersect the boundary of  $M$ . As before, extend  $\mathbf{h}^*\omega$  to all of  $H^n$  by setting it equal to zero outside  $U$ , and enclose it in a box  $B = [a_1, b_1] \times \dots \times [0, b_n]$  large enough so that  $U$  is contained in  $(a_1, b_1) \times \dots \times [0, b_n)$ . We will, as above, denote by  $B_i$  the lower-dimensional box obtained by deleting the  $i$ -th edge of  $B$ . If  $\bar{\mathbf{h}}$  denotes the restriction of  $\mathbf{h}$  to the boundary  $\partial H^n$  of  $H^n$ , then

$$\bar{\mathbf{h}}^* \iota^* \omega = f_0^n du^1 \wedge \dots \wedge du^{n-1},$$

where

$$f_0^n(u^1, \dots, u^{n-1}) = f^n(u^1, \dots, u^{n-1}, 0),$$

see also Examples and Remarks 4.1.1 (iv) regarding notation. Now, the volume form on  $\partial H^n = \mathbb{R}^{n-1} \times \{0\}$  induced by  $H^n$  is given by

$$\begin{aligned} \mu &= i(-\mathbf{D}_n) du^1 \wedge \dots \wedge du^n \\ &= \sum_{j=1}^n \mu(\mathbf{D}_1, \dots, \mathbf{D}_{j-1}, \mathbf{D}_{j+1}, \dots, \mathbf{D}_n) du^1 \wedge \dots \wedge du^{j-1} \wedge du^{j+1} \wedge \dots \wedge du^n. \end{aligned}$$

The only non-vanishing term in this sum occurs when  $j = n$ , so that

$$\begin{aligned} \mu &= \mu(\mathbf{D}_1, \dots, \mathbf{D}_{n-1}) du^1 \wedge \dots \wedge du^{n-1} \\ &= du^1 \wedge \dots \wedge du^n (-\mathbf{D}_n, \mathbf{D}_1, \dots, \mathbf{D}_{n-1}) du^1 \wedge \dots \wedge du^{n-1} \\ &= (-1)^n du^1 \wedge \dots \wedge du^{n-1}. \end{aligned}$$

Thus,

$$\int_{\partial M} \iota^* \omega = \int_{B_n} (-1)^n f^n(u^1, \dots, u^{n-1}, 0) du^1 \dots du^{n-1}. \tag{5.6.1}$$

On the other hand, we have, as before,

$$\int_M d\omega = \sum_{i=1}^n (-1)^{i+1} \int_B D_i f^i = \sum_{i=1}^n (-1)^{i+1} \int_{B_i} g^i.$$

Since  $g^i = 0$  for  $i < n$ ,

$$\begin{aligned} \int_M d\omega &= (-1)^{n+1} \int_{B_n} (f^n(u^1, \dots, u^{n-1}, b_n) \\ &\quad - f^n(u^1, \dots, u^{n-1}, 0)) du^1 \dots du^{n-1} \\ &= \int_{B_n} (-1)^n f^n(u^1, \dots, u^{n-1}, 0) du^1 \dots du^{n-1}. \end{aligned}$$

Comparing with (5.6.1) now completes the proof.  $\square$

**Examples and Remarks 5.6.1.** (i) If  $M^n$  is a manifold without boundary, then

$$\int_M d\omega = 0 \text{ for any } n\text{-form with compact support in } M.$$

(ii) The compact support hypothesis in Stokes' theorem cannot be removed: consider for example the interior of a unit disk  $M^2 = \{(a, b) \in \mathbb{R}^2 \mid a^2 + b^2 < 1\}$ , and the 2-form  $\alpha = du^1 \wedge du^2$  on  $M$ . We have  $\alpha = d(u^1 du^2)$ , so if Stokes' theorem were to hold, then by (i), the integral of  $\alpha$  over  $M$  would vanish. By the definition of integral, however,  $\int_M \alpha$  equals the area  $\pi$  of  $M^2$ .

(iii) Consider the one-form

$$\omega = \frac{-u^2}{(u^1)^2 + (u^2)^2} du^1 + \frac{u^1}{(u^1)^2 + (u^2)^2} du^2$$

on  $\mathbb{R}^2 \setminus \{(0, 0)\}$ . It is related to the polar angle function  $\theta$  which is defined on  $\mathbb{R}^2 \setminus \{(a, 0) \mid a \geq 0\}$ . In fact, in each of the four open quadrants in the plane,  $\theta$  can be written as  $\arctan(u^2/u^1) + c$ , where  $c$  is some constant depending on the quadrant, cf. Section 4.6.1. A direct computation shows that  $d(\arctan(u^2/u^1)) = \omega$ . There is, however, no zero-form  $f$  defined on *all* of  $\mathbb{R}^2$  such that  $df = \omega$ : otherwise, with  $\iota : S^1 \hookrightarrow \mathbb{R}^2$  denoting the inclusion of the circle  $S^1$  with radius 1 around the origin, we would have that

$$\int_{S^1} \iota^* \omega = \int_{S^1} \iota^*(df) = \int_{S^1} d(\iota^* f) = 0$$

by (i). On the unit circle, however, this one-form coincides with that from Examples 5.4.1 (iii), so that  $\int_{S^1} \iota^* \omega = 2\pi$ . A differential form  $\omega$  is said to be *closed* if  $d\omega = 0$ , and *exact* if  $\omega = d\alpha$  for some form  $\alpha$  of degree one less. Since  $d \circ d = 0$ , any exact form is closed, but this example shows the converse is false in general. It also illustrates yet again how Stokes' theorem can fail without the compactness assumption: if  $M$  denotes the manifold with boundary that consists of the closed unit disk about the origin with the origin deleted, then  $\int_{\partial M} \iota^* \omega = 2\pi$ , but  $\int_M d\omega = 0$ , since  $\omega$  is closed.

(iv) The form  $\iota^* \omega$  on  $\partial M$  is by definition just the restriction of  $\omega$  to the boundary. For this reason, Stokes' theorem is often stated simply as  $\int_M d\omega = \int_{\partial M} \omega$ .

## 5.7 Classical versions of Stokes' theorem

We present some classical theorems found in multivariable calculus texts that are all special cases of Stokes' theorem. Since most books do not use differential forms, each statement needs some degree of reinterpretation.

**Theorem 5.7.1** (Green's Theorem). *Suppose  $M \subset \mathbb{R}^2$  is a compact 2-dimensional manifold with boundary, with the usual orientation. If  $P$  and  $Q$  are differentiable, then*

$$\int_{\partial M} P dx + Q dy = \iint_M (D_1 Q - D_2 P) dA.$$

*Proof.* As observed in Examples and Remarks 5.4.1 (ii), the left side is interpreted in physics as the line integral along  $\partial M$  of the vector field  $\mathbf{F} = P\mathbf{D}_1 + Q\mathbf{D}_2$ . In the present context, the identity means

$$\int_{\partial M} i^*(P du^1 + Q du^2) = \int_M (D_1 Q - D_2 P) du^1 \wedge du^2.$$

This is precisely Stokes' theorem for the 1-form  $\omega = \mathbf{F}^\flat = P du^1 + Q du^2$ .  $\square$

Before stating the second theorem, two more concepts must be defined: first of all, if  $M$  is an oriented 2-dimensional manifold in  $\mathbb{R}^3$ , the *positive unit normal field*  $\mathbf{N}$  of  $M$  is the unit normal field (out of two) such that the restriction of  $i(\mathbf{N}) du^1 \wedge du^2 \wedge du^3$  to  $M$  equals the volume form  $\eta$  of  $M$ . Extend the determinant and cross product to each tangent space via the canonical isomorphism  $\mathcal{I}_u : \mathbb{R}^3 \rightarrow \mathbb{R}_u^3$ . If  $\mathbf{N} = \sum N_i \mathbf{D}_i$ , then for a positive local orthonormal basis  $\mathbf{X}_1, \mathbf{X}_2$  of vector fields on  $M$ ,

$$\begin{aligned} 1 &= \eta(\mathbf{X}_1, \mathbf{X}_2) = du^1 \wedge du^2 \wedge du^3(\mathbf{N}, \mathbf{X}_1, \mathbf{X}_2) = \det \begin{bmatrix} \mathbf{N} & \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix} \\ &= \langle \mathbf{N}, \mathbf{X}_1 \times \mathbf{X}_2 \rangle \\ &= (N^1 du^2 \wedge du^3 + N^2 du^3 \wedge du^1 + N^3 du^1 \wedge du^2)(\mathbf{X}_1, \mathbf{X}_2), \end{aligned}$$

where the last equality was obtained by expanding the determinant along the first row. Thus, the volume form of an oriented 2-dimensional submanifold of  $\mathbb{R}^3$  with positive unit normal  $\mathbf{N} = \sum_i N^i \mathbf{D}_i$  is

$$\eta = N^1 du^2 \wedge du^3 + N^2 du^3 \wedge du^1 + N^3 du^1 \wedge du^2.$$

The second concept is that of the *curl* of a vector field. The curl of  $\mathbf{F} = \sum_i F^i \mathbf{D}_i$  is the vector field

$$\text{curl } \mathbf{F} = (D_2 F^3 - D_3 F^2) \mathbf{D}_1 + (D_3 F^1 - D_1 F^3) \mathbf{D}_2 + (D_1 F^2 - D_2 F^1) \mathbf{D}_3,$$

cf. also the exercises for a more motivating definition.

With this in mind, the classical Stokes' theorem may now be stated as follows:

**Theorem 5.7.2** (Stokes' theorem). *Let  $M^2 \subset \mathbb{R}^3$  denote a compact oriented 2-dimensional submanifold of  $\mathbb{R}^3$  with boundary and positive unit normal field  $\mathbf{N}$ . If  $\mathbf{F}$  is a differentiable vector field on an open set in  $\mathbb{R}^3$  containing  $M$ , then*

$$\int_M \langle \operatorname{curl} \mathbf{F}, \mathbf{N} \rangle dS = \int_{\partial M} \langle \mathbf{F}, d\mathbf{c} \rangle.$$

*Proof.* The right side is standard notation for the line integral of  $\mathbf{F}$  along the boundary. Thus, if  $\mathbf{F} = \sum_i F^i \mathbf{D}_i$ , and  $\omega = \mathbf{F}^\flat = \sum_i F^i du^i$  is the dual 1-form, then the identity that needs to be established is

$$\int_M \langle \operatorname{curl} \mathbf{F}, \mathbf{N} \rangle \eta = \int_{\partial M} \iota^* \omega,$$

where  $\eta$  is the volume form of  $M$ . In view of our version of Stokes', it must then be shown that  $\langle \operatorname{curl} \mathbf{F}, \mathbf{N} \rangle \eta = d\omega$ . Notice that for orthonormal vector fields  $\mathbf{X}_i$  on  $M$ ,

$$\det [\langle \mathbf{F}, \mathbf{N} \rangle \mathbf{N} \quad \mathbf{X}_1 \quad \mathbf{X}_2] = \det [\mathbf{F} \quad \mathbf{X}_1 \quad \mathbf{X}_2] \quad (5.7.1)$$

since  $\langle \mathbf{F}, \mathbf{N} \rangle \mathbf{N}$  is the component of  $\mathbf{F}$  orthogonal to  $TM = \operatorname{span}\{\mathbf{X}_1, \mathbf{X}_2\}$ . Thus,

$$\begin{aligned} \langle \operatorname{curl} \mathbf{F}, \mathbf{N} \rangle \eta(\mathbf{X}_1, \mathbf{X}_2) &= \det [\langle \operatorname{curl} \mathbf{F}, \mathbf{N} \rangle \mathbf{N} \quad \mathbf{X}_1 \quad \mathbf{X}_2] \\ &= \det [\operatorname{curl} \mathbf{F} \quad \mathbf{X}_1 \quad \mathbf{X}_2]. \end{aligned}$$

Writing out this last expression yields

$$\begin{aligned} ((D_2 F^3 - D_3 F^2) du^2 \wedge du^3 + (D_3 F^1 - D_1 F^3) du^3 \wedge du^1 \\ + (D_1 F^2 - D_2 F^1) du^1 \wedge du^2)(\mathbf{X}_1, \mathbf{X}_2), \end{aligned}$$

which is precisely  $d\omega(\mathbf{X}_1, \mathbf{X}_2)$ .  $\square$

The next result involves yet another concept commonly used in physics: the *divergence* of a vector field  $\mathbf{X} = \sum_i X^i \mathbf{D}_i$  in  $\mathbb{R}^3$  is the function  $\operatorname{div} \mathbf{X} = D_1 X^1 + D_2 X^2 + D_3 X^3$ . It can actually be defined for vector fields on a manifold of arbitrary dimension, and is explored further in the exercises.

**Theorem 5.7.3** (Divergence theorem). *Let  $M$  be a compact 3-dimensional manifold with boundary in  $\mathbb{R}^3$ , with positive unit normal field  $\mathbf{N}$ . If  $\mathbf{F}$  is a differentiable vector field on  $M$ , then*

$$\int_M \operatorname{div} \mathbf{F} dV = \int_{\partial M} \langle \mathbf{F}, \mathbf{N} \rangle dS.$$

*Proof.* The identity to be established is

$$\int_M \operatorname{div} \mathbf{F} du^1 \wedge du^2 \wedge du^3 = \int_{\partial M} \langle \mathbf{F}, \mathbf{N} \rangle \eta,$$

where  $\eta$  denotes the volume form of the boundary. If  $\mathbf{F} = \sum_i F^i \mathbf{D}_i$ , let

$$\omega = F^1 du^2 \wedge du^3 + F^2 du^3 \wedge du^1 + F^3 du^1 \wedge du^2,$$

so that  $d\omega = \operatorname{div} \mathbf{F} du^1 \wedge du^2 \wedge du^3$ . Now, by (5.71), given vector fields  $\mathbf{X}_i$  on  $\partial M$ ,

$$\begin{aligned} \langle \mathbf{F}, \mathbf{N} \rangle \eta(\mathbf{X}_1, \mathbf{X}_2) &= \det \left[ \langle \mathbf{F}, \mathbf{N} \rangle \mathbf{N} \quad \mathbf{X}_1 \quad \mathbf{X}_2 \right] = \det \left[ \mathbf{F} \quad \mathbf{X}_1 \quad \mathbf{X}_2 \right] \\ &= (F^1 du^1 \wedge du^2 + F^2 du^3 \wedge du^1 \\ &\quad + F^3 du^1 \wedge du^2)(\mathbf{X}_1, \mathbf{X}_2) \\ &= \iota^* \omega(\mathbf{X}_1, \mathbf{X}_2), \end{aligned}$$

so that  $\langle \mathbf{F}, \mathbf{N} \rangle \eta = \iota^* \omega$ , and the identity is once again that from Stokes' theorem.  $\square$

The above theorems provide physical interpretations for the notions of curl and divergence. Consider for example a fluid with constant density  $\rho$  in a region of 3-space. If  $\mathbf{v}$  is the velocity field of the fluid,  $\mathbf{F} = \rho \mathbf{v}$ , and  $\mathbf{N}$  a unit field, then  $\langle \mathbf{F}, \mathbf{N} \rangle$  represents the rate of flow (in the physical sense, not to be confused with the flow of a vector field) per unit area in direction  $\mathbf{N}$ . If  $\mathbf{a}$  is a point in the fluid,  $B_r(\mathbf{a})$  the set of points at distance less than  $r$  from  $\mathbf{a}$ , let

$$m_r = \inf\{\operatorname{div} \mathbf{F}(\mathbf{b}) \mid \mathbf{b} \in \overline{B_r(\mathbf{a})}\}, \quad M_r = \sup\{\operatorname{div} \mathbf{F}(\mathbf{b}) \mid \mathbf{b} \in \overline{B_r(\mathbf{a})}\}.$$

Denote by  $S_r(\mathbf{a})$  the sphere of radius  $r$  centered at  $\mathbf{a}$ . By the divergence theorem,

$$m_r \leq \frac{1}{\operatorname{vol}(B_r(\mathbf{a}))} \int_{B_r(\mathbf{a})} \operatorname{div} \mathbf{F} = \frac{1}{\operatorname{vol}(B_r(\mathbf{a}))} \int_{S_r(\mathbf{a})} \langle \mathbf{F}, \mathbf{N} \rangle \eta \leq M_r,$$

so that

$$\operatorname{div} \mathbf{F}(\mathbf{a}) = \lim_{r \rightarrow 0} \frac{1}{\operatorname{vol}(B_r(\mathbf{a}))} \int_{S_r(\mathbf{a})} \langle \mathbf{F}, \mathbf{N} \rangle \eta.$$

In this case, the divergence may be interpreted as the net rate of outward flow per unit volume at  $\mathbf{a}$ .  $\mathbf{a}$  is called a *source* if  $\operatorname{div} \mathbf{F}(\mathbf{a}) > 0$ , and a *sink* if it is negative. A fluid is said to be *incompressible* if it has everywhere vanishing divergence. More generally, the reader is asked to show in the exercises that a vector field on  $\mathbb{R}^n$  has everywhere vanishing divergence if and only if the (mathematical) flow  $\Phi_t$  of the vector field preserves volumes; i.e.,  $\Phi_t^* \omega = \omega$  for all  $t$ , where  $\omega$  is the volume form of  $\mathbb{R}^n$ .

The curl has a similar interpretation, but in terms of rotation of a fluid. This is also discussed in the exercises.

**Examples 5.7.1.** (i) One of Maxwell's equations in electromagnetism is

$$\operatorname{curl} \mathbf{H} = \mathbf{J},$$

where  $\mathbf{H}$  denotes the magnetizing field and  $\mathbf{J}$  the current density. It can be used to prove Ampère's law, which states that if  $C$  is a closed curve, then the current  $I$



crossing any surface  $M$  bounded by  $C$  is

$$I = \int_C \langle \mathbf{H}, d\mathbf{c} \rangle.$$

To see this, let  $M$  be a 2-dimensional manifold with boundary  $\partial M = C$ . As in the case of a fluid,  $\langle \mathbf{J}, \mathbf{N} \rangle$  represents the current per unit area that crosses a surface orthogonal to  $\mathbf{N}$ . Then, by Stokes' theorem, the current crossing  $M$  is

$$I = \int_M \langle \mathbf{J}, \mathbf{N} \rangle dS = \int_M \langle \text{curl } \mathbf{H}, \mathbf{N} \rangle dS = \int_C \langle \mathbf{H}, d\mathbf{c} \rangle.$$

- (ii) Yet another of Maxwell's equations is that if  $\mathbf{E}$  is the electric field created by a charge distribution with charge density  $\rho$ , then  $\text{div } \mathbf{E} = \rho / \epsilon_0$ , where  $\epsilon_0$  is a universal constant. One can use it to prove Gauss's law, which says that if  $M$  is a compact 3-dimensional manifold with boundary  $\partial M$ , then the total electric charge inside  $\partial M$  (i.e., in  $M$ ) is

$$Q = \epsilon_0 \int_{\partial M} \langle \mathbf{E}, \mathbf{N} \rangle dS.$$

Indeed, by the divergence theorem and the definition of  $Q$ ,

$$Q = \int_M \rho dV = \epsilon_0 \int_M \text{div } \mathbf{E} = \epsilon_0 \int_{\partial M} \langle \mathbf{E}, \mathbf{N} \rangle dS.$$

### 5.7.1 An application: the polar planimeter

The fundamental theorem of Calculus and its generalization, Stokes' theorem, illustrate how the behavior of a function on an open set is controlled by its behavior on the boundary of that set. This has some important practical applications: One such is an instrument, called a *planimeter*, that calculates the area of the region bounded by a simple closed curve in the plane. It does so by merely tracing the boundary curve. There are several types of planimeters. The first one was invented in 1814 in Bavaria, and the one we discuss here, the polar planimeter, was created by Jakob Amsler in 1856. Although the underlying principle relies on Green's theorem, it seems that this was only established almost a century later, cf. [10, 3].

Our planimeter consists of a pole that is fixed at the origin in the plane. A pole arm which can rotate about the pole is attached to it, and at the end of the pole arm is a pivot connected to the tracer arm. The tracer arm terminates in a wheel that is perpendicular to the arm, so that the wheel can only roll about the pivot. It can, however, also slide if moved parallel to the arm. The wheel then traces out the boundary curve by alternatively rolling and sliding, but only the rolling part is recorded.

In order to simplify the computations, our planimeter has both arms of equal length  $R$ . Thus, the curve that is traced has to lie within a disk of radius  $2R$  about the pole. We will see that if  $l$  is the total distance the wheel has rolled after tracing the

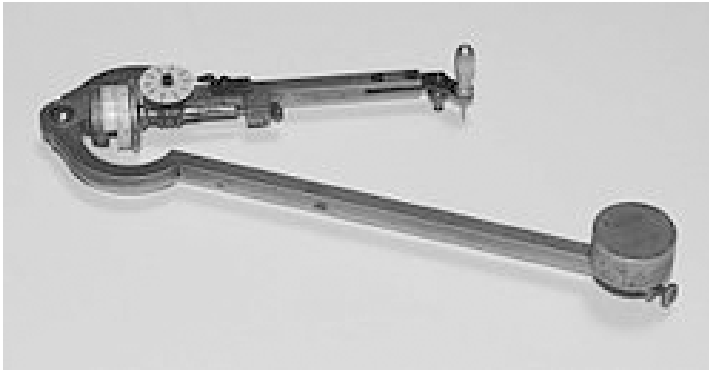
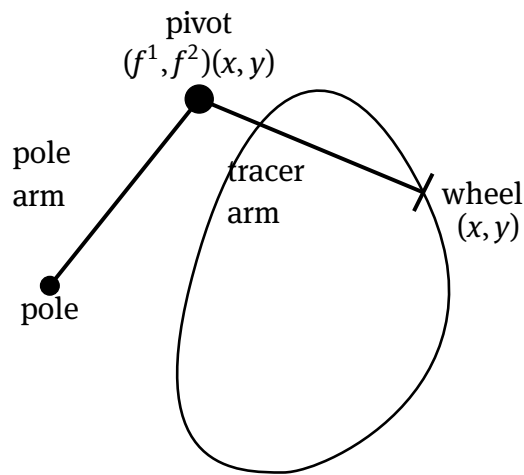


Fig. 5.1: A polar planimeter [1]

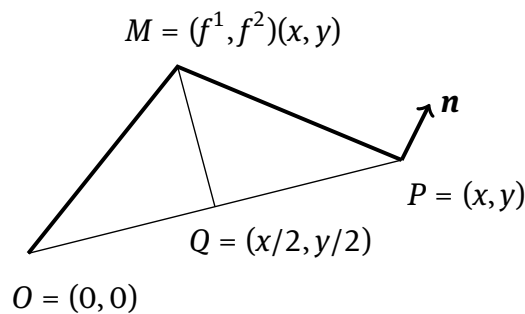
boundary curve, then the area  $A$  of the region enclosed by the curve is

$$A = IR.$$



Simple geometry dictates that if the wheel is at a given point  $(x, y)$ , there are exactly two possible locations for the pivot; these are mutual reflections in the line connecting the pole to the wheel. One way to see this is that the pivot must lie at one of the two intersections of circles with radius  $R$ , the first one centered at the origin and the other at  $(x, y)$ .

In order to find the coordinates  $(f^1, f^2)$  of the pivot  $M$ , we use the figure below, and denote by  $\vec{QM}$  the vector version of the point  $M - Q$ . Thus,  $|\vec{QM}|^2 = R^2 - (x^2 + y^2)/4$ , and  $\vec{QM}$  is orthogonal to  $\vec{OP} = [x \ y]$ .



There are two possible directions for this orthogonal vector, and we choose a counterclockwise rotation by  $\pi/2$ . This transforms  $[x \ y]$  into  $[-y \ x]$ , so that

$$\overrightarrow{QM} = \left( \frac{R^2 - (x^2 + y^2)/4}{(x^2 + y^2)} \right)^{1/2} [-y \ x],$$

and since  $\overrightarrow{OM} = \overrightarrow{OQ} + \overrightarrow{QM}$ , the coordinates of the pivot  $M$  are

$$f^1(x, y) = \frac{1}{2} \left[ x - y \left( \frac{4R^2 - x^2 - y^2}{x^2 + y^2} \right)^{1/2} \right],$$

$$f^2(x, y) = \frac{1}{2} \left[ y + x \left( \frac{4R^2 - x^2 - y^2}{x^2 + y^2} \right)^{1/2} \right].$$

It immediately follows that

$$\overrightarrow{MP} = \frac{1}{2} \left[ x + y \left( \frac{4R^2 - x^2 - y^2}{x^2 + y^2} \right)^{1/2} \quad y - x \left( \frac{4R^2 - x^2 - y^2}{x^2 + y^2} \right)^{1/2} \right],$$

and performing once again a counterclockwise rotation by  $\pi/2$ , we see that the unit normal vector field orthogonal to the tracer arm  $\overrightarrow{MP}$  is  $\mathbf{n} = PD_1 + D_2$ , where

$$P(x, y) = \frac{1}{2R} \left[ -y + x \left( \frac{4R^2 - x^2 - y^2}{x^2 + y^2} \right)^{1/2} \right],$$

$$Q(x, y) = \frac{1}{2R} \left[ x + y \left( \frac{4R^2 - x^2 - y^2}{x^2 + y^2} \right)^{1/2} \right].$$

Now, the total distance that the wheel has rolled is the line integral of  $\mathbf{n}$  along the boundary curve of the region  $M$ , and according to Green's theorem, this integral equals

$$\int_M (D_1Q - D_2P) du^1 \wedge du^2.$$

But

$$D_1Q(x, y) = \frac{1}{2R} + yD_1g(x, y), \quad D_2P(x, y) = -\frac{1}{2R} + xD_2g(x, y), \quad (5.7.2)$$

where

$$g(x, y) = \left( \frac{4R^2 - x^2 - y^2}{x^2 + y^2} \right)^{1/2}.$$

A straightforward calculation yields that

$$D_1g(x, y) = -\frac{4R^2x}{(x^2 + y^2)^{3/2}(4R^2 - x^2 - y^2)^{1/2}}.$$

Since the formula for  $g$  is symmetric in  $x$  and  $y$ ,  $D_2g$  is obtained by interchanging  $x$  and  $y$  in the formula for  $D_1g$ . Thus  $yD_1g(x, y) = xD_2g(x, y)$ , and (5.7.2) implies that the integrand in Green's theorem is  $D_1Q - D_2P = 1/R$ . This means that the distance the wheel has rolled equals the area of  $M$  divided by  $R$ , as claimed.

## 5.8 Closed forms and exact forms

Recall that a differential form  $\alpha$  on  $M$  is said to be closed if  $d\alpha = 0$ , and exact if there is some form  $\beta$  such that  $d\beta = \alpha$ . Every exact form is closed, but not every closed form is exact, cf. Examples and Remarks 5.6.1 (iii). We will show, however, that a closed form is always locally exact.

**Remark 5.8.1.** A one-form  $\alpha$  on a Riemannian manifold is exact if and only if its dual vector field is the gradient of some function: the equation  $\alpha = df$  is equivalent to  $\alpha^\sharp = \nabla f$ . Alternatively, for a vector field  $\mathbf{X}$ ,  $\mathbf{X} = \nabla f$  if and only if  $\mathbf{X}^\flat = df$ .

This concept is important enough in physics to warrant its own terminology: a force (i.e., a vector field)  $\mathbf{F}$  is said to be *conservative* if there is some function  $f$  such that  $\mathbf{F} = \nabla f$ . In this case, the function  $P = -f$  is called a *potential energy* function for  $f$ . Such a function is only defined up to an additive constant (if  $P$  is a potential energy, then so is  $P + c$  for any constant  $c$ ), but this is usually irrelevant, since one is mostly interested in the change in potential energy. In fact, the reason such forces are called conservative is that they satisfy the *conservation of energy law*: the sum of potential and kinetic energy is constant.

To see why, recall from Examples and Remarks 5.4.1 (ii) that if an object of mass  $m$  moves along a curve  $\mathbf{c} : [a, b] \rightarrow M$  under the action of a (not necessarily conservative) force  $\mathbf{F}$ , then its kinetic energy at  $\mathbf{c}(t)$  is  $K(\mathbf{c}(t)) = (m|\dot{\mathbf{c}}(t)|^2)/2$  and the work done by  $\mathbf{F}$  in moving the object from  $\mathbf{c}(a)$  to  $\mathbf{c}(b)$  is

$$W = \int_{\mathbf{c}} \mathbf{F}^\flat = K(\mathbf{c}(b)) - K(\mathbf{c}(a)).$$

If in addition  $\mathbf{F}$  is conservative, then the work also equals

$$W = \int_{\mathbf{c}} \mathbf{F}^\flat = \int_{\mathbf{c}} -dP = P(\mathbf{c}(a)) - P(\mathbf{c}(b)),$$

so that  $(P + K)(\mathbf{c}(a)) = (P + K)(\mathbf{c}(b))$ , as claimed.

A typical example of a conservative force is the gravitational force: an object of mass  $M$  that is placed at the origin in 3-space exerts a gravitational force on nearby objects. The force acting on an object of mass  $m$  located at  $\mathbf{a} \in \mathbb{R}^3 \setminus \{\mathbf{0}\}$  is

$$\mathbf{F}(\mathbf{a}) = -\frac{GMm}{|\mathbf{a}|^3} \mathbf{P}(\mathbf{a}),$$

where  $G$  is a universal constant and  $\mathbf{P}$  denotes the position vector field  $\mathbf{P} = \sum_i u^i \mathbf{D}_i$ . Thus,

$$\mathbf{F} = \nabla \left( \frac{GMm}{|\mathbf{P}|} \right),$$

as a straightforward calculation shows.

Denote by  $Z^k(M)$  the space of all closed  $k$ -forms on  $M$ , and by  $B^k(M)$  that of exact  $k$ -forms.  $B^k(M)$  is a subspace of  $Z^k(M)$ . In general, if  $W$  is a subspace of a vector space  $V$ , the *quotient space*  $V/W$  is the collection of all elements of the form  $\mathbf{v} + W$ ,  $\mathbf{v} \in V$ . Two elements  $\mathbf{v}_1 + W$  and  $\mathbf{v}_2 + W$  are said to be equal if  $\mathbf{v}_1 - \mathbf{v}_2 \in W$ .  $V/W$  inherits a vector space structure from the operations in  $V$ , if we define

$$(\mathbf{v}_1 + W) + (\mathbf{v}_2 + W) = (\mathbf{v}_1 + \mathbf{v}_2) + W, \quad c(\mathbf{v} + W) = (c\mathbf{v}) + W,$$

for  $\mathbf{v}_i \in V$ ,  $c \in \mathbb{R}$ . It must be checked that these operations are well defined: for example, if  $\mathbf{v}_1 + W = \mathbf{v}_2 + W$ , then for any  $\mathbf{v}_3$ , the identity

$$(\mathbf{v}_1 + W) + (\mathbf{v}_3 + W) = (\mathbf{v}_2 + W) + (\mathbf{v}_3 + W)$$

must hold. The left side is  $\mathbf{v}_1 + \mathbf{v}_3 + W$ , whereas the right side equals  $\mathbf{v}_2 + \mathbf{v}_3 + W$ . Since  $\mathbf{v}_1 + \mathbf{v}_3 - (\mathbf{v}_2 + \mathbf{v}_3) = \mathbf{v}_1 - \mathbf{v}_2$  lies in  $W$ , both sides indeed agree. A similar argument works for scalar multiplication. In fact, with these operations, the projection  $\mathbf{V} \rightarrow \mathbf{V}/\mathbf{W}$  which maps  $\mathbf{v}$  to  $\mathbf{v} + W$  is a linear transformation that is by construction surjective and has  $W$  as kernel, so that if  $V$  is finite-dimensional, then  $V/W$  has dimension  $\dim V - \dim W$ .

**Definition 5.8.1.** Given a positive integer  $k$ , the  $k$ -th de Rham cohomology space of  $M$  is  $H^k(M) = Z^k(M)/B^k(M)$ .  $H^0(M)$  is defined as  $Z^0(M)$ .

The reader is invited to verify that if  $M$  is connected, then  $H^0(M) \cong \mathbb{R}$ , see Exercise 5.22. Even though determining cohomology spaces can, in general, be challenging, there is one case that can already be settled:

**Proposition 5.8.1.** *If  $M$  is a compact, oriented  $n$ -dimensional manifold without boundary,  $n > 0$ , then  $H^n(M) \neq 0$ .*

*Proof.* Let  $\omega$  be the volume form of  $M$ .  $\omega$  is closed because there are no  $(n + 1)$ -forms. By definition, if  $(U, \mathbf{x})$  is any positive chart of  $M$ , then  $\int_{\mathbf{x}(U)} (\mathbf{x}^{-1})^* \omega > 0$ . Thus,  $\int_M \omega > 0$ , so that  $\omega$  cannot be exact (if  $\omega = d\alpha$ , then  $\int_M \omega = \int_{\partial M} \alpha = 0$  because the boundary is empty).  $\square$

Our next goal is to show that the pullback of  $\mathbf{f} : M \rightarrow N$  induces a linear transformation  $\mathbf{f}^* : H^k(N) \rightarrow H^k(M)$ . In general, if  $W_i$  is a subspace of  $V_i$ ,  $i = 1, 2$ , then any linear transformation  $L : V_1 \rightarrow V_2$  that maps  $W_1$  to  $W_2$  induces a linear map  $L : V_1/W_1 \rightarrow V_2/W_2$ : define  $L(\mathbf{v} + W_1) = (L\mathbf{v}) + W_2$ . One only needs to check that this definition makes sense, because linearity follows from linearity of the original map. But if  $\mathbf{u} + W_1 = \mathbf{v} + W_1$ , then  $\mathbf{u} - \mathbf{v} \in W_1$ , so that  $L\mathbf{u} - L\mathbf{v} = L(\mathbf{u} - \mathbf{v}) \in W_2$ ; i.e.,  $L\mathbf{u} + W_2 = L\mathbf{v} + W_2$ , and the map is well defined.

**Proposition 5.8.2.** *A differentiable map  $\mathbf{f} : M \rightarrow N$  induces linear transformations  $\mathbf{f}^* : H^k(N) \rightarrow H^k(M)$  for all  $k$ .*

*Proof.* By Theorem 5.3.3, the pullback  $f^* : \Lambda_k(N) \rightarrow \Lambda_k(M)$  maps closed forms to closed forms, and exact forms to exact forms. The first property means that  $f^*$  restricts to a map  $f^* : Z^k(N) \rightarrow Z^k(M)$ , and the second one implies the claim by the remark before the Proposition.  $\square$

By definition of the pullback, if  $f : M_1 \rightarrow M_2$  and  $g : M_2 \rightarrow M_3$ , then  $g \circ f : M_1 \rightarrow M_3$  has as pullback

$$(g \circ f)^* = f^* \circ g^* : H^k(M_3) \rightarrow H^k(M_1).$$

One consequence is that diffeomorphic manifolds have isomorphic cohomology spaces: the pullback  $1_M^*$  of the identity map  $1_M : M \rightarrow M$  is the identity on the cohomology of  $M$ , so that if  $f : M \rightarrow N$  is a diffeomorphism, then  $f^* \circ f^{-1*} = 1_{H^k(N)}$  and  $f^{-1*} \circ f^* = 1_{H^k(M)}$ . Thus,  $f^*$  is an isomorphism with inverse  $f^{-1*}$ .

An equivalent formulation is that if two manifolds have different cohomology, then they are not diffeomorphic. It turns out that a much stronger property holds. In order to state it, we will need the following:

**Definition 5.8.2.** Two maps  $f, g : M \rightarrow N$  are said to be *homotopic* if there exists a map  $H : M \times [0, 1] \rightarrow N$  such that  $H \circ \iota_0 = f$  and  $H \circ \iota_1 = g$ , where  $\iota_t : M \rightarrow M \times [0, 1]$  maps  $p \in M$  to  $(p, t)$ ,  $0 \leq t \leq 1$ . In this case,  $H$  is called a *homotopy* between  $f$  and  $g$ .

Roughly speaking, a homotopy smoothly deforms  $f$  into  $g$  by a “curve of maps”  $t \mapsto H \circ \iota_t$  joining  $f$  to  $g$ . Two spaces  $M$  and  $N$  are said to be *homotopy equivalent* if there exist maps  $f : M \rightarrow N$  and  $g : N \rightarrow M$  such that  $g \circ f$  and  $f \circ g$  are homotopic to the identity maps on the two spaces. Spaces that are diffeomorphic are of course homotopy equivalent, but the latter concept is much weaker: for example,  $\mathbb{R}^n$  is homotopy equivalent to the one point subset consisting of the origin. In fact, if  $f : \mathbb{R}^n \rightarrow \{\mathbf{0}\}$  maps everything to the origin, and  $g : \{\mathbf{0}\} \rightarrow \mathbb{R}^n$  is inclusion, then  $f \circ g$  is the identity on  $\{\mathbf{0}\}$ , whereas  $g \circ f$  is homotopic to the identity on  $\mathbb{R}^n$  via  $H : \mathbb{R}^n \times [0, 1] \rightarrow \mathbb{R}^n$ , where  $H(p, t) = (1 - t)p$ . In general, a manifold that is homotopy equivalent to a one point subset is said to be *contractible*; alternatively,  $M$  is contractible if the identity map on  $M$  is homotopic to a constant map (i.e., to one that maps all of  $M$  to a single point).

Let  $\pi_M$  and  $t$  denote the projections of the product  $M \times [0, 1]$  onto the factors, so that  $(\pi_{M*}, t_*) : (M \times [0, 1])_{(p, t_0)} \rightarrow M_p \times [0, 1]_{t_0}$  is an isomorphism, cf. Exercise 3.14. Denote by  $\tilde{D}$  the vector field on  $M \times [0, 1]$  given by  $\tilde{D}(p, t_0) = (\pi_{M*}, t_*)_{(p, t_0)}^{-1}(\mathbf{0}, D(t_0))$ . Alternatively,  $\tilde{D}(p, t_0) = J_{p*} D(t_0)$ , where  $J_p : [0, 1] \rightarrow M \times [0, 1]$  maps  $t_0$  to  $(p, t_0)$ .

**Lemma 5.8.1.** Any  $\omega \in \Lambda_k(M \times [0, 1])$  can be uniquely written as  $\omega = \omega_1 + dt \wedge \chi$ , where  $\omega_1$  and  $\chi$  are  $k$  and  $(k - 1)$ -forms on  $M \times [0, 1]$  respectively such that  $i(\tilde{D})\omega_1 = i(\tilde{D})\chi = 0$ .

*Proof.* Because of the isomorphism  $(\pi_{M*}, t_*)$ , it suffices to check that if a vector space  $V$  decomposes as  $W \times \mathbb{R}$ , then any  $\omega \in \Lambda_k(V)$  can be uniquely written as  $\omega_1 + \alpha \wedge \eta$ , where  $i(t)\omega_1, i(t)\eta = 0$  for  $t \in \mathbb{R}$  and  $\alpha \in \mathbb{R}^* \setminus \{0\}$ . But this is clear, since if  $\alpha_1, \dots, \alpha_n$  is a basis of  $W^*$ , then  $\alpha_1, \dots, \alpha_n, \alpha$  is a basis of  $V^*$ .  $\square$

Writing  $\omega = \omega_1 + dt \wedge \chi \in \Lambda_k(M \times [0, 1])$  as in Lemma 5.8.1, define a linear operator  $I : A_k(M \times [0, 1]) \rightarrow A_{k-1}(M)$  by

$$I\omega(\mathbf{p})(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) = \int_0^1 \chi(\mathbf{p}, t)(\iota_{t*} \mathbf{v}_1, \dots, \iota_{t*} \mathbf{v}_{k-1}) dt. \quad (5.8.1)$$

**Proposition 5.8.3.** *If  $\omega$  is a  $k$ -form on  $M \times [0, 1]$ , then  $\iota_1^* \omega - \iota_0^* \omega = d(I\omega) + I(d\omega)$ .*

*Proof.* Let  $(U, \mathbf{x})$  be a chart on  $M$ , so that  $(U \times [0, 1])(\bar{\mathbf{x}}, t)$  is a chart on  $M \times [0, 1]$ , where  $\bar{\mathbf{x}} = \mathbf{x} \circ \pi_M$ . By Lemma 5.8.1, the restriction of  $\omega$  can be written as a sum of terms of two types:

- (1)  $f d\bar{x}^{i_1} \wedge \dots \wedge d\bar{x}^{i_k}$  for some function  $f$  on  $U$ , and
- (2)  $f dt \wedge d\bar{x}^{i_1} \wedge \dots \wedge d\bar{x}^{i_{k-1}}$ .

Since  $I$  is linear, it suffices to consider the case when  $\omega$  is one of the above types.

If  $\omega$  is of type 1, then

$$d\omega = (\text{terms not involving } dt) + (\partial f / \partial t) dt \wedge d\bar{x}^{i_1} \wedge \dots \wedge d\bar{x}^{i_k}.$$

Now, for any  $t \in [0, 1]$ ,  $\pi_M \circ \iota_t = 1_M$ , so that  $\iota_t^* d\bar{x}^i = d(\bar{x}^i \circ \iota_t) = dx^i$ , and

$$\begin{aligned} I(d\omega)(\mathbf{p}) \left( \frac{\partial}{\partial x^{j_1}}, \dots, \frac{\partial}{\partial x^{j_k}} \right) &= \int_0^1 \frac{\partial f}{\partial t}(\mathbf{p}, t) (d\bar{x}^{i_1} \wedge \dots \wedge d\bar{x}^{i_k})(\mathbf{p}, t) \left( \iota_{t*} \frac{\partial}{\partial x^{j_1}}, \dots, \iota_{t*} \frac{\partial}{\partial x^{j_k}} \right) dt \\ &= \int_0^1 \frac{\partial f}{\partial t}(\mathbf{p}, t) (dx^{i_1} \wedge \dots \wedge dx^{i_k})(\mathbf{p}) \left( \frac{\partial}{\partial x^{j_1}}, \dots, \frac{\partial}{\partial x^{j_k}} \right) dt \\ &= \left( \int_0^1 \frac{\partial f}{\partial t}(\mathbf{p}, t) dt \right) (dx^{i_1} \wedge \dots \wedge dx^{i_k})(\mathbf{p}) \left( \frac{\partial}{\partial x^{j_1}}, \dots, \frac{\partial}{\partial x^{j_k}} \right). \end{aligned}$$

Thus,

$$\begin{aligned} I(d\omega)(\mathbf{p}) &= (f(\mathbf{p}, 1) - f(\mathbf{p}, 0))(dx^{i_1} \wedge \dots \wedge dx^{i_k})(\mathbf{p}) \\ &= \iota_1^* \omega(\mathbf{p}) - \iota_0^* \omega(\mathbf{p}), \end{aligned}$$

which proves the result in this case, since  $I\omega = 0$ .

Suppose now that  $\omega$  is of type 2; then  $\iota_1^* \omega = \iota_0^* \omega = 0$  because  $\iota_{t_0}^* dt = 0$  for any  $t_0$ . On

the other hand,

$$\begin{aligned} I(d\omega)(\mathbf{p}) &= I\left(\sum_{l=1}^n \frac{\partial f}{\partial \bar{x}^l} d\bar{x}^l \wedge dt \wedge d\bar{x}^{i_1} \wedge \cdots \wedge d\bar{x}^{i_k}\right)(\mathbf{p}) \\ &= I\left(-\sum_{l=1}^n \frac{\partial f}{\partial \bar{x}^l} dt \wedge d\bar{x}^l \wedge d\bar{x}^{i_1} \wedge \cdots \wedge d\bar{x}^{i_k}\right)(\mathbf{p}) \\ &= -\sum_l \left(\int_0^1 \frac{\partial f}{\partial \bar{x}^l}(\mathbf{p}, t) dt\right) (dx^l \wedge dx^{i_1} \wedge \cdots \wedge dx^{i_k})(\mathbf{p}), \end{aligned}$$

while

$$\begin{aligned} d(I\omega)(\mathbf{p}) &= d\left(\left[\int_0^1 f(\mathbf{p}, t) dt\right] dx^{i_1} \wedge \cdots \wedge dx^{i_k}\right)(\mathbf{p}) \\ &= \sum_l \frac{\partial}{\partial x^l} \left(\int_0^1 f(\mathbf{p}, t) dt\right) (dx^l \wedge dx^{i_1} \wedge \cdots \wedge dx^{i_k})(\mathbf{p}) \\ &= \sum_l \left(\int_0^1 \frac{\partial f}{\partial \bar{x}^l}(\mathbf{p}, t) dt\right) (dx^l \wedge dx^{i_1} \wedge \cdots \wedge dx^{i_k})(\mathbf{p}), \end{aligned}$$

so that  $I(d\omega) + d(I\omega) = 0$ . □

**Theorem 5.8.1.** *If  $f_0, f_1 : M \rightarrow N$  are homotopic, then the induced linear maps  $f_0^*, f_1^* : H^k(N) \rightarrow H^k(M)$  are equal for every  $k$ .*

*Proof.* Let  $\mathbf{H} : M \times [0, 1] \rightarrow N$  be a homotopy, with  $\mathbf{H} \circ \iota_j = f_j$ ,  $j = 0, 1$ . Then for any closed  $k$ -form  $\omega$  on  $N$ ,

$$f_1^* \omega - f_0^* \omega = (\iota_1^* - \iota_0^*) \mathbf{H}^* \omega = (dI + Id) \mathbf{H}^* \omega = dI \mathbf{H}^* \omega + I \mathbf{H}^* d\omega = dI \mathbf{H}^* \omega$$

is exact, so that  $f_1^* \omega + B^k(M) = f_0^* \omega + B^k(M)$ . □

**Corollary 5.8.1.** *If  $M$  is contractible, then  $H^k(M) = \{\mathbf{0}\}$  for all  $k \geq 1$ .*

*Proof.* If  $M$  is contractible, then there is a point  $\mathbf{p}_0 \in M$  such that the constant map  $\mathbf{f} : M \rightarrow M$ ,  $\mathbf{f}(\mathbf{p}) = \mathbf{p}_0$ , is homotopic to the identity map  $1_M$  of  $M$ . By Theorem 5.8.1,  $1_M^*, \mathbf{f}^* : H^k(M) \rightarrow H^k(M)$  are equal. But  $1_M^*$  is the identity map on  $H^k(M)$ , whereas  $\mathbf{f}^* \omega = 0$  for any  $\omega \in \Lambda_k(M)$  if  $k \geq 1$  because  $\mathbf{f}_* = \mathbf{0}$ . □

In particular, a compact oriented manifold of positive dimension  $n$  cannot be contractible, since its  $n$ -th cohomology is nontrivial.

**Corollary 5.8.2 (Poincaré Lemma).** *Let  $\alpha$  be a closed  $k$ -form on a manifold  $M$ ,  $k \geq 1$ . Then any  $\mathbf{p} \in M$  has a neighborhood on which the restriction of  $\alpha$  is exact.*

*Proof.* Let  $(V, \mathbf{h})$  be a local parametrization of  $M$  with  $\mathbf{h}(\mathbf{0}) = \mathbf{p}$ , and  $\varepsilon > 0$  small enough that  $B_\varepsilon(\mathbf{0}) \subset V$ . Then  $U = \mathbf{h}(B_\varepsilon(\mathbf{0}))$  is a contractible neighborhood of  $\mathbf{p}$ . If  $\iota : U \hookrightarrow M$



denotes inclusion, then the restriction  $\iota^* \alpha$  of  $\alpha$  to  $U$  is a closed form on a contractible manifold, and is therefore exact by the previous corollary.  $\square$

We emphasize again that the conclusion of the Poincaré Lemma is local: the 2-form  $\omega$  from Examples and Remarks 5.6.1 (iii) is the differential of the polar angle function  $\theta$  in a neighborhood of any point, even though it is not exact.

## 5.9 Exercises

**5.1.** Prove that a  $k$ -tensor ( $k > 1$ )  $T$  on  $V$  is alternating if and only if

$$T(\dots, \mathbf{v}, \dots, \mathbf{v}, \dots) = 0, \quad \mathbf{v} \in V.$$

**5.2.** Given a  $k$ -tensor  $T$  on  $V$ , define  $T_s$  by

$$T_s(\mathbf{v}_1, \dots, \mathbf{v}_k) = \frac{1}{k!} \sum_{\sigma \in S_k} T(\mathbf{v}_{\sigma(1)}, \dots, \mathbf{v}_{\sigma(k)}),$$

for  $\mathbf{v}_i \in V$ . Show that  $T_s$  is a symmetric  $k$ -tensor. Is it true that  $T = T_s + T_a$ ?

**5.3.** Prove Proposition 5.1.1.

**5.4.** Prove that if  $\alpha$  is alternating, then  $\alpha_a = \alpha$ , and conclude that for any  $k$ -tensor  $T$ ,  $(T_a)_a = T_a$ .

**5.5.** Show that  $\alpha \wedge \beta = (-1)^{kl} \beta \wedge \alpha$  for  $\alpha \in \Lambda_k(V)$  and  $\beta \in \Lambda_l(V)$ . *Hint:* It is enough to show this for  $\alpha = \alpha^1 \wedge \dots \wedge \alpha^k$  and  $\beta = \beta^1 \wedge \dots \wedge \beta^l$ , where  $\alpha^i, \beta^j \in V^*$ .

**5.6.** Prove that any form on  $V$  is decomposable if  $V$  has dimension  $\leq 3$ . Show that this is false if  $\dim V \geq 4$ . *Hint:* Use Proposition 5.2.2.

**5.7.** Prove that  $\alpha \in \Lambda_2(V)$  is decomposable if and only if  $\alpha \wedge \alpha = 0$ . *Hint:* Use Proposition 5.2.2.

**5.8.** Let  $V$  be an  $n$ -dimensional vector space. Show that any  $n$ -form on  $V$  is the volume form of some inner product and some orientation on  $V$ . Are the inner product and orientation unique?

**5.9.** Show that if  $W$  is a subspace of  $V$ , then the operation

$$c(\mathbf{v} + W) = (c\mathbf{v}) + W, \quad c \in \mathbb{R}, \quad \mathbf{v} \in V,$$

on the set  $V/W = \{\mathbf{v} + W \mid \mathbf{v} \in V\}$  is well defined; i.e., if  $\mathbf{u}_i + W = \mathbf{v}_i + W$ , then  $(c\mathbf{u}_i) + W = (c\mathbf{v}_i) + W$ . What is the zero vector in the space  $V/W$ ?

**5.10.** Let  $V$  be an  $n$ -dimensional oriented inner product space. This exercise constructs an isomorphism  $\star : \Lambda_k(V) \rightarrow \Lambda_{n-k}(V)$ ,  $k = 1, \dots, n-1$ , called the *Hodge star*

operator, that preserves decomposable elements. In particular, it shows that every element in  $\Lambda_{n-1}(V) \cong \Lambda_1(V) = V^*$  is decomposable.

Recall that a linear transformation is entirely determined by its values on a basis. So let  $\mathbf{v}_1, \dots, \mathbf{v}_n$  be a positively oriented orthonormal basis,  $\alpha^1, \dots, \alpha^n$  its dual basis, and  $\omega = \alpha^1 \wedge \dots \wedge \alpha^n \in \Lambda_n(V)$  the volume form of  $V$ .

(a) Show that for any basis element

$$\alpha \in \mathcal{B} = \{\alpha^{i_1} \wedge \dots \wedge \alpha^{i_k} \mid 1 \leq i_1 < \dots < i_k \leq n\}$$

of the induced basis  $\mathcal{B}$  of  $\Lambda_k(V)$ , there exists a unique  $\beta \in \Lambda_{n-k}(V)$  such that  $\alpha \wedge \beta = \omega$ . Define  $\star\alpha = \beta$  and extend linearly to all  $\Lambda_k(V)$ . *Hint:*  $\beta$  or  $-\beta$  belongs to the induced basis of  $\Lambda_{n-k}(V)$ .

(b) Prove that  $\star \circ \star = (-1)^{k(n-k)} 1_{\Lambda_k(V)}$ . In particular,  $\star$  is an isomorphism.

(c) Prove that if  $\alpha$  is decomposable, then so is  $\star\alpha$ . *Hint:* a decomposable  $\alpha$  corresponds to a  $k$ -dimensional subspace  $W$  of  $V$ . Consider its orthogonal complement.

(d) Show that  $\star$  does not depend on the particular choice of basis.

**5.11.** Let  $V$  be a vector space with basis  $\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3\}$ . According to Exercise 5.10, every element of  $\Lambda_2(V)$  is decomposable. If  $\alpha_i$  is the basis dual to  $\mathbf{v}_i$ , write  $\alpha_1 \wedge \alpha_2 + \alpha_2 \wedge \alpha_3 + \alpha_1 \wedge \alpha_3$  as a decomposable element.

**5.12.** Let  $V$  be an oriented  $n$ -dimensional inner product space with volume form  $\omega$  and musical isomorphisms  $\sharp, \flat$ . Given  $\mathbf{v}_1, \dots, \mathbf{v}_{n-1} \in V$ , define their *cross product* to be the vector

$$\mathbf{v}_1 \times \dots \times \mathbf{v}_{n-1} := \left( \star(\mathbf{v}_1^\flat \wedge \dots \wedge \mathbf{v}_{n-1}^\flat) \right)^\sharp \in V.$$

Here,  $\star$  is the Hodge star operator from Exercise 5.10.

(a) Show that  $\langle \mathbf{v}_1 \times \dots \times \mathbf{v}_{n-1}, \mathbf{u} \rangle = \omega(\mathbf{v}_1, \dots, \mathbf{v}_{n-1}, \mathbf{u})$ ,  $\mathbf{u} \in V$ .

(b) Prove that when  $V = \mathbb{R}^3$  with the standard inner product and orientation, this definition coincides with the usual cross product.

**5.13.** Recall that the musical isomorphisms in  $\mathbb{R}^n$  extend naturally to each tangent space and its dual, and thus establish a bijective correspondence between vector fields and 1-forms. The same is true for vector fields and one-forms on a manifold. Similarly, if  $M^n$  is an oriented manifold, then the Hodge star operator from Exercise 5.10 on each tangent space induces a bijective correspondence  $\star : \Lambda_k(M) \rightarrow \Lambda_{n-k}(M)$ .

(a) The *gradient* of  $f : M \rightarrow \mathbb{R}$  is the vector field  $\nabla f$  on  $M$  given by  $(df)^\sharp$ . Show that if  $M = \mathbb{R}^n$ , this coincides with the usual definition of gradient.

(b) Given a vector field  $\mathbf{X}$  on an oriented manifold  $M$ , the *divergence* of  $\mathbf{X}$  is the function  $\operatorname{div} \mathbf{X} = \star d \star \mathbf{X}^\flat$ . Prove that when  $M = \mathbb{R}^n$  and  $\mathbf{X} = \sum_i X^i \mathbf{D}_i$ , then

$$\operatorname{div} \mathbf{X} = \sum_j D_j(X^j).$$

- (c) Suppose  $\mathbf{X}$  is a vector field on  $\mathbb{R}^3$ . The *curl* of  $\mathbf{X}$  is the vector field  $\text{curl } \mathbf{X} = (*d\mathbf{X}^\flat)^\sharp$ . Show that if  $X^i = \langle \mathbf{X}, \mathbf{D}_i \rangle$ , then

$$\text{curl } \mathbf{X} = (D_2X^3 - D_3X^2)\mathbf{D}_1 + (D_3X^1 - D_1X^3)\mathbf{D}_2 + (D_1X^2 - D_2X^1)\mathbf{D}_3.$$

**5.14.** Let  $\mathbf{X}, \mathbf{Y}$  denote vector fields on  $\mathbb{R}^3$ ,  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ . Prove the following identities:

- $\text{curl}(\mathbf{X} + \mathbf{Y}) = \text{curl}(\mathbf{X}) + \text{curl}(\mathbf{Y})$  (see previous exercise for the definition of curl).
- $\text{curl}(f\mathbf{X}) = f \text{curl } \mathbf{X} + \nabla f \times \mathbf{X}$ .
- $\text{div}(f\mathbf{X}) = f \text{div } \mathbf{X} + \langle \nabla f, \mathbf{X} \rangle$ .
- $\text{div}(\mathbf{X} \times \mathbf{Y}) = \langle \text{curl } \mathbf{X}, \mathbf{Y} \rangle - \langle \text{curl } \mathbf{Y}, \mathbf{X} \rangle$ .
- $\text{div } \text{curl } \mathbf{X} = 0$ .
- $\text{curl } \nabla f = \mathbf{0}$ .

**5.15.** The next two problems examine in more detail the concept of divergence; they show, in particular, that a vector field  $\mathbf{X}$  has vanishing divergence if and only if its flow  $\Phi_t$  preserves the volume form  $\omega$ ; i.e.,  $\Phi_t^* \omega = \omega$  for all  $t$ . For simplicity, we limit ourselves to the case  $M = \mathbb{R}^n$ .

The *Lie derivative of a  $k$ -form  $\alpha$  with respect to a vector field  $\mathbf{X}$*  is the  $k$ -form  $L_{\mathbf{X}}\alpha$  defined by

$$L_{\mathbf{X}}\alpha(\mathbf{p}) = \lim_{t \rightarrow 0} \frac{1}{t} (\Phi_t^* \omega(\mathbf{p}) - \omega(\mathbf{p})),$$

where  $\Phi_t$  denotes the flow of  $\mathbf{X}$ .

This problem shows that the divergence of  $\mathbf{X}$  satisfies  $L_{\mathbf{X}}\omega = (\text{div } \mathbf{X})\omega$ , where  $\omega$  is the volume form of  $\mathbb{R}^n$ . The next problem deals with the statement about vanishing divergence.

- (a) Show that  $\Phi_t^* \omega = (\det D\Phi_t)\omega$ . Thus,  $L_{\mathbf{X}}\omega(\mathbf{p}) = f'(0)\omega(\mathbf{p})$ , where

$$f(t) = \det M(t), \quad M(t) = D\Phi_t(\mathbf{p}).$$

- (b) Use the chain rule to prove that  $f'(t) = \sum_i \langle M'_i(t), \mathbf{e}_i \rangle$ , where  $M_i$  denotes the  $i$ -th column of  $M$ .
- (c) Let  $\Phi : \mathbb{R}^{n+1} \supset U \rightarrow \mathbb{R}^n$  be given by  $\Phi(t, \mathbf{q}) = \Phi_t(\mathbf{q})$ . Show that

$$M'_i(t) = D(D\Phi(t, \mathbf{p})\mathbf{e}_i)\mathbf{e}_1 = D(D\Phi(t, \mathbf{p})\mathbf{e}_1)\mathbf{e}_i,$$

so that

$$M'_i(t) = D(\pi_2(\mathbf{X} \circ \Phi(t, \mathbf{p})))\mathbf{e}_i,$$

where  $\pi_2 : T\mathbb{R}^n = \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n$  denotes projection. Conclude that  $L_{\mathbf{X}}\omega = (\text{div } \mathbf{X})\omega$ .

**5.16.** This problem uses notation and results from the previous one.

- Show that if the volume form  $\omega$  of  $\mathbb{R}^n$  is invariant under the flow  $\Phi_t$  of a vector field  $\mathbf{X}$  (in the sense that  $\Phi_t^* \omega = \omega$ ), then the divergence of  $\mathbf{X}$  is identically zero.
- Conversely, suppose that a vector field with flow  $\Phi_t$  has zero divergence. Given  $\mathbf{p} \in \mathbb{R}^n$ , consider as before  $f(t) = \det(D\Phi_t)(\mathbf{p})$ , so that  $f'(0) = 0$ . Prove that  $f'(t) = 0$

for all  $t$ , and conclude that  $\Phi_t^* \omega = \omega$ . *Hint:* Use an argument similar to that in the last part of the proof of Theorem 2.9.4.

**5.17.** The spectral theorem describes a canonical form for self-adjoint operators on an inner product space  $V$ . This problem describes one for skew-adjoint operators; i.e., operators  $L$  satisfying  $\langle L\mathbf{u}, \mathbf{v} \rangle = -\langle L\mathbf{v}, \mathbf{u} \rangle$  for any  $\mathbf{u}, \mathbf{v} \in V$ .

Let  $V$  be an inner product space, and  $\mathfrak{o}(V)$  the space of skew-adjoint operators on  $V$ .

- (a) Prove that the map  $\mathfrak{o}(V) \rightarrow \Lambda_2(V)$  that assigns to  $L \in \mathfrak{o}(V)$  the element  $\alpha \in \Lambda_2(V)$  given by  $\alpha(\mathbf{v}, \mathbf{w}) = \langle L\mathbf{v}, \mathbf{w} \rangle$ , is an isomorphism.
- (b) Use Proposition 5.2.2 to show that for any  $L \in \mathfrak{o}(V)$ , there exists an orthonormal basis in which the matrix of  $L$  has the form

$$\begin{bmatrix} 0 & -\lambda_1 & 0 & 0 & \cdots \\ \lambda_1 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 0 & -\lambda_2 & \cdots \\ 0 & 0 & \lambda_2 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

with the last row and column consisting of zeros if the dimension of  $V$  is odd.

**5.18.** An (orthonormal) *moving frame* on  $\mathbb{R}^n$  is an orthonormal basis of vector fields  $\mathbf{X}_1, \dots, \mathbf{X}_n$  on  $\mathbb{R}^n$ . The term “basis” is of course an abuse of notation; what is really meant is that when evaluated at any point, these vector fields form a basis of the tangent space at that point. Moving frames were introduced by E. Cartan to study the geometry of manifolds.

Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  denote a moving frame on  $\mathbb{R}^n$ . Given  $\mathbf{p} \in \mathbb{R}^n$  and  $1 \leq i, j \leq n$ , the map

$$\begin{aligned} \mathbb{R}_\mathbf{p}^n &\longrightarrow \mathbb{R} \\ \mathbf{u} &\longmapsto \langle \nabla_\mathbf{u} \mathbf{X}_i, \mathbf{X}_j(\mathbf{p}) \rangle \end{aligned}$$

is linear, and thus defines a one-form  $\omega_i^j(\mathbf{p})$  on  $\mathbb{R}_\mathbf{p}^n$ . The  $n^2$  forms  $\omega_i^j$  are called the *connection forms* of  $\mathbb{R}^n$  in the moving frame.

- (a) Show that  $\omega_i^j = -\omega_j^i$ .
- (b) Let  $X_i^j = \langle \mathbf{X}_i, \mathbf{D}_j \rangle$ . By definition,

$$\nabla_\mathbf{X} \mathbf{X}_i = \sum_k \omega_i^k(\mathbf{X}) \mathbf{X}_k = \sum_{k,j} \omega_i^k(\mathbf{X}) X_k^j \mathbf{D}_j.$$

Use this to prove that  $dX_i^j = \sum_k \omega_i^k X_k^j$ .

- (c) Let  $\omega_1, \dots, \omega_n$  denote the one-forms that are dual to the moving frame:  $\omega_i = \mathbf{X}_i^\flat$ , so that  $\omega_i = \sum_j X_i^j du^j$ . Prove the *structural equations of  $\mathbb{R}^n$* :

$$d\omega_i = -\sum_k \omega_i^k \wedge \omega_k \tag{5.9.1}$$

$$d\omega_i^j = \sum_k \omega_i^k \wedge \omega_k^j. \tag{5.9.2}$$

**5.19.** A Riemannian manifold is a pair  $(M, \mathbf{g})$ , where  $M$  is a manifold and  $\mathbf{g}$  a Riemannian metric on  $M$ , cf. Example 5.1.1. Although we have mostly been using the metric induced from the ambient Euclidean space, there are many other interesting examples: given  $\kappa \in \mathbb{R}$ , denote by  $U_\kappa$  the open set

$$U_\kappa = \begin{cases} \mathbb{R}^n, & \text{if } \kappa \geq 0; \\ \{\mathbf{a} \in \mathbb{R}^n \mid |\mathbf{a}|^2 < \frac{1}{|\kappa|}\}, & \text{if } \kappa < 0. \end{cases}$$

Show that the tensor field  $\mathbf{g}_\kappa$  on  $U_\kappa$ , where

$$\mathbf{g}_\kappa(\mathbf{a}) = \frac{4}{(1 + \kappa|\mathbf{a}|^2)^2} \sum_{i=1}^n du^i \otimes du^i, \quad \mathbf{a} \in U_\kappa,$$

is a Riemannian metric.

Notice that when  $\kappa = 0$ , we recover standard Euclidean space, which is flat. In general, one can define the curvature tensor of a Riemannian metric in the same way we did for the metric induced from Euclidean space. The same is true for isometries, and as in Euclidean space, the curvature tensor is preserved under isometries. The next exercise shows that when  $\kappa > 0$ ,  $(U_\kappa, \mathbf{g}_\kappa)$  has constant curvature  $\kappa$ . This turns out to be true also when  $\kappa < 0$ .

**5.20.** If  $(M_1, \mathbf{g}_1)$  and  $(M_2, \mathbf{g}_2)$  are Riemannian manifolds (see Exercise 5.19), an isometry between them is a diffeomorphism  $f : M_1 \rightarrow M_2$  such that  $f^* \mathbf{g}_2 = \mathbf{g}_1$ ; i.e.,

$$\mathbf{g}_2(f_{*p} \mathbf{u}, f_{*p} \mathbf{v}) = \mathbf{g}_1(\mathbf{u}, \mathbf{v}), \quad p \in M_1, \quad \mathbf{u}, \mathbf{v} \in M_{1p}.$$

When this is the case, the two Riemannian manifolds are said to be isometric. Notice that when the metrics are those induced by the ambient Euclidean space, this definition coincides with the usual one.

Prove that when  $\kappa > 0$ , the Riemannian manifold  $(U_\kappa, \mathbf{g}_\kappa)$  from Exercise 5.19 is isometric to a sphere of radius  $1/\sqrt{\kappa}$  with a point removed (and therefore with constant curvature  $\kappa$ ) in  $\mathbb{R}^{n+1}$  with the usual metric. *Hint:* Consider stereographic projection.

**5.21.** In Examples and Remarks 5.4.1 (ii), we defined the integral of a 1-form  $\omega \in \Lambda_k(\mathbb{R}^n)$  along a curve  $\mathbf{c} : I \rightarrow \mathbb{R}^n$  to be

$$\int_{\mathbf{c}} \omega = \int_I \mathbf{c}^* \omega.$$

Let  $\omega$  be a 1-form on a connected open set  $U$  in  $\mathbb{R}^n$ , and consider the following three statements:

- (1)  $\omega$  is exact in  $U$ .
- (2) If  $\mathbf{c}_1$  and  $\mathbf{c}_2$  are curves in  $U$  with the same beginning and endpoints, then  $\int_{\mathbf{c}_1} \omega = \int_{\mathbf{c}_2} \omega$ .
- (3) If  $\mathbf{c}$  is a closed curve (i.e., beginning and endpoints coincide) in  $U$ , then  $\int_{\mathbf{c}} \omega = 0$ .

- (a) Prove that 1 implies 2.  
 (b) Conclude that all three statements are equivalent. *Hint:* to show that 2 implies 1, fix any  $\mathbf{p} \in U$ , and define a function  $f$  on  $U$  by  $f(\mathbf{a}) = \int_{\mathbf{c}} \omega$ , where  $\mathbf{c}$  is any curve in  $U$  beginning at  $\mathbf{p}$  and ending at  $\mathbf{a}$ . Show that  $\omega = df$ .

**5.22.** Recall that the 0-th cohomology space of  $M$  is

$$H^0(M) = Z^0(M) = \{f : M \rightarrow \mathbb{R} \mid df = 0\}.$$

Prove that if  $M$  is connected, then  $H^0(M) \cong \mathbb{R}$ .

**5.23.** Let  $\mathbf{a} \in \mathbb{R}^n$ ,  $r > 0$ , and  $M$  be the  $n$ -dimensional manifold with boundary  $\overline{B_r(\mathbf{a})}$ . Show that every closed  $k$ -form on  $M$  is exact.

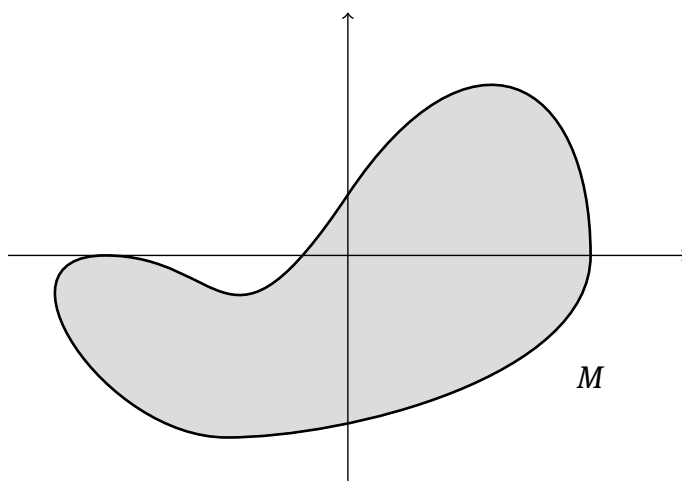
**5.24.** Give examples of manifolds with boundary where the manifold boundary does not coincide with the topological boundary of the set  $M$  from Chapter 1. Is one of these boundaries always contained in the other?

**5.25.** Let  $M, N$  be open sets in  $\mathbb{R}^n$  whose topological boundaries are smooth  $(n - 1)$ -dimensional manifolds, so that  $\bar{M}$  and  $\bar{N}$  are manifolds with boundaries  $\partial M$  and  $\partial N$  respectively.

- (a) Show that if  $\bar{N} \subset M$ , then  $\int_{\partial N} \omega = \int_{\partial M} \omega$  for any closed  $(n - 1)$ -form  $\omega$ . (It should really be the pullback of  $\omega$  by the respective inclusion maps, of course). *Hint:* apply Stokes' theorem to  $\bar{M} \setminus N$ .  
 (b) Show that if

$$\omega = \frac{-u^2}{(u^1)^2 + (u^2)^2} du^1 + \frac{u^1}{(u^1)^2 + (u^2)^2} du^2,$$

then  $\int_{\partial M} \omega = 2\pi$ , with  $M$  the region pictured below. *Hint:* see Examples and Remarks 5.6.1 (iii).



- (c) Prove that if  $M$  is the region above but translated so that it no longer contains the origin, then  $\int_{\partial M} \omega = 0$ .

**5.26.** The higher-dimensional version of the 1-form  $\omega$  on  $\mathbb{R}^2 \setminus \{\mathbf{0}\}$  from part (b) in the previous problem is the 2-form

$$\Omega = \frac{1}{((u^1)^2 + (u^2)^2 + (u^3)^2)} (u^1 du^2 \wedge du^3 + u^2 du^3 \wedge du^1 + u^3 du^1 \wedge du^2)$$

on  $\mathbb{R}^3 \setminus \{\mathbf{0}\}$ .

- (a) Show that  $\Omega$  is closed.  
 (b) Extend the concepts of determinant and cross product to the tangent space of  $\mathbb{R}^3$  at any point  $\mathbf{p}$  by means of the canonical isomorphism  $\mathcal{I}_{\mathbf{p}} : \mathbb{R}^3 \rightarrow \mathbb{R}_{\mathbf{p}}^3$ . Show that for  $\alpha = u^1 du^2 \wedge du^3 + u^2 du^3 \wedge du^1 + u^3 du^1 \wedge du^2$ ,

$$\alpha(\mathbf{p})(\mathbf{u}, \mathbf{v}) = \det [\mathbf{P}(\mathbf{p}) \quad \mathbf{u} \quad \mathbf{v}] = \langle \mathbf{P}(\mathbf{p}), \mathbf{u} \times \mathbf{v} \rangle, \quad \mathbf{p} \in \mathbb{R}^3 \setminus \{\mathbf{0}\}, \quad \mathbf{u}, \mathbf{v} \in \mathbb{R}_{\mathbf{p}}^3,$$

where  $\mathbf{P} = \sum_i u^i \mathbf{D}_i$  is the position vector field.

- (c) Conclude that if  $\mu$  is the volume form of  $S^2(r)$  and  $\iota : S^2(r) \hookrightarrow \mathbb{R}^3 \setminus \{\mathbf{0}\}$  denotes the inclusion map, then

$$\mu = r^2 \iota^* \Omega, \text{ so that } \int_{S^2(r)} \iota^* \Omega = 4\pi.$$

Why does this imply that  $\Omega$  is not exact?

- (d) Let  $M$  be a bounded, open set in  $\mathbb{R}^3$  whose topological boundary  $\partial M$  is a smooth manifold. Show that the integral of  $\Omega$  over  $\partial M$  equals zero if  $\bar{M}$  does not contain the origin, and  $4\pi$  otherwise.

**5.27.** Let  $\alpha$  be the 1-form on  $\mathbb{R}^3$  given by

$$\alpha = \ln((u^2)^2 + 1) du^1 + u^1 u^2 u^3 du^2 - (u^2)^2 e^{u^1 u^2} du^3,$$

and  $M$  the manifold with boundary consisting of the upper hemisphere of the unit sphere centered at the origin, oriented so that the position vector field is the positive unit normal field. If  $\iota : M \hookrightarrow \mathbb{R}^3$  denotes inclusion, determine  $\int_M \iota^* \alpha$ . *Hint:* Direct computation is difficult, and applying Stokes' theorem isn't easy either. The integral of  $\alpha$  along the boundary  $\partial M$  may, however, be computed by applying Stokes' theorem to any manifold that has  $\partial M$  as boundary.

**5.28.** Suppose that  $a$  is a regular value of  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ , and that  $M = f^{-1}(a) \neq \emptyset$ , so that  $M$  is a 2-dimensional manifold. Let  $\alpha$  be the 2-form on  $\mathbb{R}^3$  given by

$$\alpha = D_1 f du^2 \wedge du^3 + D_2 f du^3 \wedge du^1 + D_3 f du^1 \wedge du^2.$$

Show that if  $\iota : M \hookrightarrow \mathbb{R}^3$  is inclusion, then  $\iota^* \alpha$  is a nowhere-zero 2-form on  $M$ , so that the latter is orientable. In fact, prove that  $(1/|\nabla f|) \iota^* \alpha$  is the volume form of this orientation. *Hint:* the 2-form  $\alpha$  from Exercise 5.26 (b) is the special case when  $f(\mathbf{x}) = (1/2)|\mathbf{x}|^2$ .

**5.29.** The divergence of a vector field on  $\mathbb{R}^n$  was defined in Exercise 5.13 for arbitrary  $n$ . Formulate and prove a divergence theorem in  $\mathbb{R}^n$ .

**5.30.** The *Laplacian* of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is the function  $\Delta f = \operatorname{div} \nabla f$ .  $f$  is said to be *harmonic* if its Laplacian is identically zero.

Let  $M$  be a connected, compact 3-dimensional manifold with boundary in  $\mathbb{R}^3$ ,  $f : M \rightarrow \mathbb{R}$  a harmonic function.

(a) Show that  $\int_{\partial M} \langle \nabla f, \mathbf{N} \rangle \eta = 0$ .

(b) Prove that

$$\int_{\partial M} \langle f \nabla f, \mathbf{N} \rangle \eta = \int_M |\nabla f|^2 du^1 \wedge du^2 \wedge du^3.$$

Conclude that the only harmonic functions  $f$  on  $M$  that satisfy

$$\int_{\partial M} \langle f \nabla f, \mathbf{N} \rangle \eta = 0$$

are the constant ones.

**5.31.** Let  $f, g : \mathbb{R}^3 \rightarrow \mathbb{R}$  be differentiable, and  $M \subset \mathbb{R}^3$  a compact 3-dimensional manifold with boundary. Prove *Green's identities*:

$$\int_M (\langle \nabla f, \nabla g \rangle + f \Delta g) du^1 \wedge du^2 \wedge du^3 = \int_{\partial M} f \langle \nabla g, \mathbf{N} \rangle \eta;$$

$$\int_M (f \Delta g - g \Delta f) du^1 \wedge du^2 \wedge du^3 = \int_{\partial M} \langle f \nabla g - g \nabla f, \mathbf{N} \rangle \eta.$$

**5.32.** (a) Prove that  $g : \mathbb{R}^3 \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}$ , where  $g(\mathbf{p}) = 1/|\mathbf{p}|$ , is harmonic.

(b) Suppose  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  is harmonic on some open set  $U$ , and consider  $\mathbf{p}_0, r > 0$  such that  $B_r(\mathbf{p}_0) \subset U$ . Show that for any  $0 < r' < r$ ,

$$\frac{1}{4\pi r'^2} \int_{\partial B_{r'}(\mathbf{p}_0)} f \eta = \frac{1}{4\pi r^2} \int_{\partial B_r(\mathbf{p}_0)} f \eta,$$

by applying Green's second identity (Exercise 5.31) to  $f$  and  $g$  on the manifold with boundary  $\overline{B_r(\mathbf{p}_0)} \setminus B_{r'}(\mathbf{p}_0)$ . Conclude that

$$f(\mathbf{p}_0) = \frac{1}{4\pi r^2} \int_{\partial B_r(\mathbf{p}_0)} f \eta.$$

This identity is known as the *mean value theorem for harmonic functions*.

(c) Prove the *maximum principle for harmonic functions*: If  $f$  is harmonic on a closed and bounded region  $K \subset \mathbb{R}^3$ , then the maximum and minimum of  $f$  occur on the (topological) boundary of  $K$ . *Hint*: We may assume  $f$  is not constant. If, say, the maximum were to occur at an interior point  $\mathbf{p}_0$ , this would contradict the mean value theorem.



**5.33.** Let  $\mathbf{X}$  be a vector field on  $\mathbb{R}^3$ , and  $\omega = \mathbf{X}^\flat$  the dual 1-form. Given  $\mathbf{a} \in \mathbb{R}^3$ , consider any plane  $P$  through  $\mathbf{a}$  and a unit vector  $\mathbf{n}$  perpendicular to  $P$ . As usual, we identify  $\mathbf{n}$  with a vector in the tangent space at  $\mathbf{a}$ . Show that if  $S_r$  denotes the circle of radius  $r$  centered at  $\mathbf{a}$  that lies in  $P$ , then

$$\langle \text{curl } \mathbf{X}(\mathbf{a}), \mathbf{n} \rangle = \lim_{r \rightarrow 0} \frac{1}{\pi r^2} \int_{S_r} \iota^* \omega.$$

Here  $\iota$  denotes the inclusion map of the circle in the disk of radius  $r$  in  $P$ , and the disk is oriented so that  $\mathbf{n}$  is the positive unit normal.

If  $\mathbf{X}$  represents the velocity field of a fluid in  $\mathbb{R}^3$ , interpret the above identity in terms of rotation of the fluid about an axis. Which axis direction measures maximal rotation? A fluid with zero curl everywhere is said to be *irrotational*.



## 6 Manifolds as metric spaces

In this chapter, we endow any connected manifold  $M$  with a distance function (see Section 1.8). This distance is not the restriction of the standard distance in the ambient Euclidean space, but rather the length of the shortest path in the space between two points, provided such a path exists. It turns out that the open sets this distance generates coincide with the open sets of  $M$  induced by the ambient space: i.e., those of the form  $U \cap M$ , where  $U$  is open in the ambient space. We begin by taking a more detailed look at geodesics, since they represent the kind of curves having this property.

### 6.1 Extremal properties of geodesics

One aim of this section is to establish the fact that geodesics are locally length-minimizing; i.e., if  $\mathbf{c}$  is a geodesic that is short enough, then it is (modulo reparametrization) the shortest curve between its endpoints. This is clearly not true in general for arbitrarily long geodesics, since for example on a sphere, any geodesic only minimizes up to the antipodal point. In order to state this more precisely, we need to introduce the following concept: If  $\mathbf{c} : [0, a] \rightarrow M$  is a curve in  $M \subset \mathbb{R}^n$ , a *variation of  $\mathbf{c}$*  is a map  $\mathbf{V} : [0, a] \times [-\varepsilon, \varepsilon] \rightarrow M$  such that if  $\mathbf{V}_s : [0, a] \rightarrow M$  denotes the curve  $\mathbf{V}_s(t) = \mathbf{V}(t, s)$  for  $|s| \leq \varepsilon$ , then  $\mathbf{V}_0 = \mathbf{c}$ .

Notice that  $\nabla_{\mathbf{D}} \dot{\mathbf{V}}_s(t) = \nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_1(t, s)$ . In fact, if  $\iota_s : \mathbb{R} \rightarrow \mathbb{R}^2$  is given by  $\iota_s(t) = (t, s)$ , then  $\mathbf{V}_s = \mathbf{V} \circ \iota_s$ , so that

$$\dot{\mathbf{V}}_s = \mathbf{V}_* \iota_{s*} \mathbf{D} = \mathbf{V}_* \mathbf{D}_1 \circ \iota_s,$$

and

$$\nabla_{\mathbf{D}} \dot{\mathbf{V}}_s = \nabla_{\mathbf{D}} (\mathbf{V}_* \mathbf{D}_1 \circ \iota_s) = \nabla_{\iota_{s*} \mathbf{D}} \mathbf{V}_* \mathbf{D}_1 = \nabla_{\mathbf{D}_1 \circ \iota_s} \mathbf{V}_* \mathbf{D}_1 = (\nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_1) \circ \iota_s.$$

**Lemma 6.1.1.** *If  $\mathbf{V}$  is a variation, then  $\nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_2 = \nabla_{\mathbf{D}_2} \mathbf{V}_* \mathbf{D}_1$ .*

*Proof.* This is, after deciphering notation, just the fact that mixed partial derivatives are equal. First of all, it is enough to show that  $D_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_2 = D_{\mathbf{D}_2} \mathbf{V}_* \mathbf{D}_1$ . Now,

$$\mathbf{V}_* \mathbf{D}_2(t_0, s_0) = \mathcal{I}_{\mathbf{V}(t_0, s_0)} D\mathbf{V}(t_0, s_0) \mathbf{e}_2 = \mathcal{I}_{\mathbf{V}(t_0, s_0)} \begin{bmatrix} D_2 \mathbf{V}^1 \\ \vdots \\ D_2 \mathbf{V}^n \end{bmatrix} (t_0, s_0),$$

so that

$$D_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_2 = \mathcal{I}_{\mathbf{V}} \begin{bmatrix} D_1 D_2 \mathbf{V}^1 \\ \vdots \\ D_1 D_2 \mathbf{V}^n \end{bmatrix} = \mathcal{I}_{\mathbf{V}} \begin{bmatrix} D_2 D_1 \mathbf{V}^1 \\ \vdots \\ D_2 D_1 \mathbf{V}^n \end{bmatrix} = D_{\mathbf{D}_2} \mathbf{V}_* \mathbf{D}_1. \quad \square$$

Recall that we defined the inner product on the tangent space of  $\mathbb{R}^n$  at a point  $\mathbf{p}$  to be the one for which the canonical isomorphism  $\mathcal{I}_{\mathbf{p}} : \mathbb{R}^n \rightarrow \mathbb{R}_{\mathbf{p}}^n$  becomes a linear isometry.

Given a manifold  $M$  and  $\mathbf{p} \in M$ , the tangent space  $M_{\mathbf{p}}$  is also an inner product space, and we can, in the same way, define for any  $\mathbf{v} \in M_{\mathbf{p}}$  a *canonical inner product on  $(M_{\mathbf{p}})_{\mathbf{v}}$*  by requiring that  $\mathcal{I}_{\mathbf{v}} : M_{\mathbf{p}} \rightarrow (M_{\mathbf{p}})_{\mathbf{v}}$  be a linear isometry. The following result may then be interpreted as saying that the derivative of the exponential map at a point is “radially” isometric. Its kernel is orthogonal to the radial direction.

**Lemma 6.1.2** (The Gauss lemma). *Let  $\mathbf{p} \in M$ ,  $\mathbf{v} \in \tilde{M}_{\mathbf{p}} \subset M_{\mathbf{p}}$  in the domain  $\tilde{M}_{\mathbf{p}}$  of  $\exp_{\mathbf{p}}$ . Consider the ray  $\boldsymbol{\varphi}_{\mathbf{v}} : \mathbb{R} \rightarrow M_{\mathbf{p}}$  in direction  $\mathbf{v}$ ,  $\boldsymbol{\varphi}_{\mathbf{v}}(t) = t\mathbf{v}$ , and let  $\mathbf{u} = \dot{\boldsymbol{\varphi}}_{\mathbf{v}}(1) \in (M_{\mathbf{p}})_{\mathbf{v}}$ . If  $\mathbf{w} \in (M_{\mathbf{p}})_{\mathbf{v}}$ , then*

$$\langle \exp_{\mathbf{p}*} \mathbf{u}, \exp_{\mathbf{p}*} \mathbf{w} \rangle = \langle \mathbf{u}, \mathbf{w} \rangle.$$

*Proof.* We may assume  $\mathbf{v} \neq \mathbf{0}$ , for otherwise  $\mathbf{u} = \mathbf{0}$ , and the statement is trivial. Consider the variation

$$\begin{aligned} \mathbf{V} : [0, 1] \times [-\varepsilon, \varepsilon] &\rightarrow M, \\ (t, s) &\mapsto \exp_{\mathbf{p}}(t(\mathbf{v} + s\mathcal{I}_{\mathbf{v}}^{-1}\mathbf{w})) \end{aligned}$$

of  $\mathbf{c}_{\mathbf{v}}$ , where  $\varepsilon$  is chosen small enough that  $\mathbf{V}$  is defined. Each  $\mathbf{V}_s$  is a geodesic, so that

$$(\nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_1)(t, s) = (\nabla_{\mathbf{D}} \dot{\mathbf{V}}_s)(t) = 0.$$

Furthermore, the tangent vector field of  $\mathbf{V}_s$  has constant norm equal to

$$|\dot{\mathbf{V}}_s(0)| = |\mathbf{V}_* \mathbf{D}_1(0, s)| = |\mathbf{v} + s\mathcal{I}_{\mathbf{v}}^{-1}\mathbf{w}|.$$

Since the partial derivative of  $\langle \mathbf{V}_* \mathbf{D}_1, \mathbf{V}_* \mathbf{D}_2 \rangle$  with respect to  $t$  equals

$$\begin{aligned} D_1 \langle \mathbf{V}_* \mathbf{D}_1, \mathbf{V}_* \mathbf{D}_2 \rangle &= \langle \nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_1, \mathbf{V}_* \mathbf{D}_2 \rangle + \langle \mathbf{V}_* \mathbf{D}_1, \nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_2 \rangle \\ &= \langle \mathbf{V}_* \mathbf{D}_1, \nabla_{\mathbf{D}_1} \mathbf{V}_* \mathbf{D}_2 \rangle = \langle \mathbf{V}_* \mathbf{D}_1, \nabla_{\mathbf{D}_2} \mathbf{V}_* \mathbf{D}_1 \rangle \\ &= \frac{1}{2} D_2 \langle \mathbf{V}_* \mathbf{D}_1, \mathbf{V}_* \mathbf{D}_1 \rangle, \end{aligned}$$

we deduce that

$$\begin{aligned} D_1 \langle \mathbf{V}_* \mathbf{D}_1, \mathbf{V}_* \mathbf{D}_2 \rangle(t, s) &= \frac{1}{2} D_2 (|\mathbf{v}|^2 + 2s\langle \mathbf{v}, \mathcal{I}_{\mathbf{v}}^{-1}\mathbf{w} \rangle + s^2|\mathcal{I}_{\mathbf{v}}^{-1}\mathbf{w}|) \\ &= \langle \mathbf{v}, \mathcal{I}_{\mathbf{v}}^{-1}\mathbf{w} \rangle + 2s|\mathcal{I}_{\mathbf{v}}^{-1}\mathbf{w}|, \end{aligned}$$

which becomes  $\langle \mathbf{v}, \mathcal{I}_{\mathbf{v}}^{-1}\mathbf{w} \rangle$  when  $s = 0$ . Now,

$$\langle \mathbf{v}, \mathcal{I}_{\mathbf{v}}^{-1}\mathbf{w} \rangle = \langle \mathcal{I}_{\mathbf{v}}\mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{u}, \mathbf{w} \rangle,$$

so that the function

$$t \mapsto \langle \mathbf{V}_* \mathbf{D}_1, \mathbf{V}_* \mathbf{D}_2 \rangle(t, 0)$$

has constant derivative  $\langle \mathbf{u}, \mathbf{w} \rangle$  on  $[0, 1]$ . But it vanishes at the origin (because  $\mathbf{V}_* \mathbf{D}_2(0, 0) = \mathbf{0}$ ), so that

$$\langle \mathbf{V}_* \mathbf{D}_1, \mathbf{V}_* \mathbf{D}_2 \rangle(t, 0) = t\langle \mathbf{u}, \mathbf{w} \rangle. \quad (6.1.1)$$

On the other hand, the very definition of  $\mathbf{V}$  implies that

$$\langle \mathbf{V}_* \mathbf{D}_1, \mathbf{V}_* \mathbf{D}_2 \rangle(1, 0) = \langle \exp_{\mathbf{p}^*} \mathbf{u}, \exp_{\mathbf{p}^*} \mathbf{w} \rangle.$$

Comparing this expression with (6.1.1) now yields the statement.  $\square$

Consider a curve  $\gamma : I \rightarrow M_{\mathbf{p}}$ . If  $\gamma(t) \neq \mathbf{0}$ , we may write  $\dot{\gamma}(t) = \dot{\gamma}_r(t) + \dot{\gamma}_\theta(t)$ , with

$$\dot{\gamma}_r = \frac{1}{|\gamma|^2} \langle \dot{\gamma}, \mathcal{I}_\gamma \gamma \rangle \mathcal{I}_\gamma \gamma \quad \text{and} \quad \dot{\gamma}_\theta = \dot{\gamma} - \dot{\gamma}_r.$$

$\dot{\gamma}_r$  and  $\dot{\gamma}_\theta$  are called the *radial* and *polar* components of  $\dot{\gamma}$  respectively. Notice that they are mutually orthogonal, because

$$\begin{aligned} \langle \dot{\gamma}_\theta, \dot{\gamma}_r \rangle &= \langle \dot{\gamma}, \dot{\gamma}_r \rangle - \langle \dot{\gamma}_r, \dot{\gamma}_r \rangle \\ &= \frac{1}{|\gamma|^2} \langle \dot{\gamma}, \mathcal{I}_\gamma \gamma \rangle^2 - \frac{1}{|\gamma|^4} \langle \dot{\gamma}, \mathcal{I}_\gamma \gamma \rangle^2 \langle \mathcal{I}_\gamma \gamma, \mathcal{I}_\gamma \gamma \rangle \\ &= 0, \end{aligned}$$

since  $\langle \mathcal{I}_\gamma \gamma, \mathcal{I}_\gamma \gamma \rangle = |\gamma|^2$  by definition of the inner product on the “double” tangent space  $(M_{\mathbf{p}})_{\gamma(t)}$ . When  $\gamma$  is a ray  $\boldsymbol{\varphi}_v$ , i.e.,  $\gamma(t) = t\mathbf{v}$  for some  $\mathbf{v} \in M_{\mathbf{p}}$ , then  $\dot{\gamma} = \dot{\gamma}_r$ , and  $\gamma$  is length-minimizing. The following lemma says that this property is preserved under the exponential map:

**Lemma 6.1.3.** *Let  $\mathbf{p} \in M$ , and consider a vector  $\mathbf{v}$  in the domain  $\tilde{M}_{\mathbf{p}}$  of  $\exp_{\mathbf{p}}$ . Denote by  $\boldsymbol{\varphi}_v : [0, 1] \rightarrow \tilde{M}_{\mathbf{p}}$  the ray from  $\mathbf{0}$  to  $\mathbf{v}$ ,  $\boldsymbol{\varphi}_v(t) = t\mathbf{v}$ . If  $\gamma : [0, 1] \rightarrow \tilde{M}_{\mathbf{p}}$  is any (piecewise-smooth) curve with  $\gamma(0) = \boldsymbol{\varphi}_v(0) = \mathbf{0}$  and  $\gamma(1) = \boldsymbol{\varphi}_v(1) = \mathbf{v}$ , then*

$$L(\exp \circ \gamma) \geq L(\exp \circ \boldsymbol{\varphi}_v).$$

*Inequality is strict if there is some  $t_0 \in [0, 1]$  for which the polar component of  $\dot{\gamma}$  at  $t_0$  does not vanish under the exponential map; i.e., if  $\exp_{\mathbf{p}^*} \dot{\gamma}_\theta(t_0) \neq \mathbf{0}$ .*

*Proof.* We may assume, by Proposition 2.4.1, that  $\gamma$  is differentiable, and that both  $\mathbf{v}$  and  $\gamma(t)$  are nonzero for  $t \in [0, 1]$ . By the Gauss Lemma,

$$|\exp_{\mathbf{p}^*} \dot{\gamma}|^2 = |\exp_{\mathbf{p}^*} \dot{\gamma}_r|^2 + |\exp_{\mathbf{p}^*} \dot{\gamma}_\theta|^2 \geq |\exp_{\mathbf{p}^*} \dot{\gamma}_r|^2 = |\dot{\gamma}_r|^2. \quad (6.1.2)$$

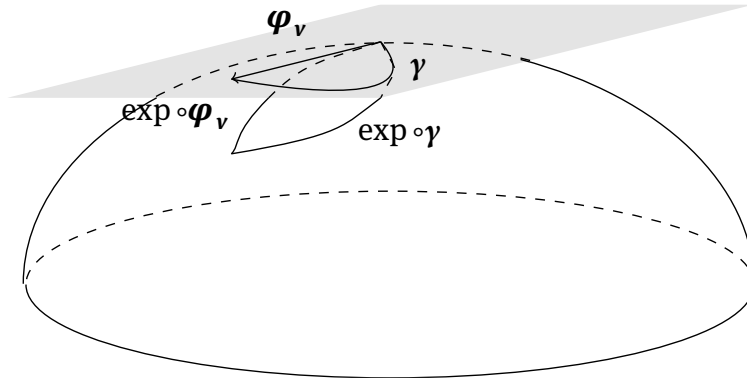
Now,  $|\gamma|' = |\dot{\gamma}_r|$ , because

$$|\gamma|' = \langle \gamma, \gamma \rangle^{1/2'} = \frac{\langle \gamma, \gamma' \rangle}{|\gamma|} = \frac{\langle \mathcal{I}_\gamma \gamma, \dot{\gamma} \rangle}{|\gamma|} = |\dot{\gamma}_r|.$$

Thus,

$$\begin{aligned} L(\exp_{\mathbf{p}} \circ \gamma) &= \int_0^1 |\exp_{\mathbf{p}^*} \dot{\gamma}| \geq \int_0^1 |\dot{\gamma}_r| = \int_0^1 |\gamma|' = |\gamma(1)| = |\mathbf{v}| \\ &= L(\exp_{\mathbf{p}} \circ \boldsymbol{\varphi}_v). \end{aligned}$$

The last assertion of the lemma is clear, since the inequality in (6.1.2) is strict on some interval around  $t_0$  if  $\exp_{\mathbf{p}^*} \dot{\gamma}_\theta(t_0) \neq \mathbf{0}$ .  $\square$



The following theorem is a mathematical formulation of the first sentence in this section:

**Theorem 6.1.1.** *Let  $\mathbf{p} \in M$ , and choose  $\varepsilon > 0$  so that  $\exp_{\mathbf{p}} : U_{\varepsilon} \rightarrow \exp_{\mathbf{p}}(U_{\varepsilon})$  is a diffeomorphism, where  $U_{\varepsilon} = \{\mathbf{v} \in M_{\mathbf{p}} \mid |\mathbf{v}| < \varepsilon\}$  is the open ball of radius  $\varepsilon$  centered at  $\mathbf{0} \in M_{\mathbf{p}}$ . For  $\mathbf{v} \in U_{\varepsilon}$ , denote as usual by  $\mathbf{c}_{\mathbf{v}} : [0, 1] \rightarrow M$  the geodesic in direction  $\mathbf{v}$ ,  $\mathbf{c}_{\mathbf{v}}(t) = \exp_{\mathbf{p}}(t\mathbf{v})$ . Then for any piecewise-smooth curve  $\mathbf{c} : [0, 1] \rightarrow M$  with  $\mathbf{c}(0) = \mathbf{c}_{\mathbf{v}}(0) = \mathbf{p}$  and  $\mathbf{c}(1) = \mathbf{c}_{\mathbf{v}}(1) = \mathbf{q}$ , the length of  $\mathbf{c}$  is at least as great as that of  $\mathbf{c}_{\mathbf{v}}$ , and is strictly greater unless  $\mathbf{c}$  equals  $\mathbf{c}_{\mathbf{v}}$  up to reparametrization.*

*Proof.* For  $\mathbf{u} \in U_{\varepsilon}$ , let  $\boldsymbol{\varphi}_{\mathbf{u}}$  be the ray  $t \mapsto t\mathbf{u}$ . Suppose first that the image of  $\mathbf{c}$  is contained inside  $\exp_{\mathbf{p}}(U_{\varepsilon})$ , so that there exists a lift  $\gamma$  of  $\mathbf{c}$  in  $U_{\varepsilon}$ ; i.e.,  $\mathbf{c} = \exp_{\mathbf{p}} \circ \gamma$ ,  $\gamma(0) = \mathbf{0}$ ,  $\gamma(1) = \mathbf{v}$ . By Lemma 6.1.3,  $L(\mathbf{c}) \geq L(\mathbf{c}_{\mathbf{v}})$ . We claim that if  $\mathbf{c}$  is not a reparametrization of  $\mathbf{c}_{\mathbf{v}}$ , then for some  $t \in [0, 1]$ ,  $\dot{\gamma}$  is not radial (and therefore  $L(\mathbf{c}) > L(\mathbf{c}_{\mathbf{v}})$  by the same lemma): Otherwise,

$$L(\gamma) = \int_0^1 |\dot{\gamma}| = \int_0^1 |\dot{\gamma}_r| = \int_0^1 |\gamma'| = |\gamma(1)| \tag{6.1.3}$$

as in the proof of Lemma 6.1.3. On the other hand, there exists  $t_0 \in (0, 1)$  such that  $\gamma(t_0) \notin \{s\mathbf{v} \mid s \in [0, 1]\}$ . Then

$$L(\gamma) = L(\gamma|_{[0, t_0]}) + L(\gamma|_{[t_0, 1]}) \geq |\gamma(t_0)| + |\gamma(1) - \gamma(t_0)| > |\gamma(1)|,$$

which contradicts (6.1.3). This establishes the result if the image of  $\mathbf{c}$  lies in  $\exp_{\mathbf{p}}(U_{\varepsilon})$ . Next, suppose  $\mathbf{c}$  is not entirely contained inside  $\exp_{\mathbf{p}}(U_{\varepsilon})$ , and let  $b = \sup\{t \mid \mathbf{c}[0, t] \subset \exp_{\mathbf{p}}(U_{\varepsilon})\}$ . There must exist some  $t_0 \in (0, b)$  such that  $\mathbf{v}_0 := (\exp_{\mathbf{p}}|_{U_{\varepsilon}})^{-1}\mathbf{c}(t_0)$  has norm greater than that of  $\mathbf{v}$ . Then

$$L(\mathbf{c}) \geq L(\mathbf{c}|_{[0, t_0]}) \geq |\mathbf{v}_0| > |\mathbf{v}| = L(\mathbf{c}_{\mathbf{v}}). \quad \square$$

A geodesic  $\mathbf{c}$  is said to be *minimal* if its length is equal to the distance between its endpoints. Theorem 6.1.1 asserts that for each  $\mathbf{p} \in M$ , there exists an  $\varepsilon > 0$  such that all geodesics of length less than  $\varepsilon$  emanating from  $\mathbf{p}$  are minimal.

It is also worth noting that any curve  $\mathbf{c}$  that is minimal in the above sense is necessarily a geodesic: if  $t_0$  is a point in the domain of  $\mathbf{c}$ , then the restriction of  $\mathbf{c}$  to  $[t_0, t_0 + \varepsilon]$

is also minimal whenever  $t_0 + \varepsilon$  lies in the domain of  $\mathbf{c}$ , and by Theorem 6.1.1, it is a geodesic. Since  $t_0$  is arbitrary, the claim follows.

## 6.2 Jacobi fields

The previous section indicates that the exponential map  $\exp_p$  at  $\mathbf{p} \in M$  plays a fundamental role in the geometry of  $M$ . It turns out that its derivative can be conveniently expressed in terms of certain vector fields along geodesics emanating from  $\mathbf{p}$ . But first a word on notation: if  $\mathbf{c}$  is a curve in  $M$  and  $\mathbf{Y}$  a vector field along  $\mathbf{c}$ , we have in the past used the notation  $\mathbf{Y}'$  to denote the covariant derivative of  $\mathbf{Y}$  along  $\mathbf{c}$  in the ambient Euclidean space, and  $\mathbf{Y}'^\top = \nabla_{\mathbf{D}}\mathbf{Y}$  to denote the corresponding covariant derivative in  $M$ . This notation becomes cumbersome when taking second derivatives. Furthermore, we also wish to use the superscripts  $^\top$  and  $^\perp$  to describe components of vector fields tangent and orthogonal to  $\mathbf{c}$  respectively rather than to  $M$ . With this in mind, we introduce the following

**Notation:** Until further notice, the covariant derivative  $\nabla_{\mathbf{D}}\mathbf{Y}$  of a vector field  $\mathbf{Y}$  along a curve  $\mathbf{c}$  in  $M$  will also be denoted by  $\mathbf{Y}'$ .

**Definition 6.2.1.** Let  $\mathbf{c}$  be a geodesic in  $M$ . A vector field  $\mathbf{Y}$  along  $\mathbf{c}$  is called a *Jacobi field* along  $\mathbf{c}$  if

$$\mathbf{Y}'' + R(\mathbf{Y}, \dot{\mathbf{c}})\dot{\mathbf{c}} = 0.$$

Notice that the collection  $\mathcal{J}_{\mathbf{c}}$  of Jacobi fields along  $\mathbf{c}$  is a vector space that contains  $\dot{\mathbf{c}}$ . It turns out that the subspace of Jacobi fields orthogonal to  $\dot{\mathbf{c}}$  is the one of interest to us: If  $\mathbf{X}$  and  $\mathbf{Y}$  are Jacobi, then

$$\langle \mathbf{Y}'', \mathbf{X} \rangle = -\langle R(\mathbf{Y}, \dot{\mathbf{c}})\dot{\mathbf{c}}, \mathbf{X} \rangle = -\langle R(\mathbf{X}, \dot{\mathbf{c}})\dot{\mathbf{c}}, \mathbf{Y} \rangle = \langle \mathbf{X}'', \mathbf{Y} \rangle.$$

Thus,  $\langle \mathbf{X}', \mathbf{Y} \rangle - \langle \mathbf{Y}', \mathbf{X} \rangle$  must be constant, since its derivative is  $\langle \mathbf{X}'', \mathbf{Y} \rangle - \langle \mathbf{Y}'', \mathbf{X} \rangle = 0$ . In particular,  $\langle \mathbf{Y}, \dot{\mathbf{c}} \rangle' = \langle \mathbf{Y}', \dot{\mathbf{c}} \rangle = \langle \mathbf{Y}', \dot{\mathbf{c}} \rangle - \langle \mathbf{Y}, \dot{\mathbf{c}}' \rangle$  is constant, so that for a normal geodesic, the tangential component  $\mathbf{Y}^\top$  of  $\mathbf{Y}$  is given by

$$\mathbf{Y}^\top = \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle \dot{\mathbf{c}} = (a + bt)\dot{\mathbf{c}}, \quad a = \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle(0), \quad b = \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle'(0),$$

and satisfies the Jacobi equation. It follows that the component  $\mathbf{Y}^\perp = \mathbf{Y} - \mathbf{Y}^\top$  of  $\mathbf{Y}$  orthogonal to  $\dot{\mathbf{c}}$  is also a Jacobi field.

**Proposition 6.2.1.** Let  $\mathbf{c} : I \rightarrow M$  be a geodesic,  $t_0 \in I$ . For any  $\mathbf{v}, \mathbf{w} \in M_{\mathbf{c}(t_0)}$  there exists a unique Jacobi field  $\mathbf{Y}$  along  $\mathbf{c}$  with  $\mathbf{Y}(t_0) = \mathbf{v}$  and  $\mathbf{Y}'(t_0) = \mathbf{w}$ .

*Proof.* Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be parallel fields along  $\mathbf{c}$  with  $\mathbf{X}_n = \dot{\mathbf{c}}$ , and such that  $\mathbf{X}_1(t_0), \dots, \mathbf{X}_{n-1}(t_0)$  form an orthonormal basis of  $\dot{\mathbf{c}}(t_0)^\perp$ . Any vector field  $\mathbf{Y}$  along  $\mathbf{c}$  can then be

expressed as

$$Y = \sum_i f^i X_i, \quad f^i = \begin{cases} \langle Y, X_i \rangle, & \text{for } i \leq n-1, \\ \langle Y, \frac{X_n}{|X_n|^2} \rangle, & \text{for } i = n. \end{cases}$$

Since  $X_i$  is parallel,  $Y'' = \sum f^{i''} X_i$ . Furthermore,  $R(X_i, \dot{c})\dot{c} = \sum_{j=1}^{n-1} h_i^j X_j$ , where  $h_i^j = \langle R(X_i, \dot{c})\dot{c}, X_j \rangle$ , so that  $R(Y, \dot{c})\dot{c} = \sum_{i,j=1}^{n-1} f^i h_i^j X_j$ . The Jacobi equation then translates into

$$\sum_{j=1}^{n-1} \left( f^{j''} + \sum_{i=1}^{n-1} f^i h_i^j \right) X_j = 0, \quad f^{n''} = 0,$$

or equivalently, since the  $X_i$  are linearly independent everywhere and  $h_i^n = \langle R(X_i, \dot{c})\dot{c}, \dot{c} \rangle = 0$ ,

$$f^{j''} + \sum_{i=1}^{n-1} f^i h_i^j = 0, \quad j = 1, \dots, n.$$

This is a homogeneous system of  $n$  linear second-order equations, which has a unique solution given initial values for  $f^j$  and  $f^{j'}$  at  $t_0$ ; in our case, we have  $f^j(t_0) = \langle v, X_j(t_0) \rangle$ ,  $f^{j'}(t_0) = \langle w, X_j(t_0) \rangle$  ( $j < n$ ),  $f^n(t_0) = \langle v, (\dot{c}/|\dot{c}|^2)(t_0) \rangle$ , and  $f^{n'}(t_0) = \langle w, (\dot{c}/|\dot{c}|^2)(t_0) \rangle$ , thereby establishing the result.  $\square$

The existence part of Proposition 6.2.1 implies that the linear map

$$\begin{aligned} \mathcal{J}_c &\rightarrow M_{c(t_0)} \times M_{c(t_0)}, \\ Y &\mapsto (Y(t_0), Y'(t_0)) \end{aligned}$$

is onto. The uniqueness part implies it has trivial kernel. In other words, it is an isomorphism, and the space  $\mathcal{J}_c$  of Jacobi fields along  $c$  has dimension  $2n$ .

**Example 6.2.1.** Let  $M^n$  be a space of constant curvature  $\kappa$ , and let  $c_\kappa, s_\kappa$  denote the solutions of the differential equation

$$f'' + \kappa f = 0$$

with  $c_\kappa(0) = 1, c'_\kappa(0) = 0, s_\kappa(0) = 0, s'_\kappa(0) = 1$ . Thus,  $c_\kappa(t) = \cos \sqrt{\kappa}t$  if  $\kappa > 0$ ,  $1$  if  $\kappa = 0$ ,  $\cosh \sqrt{-\kappa}t$  if  $\kappa$  is negative, whereas  $s_\kappa$  is obtained by replacing  $\cos, 1$ , and  $\cosh$  by  $\sin, 1_{\mathbb{R}}$ , and  $\sinh$  respectively. Consider a normal geodesic  $c : [0, b] \rightarrow M$ . Given  $v, w \in M_{c(0)}$  orthogonal to  $\dot{c}(0)$ , the Jacobi field  $Y$  along  $c$  with  $Y(0) = v$  and  $Y'(0) = w$  is given by

$$Y = c_\kappa E + s_\kappa F,$$

where  $E$  and  $F$  are the parallel fields along  $c$  with  $E(0) = v$  and  $F(0) = w$ : Indeed,  $Y'' = c''_\kappa E + s''_\kappa F = -\kappa Y = -R(Y, \dot{c})\dot{c}$ , so that  $Y$  is a Jacobi field, and clearly satisfies the initial conditions at 0.

Jacobi fields essentially arise out of variations where all curves are geodesics: If  $V : [0, a] \times I \rightarrow M$  is a variation of  $c$  with  $V(t, 0) = c(t)$  for  $t \in [0, a]$ , the *variational vector field* of  $V$  is the vector field  $Y$  along  $c$  given by  $Y(t) = V_* D_2(t, 0)$ .



**Proposition 6.2.2.** *Let  $c : [0, a] \rightarrow M$  be a geodesic. If  $V$  is a variation of  $c$  through geodesics – meaning that  $V_s$  is a geodesic for each  $s$ , then the variational vector field of  $V$  is Jacobi along  $c$ . Conversely, let  $Y$  be a Jacobi field along  $c$ . Then there exists a variation  $V$  of  $c$  through geodesics whose variational vector field equals  $Y$ .*

*Proof.* Given a variation  $V$  of  $c$  through geodesics, define vector fields  $\tilde{X}$  and  $\tilde{Y}$  along  $V$  by  $\tilde{X} = V_*D_1$ ,  $\tilde{Y} = V_*D_2$ . By assumption,  $\nabla_{D_1}\tilde{X} = \mathbf{0}$ , so that

$$\begin{aligned} R(\tilde{Y}, \tilde{X})\tilde{X} &= \nabla_{D_2}\nabla_{D_1}\tilde{X} - \nabla_{D_1}\nabla_{D_2}\tilde{X} = -\nabla_{D_1}\nabla_{D_2}\tilde{X} \\ &= -\nabla_{D_1}\nabla_{D_2}V_*D_1 = -\nabla_{D_1}\nabla_{D_1}V_*D_2 \\ &= -\nabla_{D_1}\nabla_{D_1}\tilde{Y}. \end{aligned}$$

When  $s = 0$ , the above expression becomes  $R(Y, \dot{c})\dot{c} = -Y''$ , and  $Y$  is Jacobi.

Conversely, suppose  $Y$  is a Jacobi field along  $c$ , and set  $v := Y(0)$ ,  $w := Y'(0)$ . Let  $\gamma$  be a curve with  $\dot{\gamma}(0) = v$ , and  $X, W$  parallel fields along  $\gamma$  with  $X(0) = \dot{c}(0)$ ,  $W(0) = w$ . Since  $t(X(0) + sW(0))$  belongs to the domain of  $\exp_{\gamma(0)}$  for  $0 \leq t \leq a$  and small enough  $s$ , there must exist  $\varepsilon > 0$  small enough so that  $t(X(s) + sW(s))$  belongs to the domain of  $\exp_{\gamma(s)}$  for  $(t, s) \in [0, a] \times (-\varepsilon, \varepsilon)$ . Consider the variation

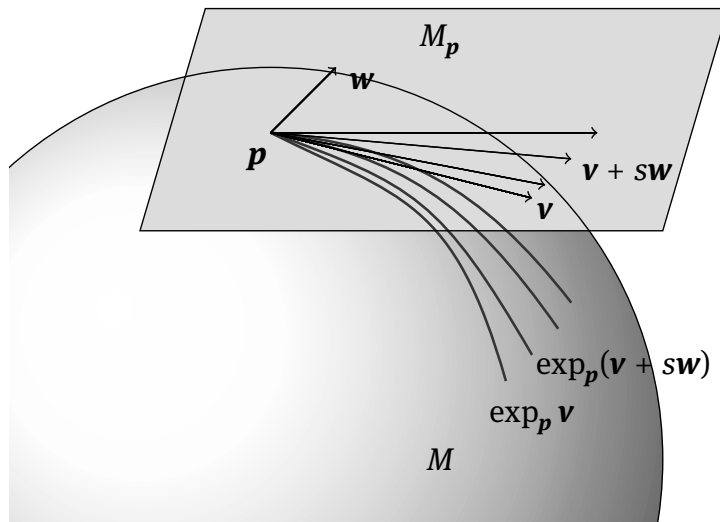
$$\begin{aligned} V : [0, a] \times (-\varepsilon, \varepsilon) &\rightarrow M, \\ (t, s) &\mapsto \exp_{\gamma(s)} t(X(s) + sW(s)) \end{aligned}$$

of  $c$ . Now, the curves  $t \mapsto V(t, s)$  are geodesics, so the variational vector field  $Z$  is Jacobi along  $c$ . Moreover,  $V(0, s) = \gamma(s)$ , so that  $Z(0) = \dot{\gamma}(0) = v$ . Finally,

$$Z'(0) = \nabla_{D_1(0,0)}V_*D_2 = \nabla_{D_2(0,0)}V_*D_1 = W(0) = w,$$

because  $V_*D_1(0, s) = X(s) + sW(s)$ , and  $X, W$  are parallel along  $\gamma$ . By Proposition 6.2.1,  $Z = Y$ . □

A geodesic variation of  $t \mapsto \exp_p tv$  with fixed initial point. The corresponding Jacobi field  $Y$  vanishes at  $t = 0$



In the special case when  $\mathbf{Y}(0) = \mathbf{0}$ , the variation from Proposition 6.2.2 becomes  $\mathbf{V}(t, s) = \exp_{c(0)} t(\dot{c}(0) + s\mathbf{w})$ . It is the image via the exponential map of a variation of the line segment  $t \mapsto t\dot{c}(0)$  in  $M_p$  by rays from the origin. The Jacobi field  $\mathbf{Y}$  with initial conditions  $\mathbf{Y}(0) = \mathbf{0}$ ,  $\mathbf{Y}'(0) = \mathbf{w}$  is given by

$$\mathbf{Y}(t) = \exp_{c(0)*}(t\mathcal{I}_{t\dot{c}(0)}\mathbf{w}). \tag{6.2.1}$$

**Definition 6.2.2.** If  $c : [a, b] \rightarrow M$  is a geodesic,  $t_0 \in (a, b)$  is said to be a *conjugate point* of  $c$  if there exists a nontrivial Jacobi field  $\mathbf{Y}$  along  $c$  that vanishes at  $a$  and  $t_0$ .

Conjugate points correspond to critical points of the exponential map: our next result implies that for  $p \in M$  and  $\mathbf{u} \in M_p$ ,  $\mathbf{u}$  is a critical point of  $\exp_p$  if and only if 1 is a conjugate point of the geodesic  $t \mapsto \exp_p(t\mathbf{u})$ .

**Theorem 6.2.1.** Let  $p \in M$ ,  $\mathbf{u} \in M_p$ ,  $c : [0, a] \rightarrow M$  the geodesic  $t \mapsto \exp_p(t\mathbf{u})$ , and  $t_0 \in (0, a)$ . The vector space  $\mathcal{J}_c^{t_0}$  of all Jacobi fields  $\mathbf{Y}$  along  $c$  that vanish at 0 and  $t_0$  has the same dimension as the kernel of the derivative of  $\exp_p$  at  $t_0\mathbf{u}$ .

*Proof.* Any Jacobi field  $\mathbf{Y}$  along  $c$  with  $\mathbf{Y}(0) = \mathbf{0}$  has the form

$$\mathbf{Y}(t) = \exp_{p*(t\mathbf{u})}(t\mathcal{I}_{t\mathbf{u}}\mathbf{Y}'(0))$$

by (6.2.1). Thus, if  $\mathbf{Y} \in \mathcal{J}_c^{t_0}$ , then

$$\mathbf{Y}(t_0) = \exp_{p*(t_0\mathcal{I}_{t_0\mathbf{u}}\mathbf{Y}'(0))} = \mathbf{0},$$

and  $\mathcal{I}_{t_0\mathbf{u}}\mathbf{Y}'(0) \in \ker \exp_{p*(t_0\mathbf{u})}$ . The map

$$\begin{aligned} L : \mathcal{J}_c^{t_0} &\longrightarrow \ker \exp_{p*(t_0\mathbf{u})}, \\ \mathbf{Y} &\longmapsto \mathcal{I}_{t_0\mathbf{u}}\mathbf{Y}'(0) \end{aligned}$$

is linear and has trivial kernel since any Jacobi field  $\mathbf{Y}$  with  $\mathbf{Y}(0) = \mathbf{Y}'(0) = \mathbf{0}$  is trivial. It is also onto, for if  $\mathbf{v}$  lies in the kernel of the derivative of  $\exp_p$  at  $t_0\mathbf{u}$ , then the Jacobi field  $\mathbf{Y}$  along  $c$  with  $\mathbf{Y}(0) = \mathbf{0}$  and  $\mathbf{Y}'(0) = \mathcal{I}_{t_0\mathbf{u}}^{-1}\mathbf{v}$  vanishes at  $t_0$  by (6.2.1), so that  $\mathbf{Y} \in \mathcal{J}_c^{t_0}$  and  $L\mathbf{Y} = \mathbf{v}$ .  $L$  is therefore an isomorphism, and the theorem follows.  $\square$

**Example 6.2.2.** If  $M^n$  is a space of constant curvature  $\kappa$  and  $c : [0, \infty) \rightarrow M$  is a normal geodesic in  $M$ , then by Example 6.2.1, any Jacobi field  $\mathbf{Y}$  along  $c$  with  $\mathbf{Y}(0) = \mathbf{0}$  can be written as  $\mathbf{Y} = s_\kappa\mathbf{E}$ , where  $\mathbf{E}$  is a parallel field and  $s_\kappa(t) = t$  if  $\kappa = 0$ ,  $\sin \sqrt{\kappa}t$  if  $\kappa > 0$ , and  $\sinh \sqrt{-\kappa}t$  if  $\kappa < 0$ . Normal geodesics therefore have no conjugate points if  $\kappa \leq 0$ , and have  $k\pi/\sqrt{\kappa}$ ,  $k \in \mathbb{N}$ , as conjugate points when  $\kappa > 0$ . In the latter case, the images via  $c$  of the critical points alternate between  $c(0)$  and its antipode.

More generally, we have the following:

**Theorem 6.2.2.** In a manifold  $M$  with nonpositive curvature, geodesics have no conjugate points.

*Proof.* Let  $\mathbf{Y}$  denote a Jacobi field along a geodesic in  $M$ . We will show that the norm of  $\mathbf{Y}$  is a convex function. Since a convex nonnegative function that vanishes at two points vanishes everywhere, the statement will follow.

Specifically, we claim that if  $\mathbf{Y}(t) \neq \mathbf{0}$ , then  $|\mathbf{Y}|''(t) \geq 0$ . To see this, observe that  $|\mathbf{Y}'| = \langle \mathbf{Y}, \mathbf{Y} \rangle^{1/2} = \langle \mathbf{Y}, \mathbf{Y}' \rangle / |\mathbf{Y}|$ . If  $K(t)$  denotes the curvature of the plane spanned by the tangent vector of the geodesic at  $t$  and by  $\mathbf{Y}(t)$ , then  $K \leq 0$  so that

$$\begin{aligned} |\mathbf{Y}|'' &= \frac{|\mathbf{Y}|(|\mathbf{Y}'|^2 + \langle \mathbf{Y}, \mathbf{Y}'' \rangle) - \langle \mathbf{Y}, \mathbf{Y}' \rangle^2 / |\mathbf{Y}|}{|\mathbf{Y}|^2} \\ &= \frac{1}{|\mathbf{Y}|^3} (|\mathbf{Y}|^2 |\mathbf{Y}'|^2 - K |\mathbf{Y}|^4 - \langle \mathbf{Y}, \mathbf{Y}' \rangle^2) \\ &\geq \frac{1}{|\mathbf{Y}|^3} (|\mathbf{Y}|^2 |\mathbf{Y}'|^2 - \langle \mathbf{Y}, \mathbf{Y}' \rangle^2), \end{aligned}$$

and the last expression is nonnegative by the Cauchy-Schwartz inequality.  $\square$

### 6.3 The length function of a variation

Let  $\mathbf{c} : [0, a] \rightarrow M$  be a normal geodesic, and  $V : [0, a] \times I \rightarrow M$  a variation of  $\mathbf{c}$ , where  $I$  is an open interval containing 0. In this section, we discuss the length function  $L : I \rightarrow \mathbb{R}$  of the variation, with  $L(s)$  denoting the length of the curve  $V_s$ ,  $V_s(t) = V(t, s)$ . It turns out that when the variation has fixed endpoints, then  $\mathbf{c}$  is shorter than nearby curves in the variation, provided it has no conjugate points.  $\mathbf{Y} : [0, a] \rightarrow TM$  will denote the variational vector field of  $V$ ,  $\mathbf{Y}(t) = V_* \mathbf{D}_2(t, 0)$ , and  $\mathbf{Y}_\perp = \mathbf{Y} - \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle \dot{\mathbf{c}}$  its component orthogonal to the geodesic.

**Lemma 6.3.1.** *With notation as above,*

$$L'(0) = \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle|_0^a, \quad (6.3.1)$$

$$L''(0) = \left( \int_0^a |\mathbf{Y}'_\perp|^2 - \langle R(\mathbf{Y}_\perp, \dot{\mathbf{c}}) \dot{\mathbf{c}}, \mathbf{Y}_\perp \rangle \right) + \langle \nabla_{\mathbf{D}_2} V_* \mathbf{D}_2(t, 0), \dot{\mathbf{c}}(t) \rangle|_0^a. \quad (6.3.2)$$

*In particular, if  $V$  has fixed endpoints (meaning  $V(0, s) = \mathbf{c}(0)$ ,  $V(a, s) = \mathbf{c}(a)$  for all  $s$ ), then*

$$L'(0) = 0, \quad (6.3.3)$$

$$L''(0) = \int_0^a |\mathbf{Y}'_\perp|^2 - \langle R(\mathbf{Y}_\perp, \dot{\mathbf{c}}) \dot{\mathbf{c}}, \mathbf{Y}_\perp \rangle. \quad (6.3.4)$$

*Proof.* By definition,  $L(s) = \int_0^a |V_* \mathbf{D}_1|(t, s) dt$ . Since  $V_* \mathbf{D}_1$  is continuous and has constant norm 1 when  $s = 0$ , compactness of  $[0, a]$  implies that  $V_* \mathbf{D}_1$  is nonzero on  $[0, a] \times$

$(-\varepsilon, \varepsilon)$  for some  $\varepsilon > 0$ . This means that  $L$  is differentiable at 0, and by Theorem 2.3.1,

$$\begin{aligned} L'(s) &= \int_0^a \mathbf{D}_2 |V_* \mathbf{D}_1|(t, s) dt = \int_0^a \frac{\mathbf{D}_2 \langle V_* \mathbf{D}_1, V_* \mathbf{D}_1 \rangle}{2|V_* \mathbf{D}_1|}(t, s) dt \\ &= \int_0^a \frac{\langle \nabla_{\mathbf{D}_2} V_* \mathbf{D}_1, V_* \mathbf{D}_1 \rangle}{|V_* \mathbf{D}_1|}(t, s) dt = \int_0^a \frac{\langle \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, V_* \mathbf{D}_1 \rangle}{|V_* \mathbf{D}_1|}(t, s) dt. \end{aligned}$$

Note that for the second equality above, we used the fact  $|\mathbf{a}| = \langle \mathbf{a}, \mathbf{a} \rangle^{1/2}$  and the chain rule. When  $s = 0$ ,  $|V_* \mathbf{D}_1| \equiv 1$ , and

$$\langle \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, V_* \mathbf{D}_1 \rangle(t, 0) = \langle \mathbf{Y}', \dot{\mathbf{c}} \rangle(t) = \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle'(t),$$

so that  $L'(0) = \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle|_0^a$ , which establishes (6.3.1). For (6.3.2), we have

$$L''(s) = \int_0^a \mathbf{D}_2 \left( \frac{\langle \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, V_* \mathbf{D}_1 \rangle}{|V_* \mathbf{D}_1|} \right) (t, s) dt,$$

and the integrand equals

$$\begin{aligned} \frac{\langle \nabla_{\mathbf{D}_2} \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, V_* \mathbf{D}_1 \rangle + \langle \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, \nabla_{\mathbf{D}_2} V_* \mathbf{D}_1 \rangle}{|V_* \mathbf{D}_1|} - \frac{\langle \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, \mathbf{D}_1 \rangle^2}{|V_* \mathbf{D}_1|^3} &= \\ \frac{\langle \nabla_{\mathbf{D}_2} \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, V_* \mathbf{D}_1 \rangle + |\nabla_{\mathbf{D}_1} V_* \mathbf{D}_2|^2}{|V_* \mathbf{D}_1|} - \frac{\langle \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, \mathbf{D}_1 \rangle^2}{|V_* \mathbf{D}_1|^3}. \end{aligned}$$

When  $s = 0$ , this expression becomes

$$\langle \nabla_{\mathbf{D}_2} \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2(t, 0), \dot{\mathbf{c}}(t) \rangle + |\nabla_{\mathbf{D}_1} \mathbf{Y}|^2(t) - \langle \nabla_{\mathbf{D}_1} \mathbf{Y}, \dot{\mathbf{c}} \rangle^2(t).$$

Now,

$$\begin{aligned} |\nabla_{\mathbf{D}_1} \mathbf{Y}|^2 &= |\nabla_{\mathbf{D}_1} (\mathbf{Y}_\perp + \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle \dot{\mathbf{c}})|^2 = |\nabla_{\mathbf{D}_1} \mathbf{Y}_\perp + \mathbf{D}_1(\langle \mathbf{Y}, \dot{\mathbf{c}} \rangle) \dot{\mathbf{c}}|^2 \\ &= |\nabla_{\mathbf{D}_1} \mathbf{Y}_\perp|^2 + (\mathbf{D}_1 \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle)^2 \quad (\text{since } \langle \nabla_{\mathbf{D}_1} \mathbf{Y}_\perp, \dot{\mathbf{c}} \rangle = \mathbf{D}_1 \langle \mathbf{Y}_\perp, \dot{\mathbf{c}} \rangle = 0) \\ &= |\nabla_{\mathbf{D}_1} \mathbf{Y}_\perp|^2 + \langle \nabla_{\mathbf{D}_1} \mathbf{Y}, \dot{\mathbf{c}} \rangle^2, \end{aligned}$$

so that

$$L''(0) = \int_0^a \left( \langle \nabla_{\mathbf{D}_2} \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2, \dot{\mathbf{c}} \rangle + |\mathbf{Y}'_\perp|^2 \right) dt.$$

Furthermore,

$$\nabla_{\mathbf{D}_2} \nabla_{\mathbf{D}_1} V_* \mathbf{D}_2(t, 0) = \nabla_{\mathbf{D}_1} \nabla_{\mathbf{D}_2} V_* \mathbf{D}_2(t, 0) - (R(\dot{\mathbf{c}}, \mathbf{Y})\mathbf{Y})(t),$$

and since  $\langle \nabla_{\mathbf{D}_1} \nabla_{\mathbf{D}_2} V_* \mathbf{D}_2(t, 0), \dot{\mathbf{c}}(t) \rangle = \mathbf{D}_1 \langle \nabla_{\mathbf{D}_2} V_* \mathbf{D}_2(t, 0), \dot{\mathbf{c}}(t) \rangle$ , we obtain

$$L''(0) = \left( \int_0^a |\mathbf{Y}'_\perp|^2 - \langle R(\dot{\mathbf{c}}, \mathbf{Y})\mathbf{Y}, \dot{\mathbf{c}} \rangle \right) + \langle \nabla_{\mathbf{D}_2} V_* \mathbf{D}_2(t, 0), \dot{\mathbf{c}}(t) \rangle|_0^a.$$

The symmetries of the curvature tensor and the fact that  $\langle R(\dot{\mathbf{c}}, \mathbf{Y})\mathbf{Y}, \dot{\mathbf{c}} \rangle = \langle R(\dot{\mathbf{c}}, \mathbf{Y}_\perp)\mathbf{Y}_\perp, \dot{\mathbf{c}} \rangle$  now yield (6.3.2).

If  $V$  has fixed endpoints, then  $V_*\mathbf{D}_2(0, s)$ ,  $V_*\mathbf{D}_2(a, s)$ ,  $\mathbf{Y}(0)$  and  $\mathbf{Y}(a)$  all vanish, thereby establishing (6.3.3) and (6.3.4).  $\square$

We will need a more general version of the lemma. Let  $I$  be an open interval containing 0. A *piecewise smooth variation* of a smooth curve  $\mathbf{c} : [0, a] \rightarrow M$  is a continuous map  $V : [0, a] \times I \rightarrow M$  for which there exists a partition  $0 = t_0 < t_1 < \dots < t_l = a$  of  $[0, a]$  such that each  $V_i := V|_{[t_{i-1}, t_i] \times I}$  is a variation of  $\mathbf{c}|_{[t_{i-1}, t_i]}$ ,  $i = 1, \dots, l$ . A *piecewise smooth vector field* along  $\mathbf{c}$  is a piecewise smooth curve  $\mathbf{Y} : [0, a] \rightarrow TM$  such that  $\mathbf{Y}(t) \in M_{\mathbf{c}(t)}$  for all  $t$ . If  $V$  is a piecewise smooth variation of  $\mathbf{c}$ , then by assumption, the map  $V_*\mathbf{D}_2 : [0, a] \times I \rightarrow TM$ , where  $V_*\mathbf{D}_2|_{[t_{i-1}, t_i] \times I} = V_{i*}\mathbf{D}_2$ , is continuous. It then induces a piecewise smooth vector field  $\mathbf{Y}$  along  $\mathbf{c}$ , with  $\mathbf{Y}(t) = V_*\mathbf{D}_2(t, 0)$ . Even though  $\mathbf{Y}$  is not, in general, differentiable at  $t_i$ ,  $\mathbf{Y}'$  may be extended to a not necessarily continuous vector field on  $[0, a]$  by defining  $\mathbf{Y}'(t_i) = \lim_{t \rightarrow t_i^+} \mathbf{Y}'(t)$ . Notice that even though  $\nabla_{\mathbf{D}_1} V_*\mathbf{D}_2(t, s)$  may not exist,  $\nabla_{\mathbf{D}_2} V_*\mathbf{D}_2(t, s)$  does. As before, set  $\mathbf{Y}_\perp = \mathbf{Y} - \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle \dot{\mathbf{c}}$ ,  $L(s) = L(V_s)$ , where  $V_s(t) = V(t, s)$ .

**Proposition 6.3.1.** *Let  $\mathbf{c} : [0, a] \rightarrow M$  be a normal geodesic, and  $V : [0, a] \times I \rightarrow M$  a piecewise smooth variation of  $\mathbf{c}$ . Then equations (6.3.1) through (6.3.4) from Lemma 6.3.1 hold.*

*Proof.* With notation as above, if  $L_i : I \rightarrow \mathbb{R}$  denotes the length function of  $V_i$ ,  $i = 1, \dots, l$ , then  $L = \sum_i L_i$ . The lemma then implies that

$$L'(0) = \sum_{i=1}^l L'_i(0) = \sum_{i=1}^l \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle|_{t_{i-1}}^{t_i} = \langle \mathbf{Y}, \dot{\mathbf{c}} \rangle|_0^a,$$

and

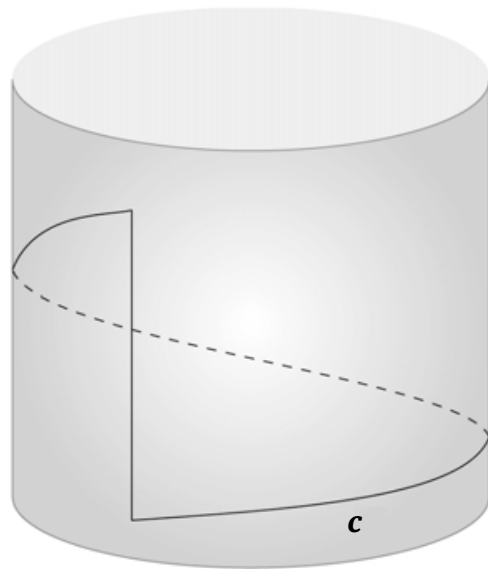
$$\begin{aligned} L''(0) &= \sum_{i=1}^l L''_i(0) \\ &= \sum_{i=1}^l \left[ \int_{t_{i-1}}^{t_i} |\mathbf{Y}'_\perp|^2 - \langle R(\mathbf{Y}_\perp, \dot{\mathbf{c}})\dot{\mathbf{c}}, \mathbf{Y}_\perp \rangle + \langle \nabla_{\mathbf{D}_2} V_*\mathbf{D}_2(t, 0), \dot{\mathbf{c}}(t) \rangle|_{t_{i-1}}^{t_i} \right] \\ &= \left( \int_0^a |\mathbf{Y}'_\perp|^2 - \langle R(\mathbf{Y}_\perp, \dot{\mathbf{c}})\dot{\mathbf{c}}, \mathbf{Y}_\perp \rangle \right) + \langle \nabla_{\mathbf{D}_2} V_*\mathbf{D}_2(t, 0), \dot{\mathbf{c}}(t) \rangle|_0^a. \end{aligned}$$

The last two identities in the statement follow just as in the Lemma.  $\square$

**Proposition 6.3.2.** *Let  $\mathbf{c} : [0, a] \rightarrow M$  be a normal geodesic without conjugate points. If  $V : [0, a] \times I \rightarrow M$  is a variation of  $\mathbf{c}$  with fixed endpoints, then for sufficiently small  $s$ ,  $L(V_s) \geq L(V_0) = L(\mathbf{c})$ , and strict inequality holds provided  $V_s$  is not a reparametrization of  $\mathbf{c}$ .*

*Proof.* Let  $\gamma : \mathbb{R} \rightarrow M_p$  denote the ray  $t \mapsto t\dot{c}(0)$  in the tangent space at  $p = c(0)$ , so that  $c = \exp \circ \gamma|_{[0,a]}$ . Since  $c$  has no conjugate points,  $\exp_p$  has maximal rank at each  $\gamma(t)$  and is therefore a diffeomorphism in a neighborhood of  $\gamma(t)$ . By compactness of  $[a, b]$ , there exists a partition  $0 = t_0 < t_1 < \dots < t_l = a$  of the interval and open sets  $U_1, \dots, U_l$  in  $M_p$  such that  $\gamma[t_{i-1}, t_i] \subset U_i$  and  $\exp_p : U_i \rightarrow \exp(U_i)$  is a diffeomorphism for each  $i$  from 1 to  $l$ . By Lemma 1.7.1, there exists for each  $i$  some  $\varepsilon_i > 0$  such that  $V([t_{i-1}, t_i] \times (-\varepsilon_i, \varepsilon_i)) \subset \exp(U_i)$ . Let  $\varepsilon = \min\{\varepsilon_1, \dots, \varepsilon_l\} > 0$ , and define  $\tilde{V} : [0, a] \times (-\varepsilon, \varepsilon) \rightarrow M_p$  by  $\tilde{V}(t, s) = (\exp_p|_{U_i})^{-1}(V(t, s))$  if  $t \in [t_{i-1}, t_i]$ . Each  $\tilde{V}_s$ , where  $\tilde{V}(s) = \tilde{V}(t, s)$ , is then a curve in  $M_p$  from  $\mathbf{0}$  to  $\dot{c}(0)$ . Since  $V_s = \exp \circ \tilde{V}_s$ , the claim follows from Lemma 6.1.3.  $\square$

It is in general not true, though, that a geodesic without conjugate points is minimal. The cylinder  $M = \{(x, y, z) \in \mathbb{R}^3 \mid x^2 + y^2 = 1\}$  provides many examples: If  $a, b > 0$  and  $a^2 + b^2 = 1$ , then by Example 3.11.1, the curve  $c : [0, 2\pi/a] \rightarrow M$ ,  $c(t) = (\cos(at), \sin(at), bt)$ , is a normal geodesic joining  $(1, 0, 0)$  to  $(1, 0, 2\pi b/a)$ , and since  $M$  is flat,  $c$  has no conjugate points by Theorem 6.2.2. The meridian  $t \mapsto (1, 0, t)$  also joins the endpoints of  $c$  and has length  $2\pi b/a < 2\pi/a$ . The proposition only guarantees that  $c$  will be shorter than sufficiently nearby curves in a variation with fixed endpoints.



The geodesic  $c$  has no conjugate points, but is not minimal

### 6.4 The index form of a geodesic

Let  $c : [0, a] \rightarrow M$  denote a normal geodesic. In this section, we introduce a symmetric bilinear form on the space  $\mathfrak{V}_c$  of piecewise smooth vector fields along  $c$  that are orthogonal to  $c$  and vanish at its endpoints. This form is closely related to the existence of conjugate points of  $c$ . As a first application, we will show that a geodesic does not minimize length past its first conjugate point; specifically, if  $c : [0, a] \rightarrow M$  has a conjugate point in  $(0, a)$ , then there exist curves in  $M$  joining  $c(0)$  to  $c(a)$  that are shorter than  $c$  and arbitrarily close to it.

**Proposition 6.4.1.** *Let  $\mathbf{c} : [0, a] \rightarrow M$  be a normal geodesic, and  $\mathfrak{V}_c$  the (infinite-dimensional) vector space of all piecewise smooth vector fields  $\mathbf{Y}$  along  $\mathbf{c}$  such that  $\langle \mathbf{Y}, \dot{\mathbf{c}} \rangle = 0$  and  $\mathbf{Y}(0) = \mathbf{0}$ ,  $\mathbf{Y}(a) = \mathbf{0}$ . Then there exists a unique symmetric bilinear form  $I$  on  $\mathfrak{V}_c$  such that if  $V$  is a variation of  $\mathbf{c}$  with fixed endpoints that has  $\mathbf{Y} \in \mathfrak{V}_c$  as variational vector field, then*

$$I(\mathbf{Y}, \mathbf{Y}) = L''(0),$$

where  $L$  denotes the length function of  $V$ .  $I$  is called the index form of  $\mathbf{c}$ .

*Proof.* Uniqueness follows from the fact that  $I$  is determined by its restriction to the diagonal  $\Delta = \{(\mathbf{Y}, \mathbf{Y}) \mid \mathbf{Y} \in \mathfrak{V}_c\}$  in  $\mathfrak{V}_c \times \mathfrak{V}_c$ , because  $I(\mathbf{Y}, \mathbf{Z}) = 1/2(I(\mathbf{Y} + \mathbf{Z}, \mathbf{Y} + \mathbf{Z}) - I(\mathbf{Y}, \mathbf{Y}) - I(\mathbf{Z}, \mathbf{Z}))$  for any  $\mathbf{Y}, \mathbf{Z} \in \mathfrak{V}_c$ . Existence follows from (6.3.4), which implies that

$$I(\mathbf{Y}, \mathbf{Z}) = \int_0^a (\langle \mathbf{Y}', \mathbf{Z}' \rangle - \langle R(\mathbf{Y}, \dot{\mathbf{c}})\dot{\mathbf{c}}, \mathbf{Z} \rangle). \quad (6.4.1)$$

This formula also shows that  $I$  is indeed symmetric and bilinear.  $\square$

Notice that any piecewise smooth vector field  $\mathbf{Y}$  along  $\mathbf{c}$  is the variational vector field of some piecewise smooth variation  $V$  of  $\mathbf{c}$ , namely

$$V(t, s) = \exp(s\mathbf{Y}(t)), \quad (6.4.2)$$

where of course,  $s$  is chosen small enough so that the exponential map is defined.

**Proposition 6.4.2.** *Let  $\mathbf{c} : [0, a] \rightarrow M$  be a normal geodesic,  $\mathbf{Y}, \mathbf{Z} \in \mathfrak{V}_c$ . If  $\mathbf{Z}$  is differentiable, then*

$$I(\mathbf{Y}, \mathbf{Z}) = - \int_0^a \langle \mathbf{Y}, \mathbf{Z}'' + R(\mathbf{Z}, \dot{\mathbf{c}})\dot{\mathbf{c}} \rangle. \quad (6.4.3)$$

*In particular, if  $\mathbf{Z}$  is a Jacobi field, then*

$$I(\mathbf{Y}, \mathbf{Z}) = 0. \quad (6.4.4)$$

*Proof.* By assumption, there exists a partition  $0 = t_0 < t_1 < \cdots < t_l = a$  of  $[0, a]$  such that  $\mathbf{Y}$  is differentiable on  $[t_{i-1}, t_i]$ ,  $1 \leq i \leq l$ . Since  $\langle \mathbf{Y}, \mathbf{Z}' \rangle'(t) = \langle \mathbf{Y}', \mathbf{Z}' \rangle(t) + \langle \mathbf{Y}, \mathbf{Z}'' \rangle(t)$  if  $t \neq t_i$ ,

$$\begin{aligned} \int_0^a \langle \mathbf{Y}', \mathbf{Z}' \rangle &= \sum_{i=1}^l \left( \int_{t_{i-1}}^{t_i} \langle \mathbf{Y}, \mathbf{Z}' \rangle' \right) - \int_0^a \langle \mathbf{Y}, \mathbf{Z}'' \rangle \\ &= \sum_{i=1}^l \langle \mathbf{Y}, \mathbf{Z}' \rangle \Big|_{t_{i-1}}^{t_i} - \int_0^a \langle \mathbf{Y}, \mathbf{Z}'' \rangle = \langle \mathbf{Y}, \mathbf{Z}' \rangle \Big|_0^a - \int_0^a \langle \mathbf{Y}, \mathbf{Z}'' \rangle \\ &= - \int_0^a \langle \mathbf{Y}, \mathbf{Z}'' \rangle, \end{aligned}$$

because  $\mathbf{Y}$  vanishes at the endpoints. Substituting this expression in (6.4.1) and using the identity  $\langle R(\mathbf{Y}, \dot{\mathbf{c}})\dot{\mathbf{c}}, \mathbf{Z} \rangle = \langle R(\mathbf{Z}, \dot{\mathbf{c}})\dot{\mathbf{c}}, \mathbf{Y} \rangle$  now yields (6.4.3); the latter immediately implies (6.4.4).  $\square$

**Remark 6.4.1.** The index form extends to the larger space of all piecewise smooth vector fields along  $\mathbf{c}$  (not necessarily zero at the endpoints) by the same formula (6.4.1). The proof of Proposition 6.4.2 shows that for vector fields  $\mathbf{Y}, \mathbf{Z}$  in this larger space,

$$I(\mathbf{Y}, \mathbf{Z}) = \langle \mathbf{Y}, \mathbf{Z}' \rangle \Big|_0^a - \int_0^a \langle \mathbf{Y}, \mathbf{Z}'' + R(\mathbf{Z}, \dot{\mathbf{c}})\dot{\mathbf{c}} \rangle, \quad (6.4.5)$$

if  $\mathbf{Z}$  is differentiable, and

$$I(\mathbf{Y}, \mathbf{Z}) = \langle \mathbf{Y}, \mathbf{Z}' \rangle \Big|_0^a \quad (6.4.6)$$

if in addition  $\mathbf{Z}$  is Jacobi. We will refer to it as the *extended index form* of  $\mathbf{c}$ .

**Theorem 6.4.1.** *If  $\mathbf{c} : [0, a] \rightarrow M$  is a normal geodesic with a conjugate point  $t_0 < a$ , then there exists a variation  $V$  of  $\mathbf{c}$  with fixed endpoints such that the curves  $V_s$ , where  $V_s(t) = V(t, s)$ , are shorter than  $\mathbf{c}$  for sufficiently small  $s$ . In particular,  $\mathbf{c}$  is not a minimal geodesic.*

*Proof.* By assumption, there exists a Jacobi field  $\mathbf{Z}$  along  $\mathbf{c}$  with  $\mathbf{Z}(0) = \mathbf{0}$ ,  $\mathbf{Z}(t_0) = \mathbf{0}$ . Let  $\mathbf{E}$  denote the parallel vector field along  $\mathbf{c}$  that equals  $-\mathbf{Z}'(t_0)$  at  $t_0$ , and define  $\mathbf{Y} := \varphi \mathbf{E} \in \mathfrak{X}_{\mathbf{c}}$ , where  $\varphi : [0, a] \rightarrow [0, 1]$  is some function satisfying  $\varphi(0) = \varphi(a) = 0$ ,  $\varphi(t_0) = 1$ , cf. Lemma 2.2.1. For each  $r > 0$ , define  $\mathbf{Z}_r \in \mathfrak{X}_{\mathbf{c}}$  by

$$\mathbf{Z}_r(t) = \begin{cases} \mathbf{Z}(t) + r\mathbf{Y}(t) & \text{if } t \leq t_0, \\ r\mathbf{Y}(t) & \text{if } t \geq t_0. \end{cases}$$

Let  $I_0$  denote the extended index form (in the sense of Remark 6.4.1) of the restriction of  $\mathbf{c}$  to  $[0, t_0]$ . Then

$$I(\mathbf{Z}_r, \mathbf{Z}_r) = I_0(\mathbf{Z}_r, \mathbf{Z}_r) + \int_{t_0}^a \left( \langle r\mathbf{Y}', r\mathbf{Y}' \rangle - \langle R(r\mathbf{Y}, \dot{\mathbf{c}})\dot{\mathbf{c}}, r\mathbf{Y} \rangle \right),$$

and recalling that  $\mathbf{Y}(t_0) = -\mathbf{Z}'(t_0)$ ,

$$\begin{aligned} I_0(\mathbf{Z}_r, \mathbf{Z}_r) &= I_0(\mathbf{Z}, \mathbf{Z}) + 2I_0(r\mathbf{Y}, \mathbf{Z}) + I_0(r\mathbf{Y}, r\mathbf{Y}) \\ &= \langle \mathbf{Z}, \mathbf{Z}' \rangle \Big|_0^{t_0} + 2r \langle \mathbf{Y}, \mathbf{Z}' \rangle \Big|_0^{t_0} + r^2 I_0(\mathbf{Y}, \mathbf{Y}) \\ &= -2r |\mathbf{Z}'|^2(t_0) + r^2 I_0(\mathbf{Y}, \mathbf{Y}), \end{aligned}$$

so that

$$I(\mathbf{Z}_r, \mathbf{Z}_r) = -2r |\mathbf{Z}'|^2(t_0) + r^2 I(\mathbf{Y}, \mathbf{Y}) = r(-2 |\mathbf{Z}'|^2(t_0) + r I(\mathbf{Y}, \mathbf{Y})).$$

$\mathbf{Z}'(t_0) \neq 0$  since  $\mathbf{Z}$  is not identically zero, so the above expression is negative for sufficiently small  $r > 0$ . By Proposition 6.4.1 and (6.4.2), the length function  $L$  of the



corresponding variation  $(t, s) \mapsto \exp(s\mathbf{Z}_r(t))$  satisfies  $L''(0) < 0$ , and the length of  $\mathbf{c} = V_0$  is then a strict local maximum.  $\square$

The *kernel* of a bilinear form  $b$  on a vector space  $E$  is defined to be the set

$$\ker b = \{\mathbf{u} \in E \mid b(\mathbf{u}, \mathbf{v}) = 0 \text{ for all } \mathbf{v} \in E\}.$$

It is clearly a subspace of  $E$ . Even though the vector space  $\mathfrak{Y}_{\mathbf{c}}$  is infinite-dimensional, the kernel of the index form has finite dimension:

**Lemma 6.4.1.** *If  $\mathbf{c} : [0, a] \rightarrow M$  is a normal geodesic, then the kernel of the index form  $I$  of  $\mathbf{c}$  consists of all Jacobi fields  $\mathbf{Y} \in \mathfrak{Y}_{\mathbf{c}}$ .*

*Proof.* By (6.4.4), every Jacobi field  $\mathbf{Y} \in \mathfrak{Y}_{\mathbf{c}}$  belongs to the kernel of  $I$ . Conversely, suppose that  $\mathbf{Z}$  lies in the kernel, and consider a partition  $0 = t_0 < t_1 < \cdots < t_l = a$  of  $[0, a]$  such that each  $\mathbf{Z}_i := \mathbf{Z}|_{[t_{i-1}, t_i]}$  is differentiable. If  $\varphi_i : [t_{i-1}, t_i] \rightarrow \mathbb{R}$  is a smooth function that vanishes at the endpoints and is positive elsewhere,  $i = 1, \dots, l$ , define a vector field  $\mathbf{Y}_i = \varphi_i(\mathbf{Z}_i'' + R(\mathbf{Z}_i, \dot{\mathbf{c}})\dot{\mathbf{c}})$  along  $\mathbf{c}_i := \mathbf{c}|_{[t_{i-1}, t_i]}$ , and  $\mathbf{Y} \in \mathfrak{Y}_{\mathbf{c}}$  by  $\mathbf{Y}|_{[t_{i-1}, t_i]} = \mathbf{Y}_i$ . Let  $I_i$  denote the extended index form of  $\mathbf{c}_i$ . Since each  $\mathbf{Z}_i$  is differentiable, (6.4.3) implies that

$$\begin{aligned} 0 = I(\mathbf{Y}, \mathbf{Z}) &= \sum_{i=1}^l I_i(\mathbf{Y}_i, \mathbf{Z}_i) = - \sum_{i=1}^l \int_{t_{i-1}}^{t_i} \langle \mathbf{Y}_i, \mathbf{Z}_i'' + R(\mathbf{Z}_i, \dot{\mathbf{c}})\dot{\mathbf{c}} \rangle \\ &= - \sum_{i=1}^l \int_{t_{i-1}}^{t_i} \varphi_i |\mathbf{Z}_i'' + R(\mathbf{Z}_i, \dot{\mathbf{c}})\dot{\mathbf{c}}|^2, \end{aligned}$$

so that each  $\mathbf{Z}_i$  is Jacobi. To see that  $\mathbf{Z}$  itself is a Jacobi field, it suffices to show that it is differentiable at  $t_i$ ,  $i = 1, \dots, l-1$ . This in turn will follow once we establish that  $\mathbf{Z}'_i(t_i) = \mathbf{Z}'_{i+1}(t_i)$ , since a Jacobi field is uniquely determined by its value and the value of its derivative at any one point. To do so, fix some  $i$ , denote by  $\mathbf{E}$  the parallel field along  $\mathbf{c}$  with  $\mathbf{E}(t_i) = \mathbf{Z}'_{i+1}(t_i) - \mathbf{Z}'_i(t_i)$ , and let  $\varphi : [0, b] \rightarrow \mathbb{R}$  be a smooth nonnegative function that equals 1 at  $t_i$  and has its support inside  $(t_{i-1}, t_{i+1})$ . If  $\mathbf{X} := \varphi\mathbf{E} \in \mathfrak{Y}_{\mathbf{c}}$ , then

$$\begin{aligned} 0 = I(\mathbf{X}, \mathbf{Z}) &= I_i(\mathbf{X}|_{[t_{i-1}, t_i]}, \mathbf{Z}_i) + I_{i+1}(\mathbf{X}|_{[t_i, t_{i+1}]}, \mathbf{Z}_{i+1}) \\ &= \langle \mathbf{Z}'_{i+1}(t_i) - \mathbf{Z}'_i(t_i), \mathbf{Z}'_i(t_i) \rangle - \langle \mathbf{Z}'_{i+1}(t_i) - \mathbf{Z}'_i(t_i), \mathbf{Z}'_{i+1}(t_i) \rangle \\ &= -|\mathbf{Z}'_{i+1}(t_i) - \mathbf{Z}'_i(t_i)|^2, \end{aligned}$$

so that  $\mathbf{Z}$  is indeed Jacobi.  $\square$

We are now able to characterize geodesics without conjugate points in terms of the index form:

**Theorem 6.4.2.** *A normal geodesic  $\mathbf{c}$  has no conjugate points if and only if its index form is positive definite on  $\mathfrak{Y}_{\mathbf{c}}$ .*

*Proof.* Notice that if  $\mathbf{c}$  has no conjugate points, then  $I(\mathbf{X}, \mathbf{X}) \geq 0$  for any  $\mathbf{X} \in \mathfrak{V}_{\mathbf{c}}$ : indeed, if  $I(\mathbf{X}, \mathbf{X}) < 0$ , then the curves  $V_s$  in the variation  $(t, s) \mapsto V(t, s) := \exp(s\mathbf{X}(t))$  would be shorter than  $\mathbf{c}$  for small  $s$ , contradicting Proposition 6.3.2. Next, suppose  $I(\mathbf{X}, \mathbf{X}) = 0$ . We must show that  $\mathbf{X}$  is identically zero. But for any  $\mathbf{Y} \in \mathfrak{V}_{\mathbf{c}}$ ,

$$0 \leq I(\mathbf{X} + t\mathbf{Y}, \mathbf{X} + t\mathbf{Y}) = t(tI(\mathbf{Y}, \mathbf{Y}) + 2I(\mathbf{X}, \mathbf{Y})),$$

which is only possible if  $I(\mathbf{X}, \mathbf{Y}) = 0$ . Thus,  $\mathbf{X}$  is Jacobi by the Lemma, and since  $\mathbf{c}$  has no conjugate points,  $\mathbf{X} \equiv \mathbf{0}$ .

Conversely, assume  $I$  is positive definite on  $\mathfrak{V}_{\mathbf{c}}$ .  $\mathbf{c}$  cannot have a conjugate point  $t_0 < b$  by Theorem 6.4.1. Nor can it have  $b$  as conjugate point, for otherwise there would exist a nonzero Jacobi field  $\mathbf{Y} \in \mathfrak{V}_{\mathbf{c}}$ , and  $I(\mathbf{Y}, \mathbf{Y}) = 0$  by the Lemma.  $\square$

**Remarks 6.4.2.** (i) Given  $\mathbf{p}, \mathbf{q} \in M$ , it is tempting to view the collection  $\Omega_{\mathbf{p}, \mathbf{q}}$  of all piecewise smooth curves from  $\mathbf{p}$  to  $\mathbf{q}$  as a manifold. We have not defined infinite-dimensional manifolds, and are therefore not in a position to formalize this, but the analogy is nevertheless suggestive: a ‘point’ in this manifold is a curve  $\mathbf{c}$  from  $\mathbf{p}$  to  $\mathbf{q}$ , a ‘curve’ through the point is a variation of  $\mathbf{c}$ , and its ‘tangent vector’ at the point is the variational vector field. Thus, the tangent space of the point is  $\mathfrak{V}_{\mathbf{c}}$ . The length function  $L : \Omega_{\mathbf{p}, \mathbf{q}} \rightarrow \mathbb{R}$  assigns to each point  $\mathbf{c}$  the length of  $\mathbf{c}$ , and (6.3.3) says that geodesics are ‘critical points’ of  $L$ . The index form then corresponds to the ‘Hessian’ of  $L$ .

(ii) Recall the extended index form of a geodesic  $\mathbf{c}$  introduced in Remark 6.4.1. In the space of all piecewise smooth vector fields along  $\mathbf{c}$ , Jacobi fields minimize this extended index form in the following sense: suppose  $\mathbf{c} : [0, a] \rightarrow M$  has no conjugate points. Notice that for any  $\mathbf{u} \in M_{\mathbf{c}(0)}$  and  $\mathbf{v} \in M_{\mathbf{c}(a)}$  there exists a unique Jacobi field  $\mathbf{Y}$  along  $\mathbf{c}$  with  $\mathbf{Y}(0) = \mathbf{u}$  and  $\mathbf{Y}(a) = \mathbf{v}$ , because the linear map

$$\begin{aligned} \mathcal{J}_{\mathbf{c}} &\longrightarrow M_{\mathbf{c}(0)} \times M_{\mathbf{c}(a)}, \\ \mathbf{Y} &\longmapsto (\mathbf{Y}(0), \mathbf{Y}(a)) \end{aligned}$$

has trivial kernel and is then an isomorphism (since both spaces have the same dimension). We claim that if  $\mathbf{X}$  is a vector field along  $\mathbf{c}$  with  $\mathbf{X}(0) = \mathbf{Y}(0)$  and  $\mathbf{X}(a) = \mathbf{Y}(a)$ , then

$$I(\mathbf{X}, \mathbf{X}) \geq I(\mathbf{Y}, \mathbf{Y}),$$

and inequality is strict provided  $\mathbf{X} \neq \mathbf{Y}$ . In fact, if  $\mathbf{X} - \mathbf{Y} \neq \mathbf{0}$ , then it is a nontrivial element of  $\mathfrak{V}_{\mathbf{c}}$ , and by Theorem 6.4.2,

$$0 < I(\mathbf{X} - \mathbf{Y}, \mathbf{X} - \mathbf{Y}) = I(\mathbf{X}, \mathbf{X}) + I(\mathbf{Y}, \mathbf{Y}) - 2I(\mathbf{X}, \mathbf{Y}).$$

The claim follows from this inequality once we observe that  $I(\mathbf{X}, \mathbf{Y}) = I(\mathbf{Y}, \mathbf{Y})$ . The latter identity, in turn, holds because by (6.4.6),

$$I(\mathbf{X}, \mathbf{Y}) = \langle \mathbf{X}, \mathbf{Y}' \rangle \Big|_0^a = \langle \mathbf{Y}, \mathbf{Y}' \rangle \Big|_0^a = I(\mathbf{Y}, \mathbf{Y}).$$

## 6.5 The distance function

A manifold  $M$  is said to be *connected* if any two points of  $M$  can be joined by a curve lying in the manifold. Notice that regardless of whether or not it is connected, a manifold is always locally connected, by which we mean that any point has a connected neighborhood, namely the domain of an appropriate chart about the point.

Unless specified otherwise, all spaces will be assumed to be connected. Given  $\mathbf{p}, \mathbf{q} \in M$ , let  $\Omega_{\mathbf{p},\mathbf{q}}$  denote the collection of all curves  $\mathbf{c} : [0, 1] \rightarrow M$  with  $\mathbf{c}(0) = \mathbf{p}$  and  $\mathbf{c}(1) = \mathbf{q}$ . As usual,  $L(\mathbf{c})$  is the length of  $\mathbf{c}$ .

**Definition 6.5.1.** The *distance* between  $\mathbf{p}$  and  $\mathbf{q}$  is the number

$$d(\mathbf{p}, \mathbf{q}) = \inf\{L(\mathbf{c}) \mid \mathbf{c} \in \Omega_{\mathbf{p},\mathbf{q}}\}.$$

The distance is well defined, since it is the infimum of a nonempty set (by connectedness of  $M$ ) bounded below (by zero). It is not true, in general, that the infimum is a minimum; i.e., there need not exist a curve from  $\mathbf{p}$  to  $\mathbf{q}$  whose length equals the distance between them (consider for example  $M = \mathbb{R}^n \setminus \{\mathbf{0}\}$ , any nonzero  $\mathbf{p}$ , and  $\mathbf{q} = -\mathbf{p}$ ). If such a curve exists, however, then the (length of the) restriction of  $\mathbf{c}$  to any subinterval  $[t_1, t_2] \subset [0, 1]$  also realizes the distance between the corresponding endpoints, so that by Theorem 6.1.1,  $\mathbf{c}$  is a geodesic.

**Theorem 6.5.1.**  $(M, d)$  is a metric space. Furthermore, the open sets in the metric space coincide with the usual open sets; i.e.,  $U \subset M$  is open in  $(M, d)$  if and only if  $U = V \cap M$  for some open set  $V$  in Euclidean space.

*Proof.* The first axiom for a metric space,  $d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p})$  for all  $\mathbf{p}$  and  $\mathbf{q}$  is clear, since  $\mathbf{c}$  is a curve from  $\mathbf{p}$  to  $\mathbf{q}$  if and only if  $-\mathbf{c}$ , where  $-\mathbf{c}(t) = \mathbf{c}(1 - t)$ , is a curve from  $\mathbf{q}$  to  $\mathbf{p}$ , and the two have the same length. The triangle inequality  $d(\mathbf{p}, \mathbf{r}) \leq d(\mathbf{p}, \mathbf{q}) + d(\mathbf{q}, \mathbf{r})$  for any three points  $\mathbf{p}, \mathbf{q}$ , and  $\mathbf{r}$  is also easy to see: given any  $\varepsilon > 0$ , there exists a curve  $\mathbf{c}_1$  from  $\mathbf{p}$  to  $\mathbf{q}$  with length smaller than  $d(\mathbf{p}, \mathbf{q}) + \varepsilon/2$ , and similarly a curve  $\mathbf{c}_2$  from  $\mathbf{q}$  to  $\mathbf{r}$  with length less than  $d(\mathbf{q}, \mathbf{r}) + \varepsilon/2$ . Then  $\mathbf{c}$ , where  $\mathbf{c}(t) = \mathbf{c}_1(2t)$  for  $0 \leq t \leq 1/2$  and  $\mathbf{c}(t) = \mathbf{c}_2(2t - 1)$  for  $1/2 \leq t \leq 1$ , is a curve from  $\mathbf{p}$  to  $\mathbf{r}$  of length less than  $d(\mathbf{p}, \mathbf{q}) + d(\mathbf{q}, \mathbf{r}) + \varepsilon$ . This implies that  $d(\mathbf{p}, \mathbf{r}) \leq d(\mathbf{p}, \mathbf{q}) + d(\mathbf{q}, \mathbf{r}) + \varepsilon$  for any  $\varepsilon > 0$ , and the triangle inequality follows.

The last condition that must be satisfied is  $d(\mathbf{p}, \mathbf{q}) \geq 0$ , and  $d(\mathbf{p}, \mathbf{q}) = 0$  if and only if  $\mathbf{p} = \mathbf{q}$ . The first part is clear, as is the fact that  $d(\mathbf{p}, \mathbf{p}) = 0$ . So assume  $d(\mathbf{p}, \mathbf{q}) = 0$ . By Theorem 6.1.1, there exists  $\varepsilon > 0$  such that  $\exp_{\mathbf{p}}$  maps the open ball  $U_\varepsilon = \{\mathbf{v} \in M_{\mathbf{p}} \mid |\mathbf{v}| < \varepsilon\}$  of radius  $\varepsilon$  centered at the origin in the tangent space at  $\mathbf{p}$  diffeomorphically onto  $V = \exp_{\mathbf{p}}(U_\varepsilon)$ . Furthermore, any curve originating at  $\mathbf{p}$  that leaves  $V$  has length greater than  $\varepsilon$ , so that  $\mathbf{q} \in V$ . But then there exists a shortest curve from  $\mathbf{p}$  to  $\mathbf{q}$  which is a geodesic of length  $|(\exp_{\mathbf{p}|U_\varepsilon})^{-1}\mathbf{q}|$ . Thus,  $(\exp_{\mathbf{p}|U_\varepsilon})^{-1}\mathbf{q} = \mathbf{0} \in M_{\mathbf{p}}$ , and  $\mathbf{q} = \mathbf{p}$  as claimed. It remains to establish that  $U \subset M$  is open in the usual sense if and only if it is open in the metric space  $(M, d)$ . The latter, we recall, means that for any  $\mathbf{p} \in U$  there exists

$\varepsilon > 0$  such that  $B_\varepsilon^M(\mathbf{p}) \subset U$ , where

$$B_\varepsilon^M(\mathbf{p}) = \{\mathbf{q} \in M \mid d(\mathbf{p}, \mathbf{q}) < \varepsilon\}.$$

We denote by  $B_\varepsilon(\mathbf{p}) = \{\mathbf{q} \in \mathbb{R}^{n+k} \mid |\mathbf{q} - \mathbf{p}| < \varepsilon\}$  the corresponding distance ball in Euclidean space. Suppose first that  $U \subset M$  is open in the usual sense. Given  $\mathbf{p} \in U$ , there exists  $r > 0$  such that  $B_r(\mathbf{p}) \cap M \subset U$ . But for  $\mathbf{q} \in M$ ,  $d(\mathbf{p}, \mathbf{q}) \geq |\mathbf{p} - \mathbf{q}|$ , so that  $B_r^M(\mathbf{p}) \subset B_r(\mathbf{p}) \cap M \subset U$ , and  $U$  is open in  $(M, d)$ . For the converse, it is enough to show that  $B_r^M(\mathbf{p})$  is open in the usual sense for any  $\mathbf{p} \in M$  and  $r > 0$ . This will in turn follow from the fact that the distance function  $d_{\mathbf{p}} : M \rightarrow [0, \infty)$  from  $\mathbf{p}$ ,  $d_{\mathbf{p}}(\mathbf{q}) := d(\mathbf{p}, \mathbf{q})$ , is continuous, since  $B_r^M(\mathbf{p}) = d_{\mathbf{p}}^{-1}[0, r)$ . To show this, we must establish that if  $\{\mathbf{q}_i\}$  is a sequence in  $M$  converging to  $\mathbf{q}$ , then  $d(\mathbf{p}, \mathbf{q}_i) \rightarrow d(\mathbf{p}, \mathbf{q})$ . Now,

$$d(\mathbf{p}, \mathbf{q}) - d(\mathbf{q}, \mathbf{q}_i) \leq d(\mathbf{p}, \mathbf{q}_i) \leq d(\mathbf{p}, \mathbf{q}) + d(\mathbf{q}, \mathbf{q}_i)$$

by the triangle inequality, so it is actually enough to show that  $d(\mathbf{p}, \mathbf{p}_i) \rightarrow 0$  if  $\mathbf{p}_i \rightarrow \mathbf{p}$ . If  $U_\varepsilon$  is as above, and  $\mathbf{p}_i \rightarrow \mathbf{p}$ , then  $\mathbf{p}_i \in \exp_{\mathbf{p}}(U_\varepsilon)$  for large enough  $i$ , and  $d(\mathbf{p}, \mathbf{p}_i) = |(\exp_{\mathbf{p}|U_\varepsilon})^{-1}\mathbf{p}_i|$ . The claim then follows from the continuity of  $|(\exp_{\mathbf{p}|U_\varepsilon})^{-1}|$ .  $\square$

The second statement in the above theorem essentially says that even though the distance on  $M$  is not the same as that of the ambient space, any map  $\mathbf{f} : M \rightarrow \mathbb{R}^k$  is continuous in the usual sense if and only if it is continuous as a map from the metric space  $(M, d)$ .

**Corollary 6.5.1.** *If  $U_\varepsilon = \{\mathbf{v} \in M_{\mathbf{p}} \mid |\mathbf{v}| < \varepsilon\}$  is a neighborhood of the origin in  $M_{\mathbf{p}}$  on which  $\exp_{\mathbf{p}}$  is a diffeomorphism, then  $\exp_{\mathbf{p}}(U_\varepsilon) = B_\varepsilon^M(\mathbf{p})$ .*

*Proof.* It is clear that the left side is contained in the right one, since any point in  $\exp_{\mathbf{p}}(U_\varepsilon)$  can be joined to  $\mathbf{p}$  by a geodesic of length smaller than  $\varepsilon$ . To show that they are equal, we only need to establish that  $\exp_{\mathbf{p}}(U_\varepsilon)$  is both open and closed, and apply Theorem 2.4.4, since  $B_\varepsilon^M(\mathbf{p})$  is connected (recall that any  $\mathbf{q} \in B_\varepsilon^M(\mathbf{p})$  can be joined to  $\mathbf{p}$  by means of a curve of length smaller than  $\varepsilon$ , so that this curve lies entirely in  $B_\varepsilon^M(\mathbf{p})$ ). By assumption,  $(\exp_{\mathbf{p}|U_\varepsilon})^{-1}$  is continuous, so that  $\exp_{\mathbf{p}}(U_\varepsilon)$  is open. To see that it is closed, consider a boundary point  $\mathbf{q}$  of  $\exp_{\mathbf{p}}(U_\varepsilon)$  in  $B_\varepsilon^M(\mathbf{p})$  and a sequence  $\{\mathbf{q}_i\}$  contained in  $\exp_{\mathbf{p}}(U_\varepsilon)$  that converges to  $\mathbf{q}$ , see also Exercise 1.27. The sequence  $\mathbf{v}_i := (\exp_{\mathbf{p}|U_\varepsilon})^{-1}(\mathbf{q}_i)$ , being bounded, contains a subsequence  $\{\mathbf{v}_{i_j}\}$  that converges to some  $\mathbf{v} \in \bar{U}_\varepsilon$ . Now,  $|\mathbf{v}_i| = d(\mathbf{p}, \mathbf{q}_i) \rightarrow d(\mathbf{p}, \mathbf{q})$ , so the subsequence also converges in norm to  $d(\mathbf{p}, \mathbf{q})$ . This means that  $|\mathbf{v}| = d(\mathbf{p}, \mathbf{q}) < \varepsilon$ , and  $\mathbf{v} \in U_\varepsilon$ . We then have that  $\mathbf{q}_{i_j} = \exp_{\mathbf{p}}(\mathbf{v}_{i_j}) \rightarrow \exp_{\mathbf{p}}(\mathbf{v})$ , and since  $\mathbf{q}_i \rightarrow \mathbf{q}$ ,  $\exp_{\mathbf{p}}(\mathbf{v}) = \mathbf{q}$ . This shows that  $\exp_{\mathbf{p}}(U_\varepsilon)$  is closed in  $B_\varepsilon^M(\mathbf{p})$  and concludes the argument.  $\square$

We will shortly see that the conclusion of the above corollary is still true under the much weaker assumption that  $\exp_{\mathbf{p}}$  is defined on  $U_\varepsilon$ . When there is no danger of confusing  $M$  with the ambient Euclidean space, we will denote  $B_\varepsilon^M(\mathbf{p})$  by  $B_\varepsilon(\mathbf{p})$ .

**Proposition 6.5.1.** *Let  $W$  denote a neighborhood of the zero section in  $TM$  such that the restriction of  $(\pi, \exp)$  to  $W$  is a diffeomorphism onto its image, cf. Remark 3.7.1. Given  $\mathbf{p} \in M$ , choose  $\varepsilon > 0$  small enough that  $B_\varepsilon(\mathbf{p}) \times B_\varepsilon(\mathbf{p}) \subset (\pi, \exp)(W)$ . Then*

$$d(\mathbf{q}_1, \mathbf{q}_2) = |((\pi, \exp)|_W)^{-1}(\mathbf{q}_1, \mathbf{q}_2)|, \quad \mathbf{q}_1, \mathbf{q}_2 \in B_{\varepsilon/3}(\mathbf{p}) \times B_{\varepsilon/3}(\mathbf{p}).$$

*In particular,  $d^2 : M \times M \rightarrow \mathbb{R}$  is differentiable on a neighborhood  $U$  of the diagonal  $\Delta = \{(\mathbf{p}, \mathbf{p}) \mid \mathbf{p} \in M\}$  in  $M \times M$ , and the distance function  $d$  itself is differentiable on  $U \setminus \Delta$ .*

*Proof.* Given  $\mathbf{p} \in M$ ,  $U_\varepsilon^{\mathbf{p}}$  will denote the open ball of radius  $\varepsilon$  centered at the origin in  $M_{\mathbf{p}}$ . Notice that if  $\mathbf{q}_1 \in B_{\varepsilon/3}(\mathbf{p})$ , then  $\exp_{\mathbf{q}_1}(U_{2\varepsilon/3}^{\mathbf{q}_1}) \subset B_{2\varepsilon/3}(\mathbf{q}_1) \subset B_\varepsilon(\mathbf{p})$ , so that the restriction of  $\exp_{\mathbf{q}_1}$  to  $U_{2\varepsilon/3}^{\mathbf{q}_1}$  is a diffeomorphism onto its image  $B_{2\varepsilon/3}(\mathbf{q}_1)$  by the corollary above. But if  $\mathbf{q}_2 \in B_{\varepsilon/3}(\mathbf{p})$ , then it also belongs to  $B_{2\varepsilon/3}(\mathbf{q}_1)$ , so that

$$d(\mathbf{q}_1, \mathbf{q}_2) = |(\exp_{\mathbf{q}_1}|_{U_{2\varepsilon/3}^{\mathbf{q}_1}})^{-1} \mathbf{q}_2| = |((\pi, \exp)|_W)^{-1}(\mathbf{q}_1, \mathbf{q}_2)|.$$

The remaining assertions then follow from the fact that in an inner product space, the norm function is differentiable away from the origin, see Exercise 2.8.  $\square$

It is in general not true, however, that the distance function is smooth everywhere outside the diagonal: consider for example  $M = S^1$ . One of the possible two normal geodesics emanating from  $\mathbf{p} = (1, 0)$  is  $\mathbf{c}$ , where  $\mathbf{c}(t) = (\cos t, \sin t)$ . Then  $\mathbf{c}$  minimizes up to time  $\pi$ , so that

$$d(\mathbf{p}, \mathbf{c}(t)) = \begin{cases} L(\mathbf{c}|_{[0,t]}) = t & \text{if } t < \pi, \\ L(\mathbf{c}|_{[t,2\pi]}) = 2\pi - t & \text{if } t > \pi. \end{cases}$$

We leave it as an exercise to show that this implies that the distance function squared  $d^2 : S^1 \times S^1 \rightarrow \mathbb{R}$  is not differentiable exactly on the anti-diagonal  $\{(\mathbf{q}, -\mathbf{q}) \mid \mathbf{q} \in S^1\}$ .

## 6.6 The Hopf-Rinow theorem

Recall that a metric space is said to be complete if every Cauchy sequence in the space converges. One of the most striking results that relates metric and differential geometric properties of a manifold is the Hopf-Rinow theorem, which roughly says that completeness as a metric space is equivalent to geodesics being defined for all time; the latter property is of course equivalent to completeness of the geodesic spray vector field introduced in Section 3.7.

Fix a point  $\mathbf{p}$  in a manifold  $M$ , and for  $\varepsilon > 0$ , denote by  $U_\varepsilon \subset M_{\mathbf{p}}$  the open ball of radius  $\varepsilon$  centered at the origin in the tangent space at  $\mathbf{p}$ .

**Definition 6.6.1.** The *injectivity radius at  $\mathbf{p}$*  is defined as

$$\text{inj}_{\mathbf{p}} = \sup\{r > 0 \mid \exp_{\mathbf{p}} : U_r \rightarrow B_r(\mathbf{p}) \text{ is a diffeomorphism}\},$$

provided the set on the right is bounded above, and as  $\infty$  if it is unbounded. The injectivity radius of a subset  $A \subset M$  is defined as

$$\text{inj}_A = \inf\{\text{inj}_p \mid p \in A\}.$$

- Examples and Remarks 6.6.1.** (i) When  $M = \mathbb{R}^n$ ,  $\text{inj}_p = \infty$  for every  $p$ .  
 (ii) On  $M = \mathbb{R}^n \setminus \{0\}$ ,  $\text{inj}_p = |p|$ . The injectivity radius of  $M$  itself is zero.  
 (iii) On the sphere of radius  $r$ , the injectivity radius is  $\pi r$  at any point, and therefore so is the injectivity radius of the whole sphere.  
 (iv) If  $q \in B_r(p)$ , where  $r \leq \text{inj}_p$ , then there exists a unique normal geodesic from  $p$  to  $q$  of length equal to  $d(p, q)$ : in fact,  $q \in B_{r'}(p)$ , where  $r' = (d(p, q) + \text{inj}_p)/2 < r$ , because  $d(p, q) \leq (d(p, q) + \text{inj}_p)/2$ . Furthermore,  $r' < r$ , so that  $\exp_p : U_{r'} \rightarrow B_{r'}(p)$  is a diffeomorphism.

**Lemma 6.6.1.** *The injectivity radius of a (nonempty) compact set is positive.*

*Proof.* Suppose, to the contrary, that  $A$  is a compact set with zero injectivity radius. Then there exists a sequence  $\{p_k\} \subset A$  such that  $\text{inj}_{p_k} \rightarrow 0$ . This in turn means that for each natural number  $k$ , there exists some  $v_k \in M_{p_k}$  with  $|v_k| \rightarrow 0$  but  $v_k$  lies outside the open set  $W$  on which  $(\pi, \exp) : W \rightarrow \pi(W) \times \exp(W)$  is a diffeomorphism. Now,  $A$  is compact, so after passing to a subsequence if necessary, it may be assumed that  $\{p_k\}$  converges to some  $p \in A$ . Then  $\{v_k\}$  converges to the zero vector in  $M_p$ , and since  $v_k$  lies in the closed set  $TM \setminus W$ , so does the zero vector. This contradicts the fact that  $W$  is a neighborhood of the zero section.  $\square$

Most of the work required in proving the Hopf-Rinow theorem is contained in the following:

**Proposition 6.6.1.** *Let  $p \in M$ . If  $\exp_p$  is defined on  $U_\varepsilon$ , then there exists a minimal geodesic joining  $p$  to any  $q \in B_\varepsilon(p)$ .*

*Proof.* Denote by  $I$  the set of all  $r \in (0, \varepsilon)$  for which there exists a minimal geodesic from  $p$  to any point in the closure of  $B_r(p)$ .  $I$  is an interval by definition, and is nonempty because the injectivity radius at  $p$  is positive. We will show that  $I = (0, \varepsilon)$  by arguing that it is both open and closed in  $(0, \varepsilon)$ . This will prove the proposition, since  $B_\varepsilon(p) = \bigcup_{r \in (0, \varepsilon)} \overline{B_r(p)}$ .

To see that it is closed, assume  $(0, \delta) \subset I$  for some  $\delta > 0$ . We must show that  $\delta \in I$ ; i.e., that for any  $q$  at distance  $\leq \delta$  from  $p$ , there exists a geodesic  $c : [0, 1] \rightarrow M$  from  $p$  to  $q$  with length  $L(c) = d(p, q)$ . So choose some sequence  $\{q_k\} \subset B_\delta(p)$  that converges to  $q$ . By assumption, there exists for each  $k$  a minimal geodesic  $c_k : [0, 1] \rightarrow M$  from  $p$  to  $q_k$ . Now,  $|\dot{c}_k(0)| = L(c_k) = d(p, q_k) < \delta$ , so that  $\{\dot{c}_k(0)\}$  is contained inside the compact set  $\overline{U_\delta}$ , and some subsequence  $\{\dot{c}_{k_i}(0)\}$  converges to, say,  $v \in \overline{U_\delta}$ . The geodesic  $c : [0, 1] \rightarrow M$ , where  $c(t) = \exp_p(tv)$ , joins  $p$  to  $q$  because

$$q = \lim_{i \rightarrow \infty} q_{k_i} = \lim_{i \rightarrow \infty} \exp(\dot{c}_{k_i}(0)) = \exp(\lim_{i \rightarrow \infty} \dot{c}_{k_i}(0)) = \exp(v).$$

Furthermore, its length satisfies

$$L(\mathbf{c}) = |\mathbf{v}| = \lim_{i \rightarrow \infty} |\dot{\mathbf{c}}_{k_i}(0)| = \lim_{i \rightarrow \infty} d(\mathbf{p}, \mathbf{q}_{k_i}) = d(\mathbf{p}, \mathbf{q})$$

by continuity of the distance function. This shows that  $I$  is closed.

To see that it is open, assume  $(0, \delta] \subset I$ . Notice that the closure of  $B_\delta(\mathbf{p})$  must then be compact: indeed, by continuity of  $\exp_{\mathbf{p}}$ ,  $\exp_{\mathbf{p}}(\overline{U_\delta})$  is contained in  $\overline{B_\delta(\mathbf{p})}$ , and being compact, is closed. But it also contains  $B_\delta(\mathbf{p})$ , and must therefore also contain its closure. Thus,  $\overline{B_\delta(\mathbf{p})} = \exp_{\mathbf{p}}(\overline{U_\delta})$  is compact. By Lemma 6.6.1, the injectivity radius  $\alpha$  of  $\overline{B_\delta(\mathbf{p})}$  is positive. We claim that if  $0 < \beta < \min\{\alpha, \varepsilon - \delta\}$ , then  $\delta + \beta \in I$ , thereby implying that  $I$  is open. To establish the claim, consider a point  $\mathbf{q} \in \overline{B_{\delta+\beta}(\mathbf{p})} \setminus \overline{B_\delta(\mathbf{p})}$ , and a sequence  $\mathbf{c}_k : [0, 1] \rightarrow M$  of curves from  $\mathbf{p}$  to  $\mathbf{q}$  with length  $\leq d(\mathbf{p}, \mathbf{q}) + 1/k$ . The intermediate value theorem guarantees the existence of a parameter value  $t_k$  for which  $d(\mathbf{p}, \mathbf{c}_k(t_k)) = \delta$ . By compactness of  $\overline{B_\delta(\mathbf{p})}$ , we may assume after passing to a subsequence if necessary that  $\mathbf{r}_k := \mathbf{c}_k(t_k)$  converges to some  $\mathbf{r}$ , which, by continuity of distance, lies at distance  $\delta$  from  $\mathbf{p}$ . Now,

$$d(\mathbf{p}, \mathbf{q}) + \frac{1}{k} \geq L(\mathbf{c}_k) = L(\mathbf{c}_k|_{[0, t_k]}) + L(\mathbf{c}_k|_{[t_k, 1]}) \geq d(\mathbf{p}, \mathbf{r}_k) + d(\mathbf{r}_k, \mathbf{q}), \quad (6.6.1)$$

so that

$$d(\mathbf{p}, \mathbf{q}) \geq d(\mathbf{p}, \mathbf{r}) + d(\mathbf{r}, \mathbf{q}). \quad (6.6.2)$$

Indeed, if  $d(\mathbf{p}, \mathbf{q})$  were smaller than  $d(\mathbf{p}, \mathbf{r}) + d(\mathbf{r}, \mathbf{q})$ , then  $d(\mathbf{p}, \mathbf{q}) < d(\mathbf{p}, \mathbf{r}) + d(\mathbf{r}, \mathbf{q}) - 2\varepsilon$  for some  $\varepsilon > 0$ , and thus, for all large enough  $k$ ,

$$d(\mathbf{p}, \mathbf{q}) + \frac{1}{k} < d(\mathbf{p}, \mathbf{r}) + d(\mathbf{r}, \mathbf{q}) - \varepsilon. \quad (6.6.3)$$

But then (6.6.1) and (6.6.3) would imply that

$$d(\mathbf{p}, \mathbf{r}_k) + d(\mathbf{r}_k, \mathbf{q}) < d(\mathbf{p}, \mathbf{r}) + d(\mathbf{r}, \mathbf{q}) - \varepsilon,$$

which contradicts the fact that  $\{\mathbf{r}_k\}$  converges to  $\mathbf{r}$ . Thus, (6.6.2) holds, and together with the triangle inequality, we obtain

$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{p}, \mathbf{r}) + d(\mathbf{r}, \mathbf{q}). \quad (6.6.4)$$

By assumption, there exists a minimal geodesic  $\mathbf{c}_1 : [0, \delta] \rightarrow M$  from  $\mathbf{p}$  to  $\mathbf{r}$ . Furthermore,  $d(\mathbf{r}, \mathbf{q}) = d(\mathbf{p}, \mathbf{q}) - d(\mathbf{p}, \mathbf{r}) \leq \delta + \beta - \delta = \beta$  is less than the injectivity radius of  $\overline{B_\delta(\mathbf{p})}$ , and since the latter set contains  $\mathbf{r}$ , there exists a minimal geodesic  $\mathbf{c}_2 : [\delta, \delta + d(\mathbf{r}, \mathbf{q})] \rightarrow M$  from  $\mathbf{r}$  to  $\mathbf{q}$ . The *a priori* only piecewise geodesic  $\mathbf{c} : [0, \delta + d(\mathbf{r}, \mathbf{q})] \rightarrow M$ , where  $\mathbf{c}(t) = \mathbf{c}_1(t)$  when  $t \leq \delta$ , and  $\mathbf{c}(t) = \mathbf{c}_2(t)$  when  $t \geq \delta$ , is a curve from  $\mathbf{p}$  to  $\mathbf{q}$  whose length realizes the distance between its endpoints by (6.6.4). It must therefore be a geodesic; i.e.,  $\dot{\mathbf{c}}_1(\delta) = \dot{\mathbf{c}}_2(\delta)$ . This completes the proof.  $\square$

**Remark 6.6.1.** We reiterate a fact observed in the proof of the proposition: if  $\exp_{\mathbf{p}}$  is defined on  $U_\varepsilon$ , then the closure  $\overline{B_r(\mathbf{p})}$  of  $B_r(\mathbf{p})$  is compact for any  $r \in (0, \varepsilon)$ : indeed, continuity of the exponential map implies that  $\exp_{\mathbf{p}}(\overline{U_r})$  is compact, hence closed, and must therefore contain  $\overline{B_r(\mathbf{p})}$ . But it must also be contained inside  $\overline{B_r(\mathbf{p})}$  by continuity of  $\exp_{\mathbf{p}}$  again, so the two sets are equal.

**Theorem 6.6.1.** *The following statements are equivalent in a connected Riemannian manifold  $M$ :*

- (1)  $M$  is complete as a metric space.
- (2) For all  $\mathbf{p} \in M$ ,  $\exp_{\mathbf{p}}$  is defined on all of  $M_{\mathbf{p}}$ .
- (3) For some  $\mathbf{p} \in M$ ,  $\exp_{\mathbf{p}}$  is defined on all of  $M_{\mathbf{p}}$ .
- (4) Any bounded set of  $M$  (with respect to the distance  $d$ ) has compact closure.

Furthermore, completeness of  $M$  implies that any two points  $\mathbf{p}$  and  $\mathbf{q}$  of  $M$  can be joined by a geodesic of length  $d(\mathbf{p}, \mathbf{q})$ .

*Proof.*  $1 \Rightarrow 2$ : Let  $S$  denote the geodesic spray on  $TM$  (see Theorem 3.7.2), and  $\gamma : I \rightarrow TM$  a maximal integral curve of  $S$ , so that  $\mathbf{c} := \pi \circ \gamma$  is a geodesic with  $\dot{\mathbf{c}} = \gamma$ . To show that  $I = \mathbb{R}$ , it is enough to show that it is closed, since it is already open. So consider a sequence  $\{t_k\} \subset I$ ,  $t_k \rightarrow t_0 \in \mathbb{R}$ . The sequence  $\{\mathbf{c}(t_k)\}$  is then a Cauchy sequence in  $M$ , because  $|\dot{\mathbf{c}}|$  is constant equal to some  $a > 0$ , so that

$$d(\mathbf{c}(t_k), \mathbf{c}(t_l)) \leq \left| \int_{t_k}^{t_l} |\dot{\mathbf{c}}| \right| = a|t_k - t_l|,$$

and the Cauchy property for  $\{\mathbf{c}(t_k)\}$  follows from that for  $\{t_k\}$ . By assumption,  $\{\mathbf{c}(t_k)\}$  converges and is therefore contained in some compact set  $K \subset M$ . But then  $\{\gamma(t_k)\}$  lies in the compact set  $\{\mathbf{v} \in TM \mid \pi(\mathbf{v}) \in K \text{ and } |\mathbf{v}| = a\}$ , and has therefore a convergent subsequence. By Theorem 3.3.3,  $t_0 \in I$ .

$2 \Rightarrow 3$ : Immediate.

$3 \Rightarrow 4$ : If  $A$  is bounded, then it is contained inside some closed metric ball of sufficiently large radius centered at  $\mathbf{p}$ . The latter is compact by Remark 6.6.1, so that  $\bar{A}$  is also compact.

$4 \Rightarrow 1$ : We already observed in Chapter 1 that any Cauchy sequence is bounded. It therefore lies inside a compact set by hypothesis, and thus admits a convergent subsequence. Then the sequence itself converges as shown in the proof of Theorem 1.8.4. To complete the proof of the theorem, it suffices to show that the second statement implies that any two points of  $M$  can be joined by a minimal geodesic. This is an immediate consequence of Proposition 6.6.1.  $\square$



### 6.7 Curvature comparison

One of the most active areas of research in Riemannian geometry is the interaction between the shape of a space and its curvature. What can one say about the manifold if its curvature is known? Although this topic is too vast and complex to allow for a comprehensive account here, we wish to establish a couple of results that hint at the beauty of this subject.

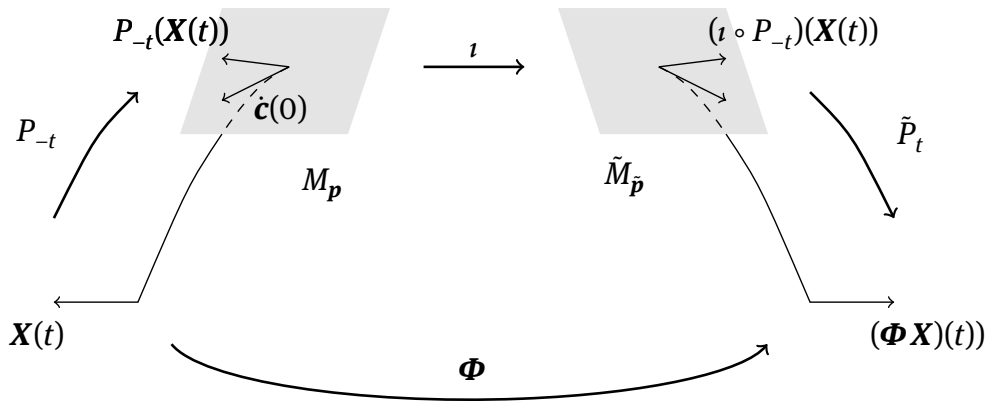
In this section we will see that the larger the curvature of a space, the earlier conjugate points appear along geodesics. This fact will be derived by comparing the index form in spaces where the curvature of one is larger than that of the other. As an application, a complete space with curvature bounded below by a positive constant must be compact. The construction used in the proof also shows how the metric is determined by the curvature.

Denote by  $\mathfrak{X}_c$  the vector space of all piecewise-smooth vector fields along a curve  $c$ . Let  $M, \tilde{M}$  denote manifolds of the same dimension  $n$ ,  $c : [0, b] \rightarrow M$  (resp.  $\tilde{c} : [0, b] \rightarrow \tilde{M}$ ) a normal geodesic in  $M$  (resp.  $\tilde{M}$ ), and  $p = c(0)$ ,  $\tilde{p} = \tilde{c}(0)$  their respective starting points. Suppose  $\iota : M_p \rightarrow \tilde{M}_{\tilde{p}}$  is a linear isometry mapping  $\dot{c}(0)$  to  $\dot{\tilde{c}}(0)$ . If  $P_t : M_p \rightarrow M_{c(t)}$  and  $\tilde{P}_t : \tilde{M}_{\tilde{p}} \rightarrow \tilde{M}_{\tilde{c}(t)}$  denote parallel translation along  $c$  and  $\tilde{c}$  respectively, define a map  $\Phi : \mathfrak{X}_c \rightarrow \mathfrak{X}_{\tilde{c}}$  by setting

$$(\Phi X)(t) = (\tilde{P}_t \circ \iota \circ P_t^{-1})X(t) = (\tilde{P}_t \circ \iota \circ P_{-t})X(t), \quad X \in \mathfrak{X}_c,$$

with  $P_{-t} = P_t^{-1}$  denoting parallel translation along  $-c$  from  $M_{c(t)}$  to  $M_p$ . Notice that  $\Phi \dot{c} = \dot{\tilde{c}}$  and that if  $X$  is parallel along  $c$ , then  $\Phi X$  is the parallel field along  $\tilde{c}$  where  $\Phi X(0) = \iota X(0)$ .

Let  $Z_1, \dots, Z_n = \dot{c}$  denote an orthonormal basis of parallel vector fields along  $c$ , so that  $\tilde{Z}_1, \dots, \tilde{Z}_n$ , with  $\tilde{Z}_i = \Phi Z_i$ , is one along  $\tilde{c}$ . Any piecewise smooth vector field  $X$  along  $c$  may then be written as  $X = \sum_i \varphi_i Z_i$ , with  $\varphi_i = \langle X, Z_i \rangle$  piecewise smooth. By construction,  $\Phi X = \sum_i \varphi_i \tilde{Z}_i$ .



It follows immediately that if  $\mathfrak{X}_c^\perp$  denotes the subspace of  $\mathfrak{X}_c$  consisting of all those orthogonal to  $\mathbf{c}$ , and  $\mathfrak{Y}_c$  the subspace consisting of all the elements in  $\mathfrak{X}_c^\perp$  that vanish at the endpoints, then  $\Phi$  maps  $\mathfrak{X}_c$  isomorphically onto  $\mathfrak{X}_{\tilde{c}}$ , and the same is true for the corresponding subspaces. Furthermore,  $\Phi$  commutes with covariant differentiation, since

$$(\Phi \mathbf{X})' = \sum_i \varphi_i' \tilde{\mathbf{Z}}_i = \Phi(\mathbf{X}'). \tag{6.7.1}$$

In order to state our next result, we introduce some terminology: for each  $t \in [0, b]$ , let  $\Phi_t : M_{\mathbf{c}(t)} \rightarrow \tilde{M}_{\tilde{\mathbf{c}}(t)}$  be the linear isometry  $\tilde{P}_t \circ \iota \circ P_t^{-1}$ . Thus,  $(\Phi \mathbf{X})(t) = \Phi_t(\mathbf{X}(t))$  for any vector field  $\mathbf{X}$  along  $\mathbf{c}$ .

**Proposition 6.7.1.** *Suppose that for any  $t \in [0, b]$  and for any 2-plane  $E_t \subset M_{\mathbf{c}(t)}$  containing  $\dot{\mathbf{c}}(t)$ , the sectional curvature  $K_{E_t}$  of  $E_t$  is greater than or equal to that of  $\Phi_t(E_t)$ . Then*

$$I(\mathbf{X}, \mathbf{X}) \leq I(\Phi \mathbf{X}, \Phi \mathbf{X}), \quad \mathbf{X} \in \mathfrak{X}_c^\perp.$$

*In particular, if  $\tilde{\mathbf{c}}$  has a conjugate point, then so does  $\mathbf{c}$ .*

*Proof.* Given  $\mathbf{X} \in \mathfrak{X}_c^\perp$ , define functions  $K, \tilde{K} : [0, b] \rightarrow \mathbb{R}$  by letting  $K(t)$  denote the sectional curvature of the plane spanned by  $\dot{\mathbf{c}}(t)$  and  $\mathbf{X}(t)$ , and similarly, letting  $\tilde{K}(t)$  denote the curvature of the image of that plane via  $\Phi_t$ . We then have

$$\begin{aligned} I(\mathbf{X}, \mathbf{X}) &= \int_0^b |\mathbf{X}'|^2 - K|\mathbf{X}| \leq \int_0^b |\mathbf{X}'|^2 - \tilde{K}|\mathbf{X}| = \int_0^b |(\Phi \mathbf{X})'|^2 - \tilde{K}|\Phi \mathbf{X}| \\ &= I(\Phi \mathbf{X}, \Phi \mathbf{X}), \end{aligned}$$

as claimed. Notice that we used both (6.7.1) and the fact that  $\Phi$  is a linear isometry.  $\square$

- Remarks 6.7.1.** (i) The same argument shows that if the curvature of each plane  $E_t$  in the statement of Proposition is less than or equal to that of  $\Phi(E_t)$ , then  $I(\mathbf{X}, \mathbf{X}) \geq I(\Phi \mathbf{X}, \Phi \mathbf{X})$  for every  $\mathbf{X} \in \mathfrak{X}_c^\perp$ .  
 (ii) The hypothesis of the Proposition is clearly satisfied if there exists a constant  $\kappa$  such that  $K_M \geq \kappa \geq K_{\tilde{M}}$ .

The *diameter* of a bounded metric space  $M$  is defined to be

$$\text{diam } M = \sup\{d(\mathbf{p}, \mathbf{q}) \mid \mathbf{p}, \mathbf{q} \in M\}.$$

For example, a sphere of constant curvature  $\kappa > 0$  has diameter  $\pi/\sqrt{\kappa}$ .

**Theorem 6.7.1.** *If  $M$  is complete and has sectional curvature  $K \geq \kappa > 0$ , then its diameter is no larger than  $\pi/\sqrt{\kappa}$ . In particular,  $M$  is compact.*

*Proof.* Recall from Example 6.2.2 that on the sphere  $\tilde{M}$  of constant curvature  $\kappa$ , any normal geodesic has  $\pi/\sqrt{\kappa}$  as conjugate point. Proposition 6.7.1 and Theorem 6.4.1 then imply that any geodesic of  $M$  with length greater than  $\pi/\sqrt{\kappa}$  is not minimal, so that

$\text{diam } M \leq \pi/\sqrt{\kappa}$ . Compactness of  $M$  follows from completeness and the Hopf-Rinow theorem.  $\square$

As a final application, we adapt the construction used at the beginning of the section to prove a result of E. Cartan that illustrates how the curvature locally determines the metric. Let  $\iota : M_{\mathbf{p}} \rightarrow \tilde{M}_{\tilde{\mathbf{p}}}$  be a linear isometry, and  $U \subset M_{\mathbf{p}}$  a neighborhood of  $\mathbf{0} \in M_{\mathbf{p}}$  on which  $\exp_{\mathbf{p}} : U \rightarrow V := \exp_{\mathbf{p}}(U)$  is a diffeomorphism. By restricting  $U$  further if necessary, we may assume that the exponential map of  $\tilde{M}$  at  $\tilde{\mathbf{p}}$  maps  $\iota(U)$  diffeomorphically onto its image. Set

$$\mathbf{f} = \exp_{\tilde{\mathbf{p}}} \circ \iota \circ (\exp_{\mathbf{p}|U})^{-1} : V \rightarrow \tilde{M}.$$

Next, we define a map  $\Phi : TV \rightarrow T\mathbf{f}(V)$  that sends each  $M_{\mathbf{q}}$  to  $\tilde{M}_{\mathbf{f}(\mathbf{q})}$ ,  $\mathbf{q} \in V$ , as follows: given any  $\mathbf{q} \in V$ , denote by  $\mathbf{c}$  the unique minimal geodesic from  $\mathbf{p}$  to  $\mathbf{q}$ , and by  $\tilde{\mathbf{c}}$  the geodesic  $t \mapsto \exp_{\tilde{\mathbf{p}}}(t \cdot \iota \dot{\mathbf{c}}(0))$ . If  $P$  and  $\tilde{P}$  are parallel translation along  $\mathbf{c}$  and  $\tilde{\mathbf{c}}$  respectively, then the restriction  $\Phi : M_{\mathbf{q}} \rightarrow \tilde{M}_{\mathbf{f}(\mathbf{q})}$  of  $\Phi$  to  $M_{\mathbf{q}}$  is given by  $\tilde{P} \circ \iota \circ P^{-1}$ . This is essentially the map  $\Phi$  used earlier, but acting on individual vectors rather than vector fields. Notice that  $\Phi$  covers  $\mathbf{f}$ , in the sense that the diagram

$$\begin{array}{ccc} TV & \xrightarrow{\Phi} & T\mathbf{f}(V) \\ \pi \downarrow & & \downarrow \tilde{\pi} \\ V & \xrightarrow{\mathbf{f}} & \mathbf{f}(V) \end{array}$$

commutes, with  $\pi$  and  $\tilde{\pi}$  denoting the respective tangent bundle projections.

**Theorem 6.7.2.** *With notation as above, suppose that*

$$\langle R(\mathbf{x}, \mathbf{y})\mathbf{u}, \mathbf{v} \rangle = \langle \tilde{R}(\Phi\mathbf{x}, \Phi\mathbf{y})\Phi\mathbf{u}, \Phi\mathbf{v} \rangle, \quad \mathbf{x}, \mathbf{y}, \mathbf{u}, \mathbf{v} \in M_{\mathbf{q}}, \quad \mathbf{q} \in V.$$

*Then  $\mathbf{f}$  is an isometry, and  $\mathbf{f}_{*\mathbf{p}} = \iota$ .*

*Proof.* By construction,  $\mathbf{f} : V \rightarrow \mathbf{f}(V)$  is a diffeomorphism and  $\mathbf{f}_{*\mathbf{p}} = \iota$ . To see that  $\mathbf{f}$  is isometric, fix an arbitrary  $\mathbf{q} \in V$  and  $\mathbf{u} \in M_{\mathbf{q}}$ . It must be shown that  $|\mathbf{f}_*\mathbf{u}| = |\mathbf{u}|$ . So consider the minimal geodesic  $\mathbf{c} : [0, t_0] \rightarrow V$  from  $\mathbf{p}$  to  $\mathbf{q}$ , and the unique Jacobi field  $\mathbf{Y}$  along  $\mathbf{c}$  with  $\mathbf{Y}(0) = \mathbf{0}$ ,  $\mathbf{Y}(t_0) = \mathbf{u}$ , cf. Exercise 6.2. Let  $\mathbf{E}_1, \dots, \mathbf{E}_n$  denote a basis of parallel orthonormal fields along  $\mathbf{c}$  with  $\mathbf{E}_n = \dot{\mathbf{c}}$ . The fields  $\tilde{\mathbf{E}}_i$ , where  $\tilde{\mathbf{E}}_i(t) = \Phi\mathbf{E}_i(t)$ , then form a basis of parallel orthonormal fields along  $\tilde{\mathbf{c}}$  by definition of  $\Phi$ , and  $\tilde{\mathbf{E}}_n = \dot{\tilde{\mathbf{c}}}$ . If  $f_i = \langle \mathbf{Y}, \mathbf{E}_i \rangle$ ,  $i = 1, \dots, n$ , so that  $\mathbf{Y} = \sum_i f_i \mathbf{E}_i$ , then the Jacobi condition for  $\mathbf{Y}$  reads

$$f_j'' + \sum_i f_i \langle R(\mathbf{E}_i, \mathbf{E}_n)\mathbf{E}_n, \mathbf{E}_j \rangle = 0, \quad j = 1, \dots, n.$$

This implies that the vector field  $\tilde{\mathbf{Y}} := \Phi\mathbf{Y}$  is then also Jacobi, because  $\tilde{\mathbf{Y}} = \sum_i f_i \tilde{\mathbf{E}}_i$ , and

$$\langle R(\mathbf{E}_i, \mathbf{E}_n)\mathbf{E}_n, \mathbf{E}_j \rangle = \langle R(\tilde{\mathbf{E}}_i, \tilde{\mathbf{E}}_n)\tilde{\mathbf{E}}_n, \tilde{\mathbf{E}}_j \rangle.$$

By definition of  $\Phi$ ,  $|\tilde{Y}(t)| = |Y(t)|$  for all  $t$ , so the theorem is proved once we establish that  $\tilde{Y}(t_0) = f_*\mathbf{u} = f_*\mathbf{q}Y(t_0)$ . But both  $Y$  and  $\tilde{Y}$  vanish at 0, so that by (6.2.1),

$$Y(t) = \exp_{\mathbf{p}^*} (t\mathcal{I}_{t\dot{\mathbf{c}}(0)}Y'(0)),$$

and similarly,

$$\tilde{Y}(t) = \exp_{\tilde{\mathbf{p}}^*} (t\mathcal{I}_{t\dot{\tilde{\mathbf{c}}}(0)}\tilde{Y}'(0)) = \exp_{\tilde{\mathbf{p}}^*} (t\mathcal{I}_{t\dot{\mathbf{c}}(0)}tY'(0)).$$

Thus,

$$\begin{aligned} \tilde{Y}(t) &= \exp_{\tilde{\mathbf{p}}^*} (\mathcal{I}_{t\dot{\mathbf{c}}(0)}tY'(0)) = (\exp_{\tilde{\mathbf{p}}^*} \circ \mathcal{I}_{t\dot{\mathbf{c}}(0)} \circ \iota \circ \mathcal{I}_{t\dot{\mathbf{c}}(0)}^{-1} \circ \exp_{\mathbf{p}^*}^{-1}) Y(t) \\ &= (\exp_{\tilde{\mathbf{p}}^*} \circ \iota_* \circ \exp_{\mathbf{p}^*}^{-1}) Y(t) \\ &= f_*\mathbf{c}(t) Y(t), \end{aligned}$$

as claimed. This concludes the proof. Notice that we have actually shown that  $\Phi$  is the derivative of  $f$ , since for every Jacobi field  $Y$  as above,  $\tilde{Y} = \Phi Y$  and  $\tilde{Y} = f_* Y$ .  $\square$

**Remark 6.7.2.** The above theorem shows that any two spaces  $M$  and  $\tilde{M}$  with the same constant curvature  $\kappa$  are locally isometric; i.e., for any  $\mathbf{p}$  in  $M$ , there exists a neighborhood  $U$  of  $\mathbf{p}$  and an isometry  $f : U \rightarrow f(U) \subset \tilde{M}$ . Such an  $f$  need not be extendable to all of  $M$ , though. For example,  $\mathbb{R}^2$  and  $S^1 \times S^1 \subset \mathbb{R}^4$  are both flat and therefore locally isometric, but not globally, since one is compact and the other is not.

## 6.8 Exercises

**6.1.** (a) Let  $Y$  be a Jacobi field along a geodesic  $\mathbf{c}$ , and  $f$  an affine function; i.e., a function of the form  $f(t) = at + b$ ,  $a, b \in \mathbb{R}$ . Show that  $\mathbf{c} \circ f$  is a geodesic, and that  $Y \circ f$  is a Jacobi field along  $\mathbf{c} \circ f$ .

(b) Notice that for any function  $f$ , not necessarily affine,  $\mathbf{c} \circ f$  has (part of) the same image as  $\mathbf{c}$ . Explain why  $\mathbf{c} \circ f$  is no longer a geodesic if  $f$  is not affine.

**6.2.** Prove that if  $t_0$  is not a conjugate point of a geodesic  $\mathbf{c} : [a, b] \rightarrow M$ , then for any  $\mathbf{v} \in M_{\mathbf{c}(a)}$  and  $\mathbf{w} \in M_{\mathbf{c}(t_0)}$ , there exists exactly one Jacobi field  $Y$  along  $\mathbf{c}$  with  $Y(a) = \mathbf{v}$  and  $Y(t_0) = \mathbf{w}$ .

**6.3.** Suppose  $V$  is a variation of  $\mathbf{c} : [0, b] \rightarrow M$  with fixed endpoints, such that all the curves  $V_s$  are geodesics. Show that the length function of  $V$  is constant and that  $b$  is a conjugate point of  $\mathbf{c}$ .

**6.4.** Recall from Chapter 3 that a Killing field on  $M$  is a vector field whose flow consists of isometries of  $M$ . Prove that the restriction of a Killing field to any geodesic is a Jacobi field along that geodesic. *Hint:* Perhaps the easiest way to see this is to consider the flow  $\Phi_s$  of the Killing field, and consider the variation  $(t, s) \mapsto \Phi_s(\mathbf{c}(t))$  of the geodesic  $\mathbf{c}$ .

**6.5.** Let  $\mathbf{c} : [0, b] \rightarrow M$  be a geodesic in a manifold with sectional curvature  $K$ , and  $V$  a (nonconstant) variation of  $\mathbf{c}$  with fixed endpoints. Prove that if  $K \leq 0$ , then the length of  $V_s$  is larger than that of  $\mathbf{c}$  if  $s$  is small enough. Show by means of an example that this does not necessarily imply that  $\mathbf{c}$  is a minimal geodesic.

**6.6.** Let  $\mathbf{E}$  denote a parallel vector field along a geodesic  $\mathbf{c} : [0, b] \rightarrow M$  that is orthogonal to  $\dot{\mathbf{c}}$ . Choose an interval  $I$  around 0 small enough that the geodesic  $s \mapsto \exp(s\mathbf{E}(t))$  is defined on  $I$  for all  $t \in [0, b]$ , and consider the variation  $V : [0, b] \times I \rightarrow M$  of  $\mathbf{c}$  given by  $V(t, s) = \exp(s\mathbf{E}(t))$ . Prove that if  $M$  has positive curvature, then  $L(V_s) < L(\mathbf{c})$  for small  $s$  and if  $M$  has negative curvature, then  $L(V_s) > L(\mathbf{c})$ . Interpret this fact in terms of the lengths of circles of latitude compared to that of the equator on a sphere.

**6.7.** Let  $M$  denote the  $n$ -dimensional sphere of radius  $r$ ,  $\mathbf{c} : [0, a] \rightarrow M$  a normal geodesic, and  $\mathfrak{E}$  the vector space of Jacobi fields along  $\mathbf{c}$  that vanish at 0.

- (a) Give an explicit formula for  $I(Y_1, Y_2)$  in terms of  $\langle \mathbf{u}_1, \mathbf{u}_2 \rangle$  if  $Y_i \in \mathfrak{E}$  with  $Y_i(0) = \mathbf{u}_i$ .
- (b) Show that the index form is negative definite on  $\mathfrak{E}$  if  $a \in (\pi r/2, \pi r)$ .

**6.8.** Show that the distance function  $d : M \times M \rightarrow \mathbb{R}$  is continuous; i.e., show that if  $\mathbf{p}_i \rightarrow \mathbf{p}$  and  $\mathbf{q}_i \rightarrow \mathbf{q}$ , then  $d(\mathbf{p}_i, \mathbf{q}_i) \rightarrow d(\mathbf{p}, \mathbf{q})$ .

**6.9.** Prove the claim made in the proof of Proposition 6.6.1: given  $\mathbf{p} \in M$ ,

$$B_\varepsilon(\mathbf{p}) = \cup_{r \in (0, \varepsilon)} \overline{B_r(\mathbf{p})}.$$

**6.10.** Let  $\tilde{M}, M$  denote manifolds with distance functions  $\tilde{d}$  and  $d$  respectively. Show that if  $\mathbf{f} : \tilde{M} \rightarrow M$  is isometric, then  $\tilde{d}(\mathbf{p}, \mathbf{q}) \geq d(\mathbf{f}(\mathbf{p}), \mathbf{f}(\mathbf{q}))$  for all  $\mathbf{p}, \mathbf{q} \in \tilde{M}$ . Deduce that if  $\mathbf{f}$  is an isometry, then it is distance-preserving; i.e.,  $\tilde{d}(\mathbf{p}, \mathbf{q}) = d(\mathbf{f}(\mathbf{p}), \mathbf{f}(\mathbf{q}))$  for all  $\mathbf{p}, \mathbf{q} \in \tilde{M}$ . R. Palais has shown that the converse is also true: any distance-preserving map between manifolds is an isometry.

**6.11.** (a) According to Palais' result mentioned in the exercise above, any distance-preserving map  $\mathbf{f} : \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$  in Euclidean space is a Euclidean motion. Prove this directly.

*Hint:* The map  $\mathbf{x} \mapsto \mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{0})$  preserves norms. Adapt the proof of Proposition 3.11.1 to conclude it is an orthogonal transformation.

(b) Let  $\mathbf{f} : S^n \rightarrow S^n$  be distance-preserving. Show that  $\mathbf{f}$  is the restriction of an orthogonal transformation.

*Hint:* Extend  $\mathbf{f}$  radially to  $\mathbb{R}^{n+1}$ .

**6.12.** Show that for  $M = S^1$ , the distance function squared  $d^2 : M \times M \rightarrow \mathbb{R}$  is not differentiable exactly on the circle  $\{(\mathbf{p}, -\mathbf{p}) \mid \mathbf{p} \in S^1\}$ .

**6.13.** Give examples showing that a manifold in which any two points can be joined by a minimal geodesic is not necessarily complete.

**6.14.** Let  $M^n$  be a manifold in  $\mathbb{R}^{n+k}$ .

- (a) Show that a Cauchy sequence in  $M$  is Cauchy as a sequence in  $\mathbb{R}^{n+k}$  (i.e., with respect to the Euclidean distance).
- (b) Prove that if  $M$  is closed in  $\mathbb{R}^{n+k}$ , then it is complete.

**6.15.** (a) Give examples showing that the converse to part (b) in the previous problem does not hold; i.e., there exist complete manifolds that are not closed as a subset of the ambient Euclidean space.

- (b) Show that if  $M$  is a subset of  $\mathbb{R}^n$  that is complete with the restriction of the distance function from  $\mathbb{R}^n$  (i.e.,  $d(\mathbf{p}, \mathbf{q}) = |\mathbf{p} - \mathbf{q}|$  for  $\mathbf{p}, \mathbf{q} \in M$ ), then  $M$  is closed in the ambient Euclidean space.

**6.16.** Suppose  $f : M \rightarrow N$  is a diffeomorphism. If  $M$  is complete, is it always true that  $N$  is also complete? What if, in addition,  $f$  is an isometry?

**6.17.**  $M$  is said to be a *Riemannian homogeneous space* if its group of isometries acts transitively on  $M$ ; i.e., given  $\mathbf{p}, \mathbf{q} \in M$ , there exists an isometry of  $M$  that maps  $\mathbf{p}$  to  $\mathbf{q}$ .

- (a) Give several examples of such spaces.
- (b) Prove that a Riemannian homogeneous space is necessarily complete.

**6.18.** A *ray* in  $M$  is a normal geodesic  $\mathbf{c} : [0, \infty) \rightarrow M$  that minimizes distance for all time; i.e.,  $d(\mathbf{c}(0), \mathbf{c}(t)) = t$  for all  $t \geq 0$ .

- (a) Prove that if  $M$  is complete and noncompact, then for any  $\mathbf{p} \in M$ , then there exists a ray  $\mathbf{c}$  with  $\mathbf{c}(0) = \mathbf{p}$ . *Hint:* Consider a sequence  $\mathbf{q}_n \in M$  with  $d(\mathbf{p}, \mathbf{q}_n) > n$ ,  $n \in \mathbb{N}$ , and minimal normal geodesics  $\mathbf{c}_n$  from  $\mathbf{p}$  to  $\mathbf{q}_n$ . The sequence  $\dot{\mathbf{c}}_n(0)$  has some convergent subsequence. If  $\mathbf{v}$  is the limit, consider the geodesic in direction  $\mathbf{v}$ .
- (b) Show by means of an example that completeness is necessary in part (a).

**6.19.** A point  $\mathbf{p} \in M$  is said to be a *pole* of  $M$  if every normal geodesic emanating from  $\mathbf{p}$  is a ray (see previous exercise). Let  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  be the function given by

$$f = u^{n+1} - \sum_{i=1}^n (u^i)^2,$$

and consider the paraboloid  $M = f^{-1}(0)$ . Show that the origin is a pole of  $M$ . Are there any other poles?

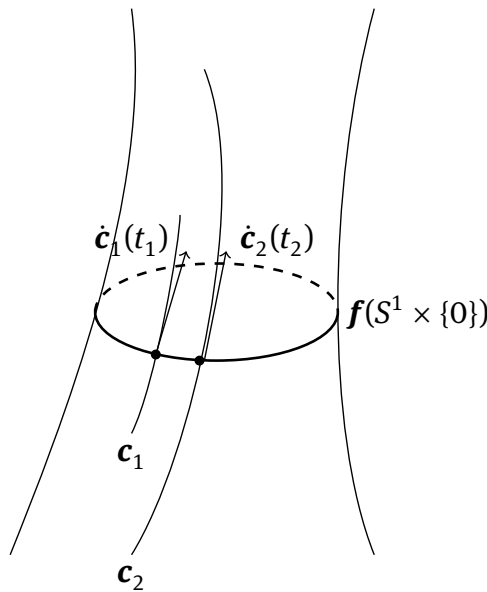
**6.20.** A *line* in  $M$  is a normal geodesic  $\mathbf{c} : (-\infty, \infty) \rightarrow M$  that minimizes distance for all time; i.e.,  $d(\mathbf{c}(t_0), \mathbf{c}(t_1)) = |t_0 - t_1|$  for all  $t_0, t_1$ .

- (a) If  $M$  admits a line, it must of course be noncompact. Show by means of examples, however, that there exist noncompact complete manifolds that admit no lines. This contrasts with rays, which exist in any noncompact, complete manifold, cf. Exercise 6.18.
- (b) Recall from Chapter 5 that a given space can have different metrics, in the sense that there are diffeomorphisms which are not isometries. For example, the right

circular cylinder

$$S^1 \times \mathbb{R} = \{(x, y, z) \mid x^2 + y^2 = 1, z \in \mathbb{R}\}$$

is diffeomorphic to the hyperboloid  $N = \{(x, y, z) \mid x^2 + y^2 - z^2 = 1\}$ , but not isometric. Show however, that *any complete* manifold  $M$  diffeomorphic to  $S^1 \times \mathbb{R}$  admits a line. *Hint:* Let  $f : S^1 \times \mathbb{R} \rightarrow M$  be a diffeomorphism, and consider a sequence of normal minimal geodesics  $c_n$  in  $M$  from  $f(1, 0, -n)$  to  $f(1, 0, n)$ ,  $n \in \mathbb{N}$ . Each  $c_n$  must intersect  $f(S^1 \times \{0\})$  for a unique value  $t_n \in \mathbb{R}$  of its parameter. Show that the sequence  $\dot{c}_n(t_n)$  has a convergent subsequence. If  $u$  is the limit of this subsequence, consider  $t \mapsto \exp(tu)$ .



**6.21.** Recall that the diameter of  $M$  is  $\text{diam}(M) = \sup\{d(p, q) \mid p, q \in M\}$  if this set is bounded above, and  $\infty$  otherwise.

- (a) Show that if  $M$  is complete, then  $\text{diam}(M) < \infty$  if and only if  $M$  is compact.
- (b) Prove that if  $M$  is compact, then there exist  $p, q \in M$  such that  $d(p, q) = \text{diam}(M)$ .

**6.22.** Suppose  $M$  is a complete manifold, and denote by  $T_1M = \{u \in TM \mid |u| = 1\}$  the *unit tangent sphere bundle* of  $M$ . If  $\pi : TM \rightarrow M$  is the tangent bundle projection, define

$$s : T_1M \rightarrow \mathbb{R}_+ \cup \{\infty\},$$

$$v \mapsto \sup\{t > 0 \mid d(\pi(v), \exp(tv)) = t\}.$$

Roughly speaking,  $s(v)$  is the largest parameter value for which the geodesic  $c_v$  in direction  $v$  is minimal.

The *tangential cut locus* of  $p \in M$  is

$$C_p = \{s(v)v \mid v \in M_p \cap T_1M, \quad s(v) < \infty\},$$

and the *cut locus* of  $p \in M$  is  $C(p) = \exp_p(C_p)$ . Determine the tangential cut locus and the cut locus at a generic point  $p \in M$ , if

- (a)  $M = \mathbb{R}^n$ ;
- (b)  $M = S^n$ ;
- (c)  $M = S^1 \times \mathbb{R}$ .

**6.23.** With the notation and terminology from Exercise 6.22, show that for  $\mathbf{p}, \mathbf{q} \in M$ ,  $\mathbf{q}$  lies in the cut locus of  $\mathbf{p}$  if and only if  $\mathbf{p}$  lies in the cut locus of  $\mathbf{q}$ .

**6.24.** With the notation and terminology from Exercise 6.22, show that the distance  $d(\mathbf{p}, C(\mathbf{p}))$  from a point to its cut locus equals the injectivity radius  $\text{inj}_{\mathbf{p}}$  at that point. *Hint:* Notice that

$$d(\mathbf{p}, C(\mathbf{p})) = \inf\{s(\mathbf{v}) \mid \mathbf{v} \in T_{\mathbf{p}}M \cap M_{\mathbf{p}}\}.$$

Argue by contradiction to rule out the possibilities  $d(\mathbf{p}, C(\mathbf{p})) > \text{inj}_{\mathbf{p}}$  and  $d(\mathbf{p}, C(\mathbf{p})) < \text{inj}_{\mathbf{p}}$ .

This fact can be used to prove that in general, the function  $\text{inj} : M \rightarrow \mathbb{R}_+ \cup \{\infty\}$  is continuous.

**6.25.** A subset  $A$  of a manifold  $M$  is said to be *convex* if any  $\mathbf{p}, \mathbf{q} \in A$  can be joined by some geodesic  $\mathbf{c}$  contained in  $A$ , with length  $L(\mathbf{c}) = d(\mathbf{p}, \mathbf{q})$ . If, in addition, this geodesic is unique, then  $A$  is said to be *strongly convex*.

- (a) Show that any convex subset of  $S^n$  is either all of  $S^n$  or else is contained in some hemisphere.
- (b) Prove that for any manifold  $M$  and  $\mathbf{p} \in M$  there exists  $\varepsilon > 0$  such that the metric ball  $B_\varepsilon(\mathbf{p})$  of radius  $\varepsilon$  about  $\mathbf{p}$  is strongly convex. *Hint:* see Proposition 6.5.1.

**6.26.** This exercise examines the gradient and Hessian of the distance function from a point.

- (a) Prove that for  $\mathbf{p} \in M$ , there exists  $\varepsilon > 0$  such that the distance function  $\mathbf{q} \mapsto d(\mathbf{p}, \mathbf{q})$  from  $\mathbf{p}$  is smooth on  $W := B_\varepsilon(\mathbf{p}) \setminus \{\mathbf{p}\}$ . Is it differentiable on all of  $B_\varepsilon(\mathbf{p})$ ?
- (b) Let  $f$  denote the restriction of this function to  $W$ ,  $f(\mathbf{q}) = d(\mathbf{p}, \mathbf{q})$ ,  $\mathbf{q} \in W$ , and  $\mathbf{c} : [0, 1] \rightarrow W$  a (necessarily minimal) geodesic from  $\mathbf{p}$  to  $\mathbf{q}$ . If  $\mathbf{u} \in M_{\mathbf{q}}$  and  $\mathbf{c}_u$  is the geodesic  $\mathbf{u} \mapsto \exp(t\mathbf{u})$ , prove that there exists a variation  $V : [0, 1] \times (-\delta, \delta) \rightarrow W$  of  $\mathbf{c}$  such that  $V_s$  is the minimal geodesic from  $\mathbf{p}$  to  $\mathbf{c}_u(s)$  if  $\delta$  is small enough.
- (c) Let  $L$  denote the length function of  $V$ ,  $L(s) = \text{length of } V_s$ . Show that

$$\langle \nabla f(\mathbf{q}), \mathbf{u} \rangle = (f \circ \mathbf{c}_u)'(0) = L'(0) = \langle \dot{\mathbf{c}}(1), \mathbf{u} \rangle,$$

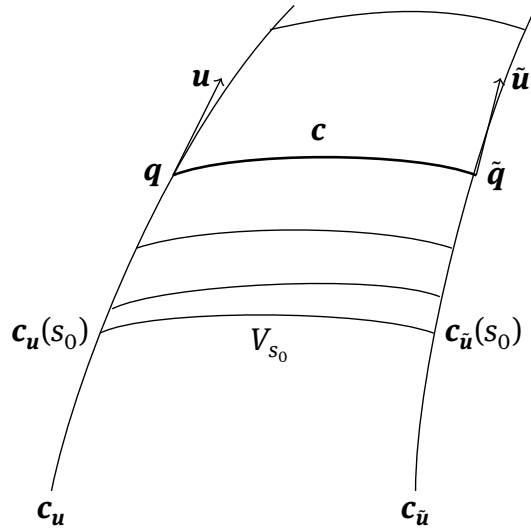
so that  $\nabla f(\mathbf{q}) = \dot{\mathbf{c}}(1)$ , and

$$h_f(\mathbf{q})(\mathbf{u}, \mathbf{u}) = (f \circ \mathbf{c}_u)''(0) = L''(0) = \langle \mathbf{Y}, \mathbf{Y}' \rangle(1),$$

where  $\mathbf{Y}$  is the Jacobi field along  $\mathbf{c}$  with  $\mathbf{Y}(0) = \mathbf{0}$ , and  $\mathbf{Y}(1) = \mathbf{u}$ .

**6.27.** Let  $\mathbf{p} \in M$ ,  $\varepsilon > 0$  such that the square of the distance function  $d^2 : M \times M \rightarrow \mathbb{R}$ , when restricted to  $B_\varepsilon(\mathbf{p}) \times B_\varepsilon(\mathbf{p})$ , is differentiable. Denote this restriction by  $f$ . The goal of this problem is to determine the gradient  $\nabla f$  and the Hessian  $H_f$  of  $f$ .





Let  $\mathbf{q}, \tilde{\mathbf{q}} \in B_\varepsilon(\mathbf{p})$ ,  $\mathbf{c} : [0, 1] \rightarrow B_\varepsilon(\mathbf{p})$  the geodesic from  $\mathbf{q}$  to  $\tilde{\mathbf{q}}$  with length  $d(\mathbf{q}, \tilde{\mathbf{q}})$ , and  $\mathbf{u} \in M_{\mathbf{q}}, \tilde{\mathbf{u}} \in M_{\tilde{\mathbf{q}}}$ . If  $\mathbf{c}_u, \mathbf{c}_{\tilde{u}}$  denote the geodesics in direction  $\mathbf{u}$  and  $\tilde{\mathbf{u}}$  respectively, strong convexity (see Exercise 6.25) guarantees the existence of a variation  $V : [0, 1] \times (-\delta, \delta) \rightarrow B_\varepsilon(\mathbf{p})$  of  $\mathbf{c}$  by geodesics, with  $V_s$  minimal from  $\mathbf{c}_u(s)$  to  $\mathbf{c}_{\tilde{u}}(s)$ ,  $|s| < \delta$ , provided  $\delta$  is small enough. Then

$$\langle \nabla f(\mathbf{q}, \tilde{\mathbf{q}}), (\mathbf{u}, \tilde{\mathbf{u}}) \rangle = D_{(\mathbf{u}, \tilde{\mathbf{u}})}f = h'(0),$$

where  $h(s) = L(V_s)$ .

- (a) Use this to show that  $\nabla f(\mathbf{q}, \tilde{\mathbf{q}}) = 2(-\dot{\mathbf{c}}(0), \dot{\mathbf{c}}(1))$ .
- (b) Use the variation from above to show that for  $\mathbf{q}, \tilde{\mathbf{q}} \in B_\varepsilon(\mathbf{p})$ ,

$$H_f(\mathbf{q}, \tilde{\mathbf{q}})(\mathbf{u}, \tilde{\mathbf{u}}) = 2(-\mathbf{Y}'(0), \mathbf{Y}'(1)), \quad \mathbf{u} \in M_{\mathbf{q}}, \quad \tilde{\mathbf{u}} \in M_{\tilde{\mathbf{q}}},$$

where  $\mathbf{Y}$  is the Jacobi field along  $\mathbf{c}$  with  $\mathbf{Y}(0) = \mathbf{u}$ ,  $\mathbf{Y}(1) = \tilde{\mathbf{u}}$ .  
Deduce that the Hessian form  $h_f$  of  $f$  satisfies

$$h_f(\mathbf{q}, \tilde{\mathbf{q}})((\mathbf{u}, \tilde{\mathbf{u}}), (\mathbf{u}, \tilde{\mathbf{u}})) = 2\langle \mathbf{Y}, \mathbf{Y}' \rangle_0^1.$$

- (c) Compute the gradient, Hessian, and Hessian form of  $g$ , where  $g$  is the square of the distance function from  $\mathbf{p}$ ,  $g(\mathbf{q}) = d^2(\mathbf{p}, \mathbf{q})$ , on  $B_\varepsilon(\mathbf{p})$ . *Hint:* An easy way to do this is to notice that  $g = f \circ \iota_{\mathbf{p}}$ , where  $\iota_{\mathbf{p}} : M \rightarrow M \times M$  maps  $\mathbf{q}$  to  $(\mathbf{p}, \mathbf{q})$ .

**6.28.** This problem, which uses concepts and results from Exercises 6.25 and 6.27, explores a result of J. H. C. Whitehead; namely, for any point  $\mathbf{p}$  in a manifold  $M$ , there exists some  $\varepsilon > 0$  such that every metric ball  $B_\delta(\mathbf{q})$  contained in  $B_\varepsilon(\mathbf{p})$  is strongly convex.

- (a) Let  $\varepsilon_0 > 0$  such that the closure  $K$  of the ball of radius  $\varepsilon_0$  centered at  $\mathbf{p}$  is compact, and consider the compact set

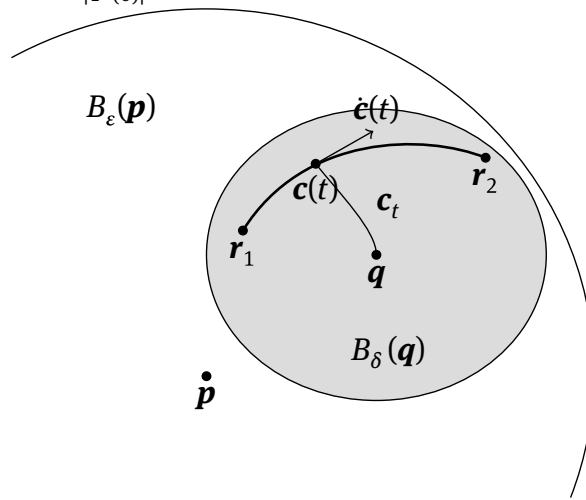
$$C = \{(\mathbf{u}, \mathbf{v}) \in TM \times TM \mid |\mathbf{u}| = |\mathbf{v}| = 1, \quad \pi(\mathbf{u}) = \pi(\mathbf{v}) \in K\},$$

where  $\pi : TM \rightarrow M$  is the bundle projection. For  $(\mathbf{u}, \mathbf{v}) \in C$ , denote by  $\mathbf{Y}_{uv}$  the Jacobi field along  $t \mapsto \exp(t\mathbf{v})$  with  $\mathbf{Y}_{uv}(0) = \mathbf{0}$ ,  $\mathbf{Y}'_{uv}(0) = \mathbf{u}$ . Use the fact that

$$\langle \mathbf{Y}_{uv}, \mathbf{Y}'_{uv} \rangle'(0) = |\mathbf{u}|^2 = 1$$

to show that there exists  $\varepsilon \in (0, \varepsilon_0)$  such that  $\langle \mathbf{Y}_{uv}, \mathbf{Y}'_{uv} \rangle|_0^t > 0$  for any  $0 < t \leq \varepsilon$  and  $(\mathbf{u}, \mathbf{v}) \in C$ .

- (b) Prove that for any  $\mathbf{p} \in M$ , there exists some positive  $\varepsilon$  such that
- (1)  $B_\varepsilon(\mathbf{p})$  is strongly convex, and the square of the distance function on  $M$ , when restricted to  $B_\varepsilon(\mathbf{p}) \times B_\varepsilon(\mathbf{p})$ , is differentiable.
  - (2) For any  $\mathbf{q} \in B_\varepsilon(\mathbf{p})$  and any geodesic  $\mathbf{c} : [0, 1] \rightarrow B_\varepsilon(\mathbf{p})$  with  $\mathbf{c}(0) = \mathbf{q}$ , if  $\mathbf{Y}$  is a Jacobi field along  $\mathbf{c}$  with  $\mathbf{Y}(0) = \mathbf{0}$ ,  $\mathbf{Y}'(0) \neq \mathbf{0}$ , then  $\langle \mathbf{Y}, \mathbf{Y}' \rangle|_0^t > 0$  for all  $t \in (0, 1]$ . *Hint:*  $\mathbf{Z} = \frac{1}{|\mathbf{Y}'(0)|} \mathbf{Y}$  is a Jacobi field satisfying  $\mathbf{Z}(0) = \mathbf{0}$ ,  $|\mathbf{Z}'(0)| = 1$ .



- (c) Show that Whitehead's result holds for the  $\varepsilon$  obtained in (b), as follows: let  $\mathbf{q} \in B_\varepsilon(\mathbf{p})$  and  $\delta > 0$  such that  $B_\delta(\mathbf{q}) \subset B_\varepsilon(\mathbf{p})$ . Since a convex subset of a strongly convex set is strongly convex, it suffices to show that the ball of radius  $\delta$  about  $\mathbf{q}$  is convex. Denote by  $h : B_\varepsilon(\mathbf{p}) \rightarrow \mathbb{R}$  the distance function from  $\mathbf{q}$  squared,  $h(\mathbf{r}) = d^2(\mathbf{q}, \mathbf{r})$ . Given  $\mathbf{r}_1, \mathbf{r}_2 \in B_\delta(\mathbf{q})$ , there is a unique geodesic  $\mathbf{c} : [0, 1] \rightarrow B_\varepsilon(\mathbf{p})$  with  $\mathbf{c}(0) = \mathbf{r}_1$ ,  $\mathbf{c}(1) = \mathbf{r}_2$ . Prove that

$$(h \circ \mathbf{c})''(t) = 2b_h(\mathbf{c}(t))(\dot{\mathbf{c}}(t), \dot{\mathbf{c}}(t)) = 2\langle \mathbf{Y}_t, \mathbf{Y}'_t \rangle|_0^1,$$

where  $b_h$  is the Hessian form of  $h$ , and  $\mathbf{Y}$  is the Jacobi field along the minimal geodesic  $\mathbf{c}_t : [0, 1] \rightarrow B_\varepsilon(\mathbf{p})$  joining  $\mathbf{q}$  to  $\mathbf{c}(t)$ , with  $\mathbf{Y}_t(0) = \mathbf{0}$ ,  $\mathbf{Y}_t(1) = \dot{\mathbf{c}}(t)$ . In particular,  $h \circ \mathbf{c}$  is a convex function. Deduce that  $(h \circ \mathbf{c})(t) < \delta^2$  for all  $t$ , and conclude that  $B_\delta(\mathbf{q})$  is convex.

**6.29.** A manifold is said to be *locally symmetric* if its curvature tensor is parallel; i.e., if for any geodesic  $\mathbf{c}$  and parallel fields  $\mathbf{X}_i$  along  $\mathbf{c}$ ,  $1 \leq i \leq 3$ ,  $R(\mathbf{X}_1, \mathbf{X}_2)\mathbf{X}_3$  is a parallel field along  $\mathbf{c}$ .

- (a) Prove that a space of constant curvature is locally symmetric.
- (b) Let  $M$  be a locally symmetric space,  $\mathbf{c}$  a normal geodesic in  $M$ . Show that there exists an orthonormal basis of parallel vector fields  $\mathbf{E}_i$  along  $\mathbf{c}$ ,  $1 \leq i \leq n$ , with

$\mathbf{E}_n = \dot{\mathbf{c}}$ , and constants  $\lambda_i$  such that  $R(\mathbf{E}_i, \dot{\mathbf{c}})\dot{\mathbf{c}} = \lambda_i \mathbf{E}_i$ . Notice that  $\lambda_i$  is the sectional curvature of the plane spanned by  $\dot{\mathbf{c}}(t)$  and  $\mathbf{E}_i(t)$ .

(c) Prove that  $\mathbf{Y}$  is a Jacobi field along  $\mathbf{c}$  if and only if

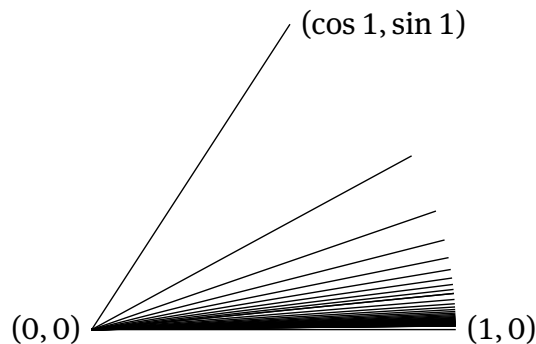
$$\mathbf{Y} = \sum_{i=0}^n f_i \mathbf{E}_i, \text{ where } f_i \text{ satisfies } f_i'' + \lambda_i f_i = 0.$$

**6.30.** The definition of distance that was used for a manifold  $M$ ,

$$d(\mathbf{p}, \mathbf{q}) = \inf\{L(\mathbf{c}) \mid \mathbf{c} : [0, 1] \rightarrow M, \mathbf{c}(0) = \mathbf{p}, \mathbf{c}(1) = \mathbf{q}\},$$

makes sense for more general connected subsets of Euclidean space that are not manifolds, see for example [6].

- (a) Show that the boundary of any square in  $\mathbb{R}^2$  is a well-defined metric space with the distance above, and that the open sets in the metric space coincide with the usual open sets as in Theorem 6.5.1.
- (b) Denote by  $L_0$  the line segment in  $\mathbb{R}^2$  connecting the origin to  $(1, 0)$ , by  $L_n$  the line segment connecting the origin to  $(\cos(1/n), \sin(1/n))$ ,  $n \in \mathbb{N}$ , and let  $M = \cup_{n=0}^{\infty} L_n$ , together with the above distance. Prove that the open sets in  $M$  are not the usual ones. *Hint:* consider metric balls centered about, say,  $(1/2, 0)$ .





# 7 Hypersurfaces

## 7.1 Hypersurfaces and orientation

A *hypersurface* is an  $n$ -dimensional submanifold of  $\mathbb{R}^{n+1}$  for some  $n$ . The case  $n = 2$ , which corresponds to surfaces in 3-dimensional space, was historically the first to be studied. Most of the concepts introduced in Chapter 3 are easier to understand and work with when the manifold's dimension is one less than that of the ambient Euclidean space.

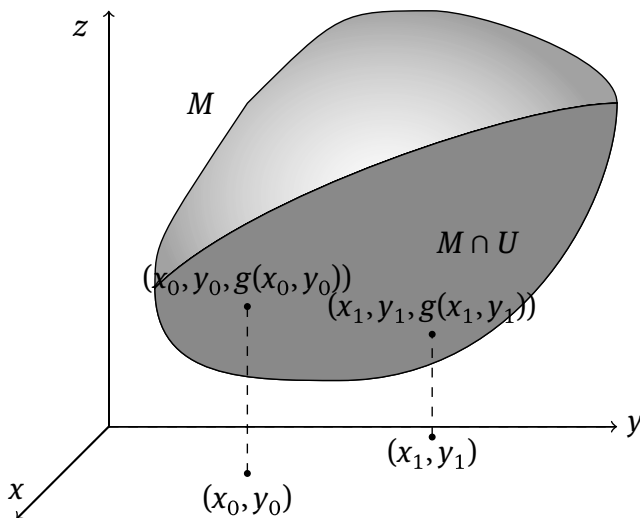
The most useful way of describing a hypersurface is by means of Corollary 3.2.1, which we recall in the present context:

**Proposition 7.1.1.**  $M^n \subset \mathbb{R}^{n+1}$  is a hypersurface if and only if any  $\mathbf{p} \in M$  admits a neighborhood  $U \subset \mathbb{R}^{n+1}$  and a function  $f : U \rightarrow \mathbb{R}$  having 0 as a regular value, such that  $M \cap U = f^{-1}(0)$ .

For example, in the above situation, the vector field  $\nabla f$  is normal – i.e., orthogonal – to  $M$  on  $M \cap U$ : by Proposition 3.1.1, given  $\mathbf{q} \in M \cap U$ ,  $M_{\mathbf{q}} = \ker f_{*\mathbf{q}}$ , and the latter is just  $\nabla f(\mathbf{q})^\perp$ .

**Proposition 7.1.2.** A hypersurface is orientable if and only if it admits a unit length normal vector field.

*Proof.* If the hypersurface is orientable, it admits an orientation. Let  $\omega$  denote the standard volume form on  $\mathbb{R}^n$ . If  $\mathbf{X}_1, \dots, \mathbf{X}_{n-1}$  is a local positive basis of vector fields on  $U \subset M$ , there is a unique unit normal field  $\mathbf{N}$  on  $U$  such that  $(i(\mathbf{N})\omega)(\mathbf{X}_1, \dots, \mathbf{X}_{n-1}) = \omega(\mathbf{N}, \mathbf{X}_1, \dots, \mathbf{X}_{n-1}) > 0$ . Since this can be done in a neighborhood of any point,  $\mathbf{N}$  is globally defined. Conversely, if  $\mathbf{N}$  is a unit normal field on  $M$ , then the map which assigns to  $\mathbf{x}_1, \dots, \mathbf{x}_{n-1} \in M_{\mathbf{p}}$  the value  $\omega(\mathbf{p})(\mathbf{N}(\mathbf{p}), \mathbf{x}_1, \dots, \mathbf{x}_{n-1})$  is a nowhere-zero  $(n-1)$ -form on  $M$ ; i.e., an orientation.  $\square$



The darker half  $M \cap U$  of  $M$  is the graph of a function  $g$  of two variables; it equals  $f^{-1}(0)$ , where  $f(x, y, z) = z - g(x, y)$ .

In light of the above, we may call the normal vector field an *orientation* of  $M$ . Clearly, if  $M$  is orientable, then it has two possible orientations, namely  $\mathbf{N}$  and  $-\mathbf{N}$ , where  $\mathbf{N}$  is any unit normal field. A large class of orientable hypersurfaces is given by the collection of functions  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  that have zero as regular value, since each  $M = f^{-1}(0)$  is a hypersurface with orientation  $\mathbf{N} = \nabla f / |\nabla f|$ .

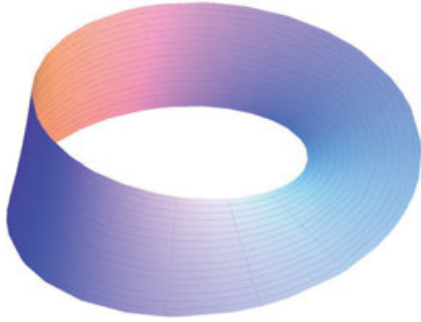


Fig. 7.1: A Möbius strip

One example of a non orientable surface is the Möbius strip. It can be described as the image of the map

$$\mathbf{h} : (-1/2, 1/2) \times [0, 2\pi] \rightarrow \mathbb{R}^3$$

$$(r, \theta) \mapsto \left( 2 \cos \theta + r \cos \frac{\theta}{2}, 2 \sin \theta + r \cos \frac{\theta}{2}, r \sin \frac{\theta}{2} \right).$$

$\mathbf{h}$  is not, strictly speaking, a parametrization, but it is one locally, since it has maximal rank everywhere. More precisely, restrict  $\mathbf{h}$  to  $(-1/2, 1/2) \times (0, 2\pi)$ , and define  $\mathbf{k} : (-1/2, 1/2) \times (-\pi, \pi) \rightarrow \mathbb{R}^3$  by the same formula.  $\mathbf{h}$  and  $\mathbf{k}$  are both parametrizations whose images cover the strip  $M$ . One obtains normal vector fields along these parametrizations by setting

$$\mathbf{N}_{\mathbf{h}} = \frac{\partial}{\partial x^1} \times \frac{\partial}{\partial x^2},$$

where  $\mathbf{x} = \mathbf{h}^{-1}$ , and using a similar formula for  $\mathbf{N}_{\mathbf{k}}$ . Notice that both fields agree on the intersection  $(-1/2, 1/2) \times (0, \pi)$  of their domains, since they are given by the same formula. They do not, however, combine to give us a well-defined normal field on  $M$ , because, for example,

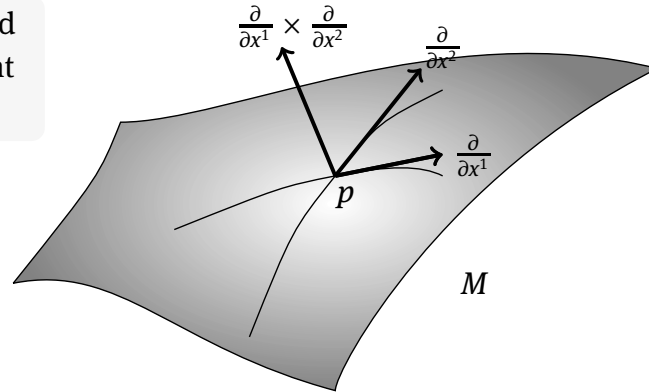
$$\mathbf{h}\left(-r, \frac{3\pi}{2}\right) = \mathbf{k}\left(r, -\frac{\pi}{2}\right), \text{ but } \mathbf{N}_{\mathbf{h}}\left(-r, \frac{3\pi}{2}\right) = -\mathbf{N}_{\mathbf{k}}\left(r, -\frac{\pi}{2}\right),$$

as is easily checked.

The Möbius strip is often realized by taking a long rectangular strip of paper. If one glues the shorter sides together so that the top vertices (and similarly the bottom ones) coincide, one obtains a cylinder. To obtain the Möbius strip instead, the shorter sides are glued in the opposite direction, so that the top left vertex is identified with the bottom right one, and similarly the bottom left is glued to the top right vertex. From a differential geometric perspective, though, this surface is different from the one we

described earlier: the paper strip was not stretched or bent, so it is flat, whereas the parametrized strip has negative curvature, as we shall soon see.

Normal vector induced by a chart  $\mathbf{x}$  at a point  $p \in M$ .



The cross product can be generalized to Euclidean spaces of arbitrary dimension, and provides further insight into the concept of orientation. Given any  $n$  vectors  $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^{n+1}$ , the map

$$\mathbf{u} \mapsto \det(\mathbf{v}_1, \dots, \mathbf{v}_n, \mathbf{u})$$

is a one-form on  $\mathbb{R}^{n+1}$ , and by Corollary 1.4.2, there exists a unique  $\mathbf{w} \in \mathbb{R}^{n+1}$  such that

$$\langle \mathbf{w}, \mathbf{u} \rangle = \det(\mathbf{v}_1, \dots, \mathbf{v}_n, \mathbf{u}).$$

This vector  $\mathbf{w}$  is denoted  $\mathbf{v}_1 \times \dots \times \mathbf{v}_n$ , and is called the *cross product* of  $\mathbf{v}_1, \dots, \mathbf{v}_n$ , cf. also Exercise 5.12. By the properties of the determinant, the cross product is nonzero if and only if the vectors are linearly independent, it changes sign whenever  $\mathbf{v}_i$  is interchanged with  $\mathbf{v}_j$  ( $i \neq j$ ), and  $\mathbf{v}_1 \times \dots \times \mathbf{v}_n \perp \mathbf{v}_i$  for  $1 \leq i \leq n$ . Given  $\mathbf{p} \in \mathbb{R}^{n+1}$ , the canonical isomorphism  $\mathcal{I}_{\mathbf{p}} : \mathbb{R}^{n+1} \rightarrow \mathbb{R}_{\mathbf{p}}^{n+1}$  extends this cross product to any tangent space.

Now, if  $M$  is a hypersurface, then any local parametrization  $\mathbf{h} = \mathbf{x}^{-1} : U \rightarrow M$  of  $M$  induces an orientation of the manifold  $\mathbf{h}(U)$ , since the vector field

$$\mathbf{n}_{\mathbf{h}} := \left( \frac{1}{\left| \frac{\partial}{\partial x^1} \times \dots \times \frac{\partial}{\partial x^n} \right|} \right) \frac{\partial}{\partial x^1} \times \dots \times \frac{\partial}{\partial x^n}$$

is a unit normal field on  $\mathbf{h}(U)$ . If  $M$  is oriented, the parametrization  $\mathbf{h}$  is said to be *consistent with the orientation* if  $\mathbf{n}_{\mathbf{h}}$  represents the given orientation. Notice that if  $\mathbf{h}$  is not consistent with the orientation, then interchanging any two component functions of  $\mathbf{h}$  yields a consistent parametrization.

Evidently,  $M$  is orientable if and only if it admits an atlas such that any two parametrizations  $\mathbf{h}$  and  $\mathbf{k}$  in the atlas satisfy  $\mathbf{n}_{\mathbf{h}} = \mathbf{n}_{\mathbf{k}}$  on the intersection of their domains. In the exercises, the reader is asked to show directly that this is equivalent to our previous criterion for orientability, namely requiring that the Jacobian matrix  $D(\mathbf{h}^{-1} \circ \mathbf{k})$  has positive determinant.

## 7.2 The Gauss map and the second fundamental form

We now associate to each oriented hypersurface  $M^n$  a map  $\gamma : M \rightarrow S^n$  which keeps track of how curved  $M$  is in  $\mathbb{R}^{n+1}$ . Let  $M^n$  be an oriented hypersurface with unit normal field  $\mathbf{n}$ , and denote by  $\pi_2 : T\mathbb{R}^{n+1} = \mathbb{R}^{n+1} \times \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$  the projection onto the second factor. Thus,  $\pi_2$  is the left inverse of  $\mathcal{I}_{\mathbf{p}}$  for any  $\mathbf{p} \in \mathbb{R}^{n+1}$ :  $\pi_2 \circ \mathcal{I}_{\mathbf{p}} = 1_{\mathbb{R}^{n+1}}$ . The *Gauss map* of the oriented hypersurface is defined to be

$$\gamma = \pi_2 \circ \mathbf{n} : M \rightarrow S^n. \quad (7.2.1)$$

Loosely speaking, the Gauss map assigns to each point of  $M$  the unit normal vector at that point, parallel translated back to the origin. The image of the Gauss map therefore measures how much the hypersurface differs from a hyperplane:

**Examples 7.2.1.** (i) A *hyperplane* in  $\mathbb{R}^{n+1}$  is a subset of the form  $M = \{\mathbf{p} \in \mathbb{R}^{n+1} \mid \langle \mathbf{p}, \mathbf{u} \rangle = a\}$ , where  $\mathbf{u}$  is any fixed nonzero element of  $\mathbb{R}^{n+1}$ , and  $a$  is some real number. When  $a = 0$ ,  $M$  is the  $n$ -dimensional subspace  $\mathbf{u}^\perp$ . When  $a \neq 0$ , it is the subspace  $\mathbf{u}^\perp$  parallel translated to any  $\mathbf{p}_0 \in M$ . Indeed,  $\langle \mathbf{p} - \mathbf{p}_0, \mathbf{u} \rangle = 0$  if and only if  $\langle \mathbf{p}, \mathbf{u} \rangle = \langle \mathbf{p}_0, \mathbf{u} \rangle$ , and the latter equals  $a$  if  $\mathbf{p}_0 \in M$ . The Gauss map of such a hyperplane has a single point as image, namely  $\mathbf{u}/|\mathbf{u}|$  (or its negative, depending on the orientation).

(ii) If  $M^2 \subset \mathbb{R}^3$  denotes the cylinder  $\{(x, y, z) \in \mathbb{R}^3 \mid x^2 + y^2 = 1\}$  with the outward orientation, then the Gauss map of  $M$  is given by  $\gamma(x, y, z) = (x, y)$ . In particular, its image consists of the equator in  $S^2$ .

(iii) Let  $M^2 \subset \mathbb{R}^3$  be the paraboloid consisting of all  $(x, y, z)$  satisfying  $z = x^2 + y^2$ .  $M = g^{-1}(0)$ , where  $g : \mathbb{R}^3 \rightarrow \mathbb{R}$ ,  $g(x, y, z) = x^2 + y^2 - z$ , has zero as regular value. By the discussion from the previous section, the normalized gradient of  $g$  yields a unit normal field

$$\mathbf{n} = \frac{1}{\sqrt{1 + 4((u^1)^2 + (u^2)^2)}} (2u^1 D_1 + 2u^2 D_2 - D_3).$$

Since  $M = \{(\mathbf{p}, |\mathbf{p}|^2) \in \mathbb{R}^2 \times \mathbb{R} \mid \mathbf{p} \in \mathbb{R}^2\}$ , the Gauss map of  $M$  for this orientation may be described as

$$\begin{aligned} \gamma : M \subset \mathbb{R}^2 \times \mathbb{R} &\rightarrow S^2 \subset \mathbb{R}^2 \times \mathbb{R}, \\ (\mathbf{p}, |\mathbf{p}|^2) &\mapsto \frac{1}{\sqrt{1 + 4|\mathbf{p}|^2}} (2\mathbf{p}, -1). \end{aligned}$$

It is clear from the formula that the image of  $\gamma$  is contained in the open southern hemisphere of  $S^2$ . The image is actually the whole open hemisphere: indeed, the function  $t \mapsto 2t/(\sqrt{1 + 4t^2})$  is easily seen to be a diffeomorphism from  $\mathbb{R}$  onto  $(-1, 1)$ : in fact, it has positive derivative everywhere so that it is strictly increasing. By the inverse function theorem, it is a diffeomorphism onto its image; finally, it



approaches 1 as  $t \rightarrow \infty$ , and is an odd function. This establishes the claim. In particular, the map

$$\mathbf{f} : \mathbb{R}^2 \rightarrow B_1(0) \subset \mathbb{R}^2,$$

$$\mathbf{p} \mapsto \frac{2}{\sqrt{1 + 4|\mathbf{p}|^2}}\mathbf{p}$$

sends the plane onto the open disk of radius 1 centered at the origin (it is, in fact, a diffeomorphism). The claim follows, since the projection of  $\boldsymbol{\gamma}(\mathbf{p}, |\mathbf{p}|^2)$  onto the  $xy$ -plane is precisely  $\mathbf{f}(\mathbf{p})$ .

(iv) For the sphere  $S^n(r) = \{\mathbf{p} \in \mathbb{R}^{n+1} \mid |\mathbf{p}| = r\}$  of radius  $r > 0$ , the Gauss map is given by

$$\boldsymbol{\gamma} : S^n(r) \rightarrow S^n(1),$$

$$\mathbf{p} \mapsto \frac{1}{r}\mathbf{p},$$

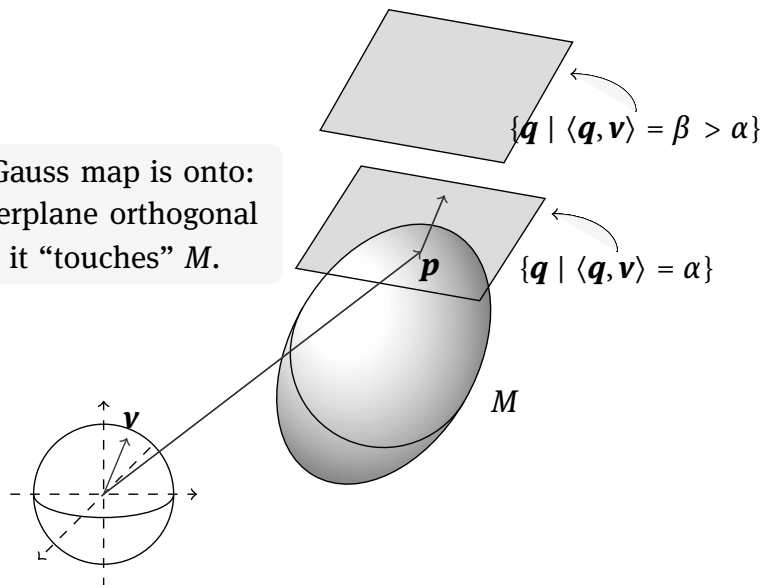
or its negative, depending on the orientation. In particular, it is onto.

The trivial fact that the Gauss map of a sphere is surjective generalizes as follows:

**Theorem 7.2.1.** *Let  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  be a function that has  $\alpha \in \mathbb{R}$  as a regular value, and suppose that  $M = f^{-1}(\alpha)$  is nonempty. If  $M$  is compact, then the image of its Gauss map is the whole sphere.*

*Proof.* Although there are some technical details that must be dealt with, the idea of the proof is simple: orient  $M$  so that its unit normal field points “outward”. Given  $\mathbf{v} \in S^n$ , consider a hyperplane  $\{\mathbf{q} \mid \langle \mathbf{q}, \mathbf{v} \rangle = \beta\}$  orthogonal to  $\mathbf{v}$  that is sufficiently far away (i.e.,  $\beta$  large) that it does not intersect the hypersurface. Parallel translate it until it hits  $M$ . At the point  $\mathbf{p}$  of intersection, this hyperplane will coincide with the tangent plane at  $\mathbf{p}$ , so that  $\boldsymbol{\gamma}(\mathbf{p}) = \mathbf{v}$ .

Showing the Gauss map is onto:  
moving a hyperplane orthogonal  
to  $\mathbf{v} \in S^n$  until it “touches”  $M$ .



Now, for the proof proper: If  $M_\alpha$  denotes the set of all  $\mathbf{p}$  such that  $f(\mathbf{p}) < \alpha$ , and  $M^\alpha$  the set of those satisfying  $f(\mathbf{p}) > \alpha$ , then the complement  $\mathbb{R}^{n+1} \setminus M$  of  $M$  equals the disjoint union of  $M_\alpha$  and  $M^\alpha$ . Notice that one of these must be compact, and the other noncompact. They cannot both be compact, since the ambient space isn't. If they were both noncompact, we could choose some ball  $B$  containing  $M$ .  $\mathbb{R}^{n+1} \setminus B$  would then contain at least one point from  $M_\alpha$  and one from  $M^\alpha$ ; being connected, it would also contain some continuous curve joining the two. This contradicts the intermediate value theorem, which guarantees that any such curve must intersect  $M$ . We may, without loss of generality, assume that  $M_\alpha$  is compact. Given  $\mathbf{v} \in S^n$ , let  $\mathbf{q} \in M$  be a point where the function  $g$ ,  $g(\mathbf{p}) = \langle \mathbf{p}, \mathbf{v} \rangle$ , takes on its maximum value when restricted to  $M$ . By the method of Lagrange multipliers,

$$\nabla f(\mathbf{q}) = \lambda \nabla g(\mathbf{q}) = \lambda \mathcal{I}_q \mathbf{v}$$

for some  $\lambda \in \mathbb{R}$ .

In order to conclude that  $(\nabla f / |\nabla f|) \mathbf{q} = \mathcal{I}_q \mathbf{v}$ , we only need to show that  $\lambda \geq 0$  (notice that  $\lambda$  is nonzero by assumption). To see this, observe that for any  $t > 0$ ,  $\langle \mathbf{q} + t\mathbf{v}, \mathbf{v} \rangle > \langle \mathbf{q}, \mathbf{v} \rangle$ , so  $\mathbf{q} + t\mathbf{v} \notin M$ , and it either belongs to  $M_\alpha$  or else to  $M^\alpha$  for all  $t > 0$ . Since  $M_\alpha$  is compact, it must belong to the latter. Therefore,

$$t \mapsto \frac{f(\mathbf{q} + t\mathbf{v}) - f(\mathbf{q})}{t} > 0, \quad t \neq 0,$$

and

$$0 \leq \lim_{t \rightarrow 0} h(t) = \langle \nabla f(\mathbf{q}), \mathcal{I}_q \mathbf{v} \rangle = \lambda,$$

which establishes the claim. □

For a hypersurface, the second fundamental tensor is unique up to sign, since there are exactly two unit normal vectors at any point, with one being the negative of the other. When the hypersurface is oriented by a unit normal vector field  $\mathbf{n}$ , any sign ambiguity disappears, and we define the second fundamental tensor to be  $S = S_n$ . The next result says that up to parallel translation, the second fundamental tensor equals minus the derivative of the Gauss map:

**Theorem 7.2.2.** *Let  $\mathbf{n}$  denote the unit normal field on  $M$  determined by some orientation of the hypersurface  $M$ , and  $\gamma : M \rightarrow S^n$  the corresponding Gauss map. If  $S$  denotes the second fundamental tensor of  $M$ , then*

$$\pi_2 \circ S = -\pi_2 \circ \gamma_*.$$

*Proof.* The argument is essentially an exercise in notation: for  $\mathbf{p} \in M$ ,  $\mathbf{x} \in M_p$ , let  $\mathbf{c}$  be a curve in  $M$  with  $\dot{\mathbf{c}}(0) = \mathbf{x}$ . By Definitions 2.8.6 and 2.8.7,

$$S\mathbf{x} = -D_x \mathbf{n} = -\mathcal{I}_p (\pi_2 \circ \mathbf{n} \circ \mathbf{c})'(0) = -\mathcal{I}_p (\gamma \circ \mathbf{c})'(0),$$

so that

$$(\pi_2 \circ S)\mathbf{x} = -(\gamma \circ \mathbf{c})'(0) = -\pi_2 (\gamma_* \dot{\mathbf{c}}(0)) = -(\pi_2 \circ \gamma_*)\mathbf{x}. \quad \square$$

The fact that the second fundamental tensor equals, up to sign, the derivative of the Gauss map, suggests that it provides a measure of how much the surface curves. For example, a manifold is said to be *totally geodesic* if its second fundamental tensor vanishes. For a hypersurface  $M$ , this means that the Gauss map has zero derivative, so that, if  $M$  is connected, the map is constant by Theorem 2.4.3. In other words, the unit normal field  $\mathbf{n}$  is parallel, and  $M$  is contained in a hyperplane: given  $\mathbf{p} \in M$ ,  $M \subset \{\mathbf{p} + \mathbf{q} \mid \mathbf{q} \perp \pi_2(\mathbf{n}(\mathbf{p}))\}$ .

One can also consider a weaker condition: a point  $\mathbf{p} \in M$  is said to be *umbilical* if the second fundamental tensor at  $\mathbf{p}$  is a multiple of the identity, and  $M$  itself is said to be *totally umbilic* if every point of  $M$  has that property; i.e., if there is a function  $f : M \rightarrow \mathbb{R}$  such that  $S\mathbf{v} = f(\mathbf{p})\mathbf{v}$  for all  $\mathbf{v} \in M_{\mathbf{p}}$  and  $\mathbf{p} \in M$ . Thus, a totally umbilic manifold is one that “curves equally in all directions”. What is remarkable is that in this case, the function  $f$  is actually constant. This is implicit in the proof of the following:

**Theorem 7.2.3.** *A totally umbilic connected hypersurface is part of a sphere or a hyperplane.*

*Proof.* We first show that the function  $f : M \rightarrow \mathbb{R}$  in the definition of umbilic manifold is constant. Denote, as usual, by  $\mathbf{n}$  the unit normal field of  $M$ . Thus,

$$D_X \mathbf{n} = -S\mathbf{X} = -f\mathbf{X}$$

for any field  $\mathbf{X}$  on  $M$ ; given fields  $\mathbf{Y}$  and  $\mathbf{Z}$ ,

$$\begin{aligned} \langle D_X D_Y \mathbf{n}, \mathbf{Z} \rangle &= \mathbf{X} \langle D_Y \mathbf{n}, \mathbf{Z} \rangle - \langle D_Y \mathbf{n}, D_X \mathbf{Z} \rangle = \mathbf{X} \langle -f\mathbf{Y}, \mathbf{Z} \rangle + f \langle \mathbf{Y}, D_X \mathbf{Z} \rangle \\ &= -(\mathbf{X}f) \langle \mathbf{Y}, \mathbf{Z} \rangle - f \langle D_X \mathbf{Y}, \mathbf{Z} \rangle - f \langle \mathbf{Y}, D_X \mathbf{Z} \rangle + f \langle \mathbf{Y}, D_X \mathbf{Z} \rangle \\ &= -(\mathbf{X}f) \langle \mathbf{Y}, \mathbf{Z} \rangle - f \langle D_X \mathbf{Y}, \mathbf{Z} \rangle. \end{aligned}$$

Using (3.9.1), we obtain

$$\begin{aligned} 0 &= \langle D_X D_Y \mathbf{n} - D_Y D_X \mathbf{n} - D_{[X,Y]} \mathbf{n}, \mathbf{Z} \rangle \\ &= -(\mathbf{X}f) \langle \mathbf{Y}, \mathbf{Z} \rangle - f \langle D_X \mathbf{Y}, \mathbf{Z} \rangle + (\mathbf{Y}f) \langle \mathbf{X}, \mathbf{Z} \rangle + f \langle D_Y \mathbf{X}, \mathbf{Z} \rangle \\ &\quad + f \langle [\mathbf{X}, \mathbf{Y}], \mathbf{Z} \rangle \\ &= \langle (\mathbf{Y}f)\mathbf{X} - (\mathbf{X}f)\mathbf{Y}, \mathbf{Z} \rangle. \end{aligned}$$

Since  $\mathbf{Z}$  is arbitrary,  $(\mathbf{Y}f)\mathbf{X} - (\mathbf{X}f)\mathbf{Y}$  must vanish. Choosing linearly independent  $\mathbf{X}$  and  $\mathbf{Y}$  then implies that  $\mathbf{X}f \equiv 0$  for any vector field  $\mathbf{X}$  on  $M$ , and since  $M$  is connected,  $f$  must equal a constant  $\lambda \in \mathbb{R}$ .

Recall from Examples 2.8.2 (ii) that the position vector field  $\mathbf{P}$  on  $\mathbb{R}^{n+1}$ ,  $\mathbf{P}(\mathbf{a}) = \mathcal{I}_a \mathbf{a}$ , satisfies  $D_X \mathbf{P} = \mathbf{X}$ . Thus, if  $\mathbf{X}$  is a vector field on  $M$ , then

$$D_X(\mathbf{n} + \lambda \mathbf{P}) = -\lambda \mathbf{X} + \lambda \mathbf{X} = 0,$$

and the vector field  $\mathbf{n} + \lambda \mathbf{P}$  is parallel. Parallel translating this field back to the origin, we conclude that the map  $\pi_2 \circ \mathbf{n} + \lambda \pi_2 \circ \mathbf{P}$  is a constant map  $\mathbf{a}$  for some  $\mathbf{a} \in \mathbb{R}^{n+1}$ .

In other words,  $\gamma(\mathbf{p}) + \lambda \mathbf{p} = \mathbf{a}$  for all  $\mathbf{p} \in M$ , where  $\gamma$  is the Gauss map of  $M$ . The case when  $\lambda$  is zero was discussed earlier, and corresponds to a hyperplane. If  $\lambda \neq 0$ , dividing by it and taking norms in the last identity yields

$$\left| \mathbf{p} - \frac{1}{\lambda} \mathbf{a} \right| = \frac{1}{|\lambda|}, \quad \mathbf{p} \in M,$$

so that  $M$  is contained in the sphere of radius  $1/|\lambda|$  centered at  $(1/\lambda)\mathbf{a}$ .  $\square$

### 7.3 Curvature of hypersurfaces

Let  $M$  denote an oriented hypersurface in  $\mathbb{R}^{n+1}$  with unit normal field  $\mathbf{n}$  and second fundamental tensor  $S = S_{\mathbf{n}}$ . In addition to those that we have already encountered, there are several other notions of curvature associated to  $M$  at each point of the hypersurface: The *principal curvatures* at  $\mathbf{p} \in M$  are the eigenvalues of  $S(\mathbf{p})$ . The *mean curvature* is defined to be  $(1/n)$  times the trace of  $S$ . The determinant of  $S(\mathbf{p})$  is called the *Gaussian curvature* at  $\mathbf{p}$ . Notice that unlike the sectional curvature, the first two (and also the third one if  $M$  is odd-dimensional) change sign if one chooses the opposite orientation. Unless otherwise specified, the word curvature by itself will always refer to the sectional curvature.

The definition of the curvature tensor  $R$  implies that

$$R(\mathbf{x}, \mathbf{y})\mathbf{z} = \langle S\mathbf{y}, \mathbf{z} \rangle S\mathbf{x} - \langle S\mathbf{x}, \mathbf{z} \rangle S\mathbf{y}, \quad \mathbf{x}, \mathbf{y}, \mathbf{z} \in M_{\mathbf{p}}, \quad \mathbf{p} \in M,$$

so that the curvature form  $k$  from Chapter 3 is given by

$$k(\mathbf{x}, \mathbf{y}) = \langle R(\mathbf{x}, \mathbf{y})\mathbf{y}, \mathbf{x} \rangle = \langle S\mathbf{x}, \mathbf{x} \rangle \cdot \langle S\mathbf{y}, \mathbf{y} \rangle - \langle S\mathbf{x}, \mathbf{y} \rangle^2. \quad (7.3.1)$$

We say  $M$  has *positive curvature* at  $\mathbf{p} \in M$  if every plane in  $M_{\mathbf{p}}$  has positive sectional curvature. When  $n = 2$ , this is the same as saying that  $M$  has positive Gaussian curvature at  $\mathbf{p}$ , since the sectional curvature of the only tangent plane equals the Gaussian curvature.  $M$  is said to be *positively curved* or to have *positive curvature* if this holds for every  $\mathbf{p} \in M$ . Obvious modifications yield the notions of *nonnegative curvature*, *negative curvature*, and *nonpositive curvature*.

One consequence of (7.3.1) is the following:

**Theorem 7.3.1.**  *$M$  has positive curvature at  $\mathbf{p}$  if and only if the second fundamental form is definite at  $\mathbf{p}$ .*

*Proof.* Recall from Section 2.6 that the second fundamental form is definite if and only if the eigenvalues of  $S$  are all positive or all negative; equivalently, the map

$$\begin{aligned} M_{\mathbf{p}} \times M_{\mathbf{p}} &\rightarrow \mathbb{R} \\ (\mathbf{x}, \mathbf{y}) &\mapsto \langle S\mathbf{x}, \mathbf{y} \rangle \end{aligned}$$

or its negative (that is, replacing  $S$  by  $-S$ ) is an inner product on  $M_p$ . So assume that  $M$  has positive curvature at  $p$ . If  $S$  has only one eigenvalue  $\lambda$  at  $p$ , then  $\lambda \neq 0$ , for otherwise  $k = 0$ . Then  $S = \lambda 1_{M_p}$  is definite. If  $S$  has more than one eigenvalue, choose unit eigenvectors  $\mathbf{x}_i$  corresponding to distinct eigenvalues  $\lambda_i$ ,  $i = 1, 2$ . (7.3.1) implies that

$$0 < k(\mathbf{x}_1, \mathbf{x}_2) = \lambda_1 \cdot \lambda_2.$$

Thus the product of any two eigenvalues is positive, which means that all eigenvalues have the same sign.

Conversely, suppose the second fundamental form is definite. Then the map  $(\mathbf{x}, \mathbf{y}) \mapsto \langle S\mathbf{x}, \mathbf{y} \rangle$  or its negative is an inner product. The Cauchy-Schwarz inequality for either case implies that the right side of (7.3.1) is positive for linearly independent  $\mathbf{x}$  and  $\mathbf{y}$ , and the claim follows.  $\square$

**Theorem 7.3.2.** *There are no compact hypersurfaces of nonpositive curvature.*

*Proof.* Let  $M$  be a compact hypersurface in  $\mathbb{R}^{n+1}$ . The square of the distance to the origin function  $f$ ,  $f(\mathbf{p}) = |\mathbf{p}|^2$ , must assume a maximum at some  $\mathbf{p}_0$  in the compact set  $M$ . We will show that  $M$  has positive curvature at  $\mathbf{p}_0$ . In view of Theorem 7.3.1, it is enough to show that the second fundamental form is definite at that point. Now, the vector  $\mathbf{n} = (1/|\mathbf{p}_0|)\mathcal{I}_{\mathbf{p}_0}\mathbf{p}_0$  is a unit normal to  $M$ , because if  $\mathbf{u} \in M_{\mathbf{p}_0}$ , and  $\mathbf{c}$  is any curve in  $M$  with  $\mathbf{c}(0) = \mathbf{p}_0$  and  $\dot{\mathbf{c}}(0) = \mathbf{u}$ , then

$$0 = (f \circ \mathbf{c})'(0) = 2\langle \mathbf{c}'(0), \mathbf{c}(0) \rangle = 2\langle \dot{\mathbf{c}}(0), \mathcal{I}_{\mathbf{p}_0}\mathbf{p}_0 \rangle = 2|\mathbf{p}_0|\langle \mathbf{u}, \mathbf{n} \rangle.$$

Furthermore,

$$0 \geq (f \circ \mathbf{c})''(0) = 2(\langle \mathbf{c}''(0), \mathbf{c}(0) \rangle + |\mathbf{c}'(0)|^2),$$

so that  $\langle \mathbf{c}''(0), \mathbf{c}(0) \rangle \leq -|\mathbf{u}|^2$ . But

$$\langle \mathbf{c}''(0), \mathbf{c}(0) \rangle = \langle \nabla_{D(0)}\dot{\mathbf{c}}, \mathcal{I}_{\mathbf{p}_0}\mathbf{p}_0 \rangle = |\mathbf{p}_0|\langle \nabla_{D(0)}\mathbf{c}, \mathbf{n} \rangle = -|\mathbf{p}_0|\langle S_n\mathbf{u}, \mathbf{u} \rangle,$$

and consequently  $\langle S_n\mathbf{u}, \mathbf{u} \rangle \geq |\mathbf{u}|^2/|\mathbf{p}_0|^2$  as claimed.  $\square$

A symmetric bilinear form  $b$  on an inner product space  $V$  is said to be *semi-definite* if either  $b(\mathbf{v}, \mathbf{v}) \geq 0$  for all  $\mathbf{v} \in V$  or  $b(\mathbf{v}, \mathbf{v}) \leq 0$  for all  $\mathbf{v} \in V$ . In the exercises, the reader is asked to show that  $M$  has nonnegative curvature at a point if and only if the second fundamental form is semi-definite at that point.

In order to formulate our next result more concisely, we introduce the following concept: the *Laplacian* of a function  $f$  is the function  $\Delta f = \text{tr } H_f$ , where  $H_f$  is the Hessian of  $f$ , and  $\text{tr}$  is the trace. To avoid the notation from becoming cumbersome, we will identify the tangent space of Euclidean space at a point with Euclidean space itself via the usual projection  $\pi_2 : T\mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$ . Thus, the gradient of  $f$  becomes the vector-valued map  $[Df]$ , and the Hessian of  $f$  satisfies

$$H_f(\mathbf{p})\mathbf{x} = D_{\mathbf{x}}\nabla f, \quad \mathbf{x} \in \mathbb{R}^{n+1},$$

where of course the  $\mathbf{x}$  on the right side is to be interpreted as  $\mathcal{I}_{\mathbf{p}}\mathbf{x}$ . To see this, it suffices to verify the identity for  $\mathbf{x} = \mathbf{e}_i$ , since both sides are linear in  $\mathbf{x}$ . Now,

$$D_{\mathbf{e}_i}\nabla f = D_i[Df](\mathbf{p}) = \mathbf{c}'(0),$$

where

$$\mathbf{c}(t) = [Df(\mathbf{p} + t\mathbf{e}_i)] = [D_1f(\mathbf{p} + t\mathbf{e}_i), \dots, D_{n+1}f(\mathbf{p} + t\mathbf{e}_i)].$$

Consequently,

$$\mathbf{c}'(0) = [D_{i1}f, \dots, D_{i(n+1)}f](\mathbf{p}) = H_f(\mathbf{p})\mathbf{e}_i,$$

as claimed.

**Theorem 7.3.3.** *Let  $M^n = f^{-1}(0) \subset \mathbb{R}^{n+1}$ , where 0 is a regular value of  $f$ , and  $\mathbf{p} \in M$ ,  $\mathbf{x}, \mathbf{y} \in M_{\mathbf{p}}$ . Denote by  $S, s, H, G$  the second fundamental tensor, second fundamental form, mean curvature, and Gaussian curvature respectively of  $M$  with respect to the standard orientation of  $M$ . Then*

- $S\mathbf{x} = -\frac{1}{|\nabla f|(\mathbf{p})}(H_f\mathbf{x})^\top$ ;
- $s(\mathbf{x}, \mathbf{y}) = -\frac{1}{|\nabla f|(\mathbf{p})}h_f(\mathbf{x}, \mathbf{y})$ ;
- $H(\mathbf{p}) = \frac{1}{n|\nabla f|(\mathbf{p})}\left(\frac{h_f(\nabla f, \nabla f)}{|\nabla f|^2}(\mathbf{p}) - \Delta f(\mathbf{p})\right)$ ;
- $G(\mathbf{p}) = \frac{1}{|\nabla f|^n(\mathbf{p})}\det H_f^\top(\mathbf{p})$ ;
- *If  $\mathbf{x}$  and  $\mathbf{y}$  form an orthonormal basis of a plane  $P \subset M_{\mathbf{p}}$ , then the sectional curvature of  $P$  is*

$$K_P = \frac{1}{|\nabla f|^2(\mathbf{p})}\det \begin{bmatrix} h_f(\mathbf{x}, \mathbf{x}) & h_f(\mathbf{x}, \mathbf{y}) \\ h_f(\mathbf{x}, \mathbf{y}) & h_f(\mathbf{y}, \mathbf{y}) \end{bmatrix}.$$

*Proof.*

$$\begin{aligned} S\mathbf{x} &= -D_{\mathbf{x}}\left(\frac{1}{|\nabla f|}\nabla f\right) = -\mathbf{x}\left(\frac{1}{|\nabla f|}\right)(\nabla f)^\top(\mathbf{p}) - \frac{1}{|\nabla f|(\mathbf{p})}(D_{\mathbf{x}}\nabla f)^\top \\ &= -\frac{1}{|\nabla f|(\mathbf{p})}(D_{\mathbf{x}}\nabla f)^\top = -\frac{1}{|\nabla f|(\mathbf{p})}(H_f\mathbf{x})^\top. \end{aligned}$$

This establishes the first identity. The others follow from the definitions of the different types of curvature, together with (7.3.1) and the identity just proved.  $\square$

There is a useful alternative formula for the Gaussian curvature in the above proposition. In order to derive it, we need the following:

**Lemma 7.3.1.** *Let  $L : V \rightarrow V$  be a linear transformation on an  $(n + 1)$ -dimensional inner product space  $V$ . Let  $W$  be an  $n$ -dimensional subspace of  $V$  and  $\pi : V \rightarrow W$  the orthogonal projection onto  $W$ ; i.e., if  $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2 \in W \oplus W^\perp$ , then  $\pi\mathbf{x} = \mathbf{x}_1$ . The restriction  $\pi \circ L|_W$  of  $\pi \circ L$  to  $W$  is then an operator on  $W$ . Its determinant is given by*

$$\det(\pi \circ L|_W) = \langle \tilde{L}\mathbf{x}, \mathbf{x} \rangle,$$

where  $\mathbf{x}$  is a unit vector in  $W^\perp$ , and  $\tilde{L}$  is the linear map adjugate to  $L$  (see Theorem 1.3.6).

*Proof.* Extend  $\mathbf{x}$  to an ordered orthonormal basis  $\mathcal{B} = \{\mathbf{x}, \mathbf{v}_1, \dots, \mathbf{v}_n\}$  of  $V$ . Then the last  $n$  vectors form an orthonormal basis  $\mathcal{B}_1$  of  $W$ , and the  $(i, j)$ -th entry of the matrix of  $\pi \circ L|_W$  with respect to this basis is

$$\langle \pi L\mathbf{v}_j, \mathbf{v}_i \rangle = \langle L\mathbf{v}_j - \langle L\mathbf{v}_j, \mathbf{x} \rangle \mathbf{x}, \mathbf{v}_i \rangle = \langle L\mathbf{v}_j, \mathbf{v}_i \rangle,$$

so that the determinant of  $\pi \circ L|_W$  equals the determinant of the matrix obtained by deleting the first row and column of  $[L]_{\mathcal{B}}$ . On the other hand,  $\langle \tilde{L}\mathbf{x}, \mathbf{x} \rangle$  is the  $(1, 1)$  entry of  $[\tilde{L}]_{\mathcal{B}}$ , which by definition of  $\tilde{L}$ , is that same determinant.  $\square$

The lemma, together with the formula for  $G$  from Theorem 7.3.3 now immediately imply the following:

**Proposition 7.3.1.** *With the hypotheses of Theorem 7.3.3, the Gaussian curvature of  $M$  is given by*

$$G = \frac{1}{|\nabla f|^{n+2}} \langle \tilde{H}_f \nabla f, \nabla f \rangle.$$

An important special case occurs when the function  $f$  is a quadratic form. As before, we have identified in the next result Euclidean space with its tangent space at a given point.

**Proposition 7.3.2.** *Suppose  $L$  is a self-adjoint operator on  $\mathbb{R}^{n+1}$ . Denote by  $b$  the associated scalar product,  $b(\mathbf{x}, \mathbf{y}) = \langle L\mathbf{x}, \mathbf{y} \rangle$ , and by  $f$  the corresponding quadratic form,  $f(\mathbf{x}) = b(\mathbf{x}, \mathbf{x})$ . Consider a regular value  $a$  of  $f$ , and the corresponding hypersurface  $M = f^{-1}(a)$ . Given  $\mathbf{p} \in M$ ,*

- (1)  $\nabla f(\mathbf{p}) = 2L\mathbf{p}$ ;
- (2)  $H_f = 2L, h_f = 2b$ ;
- (3)  $G(\mathbf{p}) = \frac{\det L}{|L\mathbf{p}|^{n+2}} a$ ;
- (4) *If  $\mathbf{x}$  and  $\mathbf{y}$  form an orthonormal basis of a plane  $P \subset M_{\mathbf{p}}$ , then the sectional curvature of  $P$  is*

$$K_P = \frac{1}{|L\mathbf{p}|^2} \det \begin{bmatrix} b(\mathbf{x}, \mathbf{x}) & b(\mathbf{x}, \mathbf{y}) \\ b(\mathbf{x}, \mathbf{y}) & b(\mathbf{y}, \mathbf{y}) \end{bmatrix}.$$

*Proof.* Let  $\mathbf{x} \in \mathbb{R}^{n+1}$ , and consider the ray  $\mathbf{c}_x, \mathbf{c}_x(t) = \mathbf{p} + t\mathbf{x}$ . Then

$$\langle \nabla f(\mathbf{p}), \mathbf{x} \rangle = \mathbf{x}f = (f \circ \mathbf{c}_x)'(0).$$

Now,

$$(f \circ \mathbf{c}_x)(t) = b(\mathbf{p} + t\mathbf{x}, \mathbf{p} + t\mathbf{x}) = b(\mathbf{p}, \mathbf{p}) + 2tb(\mathbf{p}, \mathbf{x}) + t^2b(\mathbf{x}, \mathbf{x}),$$

so that  $\langle \nabla f(\mathbf{p}), \mathbf{x} \rangle = 2b(\mathbf{p}, \mathbf{x}) = 2\langle L\mathbf{p}, \mathbf{x} \rangle$ , which establishes (1).

For (2), observe that (1) implies  $D_{ij}f(\mathbf{p}) = 2\langle L\mathbf{p}, \mathbf{e}_i \rangle$ , so that  $D_{jij}f(\mathbf{p}) = \mathbf{c}'(0)$ , where

$$\mathbf{c}(t) = 2\langle L(\mathbf{p} + t\mathbf{e}_j), \mathbf{e}_i \rangle = 2(\langle L\mathbf{p}, \mathbf{e}_i \rangle + t\langle L\mathbf{e}_j, \mathbf{e}_i \rangle).$$

Thus,  $D_{jij}f(\mathbf{p}) = 2\langle L\mathbf{e}_j, \mathbf{e}_i \rangle$ . Since these are the respective entries of the matrices of  $H_f(\mathbf{p})$  and  $2L$  in the standard basis, both transformations coincide. The formula for the

hessian form is immediate, and the one for sectional curvature follows from Theorem 7.3.3. Finally, by Proposition 7.3.1 together with (1) and (2),

$$\begin{aligned} G(\mathbf{p}) &= \frac{1}{|\nabla f(\mathbf{p})|^{n+2}} \langle \tilde{H}_f \nabla f, \nabla f \rangle(\mathbf{p}) = \frac{1}{|2L\mathbf{p}|^{n+2}} \langle \widetilde{(2L)}(2L)\mathbf{p}, 2L\mathbf{p} \rangle \\ &= \frac{\det(2L)}{2^{n+2}|L\mathbf{p}|^{n+2}} \langle 2L\mathbf{p}, \mathbf{p} \rangle = \frac{\det L}{|L\mathbf{p}|^{n+2}} a. \end{aligned} \quad \square$$

Notice that if the quadratic form  $b$  is definite, then  $M$  has positive curvature. If it is merely semi-definite, then  $M$  has nonnegative curvature by Exercise 7.17. When  $n = 2$ , the surfaces of the type considered in Proposition 7.3.2 are called *quadrics*. The reader is invited to classify them in the exercises.

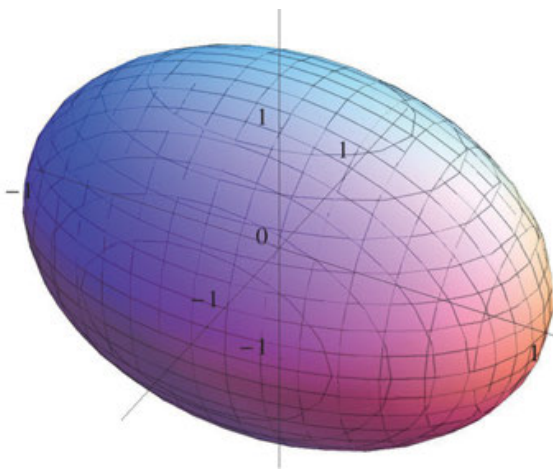


Fig. 7.2: The ellipsoid  $x^2 + \frac{2}{3}y^2 + \frac{1}{2}z^2 = 1$ .

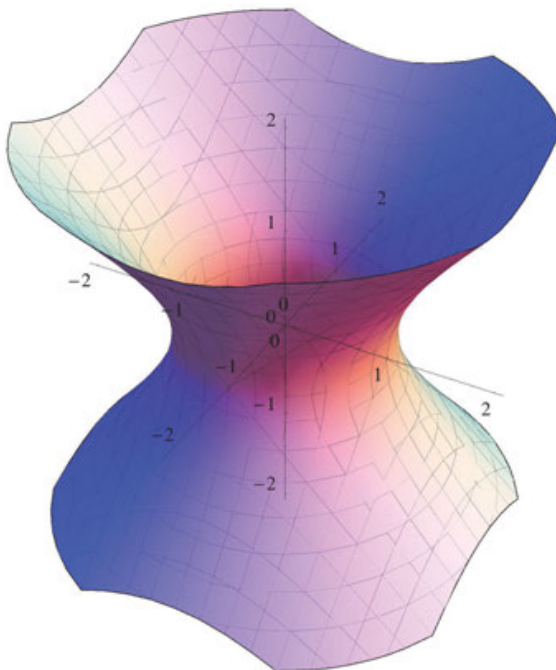


Fig. 7.3: The hyperboloid  $x^2 + y^2 - z^2 = 1$ .



**Example 7.3.1.** When  $n = 2$ , the Gaussian curvature and the sectional curvature coincide. In order to compute the latter, one can therefore use Proposition 7.3.2(3) without having to explicitly find a basis for the tangent space:

- The ellipsoid  $(x^2/a^2) + (y^2/b^2) + (z^2/c^2) = 1$  is the surface  $\{\mathbf{x} \in \mathbb{R}^3 \mid \langle L\mathbf{x}, \mathbf{x} \rangle = 1\}$ , where the matrix of  $L : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  in the standard basis is

$$[L] = \begin{bmatrix} \frac{1}{a^2} & 0 & 0 \\ 0 & \frac{1}{b^2} & 0 \\ 0 & 0 & \frac{1}{c^2} \end{bmatrix}.$$

The formula for the Gaussian curvature then implies that the curvature at  $(x, y, z)$  equals

$$K(x, y, z) = [abc(\frac{x^2}{a^4} + \frac{y^2}{b^4} + \frac{z^2}{c^4})]^{-2}.$$

- The hyperboloid  $(x^2/a^2) + (y^2/b^2) - (z^2/c^2) = 1$  differs from the ellipsoid in that the last diagonal entry of  $[L]$  has its sign reversed. The curvature is therefore given by

$$K(x, y, z) = -[abc(\frac{x^2}{a^4} + \frac{y^2}{b^4} + \frac{z^2}{c^4})]^{-2}.$$

## 7.4 The fundamental theorem for hypersurfaces

We have already remarked in Chapter 3 that if  $\mathbf{F}$  is a rigid motion of Euclidean space, and  $M$  is a submanifold, then the restriction of  $\mathbf{F}$  to  $M$  is an isometry with  $\mathbf{F}(M)$ . We also noticed that the converse is not true, that is, an isometry  $\mathbf{f} : M_1 \rightarrow M_2$  between submanifolds need not be the restriction of a rigid motion. In this section, we show it is nevertheless true for hypersurfaces, provided  $\mathbf{f}$  preserves in addition the second fundamental form (it is also true for more general submanifolds, under additional hypotheses). Specifically:

**Theorem 7.4.1.** *Let  $M_i$  denote oriented hypersurfaces in  $\mathbb{R}^{n+1}$  with unit normal fields  $\mathbf{n}_i$  and second fundamental tensors  $S_i$ ,  $i = 1, 2$ . Suppose that  $\mathbf{f} : M_1 \rightarrow M_2$  is an isometry, and that*

$$S_2 \mathbf{f}_* \mathbf{x} = \mathbf{f}_* S_1 \mathbf{x}, \quad \mathbf{x} \in M_{1\mathbf{p}}, \quad \mathbf{p} \in M_1.$$

*If  $M_1$  is connected, then there exists a rigid motion  $\mathbf{F}$  of  $\mathbb{R}^{n+1}$  such that  $\mathbf{f} = \mathbf{F} \circ \iota$ , where  $\iota : M_1 \hookrightarrow \mathbb{R}^{n+1}$  denotes inclusion.*

*Proof.* Notice first that  $\mathbf{F}$ , if it exists, is, up to sign, entirely determined by  $\mathbf{f}$ : indeed,  $\mathbf{F}(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$  for some orthogonal transformation  $A$  and  $\mathbf{b} \in \mathbb{R}^{n+1}$ . Fix any  $\mathbf{p} \in M_1$ . Then  $\mathbf{F}(\mathbf{p}) = A\mathbf{p} + \mathbf{b} = \mathbf{f}(\mathbf{p})$ , so that  $\mathbf{b}$  must equal  $\mathbf{f}(\mathbf{p}) - A\mathbf{p}$ . Furthermore,  $D\mathbf{F}(\mathbf{p}) = A$ , so that  $A$  must, up to parallel translation, equal  $\mathbf{f}_{*\mathbf{p}}$  on  $M_{1\mathbf{p}}$ . Denoting as usual by  $\pi_2$  the projection from  $T\mathbb{R}^{n+1}$  to  $\mathbb{R}^{n+1}$ , it follows that  $A(\pi_2 \mathbf{n}_1(\mathbf{p})) = \pm \pi_2 \mathbf{n}_2(\mathbf{f}(\mathbf{p}))$ . This completely determines the operator  $A$  on  $\mathbb{R}^{n+1}$ , and together with the condition on  $\mathbf{b}$ , determines  $\mathbf{F}$ .

With this in mind, define  $\mathbf{F}$  as above; i.e., consider the orthogonal transformation  $A$  determined by

$$A(\pi_2 \mathbf{n}_1(\mathbf{p})) = \pi_2 \mathbf{n}_2(\mathbf{f}(\mathbf{p})), \quad A(\pi_2 \mathbf{x}) = \pi_2 \mathbf{f}_* \mathbf{x} \text{ for } \mathbf{x} \in M_{1\mathbf{p}},$$

and the rigid motion  $\mathbf{F}$ , where

$$\mathbf{F}(\mathbf{x}) = A\mathbf{x} + \mathbf{f}(\mathbf{p}) - A\mathbf{p}.$$

The result will follow once we show that the map  $\mathbf{G} := \mathbf{F}^{-1} \circ \mathbf{f}$  is the identity on  $M_1$ . In fact, since  $M_1$  is connected, it suffices to show that  $\mathbf{G} \circ \mathbf{c} = \mathbf{c}$  for any curve  $\mathbf{c}$  in  $M_1$  with  $\mathbf{c}(0) = \mathbf{p}$ . By construction,  $\mathbf{G}(\mathbf{p}) = \mathbf{p}$  (and  $\mathbf{G}_* \mathbf{p}$  is the identity on the tangent space of  $M_1$  at  $\mathbf{p}$ ), so  $(\mathbf{G} \circ \mathbf{c})(0) = \mathbf{c}(0)$ , and we only need to show that  $(\mathbf{G} \circ \mathbf{c})' = \mathbf{c}'$ .

So choose an orthonormal basis  $\mathbf{x}_i$  of  $M_{1\mathbf{p}}$ , and let  $\mathbf{X}_i$  denote the parallel field along  $\mathbf{c}$  with  $\mathbf{X}_i(0) = \mathbf{x}_i$ . By (3.11.6), the fields  $\mathbf{Y}_i := \mathbf{G}_* \mathbf{X}_i$  are parallel and orthonormal along  $\mathbf{G} \circ \mathbf{c}$ . Moreover,  $\mathbf{Y}_i(0) = \mathbf{X}_i(0)$  because  $\mathbf{G}_* \mathbf{p}$  is the identity. Now,  $\dot{\mathbf{c}} = \sum_i \langle \dot{\mathbf{c}}, \mathbf{X}_i \rangle \mathbf{X}_i$ , and

$$\mathbf{G} \circ \dot{\mathbf{c}} = \sum_i \langle \mathbf{G}_* \dot{\mathbf{c}}, \mathbf{Y}_i \rangle \mathbf{Y}_i = \sum_i \langle \mathbf{G}_* \dot{\mathbf{c}}, \mathbf{G}_* \mathbf{X}_i \rangle \mathbf{Y}_i = \sum_i \langle \dot{\mathbf{c}}, \mathbf{X}_i \rangle \mathbf{Y}_i,$$

so that if  $X_i = \pi_2 \mathbf{X}_i$  and  $Y_i = \pi_2 \mathbf{Y}_i$ , then

$$\mathbf{c}' = \sum_i \langle \dot{\mathbf{c}}, \mathbf{X}_i \rangle X_i, \quad (\mathbf{G} \circ \mathbf{c})' = \sum_i \langle \dot{\mathbf{c}}, \mathbf{X}_i \rangle Y_i.$$

It therefore remains to show that  $X_i = Y_i$  for  $1 \leq i \leq n$ . First, observe that  $M := \mathbf{F}^{-1}(M_2)$  is a manifold isometric to  $M_2$  with unit normal  $\mathbf{n} := \mathbf{F}_*^{-1} \circ \mathbf{n}_2 \circ \mathbf{F}$ . Set  $\mathbf{X}_{n+1} = \mathbf{n}_1 \circ \mathbf{c}$ ,  $\mathbf{Y}_{n+1} = \mathbf{n} \circ \mathbf{G} \circ \mathbf{c}$ , so that  $\{\mathbf{X}_i\}_{1 \leq i \leq n+1}$  and  $\{\mathbf{Y}_i\}_{1 \leq i \leq n+1}$  are orthonormal bases of  $\mathbb{R}_c^{n+1}$  and  $\mathbb{R}_{\mathbf{G} \circ \mathbf{c}}^{n+1}$  respectively. Now,

$$X_i' = \sum_j g_{ij} X_j, \quad Y_i' = \sum_j h_{ij} Y_j, \quad \text{where } g_{ij} = \langle X_i', X_j \rangle, \quad h_{ij} = \langle Y_i', Y_j \rangle.$$

Since  $\mathbf{X}_i$  and  $\mathbf{Y}_i$  are parallel for  $i \leq n$ ,  $g_{ij} \equiv h_{ij} \equiv 0$  when  $i, j \leq n$ . In the case that  $j = n + 1$ , we use the fact that  $\mathbf{f}$  preserves the second fundamental form: The condition  $S_2 \mathbf{f}_* \mathbf{x} = \mathbf{f}_* S_1 \mathbf{x}$  means that  $D_{\mathbf{f}_* \mathbf{x}} \mathbf{n}_2 = \mathbf{f}_* D_{\mathbf{x}} \mathbf{n}_1$ . Thus,

$$\begin{aligned} D_{\mathbf{G}_* \mathbf{x}} \mathbf{n} &= D_{\mathbf{G}_* \mathbf{x}} (\mathbf{F}_*^{-1} \circ \mathbf{n}_2 \circ \mathbf{F}) = \mathbf{F}_*^{-1} D_{\mathbf{G}_* \mathbf{x}} (\mathbf{n}_2 \circ \mathbf{F}) = \mathbf{F}_*^{-1} D_{\mathbf{F}_* \circ \mathbf{G}_* \mathbf{x}} \mathbf{n}_2 \\ &= \mathbf{F}_*^{-1} D_{\mathbf{f}_* \mathbf{x}} \mathbf{n}_2 = \mathbf{F}_*^{-1} \circ \mathbf{f}_* D_{\mathbf{x}} \mathbf{n}_1 \\ &= \mathbf{G}_* D_{\mathbf{x}} \mathbf{n}_1. \end{aligned}$$

Now take  $\mathbf{x} = \dot{\mathbf{c}}$  in the above identity to conclude that  $\mathbf{Y}_{n+1}' = \mathbf{G}_* \mathbf{X}_{n+1}'$ . It follows that

$$\begin{aligned} h_{in+1} &= \langle \mathbf{Y}_i', \mathbf{Y}_{n+1} \rangle = -\langle \mathbf{Y}_i, \mathbf{Y}_{n+1}' \rangle = -\langle \mathbf{G}_* \mathbf{X}_i, \mathbf{G}_* \mathbf{X}_{n+1}' \rangle \\ &= -\langle \mathbf{X}_i, \mathbf{X}_{n+1}' \rangle = \langle \mathbf{X}_i', \mathbf{X}_{n+1} \rangle \\ &= g_{in+1} \end{aligned}$$

for all  $i \leq n$ . But this also implies equality when  $i = n + 1$ , because for any  $i$  and  $j$ ,

$$g_{ij} + g_{ji} = \langle \mathbf{X}'_i, \mathbf{X}_j \rangle + \langle \mathbf{X}_i, \mathbf{X}'_j \rangle = \langle \mathbf{X}_i, \mathbf{X}_j \rangle' = 0,$$

and similarly for  $h_{ij}$ .

Summarizing, we have established that

$$X'_i = \sum_j g_{ij} X_j, \quad Y'_i = \sum_j g_{ij} Y_j.$$

Let  $Z_i = X_i - Y_i$ . Then  $Z_i$  satisfies the system of ordinary differential equations

$$Z'_i = \sum_j g_{ij} Z_j, \quad Z_i(0) = \mathbf{0}.$$

The constant vector fields  $Z_i \equiv \mathbf{0}$  satisfy this system. By uniqueness of solutions,  $Z_i \equiv \mathbf{0}$ ; i.e.  $X_i \equiv Y_i$ .  $\square$

**Example 7.4.1.** Let  $\mathbf{c} : (a, b) \rightarrow \mathbb{R}^2$  denote a regular curve with no self-intersections (i.e.,  $\mathbf{c}$  is one-to-one, so that  $M = \mathbf{c}(a, b)$  is a – perhaps only immersed – submanifold of  $\mathbb{R}^2$ ), which we may assume is parametrized by arclength. Choose the unit normal vector field  $\mathbf{n}$  along  $\mathbf{c}$  so that  $\det(\mathbf{c}', \boldsymbol{\pi}_2 \circ \mathbf{n}) \equiv 1$ , and denote by  $S$  the corresponding second fundamental tensor. The *curvature* of  $\mathbf{c}$  is

$$\kappa = \langle S\dot{\mathbf{c}}, \dot{\mathbf{c}} \rangle.$$

The curvature is, therefore, up to sign, equal to the norm of the second fundamental tensor. It is also, up to sign, the norm of the acceleration:

$$\langle S\dot{\mathbf{c}}, \dot{\mathbf{c}} \rangle = -\langle \mathbf{n}', \dot{\mathbf{c}} \rangle = \langle \mathbf{n}, \dot{\mathbf{c}}' \rangle = \pm |\dot{\mathbf{c}}'|,$$

since  $\dot{\mathbf{c}}' \perp \dot{\mathbf{c}}$ . In the exercises, the reader is asked to show that a circle of radius  $r$  has  $|\kappa| = 1/r$ . This means that going out along the normal line to  $\mathbf{c}(t)$  at distance  $1/|\kappa(t)|$  (in the appropriate direction) and drawing a circle centered there with radius equal to that distance yields the circle that best approximates  $\mathbf{c}$  at  $t$ . It is called the *circle of curvature* at  $t$ .

Now, suppose  $\mathbf{c}_i : (a, b) \rightarrow \mathbb{R}^2$ ,  $i = 1, 2$ , are two curves parametrized by arc length with no self-intersections. Then the map  $\mathbf{f} : \mathbf{c}_1(a, b) \rightarrow \mathbf{c}_2(a, b)$ , where  $\mathbf{f} = \mathbf{c}_2 \circ \mathbf{c}_1^{-1}$ , is an isometry. The above theorem says that  $\mathbf{f}$  extends to a rigid motion of the plane if and only if the curves have the same curvature; i.e.,  $\kappa_2 \circ \mathbf{f} = \kappa_1$ .

## 7.5 Curvature in local coordinates

Although a coordinate-free approach is preferable, it is not always feasible. In this section, we compute the matrix of the second fundamental tensor of a hypersurface  $M^n$  with respect to a basis of coordinate vector fields. In the case  $n = 2$ , the formulas are simple enough to obtain a reasonable expression for the sectional curvature.

Consider a chart  $(U, \mathbf{x})$  with coordinate vector fields  $\partial/\partial x^i = \mathbf{x}_*^{-1} \mathbf{D}_i \circ \mathbf{x}$  and corresponding unit normal field

$$\mathbf{n} = \frac{\times_{i=1}^n \frac{\partial}{\partial x^i}}{\left| \times_{i=1}^n \frac{\partial}{\partial x^i} \right|}$$

to  $M$  on  $U$ . Define  $n \times n$  matrix-valued maps  $s, g$ , on  $U$  by

$$s_{ij} = \left\langle S \frac{\partial}{\partial x^i}, \frac{\partial}{\partial x^j} \right\rangle, \quad g_{ij} = \left\langle \frac{\partial}{\partial x^i}, \frac{\partial}{\partial x^j} \right\rangle.$$

As usual, let  $n = \boldsymbol{\pi}_2 \circ \mathbf{n}$ . We begin by finding an expression for  $s$ :

$$\begin{aligned} s_{ij} &= \left\langle S \frac{\partial}{\partial x^i}, \frac{\partial}{\partial x^j} \right\rangle = - \left\langle D_{\partial/\partial x^i} \mathbf{n}, \frac{\partial}{\partial x^j} \right\rangle = \left\langle D_{\partial/\partial x^i} \frac{\partial}{\partial x^j}, \mathbf{n} \right\rangle \\ &= \left\langle \left( D_{\mathbf{x}_*^{-1} \mathbf{D}_i} \frac{\partial}{\partial x^j} \right) \circ \mathbf{x}, \mathbf{n} \right\rangle = \left\langle \left( D_{\mathbf{D}_i} \left( \frac{\partial}{\partial x^j} \circ \mathbf{x}^{-1} \right) \right) \circ \mathbf{x}, \mathbf{n} \right\rangle \\ &= \left\langle (D_{\mathbf{D}_i} \mathbf{x}_*^{-1} \mathbf{D}_j) \circ \mathbf{x}, \mathbf{n} \right\rangle = \left\langle (D(D\mathbf{x}^{-1} \mathbf{e}_j) \mathbf{e}_i) \circ \mathbf{x}, \mathbf{n} \right\rangle \end{aligned}$$

by (2.8.3). By symmetry of  $s$ ,

$$s_{ij} = \langle (D^2 \mathbf{x}^{-1} \mathbf{e}_i \mathbf{e}_j) \circ \mathbf{x}, \mathbf{n} \rangle. \quad (7.5.1)$$

Alternatively, in terms of a parametrization  $\mathbf{h} = \mathbf{x}^{-1}$ ,

$$s_{ij} \circ \mathbf{h} = \langle D^2 \mathbf{h} \mathbf{e}_i \mathbf{e}_j, \mathbf{n} \circ \mathbf{h} \rangle. \quad (7.5.2)$$

We emphasize that the first term on the right side of the equality sign in (7.5.1) is the map from  $U$  to  $\mathbb{R}^{n+1}$  that sends  $\mathbf{p}$  to  $(D\mathbf{k})(\mathbf{x}(\mathbf{p}))\mathbf{e}_j$ , where  $\mathbf{k} = D\mathbf{x}^{-1} \mathbf{e}_i : \mathbf{x}(U) \rightarrow \mathbb{R}^{n+1}$ .  $s$  is not, in general, the matrix of  $S$  in the basis of coordinate vector fields, unless the latter are orthonormal. To find this matrix, we use the following:

**Lemma 7.5.1.** *Suppose  $L : V \rightarrow V$  is a linear map on an inner product space  $V$  with basis  $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ . If  $A$  denotes the matrix with  $(i, j)$  entry  $a_{ij} = \langle L\mathbf{v}_i, \mathbf{v}_j \rangle$ , and  $B$  the one with  $(i, j)$  entry  $b_{ij} = \langle \mathbf{v}_i, \mathbf{v}_j \rangle$ , then the matrix of  $L$  with respect to  $\mathcal{B}$  is  $[L]_{\mathcal{B}} = (AB^{-1})^T$ . In particular, if  $L$  is self-adjoint, then  $[L]_{\mathcal{B}} = B^{-1}A$ .*

*Proof.* If  $l_{ij}$  is the  $(i, j)$ -th entry of the matrix of  $L$  in the given basis, then

$$a_{ij} = \langle L\mathbf{v}_i, \mathbf{v}_j \rangle = \left\langle \sum_k l_{ki} \mathbf{v}_k, \mathbf{v}_j \right\rangle = \sum_k l_{kj} b_{ki}$$

so that  $A = [L]_{\mathcal{B}}^T B$ , and the claim follows. If in addition  $L$  is self-adjoint, then  $A$  is a symmetric matrix.  $B$  is always symmetric, and it is easy to see that  $B^{-1}$  must then also be symmetric. Thus,  $(AB^{-1})^T = B^{-1T} A^T = B^{-1}A$ .  $\square$

As an immediate consequence, we obtain:

**Proposition 7.5.1.** *The matrix of  $S$  in the basis  $\{\partial/\partial x^i\}$  equals  $g^{-1}s$ , where*

$$s_{ij} = \langle (D^2 \mathbf{x}^{-1} \mathbf{e}_i \mathbf{e}_j) \circ \mathbf{x}, \mathbf{n} \rangle, \quad g_{ij} = \left\langle \frac{\partial}{\partial x^i}, \frac{\partial}{\partial x^j} \right\rangle.$$

The proposition may of course be stated alternatively in terms of a parametrization  $\mathbf{h}$  using (7.5.2). This is in fact the way it is done in the following application:

**Corollary 7.5.1.** *If  $M^2 \subset \mathbb{R}^3$  is locally parametrized by  $(U, \mathbf{h})$ , then the sectional curvature of  $M$  at  $\mathbf{p} = \mathbf{h}(\mathbf{a})$  is*

$$K(\mathbf{p}) = \frac{\det \langle (D^2 \mathbf{h} \mathbf{e}_i \mathbf{e}_j)(\mathbf{a}), \mathbf{n}(\mathbf{p}) \rangle}{\det \langle D\mathbf{h} \mathbf{e}_i, D\mathbf{h} \mathbf{e}_j \rangle(\mathbf{a})} = \frac{\det \langle (D^2 \mathbf{h} \mathbf{e}_i \mathbf{e}_j)(\mathbf{a}), \mathbf{n}(\mathbf{p}) \rangle}{|D\mathbf{h} \mathbf{e}_1 \times D\mathbf{h} \mathbf{e}_2|^2(\mathbf{a})}.$$

*Proof.* By (7.3.1),  $K = \det S$ , and the first part of the identity follows from the proposition, together with (7.5.2). For the second part, if  $\theta$  denotes the angle between the two coordinate vector fields, then

$$\begin{aligned} \det \langle D\mathbf{h} \mathbf{e}_i, D\mathbf{h} \mathbf{e}_j \rangle &= |D\mathbf{h} \mathbf{e}_1|^2 |D\mathbf{h} \mathbf{e}_2|^2 - \langle D\mathbf{h} \mathbf{e}_1, D\mathbf{h} \mathbf{e}_2 \rangle^2 \\ &= |D\mathbf{h} \mathbf{e}_1|^2 |D\mathbf{h} \mathbf{e}_2|^2 - |D\mathbf{h} \mathbf{e}_1|^2 |D\mathbf{h} \mathbf{e}_2|^2 \cos^2 \theta \\ &= |D\mathbf{h} \mathbf{e}_1|^2 |D\mathbf{h} \mathbf{e}_2|^2 \sin^2 \theta \\ &= |D\mathbf{h} \mathbf{e}_1 \times D\mathbf{h} \mathbf{e}_2|^2 \end{aligned}$$

by (1.6.2). □

**Example 7.5.1.** Consider the surface  $M^2 \subset \mathbb{R}^3$  consisting of the graph of a function  $f : U \subset \mathbb{R}^2 \rightarrow \mathbb{R}$ ; i.e.,  $M = \{(x, y, f(x, y)) \mid (x, y) \in U\}$ . There is a natural parametrization  $\mathbf{h} : U \rightarrow \mathbb{R}^3$  of all of  $M$  given by  $\mathbf{h}(\mathbf{a}) = (\mathbf{a}, f(\mathbf{a}))$  for  $\mathbf{a} \in U$ .

$D\mathbf{h} \mathbf{e}_1$  and  $D\mathbf{h} \mathbf{e}_2$  are the columns of the matrix

$$[D\mathbf{h}] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ D_1 f & D_2 f \end{bmatrix},$$

so that

$$D^2 \mathbf{h} \mathbf{e}_i = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ D_{1if} & D_{2if} \end{bmatrix}.$$

It follows that

$$(D^2 \mathbf{h} \mathbf{e}_i) \mathbf{e}_j = \begin{bmatrix} 0 \\ 0 \\ D_{ijf} \end{bmatrix}, \quad \mathbf{n} \circ \mathbf{h} = \frac{1}{\sqrt{1 + (D_1 f)^2 + (D_2 f)^2}} \begin{bmatrix} -D_1 f \\ -D_2 f \\ 1 \end{bmatrix}.$$

By Proposition 7.5.1, we have

$$s \circ \mathbf{h} = \frac{1}{(1 + |\nabla f|^2)^{1/2}} \begin{bmatrix} D_{11} f & D_{12} f \\ D_{21} f & D_{22} f \end{bmatrix}, \quad g \circ \mathbf{h} = \begin{bmatrix} 1 + (D_1 f)^2 & D_1 f D_2 f \\ D_1 f D_2 f & 1 + (D_2 f)^2 \end{bmatrix},$$

and by Corollary 7.5.1, the sectional curvature of  $M$  is given by

$$K \circ \mathbf{h} = \frac{\det H_f}{(1 + |\nabla f|^2)^2}. \quad (7.5.3)$$

In particular, the curvature and the determinant of the Hessian of  $f$  have the same sign. It is worth noting that this equation can also be derived from Theorem 7.3.3 if one writes  $M = h^{-1}(0)$ , where  $h(x, y, z) = f(x, y) - z$ .

## 7.6 Convexity and curvature

We have already encountered the notion of convexity in two of its forms: convexity of a function, and convexity of a set. There is a third one, which applies to hypersurfaces and is closely related to curvature. In order to introduce it, consider a hypersurface  $M^n \subset \mathbb{R}^{n+1}$ , a point  $\mathbf{p} \in M$ , and a unit vector  $\mathbf{n} \perp M_{\mathbf{p}}$ . The two *half-spaces at  $\mathbf{p}$*  are the sets

$$\begin{aligned} H_{\mathbf{n}}^+ &= \{\mathbf{q} \in \mathbb{R}^{n+1} \mid \langle \pi_2(\mathbf{n}), \mathbf{q} - \mathbf{p} \rangle \geq 0\}, \text{ and} \\ H_{\mathbf{n}}^- &= \{\mathbf{q} \in \mathbb{R}^{n+1} \mid \langle \pi_2(\mathbf{n}), \mathbf{q} - \mathbf{p} \rangle \leq 0\}. \end{aligned}$$

Here,  $\pi_2 : T\mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$  is the usual projection. Notice that the intersection of the two half-spaces at any  $\mathbf{p} \in M$  is the “affine space”  $\mathbf{p} + \pi_2(M_{\mathbf{p}})$ , which may be construed as a visual representation of the tangent space at  $\mathbf{p}$ . The ambient space itself is the union of these half-spaces.  $M$  is said to be *convex* if it is contained in one of the half-spaces at every point of  $M$ . This notion of convexity differs from the one given in Definition 2.4.5 which applies to arbitrary sets, not just hypersurfaces. It can, however, be shown that a compact, connected hypersurface  $M$  is convex if and only if the region consisting of all points “inside” and on  $M$  is convex in the previous sense of the word.

More generally,  $M$  is said to be *convex at  $\mathbf{p} \in M$*  if there exists a neighborhood  $U$  of  $\mathbf{p}$  in  $\mathbb{R}^{n+1}$  such that  $U \cap M$  is contained in one of the half-spaces at  $\mathbf{p}$ . A convex hypersurface is therefore convex at any of its points, but the converse is not true in general: for example, the hypersurface in  $\mathbb{R}^2$  parametrized by  $\mathbf{h} : (0, 2\pi) \rightarrow \mathbb{R}^2$ , where  $\mathbf{h}(t) = (t \cos t, t \sin t)$ , is convex at every point, yet fails to be globally so, see Figure 7.4.

Convexity is closely related to curvature; to see how, recall a property of the second fundamental tensor that was touched upon previously: let  $\mathbf{n}$  be a unit vector orthogonal to  $M_{\mathbf{p}}$ , and  $S$  the second fundamental tensor with respect to  $\mathbf{n}$ . Given  $\mathbf{x} \in M_{\mathbf{p}}$ ,

$$\langle S\mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{n}, \dot{\mathbf{c}}'(0) \rangle, \quad (7.6.1)$$

where  $\mathbf{c}$  denotes any curve in  $M$  with  $\dot{\mathbf{c}}(0) = \mathbf{x}$ . This is because if  $\mathbf{n}$  is a local normal unit field to  $M$  that equals  $\mathbf{n}$  at  $\mathbf{p}$ , then

$$\langle S\mathbf{x}, \mathbf{x} \rangle = -\langle D_{\mathbf{x}}\mathbf{n}, \dot{\mathbf{c}}(0) \rangle = \langle \mathbf{n}(\mathbf{p}), \dot{\mathbf{c}}'(0) \rangle.$$

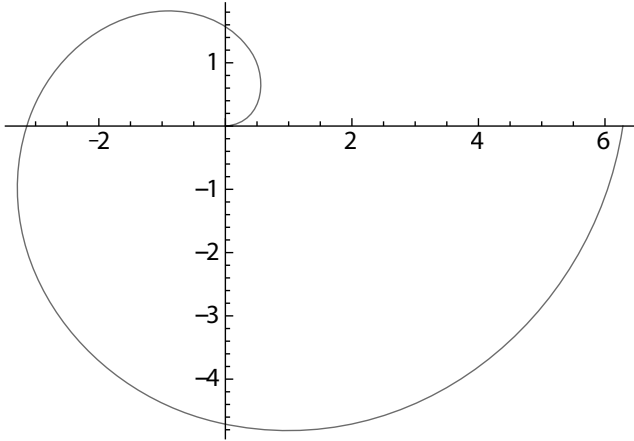


Fig. 7.4: A locally, but not globally, convex curve

**Theorem 7.6.1.** *If a hypersurface is convex at  $\mathbf{p}$ , then the sectional curvatures are non-negative at  $\mathbf{p}$ .*

*Proof.* Let  $\mathbf{n}$  denote a unit vector orthogonal to  $M_{\mathbf{p}}$ ,  $S$  the second fundamental tensor with respect to  $\mathbf{n}$ . According to Exercise 7.17, it suffices to show that  $S$  is semi-definite. By hypothesis, there exists a neighborhood  $U$  of  $\mathbf{p}$  in Euclidean space such that  $M \cap U$  is contained either in  $H_{\mathbf{n}}^+$  or in  $H_{\mathbf{n}}^-$ . Suppose the former holds. We will show that  $\langle S\mathbf{x}, \mathbf{x} \rangle \geq 0$  for any  $\mathbf{x} \in M_{\mathbf{p}}$ . To see this, let  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  be given by

$$f(\mathbf{q}) = \langle \pi_2(\mathbf{n}), \mathbf{q} - \mathbf{p} \rangle.$$

The gradient of  $f$  is the parallel vector field on  $\mathbb{R}^{n+1}$  that equals  $\mathbf{n}$  at  $\mathbf{p}$ . If  $\mathbf{c}$  is any curve in  $M$  with  $\dot{\mathbf{c}}(0) = \mathbf{x} \in M_{\mathbf{p}}$ , then  $f \circ \mathbf{c}$  is nonnegative in a neighborhood of 0, and vanishes at 0. In particular,  $f$  has a minimum at 0, so that  $(f \circ \mathbf{c})''(0) \geq 0$ . Now,  $(f \circ \mathbf{c})' = \langle \nabla f \circ \mathbf{c}, \dot{\mathbf{c}} \rangle$ , and

$$(f \circ \mathbf{c})'' = \langle \nabla f \circ \mathbf{c}, \dot{\mathbf{c}} \rangle' = \langle \nabla f \circ \mathbf{c}, \dot{\mathbf{c}}' \rangle$$

since the gradient of  $f$  is parallel. Thus, by (7.6.1),

$$0 \leq (f \circ \mathbf{c})''(0) = \langle \mathbf{n}, \dot{\mathbf{c}}'(0) \rangle = \langle S\mathbf{x}, \mathbf{x} \rangle,$$

as claimed. The case when  $M \cap U$  is contained in  $H_{\mathbf{n}}^-$  is similar, except that  $f \circ \mathbf{c}$  now has a maximum at 0, so that  $\langle S\mathbf{x}, \mathbf{x} \rangle = (f \circ \mathbf{c})''(0) \leq 0$ . □

The converse is, in general, not true: for example, the right cylinder in  $\mathbb{R}^3$  over the sine curve in  $\mathbb{R}^2$  can be parametrized by  $(s, t) \mapsto (s, \sin s, t)$ . Corollary 7.5.1 implies that the surface is flat, even though it is clearly not convex. A more striking example can be found in Exercise 7.24. The converse does hold, however, under stronger assumptions:

**Theorem 7.6.2.** *Let  $M^n$  be a hypersurface. If the sectional curvatures are positive at some point, then  $M$  is convex at that point.*

*Proof.* Let  $\mathbf{n}$  denote a unit normal field to  $M$  in a neighborhood of a point  $\mathbf{p}$  where the curvature is positive. By Theorem 7.3.1, the second fundamental form  $S$  with respect

to  $\mathbf{n}(\mathbf{p})$  is definite, and we may assume, without loss of generality, that it is positive definite. Consider the function  $f$  defined in some connected neighborhood of  $\mathbf{p}$  by  $f(\mathbf{q}) = \langle \mathbf{q} - \mathbf{p}, \pi_2(\mathbf{n}(\mathbf{p})) \rangle$ . If  $\mathbf{c}$  is any regular curve in  $M$  passing through  $\mathbf{p}$  at  $t = 0$ , then  $(f \circ \mathbf{c})'(0) = \langle \pi_2(\mathbf{n}(\mathbf{p})), \mathbf{c}'(0) \rangle = 0$ , and as in the proof of the previous theorem,  $(f \circ \mathbf{c})'' = \langle S\dot{\mathbf{c}}(0), \dot{\mathbf{c}}(0) \rangle$  which is positive by assumption. Thus,  $f \circ \mathbf{c}$  has a minimum at 0, so that all points of  $M$  close enough to  $\mathbf{p}$  will lie in the half-space  $H_{\mathbf{n}(\mathbf{p})}^+$ . This completes the argument.  $\square$

**Example 7.6.1.** If we visualize the tangent plane of a hypersurface  $M$  in  $\mathbb{R}^{n+1}$  at  $\mathbf{p}$  as the affine subspace  $\mathbf{p} + \pi_2(M_{\mathbf{p}})$  (which is the intuitive interpretation for  $n = 2$ ), then convexity at  $\mathbf{p}$  is equivalent to saying that in a neighborhood of  $\mathbf{p}$ ,  $M$  lies on one side of this affine plane. Thus, if a hypersurface does not have curvature  $\geq 0$  at some point, then it must cross its tangent plane there. This implies in particular something about the shape of minimal hypersurfaces: namely, they must cross their tangent plane at every point where not all curvatures vanish. Indeed, if the hypersurface does not cross its tangent space somewhere, then the second fundamental tensor is semi-definite there. Since the trace is zero, the tensor itself must vanish.

## 7.7 Ruled surfaces

A *ruled surface* in  $\mathbb{R}^3$  is a hypersurface parametrized by a map  $\mathbf{h}$  of the form  $\mathbf{h}(s, t) = \mathbf{c}_1(s) + t\mathbf{c}_2(s)$ , where  $\mathbf{c}_1$  and  $\mathbf{c}_2$  are curves in  $\mathbb{R}^3$ . Notice that the curves  $\mathbf{h}_s$ , where  $\mathbf{h}_s(t) = \mathbf{h}(s, t)$  are straight lines; they are called the *rulings* of the surface. Since

$$D\mathbf{h}(s, t)\mathbf{e}_1 = \mathbf{c}'_1(s) + t\mathbf{c}'_2(s), \quad D\mathbf{h}(s, t)\mathbf{e}_2 = \mathbf{c}_2(s),$$

the two vectors above must be linearly independent for all  $s, t$  if  $\mathbf{h}$  is to be a parametrization. This will be the case (for small  $t$  at least) if for example  $\mathbf{c}'_1(s)$  and  $\mathbf{c}_2(s)$  are linearly independent. An interesting feature of these surfaces is that they are always nonpositively curved:

**Proposition 7.7.1.** *The sectional curvature of a ruled surface is given by*

$$(K \circ \mathbf{h})(s, t) = -\frac{\langle \mathbf{c}'_1(s), \mathbf{c}_2(s) \times \mathbf{c}'_2(s) \rangle^2}{|(\mathbf{c}'_1(s) + t\mathbf{c}'_2(s)) \times \mathbf{c}_2(s)|^4}.$$

*Proof.* Notice that if  $\mathbf{h}_t$  denotes the map  $s \mapsto \mathbf{h}(s, t)$ , then the formula we aim to establish becomes

$$K \circ \mathbf{h}_t = -\frac{\langle \mathbf{c}'_1, \mathbf{c}_2 \times \mathbf{c}'_2 \rangle^2}{|(\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2|^4}.$$

The unit normal field determined by the parametrization is

$$\mathbf{n} \circ \mathbf{h}_t = \frac{D\mathbf{h}\mathbf{e}_1 \times D\mathbf{h}\mathbf{e}_2}{|D\mathbf{h}\mathbf{e}_1 \times D\mathbf{h}\mathbf{e}_2|} = \frac{(\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2}{|(\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2|}.$$



Now,  $D^2 \mathbf{h} \mathbf{e}_2 \mathbf{e}_2 = \mathbf{0}$ , so the matrix  $\langle D^2 \mathbf{h} \mathbf{e}_i \mathbf{e}_j, \mathbf{n} \circ \mathbf{h} \rangle$  has determinant

$$\det \langle D^2 \mathbf{h} \mathbf{e}_i \mathbf{e}_j, \mathbf{n} \circ \mathbf{h} \rangle = -\langle D^2 \mathbf{h} \mathbf{e}_1 \mathbf{e}_2, \mathbf{n} \circ \mathbf{h} \rangle^2 = -\frac{\langle \mathbf{c}'_2, (\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2 \rangle^2}{|(\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2|^2},$$

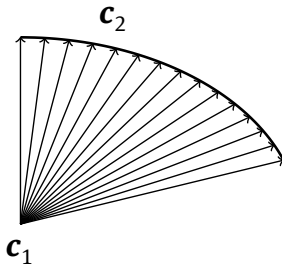
and by Corollary 7.5.1,  $K \circ \mathbf{h}_t = -\langle \mathbf{c}'_2, (\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2 \rangle^2 / |(\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2|^4$ . The definition of the cross-product implies that  $\langle \mathbf{x} \times \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{y} \times \mathbf{z}, \mathbf{x} \rangle$  for  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^3$ . Thus,

$$\langle (\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2, \mathbf{c}'_2 \rangle = \langle \mathbf{c}_2 \times \mathbf{c}'_2, \mathbf{c}'_1 + t\mathbf{c}'_2 \rangle = \langle \mathbf{c}_2 \times \mathbf{c}'_2, \mathbf{c}'_1 \rangle$$

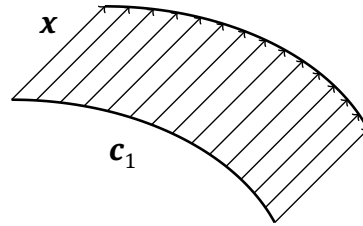
since  $\mathbf{c}'_2 \perp \mathbf{c}_2 \times \mathbf{c}'_2$ , and the result follows.  $\square$

There are two special cases worth mentioning:

- (1)  $\mathbf{c}'_1 \equiv \mathbf{0}$ .  $t$  must then be different from zero, and  $\mathbf{c}_1$  is a constant curve  $\mathbf{p}$ .  $M$  is a generalized cone over  $\mathbf{c}_2$  with vertex  $\mathbf{p}$  deleted. The cone is flat.
- (2)  $\mathbf{c}'_2 \equiv \mathbf{0}$ . Now  $\mathbf{c}_2$  is a constant  $\mathbf{x}$  and  $M$  is a generalized cylinder over  $\mathbf{c}_1$  with axis parallel to  $\mathbf{x}$ . Again,  $M$  is flat.



Ruled surface with  $\mathbf{c}_1$  constant



Ruled surface with  $\mathbf{c}_2 =$  constant  $\mathbf{x}$

These special cases having been dealt with, we will assume that  $\mathbf{c}_1$  and  $\mathbf{c}_2$  are regular and parametrized by arc length. We may also normalize  $\mathbf{c}_2$  so that it lies on the unit sphere. In this case,  $\mathbf{c}_2$  is called the *directrix*. The formula for the curvature simplifies substantially if we assume that  $\langle \mathbf{c}'_1, \mathbf{c}'_2 \rangle \equiv 0$ . Although it is not immediately obvious, it turns out this may always be done; i.e.,  $\mathbf{c}_2$  may be replaced, if necessary, by another curve  $\gamma_2$  satisfying  $\langle \mathbf{c}'_1, \gamma'_2 \rangle \equiv 0$ . Under these assumptions, we have that

$$\begin{aligned} |(\mathbf{c}'_1 + t\mathbf{c}'_2) \times \mathbf{c}_2|^2 &= |\mathbf{c}'_1 + t\mathbf{c}'_2|^2 |\mathbf{c}_2|^2 - \langle \mathbf{c}'_1 + t\mathbf{c}'_2, \mathbf{c}_2 \rangle^2 \\ &= |\mathbf{c}'_1|^2 + t^2 |\mathbf{c}'_2|^2 - \langle \mathbf{c}'_1, \mathbf{c}_2 \rangle^2 = |\mathbf{c}'_1|^2 |\mathbf{c}_2|^2 + t^2 - \langle \mathbf{c}'_1, \mathbf{c}_2 \rangle^2 \\ &= t^2 + |\mathbf{c}'_1 \times \mathbf{c}_2|^2. \end{aligned}$$

Now,  $\mathbf{c}'_2$  is orthogonal to  $\mathbf{c}'_1, \mathbf{c}_2$ , and the same is true for  $\mathbf{c}'_1 \times \mathbf{c}_2$ . Assuming linear independence of  $\mathbf{c}'_1, \mathbf{c}_2$ , it follows that  $\mathbf{c}'_1 \times \mathbf{c}_2$  and  $\mathbf{c}'_2$  are orthogonal to a common plane, and therefore linearly dependent. Thus,  $\mathbf{c}'_1 \times \mathbf{c}_2 = \langle \mathbf{c}'_1 \times \mathbf{c}_2, \mathbf{c}'_2 \rangle \mathbf{c}'_2 = \langle \mathbf{c}_2 \times \mathbf{c}'_2, \mathbf{c}'_1 \rangle \mathbf{c}'_2$ ; i.e.,

$|\mathbf{c}'_1 \times \mathbf{c}_2|^2 = \langle \mathbf{c}'_1, \mathbf{c}_2 \times \mathbf{c}'_2 \rangle^2$ , and the formula for the curvature becomes

$$K \circ \mathbf{h} = -\frac{|\mathbf{c}'_1 \times \mathbf{c}_2|^2}{(t^2 + |\mathbf{c}'_1 \times \mathbf{c}_2|^2)^2}. \quad (7.7.1)$$

Notice that since  $\mathbf{c}'_1$  and  $\mathbf{c}_2$  are assumed to be linearly independent, the curvature is strictly negative.

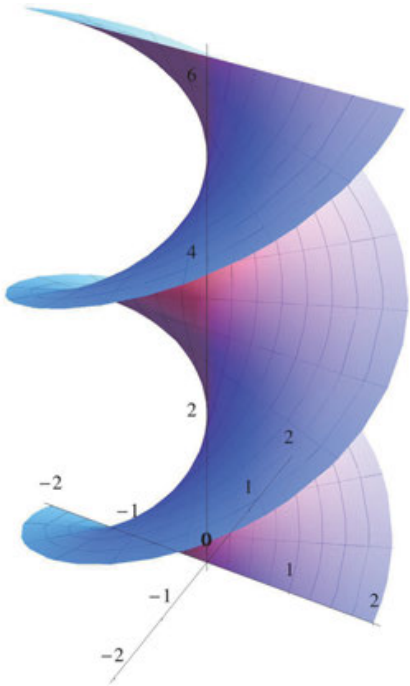


Fig. 7.5: A helicoid

One example of a ruled surface is the Möbius strip from Section 4.1, with  $\mathbf{h}(s, t) = (\cos s, \sin s, 0) + t(\cos(s/2), \sin(s/2), \sin(s/2))$ . Another is the hyperboloid  $(x/a)^2 + (y/b)^2 - (z/c)^2 = 1$  from Section 4.3, which may be parametrized by  $\mathbf{h}(s, t) = (a \cos s, b \sin s, 0) + t(-a \sin s, b \cos s, c)$ . One that we have not encountered before is the *helicoid*: For any  $s \in \mathbb{R}$ , consider the Euclidean motion  $\mathbf{k}_s$  of  $\mathbb{R}^3$ , where

$$\mathbf{k}_s(x, y, z) = \begin{bmatrix} \cos s & -\sin s & 0 \\ \sin s & \cos s & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ s \end{bmatrix}.$$

$\mathbf{k}_s$  is called a *glide rotation* as it rotates a point by angle  $s$  around the  $z$ -axis and translates it by  $s$  along that same axis. The so-called orbit  $\{\mathbf{k}_s(\mathbf{p}) \mid s \in \mathbb{R}\}$  of a point  $\mathbf{p}$  is a helix. The helicoid itself is just the orbit of the whole  $x$ -axis; i.e.,  $\mathbf{h}(s, t) = \mathbf{k}_s(t, 0, 0)$ . According to (7.7.1), the curvature of the helicoid is

$$K \circ \mathbf{h}(s, t) = -\frac{1}{(1 + t^2)^2},$$

which only depends on the distance  $|t|$  from the point  $\mathbf{h}(s, t)$  to the  $z$ -axis. The helicoid is actually related to the catenoid introduced in Exercise 3.28. Both are minimal surfaces, and either one can be smoothly and isometrically deformed into the other via  $\mathbf{F}$ , where

$$\begin{aligned} \mathbf{F}(\theta, s, t) = & (\cos \theta \sinh t \sin s + \sin \theta \cosh t \cos s, -\cos \theta \sinh t \cos s \\ & + \sin \theta \cosh t \sin s, s \cos \theta + t \sin \theta). \end{aligned}$$

When  $\theta = 0$ , one obtains the parametrization of the helicoid above rotated by  $-\pi/2$ :

$$\mathbf{F}(0, s, t)^T = \begin{bmatrix} \sinh t \sin s \\ \sinh t(-\cos s) \\ s \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{k}_s(\sinh t, 0, 0)^T.$$

Similarly,  $\theta = \pi/2$  yields the catenoid. For each fixed  $\theta_0$ , the map  $(s, t) \mapsto \mathbf{F}(\theta_0, s, t)$  parametrizes a minimal surface. Notice that

$$\mathbf{F}(\theta, s, t) = \cos \theta \mathbf{F}(0, s, t) + \sin \theta \mathbf{F}\left(\frac{\pi}{2}, s, t\right),$$

see also Exercise 7.25.

## 7.8 Surfaces of revolution

Surfaces of revolution, where the graph of a function of one variable is rotated about an axis, were introduced in Examples 3.1.1. It is convenient to allow for more general curves than just graphs. For clarity of notation in the sometimes rather complicated formulas that follow, we will use subscripts rather than superscripts to denote the component functions of the curve.

**Definition 7.8.1.** Let  $\gamma = (\gamma_1, 0, \gamma_2) : (a, b) \rightarrow \mathbb{R}^3$  denote a curve whose image lies in the  $x$ - $z$  plane. The *surface of revolution*  $M$  with *profile curve*  $\gamma$  is the surface parametrized by  $\mathbf{h} : (a, b) \times [0, 2\pi) \rightarrow \mathbb{R}^3$ , where

$$\mathbf{h}(u, v) = (\gamma_1(u) \cos v, \gamma_1(u) \sin v, \gamma_2(u)).$$

Thus, the surface is obtained by rotating the image of  $\gamma$  about the  $z$  axis. The curves  $u \mapsto \mathbf{h}(u, v)$  and  $v \mapsto \mathbf{h}(u, v)$  are called *meridians* and *parallels* respectively. The meridian  $u \mapsto \mathbf{h}(u, v_0)$  is the profile curve rotated by angle  $v_0$ , and the parallel  $v \mapsto \mathbf{h}(u_0, v)$  is the circle obtained by rotating the point  $\gamma(u_0)$  about the  $z$ -axis.

The coordinate fields  $\frac{\partial}{\partial x^i} = \mathbf{h}_* \circ \mathbf{D}_i \circ \mathbf{h}^{-1}$  are determined by the vector fields

$$\begin{aligned} \mathbf{h}_* \mathbf{D}_1(u, v) = & \gamma_1'(u) \cos v (\mathbf{D}_1 \circ \mathbf{h})(u, v) + \gamma_1'(u) \sin v (\mathbf{D}_2 \circ \mathbf{h})(u, v) \\ & + \gamma_2'(u) (\mathbf{D}_3 \circ \mathbf{h})(u, v) \end{aligned} \quad (7.8.1)$$

$$\mathbf{h}_* \mathbf{D}_2(u, v) = -\gamma_1(u) \sin v (\mathbf{D}_1 \circ \mathbf{h})(u, v) + \gamma_1(u) \cos v (\mathbf{D}_2 \circ \mathbf{h})(u, v)$$

along  $\mathbf{h}$ . The latter are just the tangent fields to the meridians and parallels. Notice that the second one is parallel to the  $x$ - $y$  plane and that they are mutually orthogonal. Their cross-product  $\mathbf{h}_* \mathbf{D}_1 \times \mathbf{h}_* \mathbf{D}_2$ , yields, after normalizing, a unit normal field  $\mathbf{N}$  along  $\mathbf{h}$  given by

$$\mathbf{N}(u, v) = \frac{1}{(\gamma_1'^2 + \gamma_2'^2)^{1/2}(u)} \left( \gamma_2'(u) \cos v \mathbf{D}_1 + \gamma_2'(u) \sin v \mathbf{D}_2 - \gamma_1'(u) \mathbf{D}_3 \right) \circ \mathbf{h}(u, v). \tag{7.8.2}$$

**Theorem 7.8.1.** *The sectional curvature of the surface of revolution parametrized by  $\mathbf{h}$  is*

$$K \circ \mathbf{h} = \frac{(\gamma_1' \gamma_2'' - \gamma_2' \gamma_1'') \gamma_2'}{\gamma_1 (\gamma_1'^2 + \gamma_2'^2)^2}.$$

*In particular, if  $\gamma$  is parametrized by arc-length, then  $K \circ \mathbf{h} = -\gamma_1'' / \gamma_1$ .*

*Proof.* Using (7.8.2), the covariant derivatives of the unit normal field along meridians and parallels are given by

$$D_{\mathbf{D}_1} \mathbf{N} = \frac{\gamma_1' \gamma_2'' - \gamma_2' \gamma_1''}{(\gamma_1'^2 + \gamma_2'^2)^{3/2}} \mathbf{h}_* \mathbf{D}_1, \quad D_{\mathbf{D}_2} \mathbf{N} = \frac{\gamma_2'}{\gamma_1 (\gamma_1'^2 + \gamma_2'^2)^{1/2}} \mathbf{h}_* \mathbf{D}_2.$$

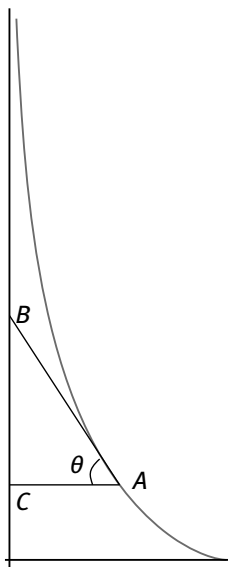
Thus, the coordinate vector fields are eigenvector fields at each point of  $M$  of the second fundamental tensor, with corresponding principal curvatures

$$\lambda_1 = \frac{\gamma_1' \gamma_2'' - \gamma_2' \gamma_1''}{(\gamma_1'^2 + \gamma_2'^2)^{3/2}}, \quad \lambda_2 = \frac{\gamma_2'}{\gamma_1 (\gamma_1'^2 + \gamma_2'^2)^{1/2}}. \tag{7.8.3}$$

Since the sectional (and Gaussian) curvature is the product of the principal curvatures, the result follows. When  $\gamma$  is parametrized by arc length,  $\sum_i \gamma_i' \gamma_i'' = 0$ , so that

$$K \circ \mathbf{h} = \frac{\gamma_1' \gamma_2' \gamma_2'' - \gamma_1'' \gamma_2'^2}{\gamma_1} = \frac{\gamma_1' (-\gamma_1' \gamma_1'') - \gamma_1'' (1 - \gamma_1'^2)}{\gamma_1} = -\frac{\gamma_1''}{\gamma_1}. \tag{7.8.4}$$

□



**Fig. 7.6:** A tractrix

**Example 7.8.1.** Consider an object located at  $(a, 0)$ ,  $a > 0$ , in the plane, attached to a rope of length  $a$  whose other end is held by someone standing initially at the origin  $(0, 0)$ . A *tractrix* (from the Latin verb “*trahere*” – to draw) is the path traced by the object as it is being pulled by the person walking along the  $y$ -axis.

Denoting by  $\gamma = (\gamma_1, \gamma_2)$  the parametrization by arc-length of the tractrix with  $\gamma(0) = (a, 0)$ , we see from the figure that

$$a \equiv AB = \frac{AC}{\cos \theta} = -\frac{AC}{\gamma_1'} = -\frac{\gamma_1}{\gamma_1'}$$

with the minus sign being due to the fact that

$$\gamma_1' = \langle \gamma', \mathbf{e}_1 \rangle = |\gamma'| \cos(\pi - \theta) = -\cos \theta.$$

Thus,  $\gamma_1' = -(1/a)\gamma_1$ , and  $\gamma_1(t) = ae^{-t/a}$ . Furthermore,

$$1 = |\gamma'|^2(t) = (ae^{-t/a})^2 + \gamma_2'^2(t),$$

so that  $\gamma_2'(t) = \pm\sqrt{1 - (ae^{-t/a})^2}$ . If the person walks along the positive portion of the  $y$ -axis, the tractrix is then given by

$$\gamma(t) = \left( ae^{-t/a}, \int_0^t \sqrt{1 - (ae^{-s/a})^2} ds \right), \quad t \geq 0.$$

The surface of revolution generated by the tractrix is called a *pseudosphere*. The name was coined in 1868 by the Italian mathematician Eugenio Beltrami who proposed it as a model of non-Euclidean geometry. By (7.8.4), this pseudosphere has constant negative curvature  $K = -1/a^2$ .

One appealing feature of surfaces of revolution is that, in contrast to general surfaces, geodesics are readily described. Notice first that any meridian is the image of a geodesic. This follows from Proposition 3.11.2, since it is (part of) the fixed point set of reflection in the plane that contains the meridian and the  $z$ -axis. For parallels, the situation is different: we claim that  $\mathbf{c}$ , where  $\mathbf{c}(t) = \mathbf{h}(u_0, t)$ , is a geodesic if and only if  $u_0$  is a critical point of  $\gamma_1$ . To see this, notice that  $\dot{\mathbf{c}}' = -\alpha \cos \mathbf{D}_1 \circ \mathbf{c} - \alpha \sin \mathbf{D}_2 \circ \mathbf{c}$ , where  $\alpha := \gamma_1(u_0)$ . On the other hand,  $\mathbf{c}$  is a geodesic if and only if  $\mathbf{N} \circ \mathbf{c}$  and  $\dot{\mathbf{c}}'$  are linearly dependent. Comparing with (7.8.2) now yields the claim.

The other geodesics can be described by means of the following

**Theorem 7.8.2 (Clairaut).** Let  $\mathbf{c} : I \rightarrow M$  denote a normal geodesic in  $M$ , and  $\boldsymbol{\rho} = \mathbf{h}^{-1} \circ \mathbf{c} : I \rightarrow \mathbb{R}^2$ . If  $\alpha(t) = \angle(\dot{\mathbf{c}}(t), \mathbf{h}_*(\mathbf{D}_2 \circ \boldsymbol{\rho})(t))$  is the angle between  $\dot{\mathbf{c}}(t)$  and the parallel through  $\mathbf{c}(t)$ , and if  $r = ((u^1)^2 + (u^2)^2)^{1/2}$  denotes the distance from a point to the  $z$ -axis, then  $(r \circ \mathbf{c}) \cos \alpha$  is constant.

Conversely, if  $\mathbf{c}$  is a curve parametrized by arc length for which  $(r \circ \mathbf{c}) \cos \alpha$  is constant, and  $\mathbf{c}$  is not a parallel, then  $\mathbf{c}$  is a geodesic.

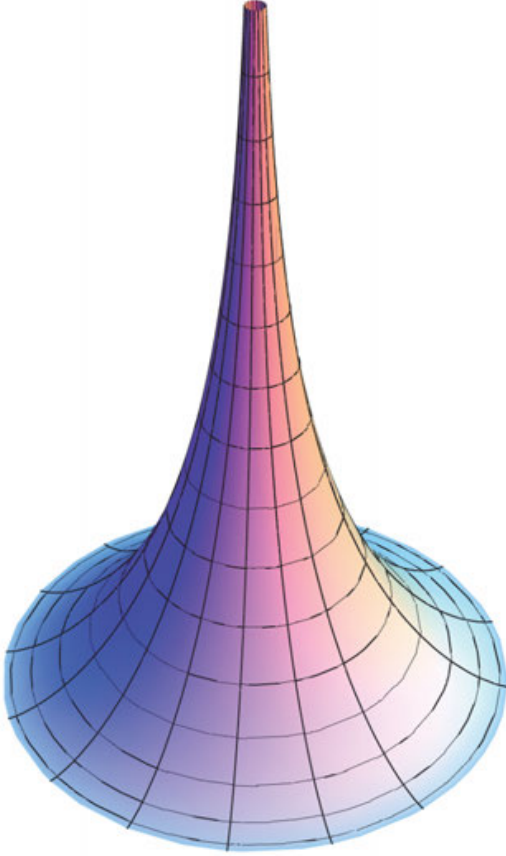


Fig. 7.7: A pseudosphere

*Proof.* We will, as we did with  $\mathbf{c}$ , denote with a subscript the component functions  $\rho_i = u^i \circ \boldsymbol{\rho}$ ,  $i = 1, 2$ , of  $\boldsymbol{\rho}$ . We will also abbreviate  $\mathbf{h}_*(\mathbf{D}_i \circ \boldsymbol{\rho})$  by  $\mathbf{h}_* \mathbf{D}_i \circ \boldsymbol{\rho}$ . Then

$$\begin{aligned} \dot{\mathbf{c}} &= \rho_1' \mathbf{h}_* \mathbf{D}_1 \circ \boldsymbol{\rho} + \rho_2' \mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho} \\ &= \left( (\gamma_1 \circ \rho_1)' \cos \rho_2 - (\gamma_1 \circ \rho_1) (\sin \rho_2) \rho_2' \right) \mathbf{D}_1 \circ \mathbf{c} \\ &\quad + \left( (\gamma_1 \circ \rho_1)' \sin \rho_2 + (\gamma_1 \circ \rho_1) (\cos \rho_2) \rho_2' \right) \mathbf{D}_2 \circ \mathbf{c} + \rho_1' (\gamma_2' \circ \rho_1) \mathbf{D}_3 \circ \mathbf{c} \\ &= \left( (\gamma_1 \circ \rho_1) \cos \rho_2 \right)' \mathbf{D}_1 \circ \mathbf{c} + \left( (\gamma_1 \circ \rho_1) \sin \rho_2 \right)' \mathbf{D}_2 \circ \mathbf{c} + (\gamma_2 \circ \rho_1)' \mathbf{D}_3 \circ \mathbf{c}. \end{aligned}$$

Now, (7.8.1) implies

$$(\mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho})' = - \left( (\gamma_1 \circ \rho_1) \sin \rho_2 \right)' \mathbf{D}_1 \circ \mathbf{c} + \left( (\gamma_1 \circ \rho_1) \cos \rho_2 \right)' \mathbf{D}_2 \circ \mathbf{c},$$

so that

$$\langle \dot{\mathbf{c}}, \mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho} \rangle' = \langle \dot{\mathbf{c}}, (\mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho})' \rangle = 0,$$

and  $\langle \dot{\mathbf{c}}, \mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho} \rangle$  is constant. But the latter also equals

$$|\dot{\mathbf{c}}| |\mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho}| \cos \alpha = |\gamma_1 \circ \rho_1| \cos \alpha = (r \circ \mathbf{c}) \cos \alpha.$$

Conversely, if  $(r \circ \mathbf{c}) \cos \alpha$  is constant, then by the above calculation

$$\langle \dot{\mathbf{c}}', \mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho} \rangle = \langle \dot{\mathbf{c}}, \mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho} \rangle' - \langle \dot{\mathbf{c}}, (\mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho})' \rangle = - \langle \dot{\mathbf{c}}, (\mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho})' \rangle = 0.$$

Now,  $\dot{\mathbf{c}}'$  is orthogonal to  $\dot{\mathbf{c}}$  because  $\mathbf{c}$  is parametrized by arc length. It follows that at any  $t_0$  where  $\dot{\mathbf{c}}(t_0)$  and  $\mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho}(t_0)$  are linearly independent,  $\dot{\mathbf{c}}'(t_0)$  is orthogonal to the surface; i.e.,  $(\nabla_{\mathbf{D}} \dot{\mathbf{c}})(t_0) = 0$ . On the other hand, if  $\dot{\mathbf{c}}(t_0)$  and  $\mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho}(t_0)$  are linearly dependent, then there exists a sequence  $t_n \rightarrow t_0$  such that  $\dot{\mathbf{c}}(t_n)$  and  $\mathbf{h}_* \mathbf{D}_2 \circ \boldsymbol{\rho}(t_n)$  are linearly independent, for otherwise  $\mathbf{c}$  would be a meridian. Then  $(\nabla_{\mathbf{D}} \dot{\mathbf{c}})(t_n) = 0$ , and by continuity, so is  $(\nabla_{\mathbf{D}} \dot{\mathbf{c}})(t_0)$ .  $\square$

Clairaut's theorem yields a fairly complete qualitative description of the geodesics on  $M$ . Let us illustrate this in the case when  $M$  is a paraboloid of revolution; specifically, the image of the curve  $\mathbf{c}$  is the set  $\{(x, 0, x^2) \mid x \geq 0\}$ . All the meridians are the (images of the) geodesics emanating from the vertex  $\mathbf{0}$ . Since none of these intersect, any such geodesic  $\mathbf{c} : [0, \infty) \rightarrow M$  is a *ray*; i.e.,  $\mathbf{c}$  is the shortest curve from the vertex to  $\mathbf{c}(t)$  for any  $t > 0$ . Assume then that  $\mathbf{c} : (-\infty, \infty) \rightarrow M$  is a geodesic that is not a meridian. Consider the lowest point on the image of  $\mathbf{c}$ . The existence of such a point is discussed in Exercise 7.10. Since the distance to the  $z$ -axis increases with height, at this lowest point,  $r \circ \mathbf{c}$  is minimal and the angle  $\alpha$  is zero. After reparametrizing  $\mathbf{c}$  if necessary, we may assume this occurs at  $t = 0$ , so that

$$(r \circ \mathbf{c})(t) \cos \alpha(t) = (r \circ \mathbf{c})(0) =: k.$$

Since no parallel is a geodesic,  $(r \circ \mathbf{c})(t) > (r \circ \mathbf{c})(0)$  for  $t > 0$ , and  $\mathbf{c}$  rises; i.e.,  $u^3 \circ \mathbf{c}$  increases. We claim  $(u^3 \circ \mathbf{c})(t) \rightarrow \infty$  as  $t \rightarrow \infty$ . To see this, let  $\mathbf{c}_0(t)$  denote the point on the (image of the) profile curve  $\boldsymbol{\gamma}$  that lies on the same parallel as  $\mathbf{c}(t)$ ; i.e.,  $\mathbf{c}_0(t)$  is obtained by rotating  $\mathbf{c}(t)$  about the  $z$ -axis by an angle  $-\gamma_2(t)$ . Thus,

$$\mathbf{c}_0 = \mathbf{h}(\rho_1, 0) = (\gamma_1 \circ \rho_1, 0, \gamma_2 \circ \rho_1) = \boldsymbol{\gamma} \circ \rho_1,$$

and  $\dot{\mathbf{c}}_0 = \rho_1' \dot{\boldsymbol{\gamma}} \circ \rho_1$ . The claim follows once we show that  $\rho_1$  is unbounded. Assuming  $\mathbf{c}$  is parametrized by arc length, the speed of  $\mathbf{c}_0$  is

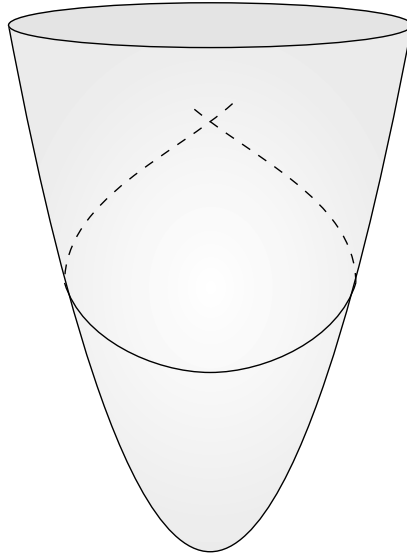
$$\begin{aligned} |\dot{\mathbf{c}}_0| &= \rho_1' = \langle \dot{\mathbf{c}}, \mathbf{h}_* \mathbf{D}_1 \circ \boldsymbol{\rho} \rangle = |\dot{\mathbf{c}}| |\mathbf{h}_* \mathbf{D}_1 \circ \boldsymbol{\rho}| \sin \alpha = \sin \alpha \\ &= \left( 1 - \frac{k^2}{(\gamma_1 \circ \rho_1)^2} \right)^{1/2}, \end{aligned} \quad (7.8.5)$$

since

$$\sin \alpha = (1 - \cos^2 \alpha)^{1/2} = \left( 1 - \frac{k^2}{(r \circ \mathbf{c})^2} \right)^{1/2} = \left( 1 - \frac{k^2}{(r \circ \mathbf{c}_0)^2} \right)^{1/2}$$

and  $r \circ \mathbf{c}_0 = \gamma_1 \circ \rho_1$ . Differentiating (7.8.5) yields

$$\rho_1'' = \left( 1 - \frac{k^2}{(\gamma_1 \circ \rho_1)^2} \right)^{-1/2} k^2 (\gamma_1 \circ \rho_1)^{-3} (\gamma_1' \circ \rho_1) \rho_1'.$$



Now,  $\gamma_1$  and  $\gamma_1'$  are both positive in our case, so  $\rho_1''$  and  $\rho_1'$  have the same sign. Furthermore,  $\rho_1'(0) = 0$  and  $\rho_1'(t) > 0$  for small  $t > 0$ . By the mean value theorem,  $\rho_1'(t) > 0$  for all  $t > 0$ : otherwise, letting  $t_0$  denote the infimum of those  $t > 0$  such that  $\rho_1(t) \leq 0$ , we have that  $\rho_1'(t_0) \leq 0$ , and  $\rho_1'(t) > 0$  for all  $t < t_0$ . This is impossible, since the mean value theorem guarantees the existence of some  $t \in (0, t_0)$  satisfying

$$\rho_1''(t) = \frac{\rho_1'(t_0) - \rho_1'(0)}{t_0} \leq 0,$$

and therefore,  $\rho_1'(t)$  is also nonpositive. Thus,  $\rho_1', \rho_1'' > 0$  always; i.e.,  $\rho_1$  is a convex function, hence unbounded. This establishes the claim.

Notice also that reflection  $R$  in the plane containing the meridian through  $\mathbf{c}(0)$  is an isometry that must leave the image of  $\mathbf{c}$  invariant:  $R \circ \mathbf{c}|_{[0, \infty)}$  is a geodesic with initial tangent vector  $-\dot{\mathbf{c}}(0)$ , so that by uniqueness, it must equal  $t \mapsto \mathbf{c}(-t)$ .

Further examples are explored in the exercises.

## 7.9 Exercises

**7.1.** Suppose  $\mathbf{h}$  and  $\mathbf{k}$  are two overlapping parametrizations – in the sense that their images have nonempty intersection – of  $M^n \subset \mathbb{R}^{n+1}$  with inverses  $\mathbf{x}$  and  $\mathbf{y}$  respectively, and consider a point  $\mathbf{p} = \mathbf{h}(\mathbf{a}) \in M$  in the common image. Show that for any  $\mathbf{u} \in \mathbb{R}_p^{n+1}$ ,

$$\begin{aligned} \det\left(\frac{\partial}{\partial x^1}(\mathbf{p}), \dots, \frac{\partial}{\partial x^n}(\mathbf{p}), \mathbf{u}\right) \\ = \det[D(\mathbf{k}^{-1} \circ \mathbf{h})](\mathbf{a}) \det\left(\frac{\partial}{\partial y^1}(\mathbf{p}), \dots, \frac{\partial}{\partial y^n}(\mathbf{p}), \mathbf{u}\right). \end{aligned}$$

Conclude that  $M$  is orientable if and only if it admits an atlas with the property that  $\det(\mathbf{k}^{-1} \circ \mathbf{h}) > 0$  for any two overlapping parametrizations  $\mathbf{h}, \mathbf{k}$  in the atlas.

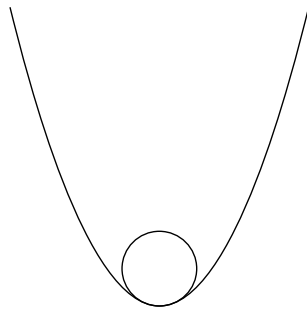


**7.2.** Show that an  $n \times n$  invertible matrix is symmetric if and only if its inverse is symmetric.

**7.3.** With the terminology from Example 7.4.1, determine the curvature of

- (a) a straight line in  $\mathbb{R}^2$ ;  
 (b) a circle of radius  $r > 0$  in  $\mathbb{R}^2$ .

**7.4.** This long exercise is meant to clarify the remark made in Example 7.4.1, which asserts that the circle of curvature of a curve  $\mathbf{c}$  at a point  $\mathbf{c}(t_0)$  is the circle that best approximates the curve at that point. The underlying idea is that there is one, and only one, circle that passes through any three non-colinear points. Letting those points approach  $\mathbf{c}(t_0)$  results in this limit circle.

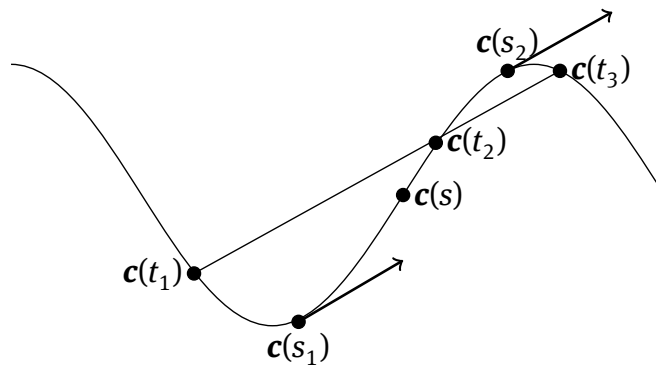


The circle of curvature at the vertex of a parabola

So let  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^2$  denote a regular curve parametrized by arc length, and  $t_0 \in (a, b)$  be a point where  $\mathbf{c}''(t_0) \neq \mathbf{0}$ .

- (a) Show that any three points on the curve that are sufficiently close to  $\mathbf{c}(t_0)$  cannot be colinear; i.e., they cannot lie on a common line.

*One outline of a proof:* Suppose, to the contrary, that  $\mathbf{c}(t_i)$  are colinear for  $t_1 < t_2 < t_3$ . The mean value theorem then implies that there exist  $s_1 \in (t_1, t_2)$  and  $s_2 \in (t_2, t_3)$  satisfying  $\mathbf{c}'(s_1) = \mathbf{c}'(s_2)$ .



Now,  $\mathbf{c}'$  has its image in  $S^1$  and cannot be onto  $S^1$  if  $s_1$  and  $s_2$  are close enough; i.e., if the  $t_i$  are close enough to  $t_0$ ,  $1 \leq i \leq 3$ . If  $s \in (s_1, s_2)$  is such that  $\mathbf{c}'(s)$  is furthest away from  $\mathbf{c}'(s_i)$ , show that  $\mathbf{c}''(s) = \mathbf{0}$ . This is impossible if the  $t_i$  are sufficiently close to  $t_0$ .

- (b) Let  $t_i$ ,  $1 \leq i \leq 3$ , be as in part (a) so that the  $\mathbf{c}(t_i)$  are not colinear, and consider the center  $\mathbf{p}_{123}$  of the unique circle that passes through these three points. Prove that there exist  $s_1, s_2 \in (t_1, t_3)$  such that

$$\langle \mathbf{c}'(s_1), \mathbf{c}(s_1) - \mathbf{p}_{123} \rangle = 0, \quad (7.9.1)$$

$$\langle \mathbf{c}''(s_2), \mathbf{c}(s_2) - \mathbf{p}_{123} \rangle = -|\mathbf{c}'(s_2)|^2 = -1. \quad (7.9.2)$$

*Hint:* the function  $t \mapsto |\mathbf{c}(t) - \mathbf{p}_{123}|^2$  has the same value at each  $t_i$ , so its derivative must vanish at some  $r_1 \in (t_1, t_2)$  and  $r_2 \in (t_2, t_3)$ . By the same reasoning, its second derivative vanishes at some point in  $(r_1, r_2)$ .

- (c) Show that there exists a unique point  $\mathbf{p} \in \mathbb{R}^2$  satisfying

$$\langle \mathbf{c}'(t_0), \mathbf{c}(t_0) - \mathbf{p} \rangle = 0,$$

$$\langle \mathbf{c}''(t_0), \mathbf{c}(t_0) - \mathbf{p} \rangle = -1.$$

Compare these equations with (7.9.1) to conclude that as  $t_i \rightarrow t_0$ ,  $1 \leq i \leq 3$ , the unique circle passing through  $\mathbf{c}(t_1)$ ,  $\mathbf{c}(t_2)$ , and  $\mathbf{c}(t_3)$  approaches a circle passing through  $\mathbf{c}(t_0)$ . Furthermore, the latter circle has curvature equal to the curvature of  $\mathbf{c}$  at  $t_0$ , and its center  $\mathbf{p}$  lies on the line through  $\mathbf{c}(t_0)$  perpendicular to the tangent line to  $\mathbf{c}$  at  $\mathbf{c}(t_0)$ .

- 7.5.** Let  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^2$  be a regular curve parametrized by arc length. Prove that the curvature  $\kappa$  of  $\mathbf{c}$  satisfies

$$\kappa = \det \begin{bmatrix} \mathbf{c}' & \mathbf{c}'' \end{bmatrix}.$$

*Hint:* The right side of the above equation equals, up to sign, the area of the parallelogram spanned by  $\mathbf{c}'$  and  $\mathbf{c}''$ .

- 7.6.** Even for simple curves, parametrization by arc length can be difficult to obtain (try using the formula from Exercise 7.5 to compute the curvature of a parabola). This problem generalizes that formula to regular curves that are not parametrized by arc length. If  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^2$  is such a curve, let  $s(t) = \int_a^t |\mathbf{c}'|$ , so that  $\tilde{\mathbf{c}} : \mathbf{c} \circ s^{-1}$  is the reparametrization of  $\mathbf{c}$  by arc length. The curvature  $\kappa(t)$  of  $\mathbf{c}$  at  $t \in (a, b)$  is defined to be the curvature of  $\tilde{\mathbf{c}}$  at  $s(t)$ . Show that

$$\kappa = \frac{1}{|\mathbf{c}'|^3} \det \begin{bmatrix} \mathbf{c}' & \mathbf{c}'' \end{bmatrix}.$$

Compute the curvature of the parabola  $y = x^2$ .

- 7.7.** The concepts discussed in the previous exercises generalize in part to curves in  $\mathbb{R}^3$ . Even though these are no longer hypersurfaces, they are important enough in classical differential geometry to warrant at least cursory mention.

So let  $\mathbf{c} : [a, b] \rightarrow \mathbb{R}^3$  be a regular curve parametrized by arc length. Define  $\mathbf{T} := \mathbf{c}'$ . Notice that  $\mathbf{T}$  is actually a map into  $S^2$ , even though we often identify it with the velocity vector field  $\dot{\mathbf{c}}$  of the curve. The *curvature* of  $\mathbf{c}$  is the function

$$\kappa = |\mathbf{T}'|.$$

Thus, unlike plane curves, the curvature is always nonnegative – essentially because we can no longer appeal to the orientation of  $\mathbb{R}^2$ . When  $\kappa(t) \neq 0$ , define the *principal normal*  $\mathbf{N}$  at  $t$  by

$$\mathbf{N}(t) = \frac{1}{|\mathbf{T}'(t)|} \mathbf{T}'(t),$$

so that

$$\mathbf{T}' = \kappa \mathbf{N}.$$

Finally, the *binormal* is by definition

$$\mathbf{B}(t) = \mathbf{T}(t) \times \mathbf{N}(t),$$

so that when  $\kappa(t) \neq 0$ , the vectors  $\mathbf{T}(t)$ ,  $\mathbf{N}(t)$ , and  $\mathbf{B}(t)$  form a positively oriented orthonormal basis of  $\mathbb{R}^3$ .

(a) Show that  $\mathbf{B}'$  is a multiple of  $\mathbf{N}$ . We may therefore define a function  $\tau$  by the equation

$$\mathbf{B}' = -\tau \mathbf{N}.$$

$\tau$  is called the *torsion* of  $\mathbf{c}$ .

(b) Prove that  $\mathbf{N}' = -\kappa \mathbf{T} + \tau \mathbf{B}$ . The three identities

$$\begin{array}{rcl} \mathbf{T}' & = & \kappa \mathbf{N} \\ \mathbf{N}' & = & -\kappa \mathbf{T} + \tau \mathbf{B} \\ \mathbf{B}' & = & -\tau \mathbf{N} \end{array}$$

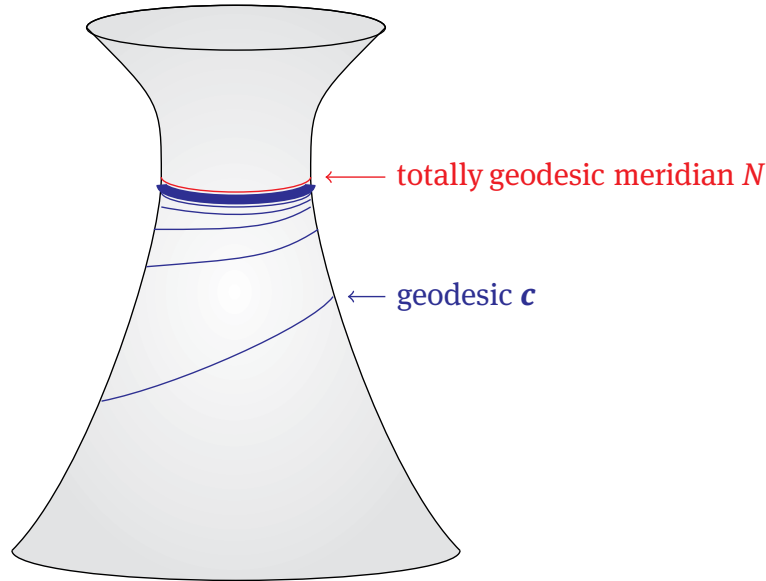
are called the *Serret-Frénet formulas*.

It can be shown that given continuous functions  $\kappa, \tau : [a, b] \rightarrow \mathbb{R}$  with  $\kappa > 0$ , there exists one and only one (up to a rigid motion of  $\mathbb{R}^3$ ) curve in 3-space parametrized by arc length that has  $\kappa$  as curvature and  $\tau$  as torsion. This should be compared with Example 7.4.1 which asserts a similar property for plane curves involving only the curvature.

**7.8.** (a) Prove that any two surfaces in  $\mathbb{R}^3$  with the same constant curvature  $\kappa$  are locally isometric. *Hint:* Use Theorem 6.7.2.

(b) Show that there exist (noncomplete) surfaces of revolution with constant curvature 1 that are not round spheres. These are not totally umbilic, however, and therefore do not contradict the fundamental theorem for submanifolds.

**7.9.** Consider a surface of revolution, and suppose that there is a point  $\mathbf{p} = \boldsymbol{\gamma}(t_0)$  on the profile curve  $\boldsymbol{\gamma}$  where the distance  $\gamma_1$  to the  $z$ -axis has a strict local minimum. Thus, the parallel  $N$  through  $\mathbf{p}$  is (the image of) a geodesic. Assume without loss of generality that  $\boldsymbol{\gamma}$  is pointing upward, and choose  $t_1 < t_0$  close enough to  $t_0$  that  $\gamma_1|_{(t_1, t_0)}$  is strictly decreasing. Set  $\mathbf{q} = \boldsymbol{\gamma}(t_1)$ , and  $\alpha \in (0, \pi/2)$  the angle with  $\cos \alpha = \gamma_1(t_0)/\gamma_1(t_1)$ . Show that the two geodesics emanating from  $\mathbf{q}$  at angle  $\alpha$  with the parallel increase in height for all time, and come arbitrarily close to  $N$  without ever reaching it.



**7.10.** As in the previous exercise, assume  $\gamma$  is pointing upward. Let  $\mathbf{p} \in M$ ,  $\alpha \in (0, \pi/2)$ , and set  $k := r(\mathbf{p}) \cos \alpha$ . Denote by  $\mathbf{c}$  one of the two geodesics from  $\mathbf{p}$  that make an angle  $\alpha$  with the parallel  $\frac{\partial}{\partial x^2}$ . The previous exercise investigated the behavior of  $\mathbf{c}$  when  $k$  is a local minimum value of  $\gamma_1$ . This problem explores the case when  $k$  is not a critical value of  $\gamma_1$ ; for instance,  $M$  could be the paraboloid discussed earlier, but inverted so that its vertex is the highest point. Let  $N$  denote the first parallel above  $\mathbf{p}$  at distance  $k$  from the  $z$ -axis. Show that  $\mathbf{c}$  hits  $N$  tangentially and then winds back down. The image of  $\mathbf{c}$  is invariant under reflection in the plane containing the meridian through the point where  $\mathbf{c}$  hits  $N$ , so that  $\mathbf{c}$  passes again through  $\mathbf{p}$  at an angle  $\alpha$  with the parallel.

**7.11.** Let  $M^2$  be a surface in  $\mathbb{R}^3$ ,  $\mathbf{p} \in M$ , and  $U$  a neighborhood of the origin in  $M_{\mathbf{p}}$  such that  $\exp_{\mathbf{p}} : U \rightarrow \exp_{\mathbf{p}}(U) \subset M$  is a diffeomorphism. By means of a linear isometry  $M_{\mathbf{p}} \cong \mathbb{R}^2$ , introduce polar coordinates  $(\tilde{r}, \tilde{\theta})$  on  $M_{\mathbf{p}} \setminus L$ , where  $L$  is some ray from the origin. If  $V$  denotes  $\exp_{\mathbf{p}}(U \setminus L)$ , then

$$\mathbf{x} = (r, \theta) := (\tilde{r}, \tilde{\theta}) \circ (\exp_{\mathbf{p}}^{-1})|_V$$

defines a chart on  $V$ . The object of this exercise is to derive a formula for the curvature  $K$  in terms of the “polar” chart  $\mathbf{x}$ .

For the sake of brevity, we denote the coordinate vector fields  $\partial/\partial r$  and  $\partial/\partial \theta$  by  $\partial_r$  and  $\partial_\theta$  respectively.

- (a) Prove that  $\nabla_{\partial_r} \partial_r = 0$ ,  $\langle \partial_r, \partial_r \rangle = 1$ , and  $\langle \partial_r, \partial_\theta \rangle = 0$ .
- (b) Let  $G = \langle \partial_\theta, \partial_\theta \rangle$ . Show that

$$G^2 K = \langle D_{\partial_r} \partial_r, \partial_r \times \partial_\theta \rangle \langle D_{\partial_\theta} \partial_\theta, \partial_r \times \partial_\theta \rangle - \langle D_{\partial_r} \partial_\theta, \partial_r \times \partial_\theta \rangle^2.$$

- (c) Identify all vector fields in the identity above with their coordinate vector fields in the standard basis; i.e., identify  $\partial_r$  with the column matrix whose transpose is

$[\langle \partial_r, D_1 \rangle \quad \langle \partial_r, D_2 \rangle \quad \langle \partial_r, D_3 \rangle]$ , etc. Use properties of the cross product to show that

$$G^2 K = \det \begin{bmatrix} (D_{\partial_r} \partial_r)^T \\ \partial_r^T \\ \partial_\theta^T \end{bmatrix} \cdot \det \begin{bmatrix} (D_{\partial_\theta} \partial_\theta)^T \\ \partial_r^T \\ \partial_\theta^T \end{bmatrix} - \det \left( \begin{bmatrix} (D_{\partial_r} \partial_\theta)^T \\ \partial_r^T \\ \partial_\theta^T \end{bmatrix} \right)^2,$$

and deduce that

$$G^2 K = \det \begin{bmatrix} \langle D_{\partial_r} \partial_r, D_{\partial_\theta} \partial_\theta \rangle & \langle D_{\partial_r} \partial_r, \partial_r \rangle & \langle D_{\partial_r} \partial_r, \partial_\theta \rangle \\ \langle \partial_r, D_{\partial_\theta} \partial_\theta \rangle & 1 & 0 \\ \langle \partial_\theta, D_{\partial_\theta} \partial_\theta \rangle & 0 & G \end{bmatrix} \\ - \det \begin{bmatrix} \langle D_{\partial_r} \partial_\theta, D_{\partial_r} \partial_\theta \rangle & \langle D_{\partial_r} \partial_\theta, \partial_r \rangle & \langle D_{\partial_r} \partial_\theta, \partial_\theta \rangle \\ \langle \partial_r, D_{\partial_r} \partial_\theta \rangle & 1 & 0 \\ \langle \partial_\theta, D_{\partial_r} \partial_\theta \rangle & 0 & G \end{bmatrix}.$$

(d) Show that

$$G^2 K = G(\langle D_{\partial_r} \partial_r, D_{\partial_\theta} \partial_\theta \rangle - \langle D_{\partial_r} \partial_\theta, D_{\partial_r} \partial_\theta \rangle) + \frac{1}{4}(D_{\partial_r} G)^2.$$

(e) Prove that

$$\frac{1}{2} D_{\partial_r} D_{\partial_r} G = \langle D_{\partial_r} D_{\partial_\theta} \partial_r, \partial_\theta \rangle + \langle D_{\partial_r} \partial_\theta, D_{\partial_r} \partial_\theta \rangle,$$

and

$$0 = D_{\partial_\theta} \langle D_{\partial_r} \partial_r, \partial_\theta \rangle = \langle D_{\partial_\theta} D_{\partial_r} \partial_r, \partial_\theta \rangle + \langle D_{\partial_r} \partial_r, D_{\partial_\theta} \partial_\theta \rangle.$$

Conclude that

$$\langle D_{\partial_r} \partial_r, D_{\partial_\theta} \partial_\theta \rangle - \langle D_{\partial_r} \partial_\theta, D_{\partial_r} \partial_\theta \rangle = -\frac{1}{2} D_{\partial_r} D_{\partial_r} G.$$

(f) Combine (d) and (e) to obtain

$$4KG^2 = -2GD_{\partial_r} D_{\partial_r} G + (D_{\partial_r} G)^2.$$

This last identity may be rewritten as

$$K = -\frac{D_{\partial_r} D_{\partial_r} \sqrt{G}}{\sqrt{G}}. \quad (7.9.3)$$

Observe that integral curves of  $\partial_r$  are the normal geodesics originating from  $\mathbf{p}$ . Let  $\mathbf{v}$  be a unit vector in  $M_{\mathbf{p}}$ ,  $\gamma_{\mathbf{v}}$  the ray  $t \mapsto t\mathbf{v}$  in  $M_{\mathbf{p}}$ , and  $\mathbf{c}_{\mathbf{v}} = \exp \circ \gamma_{\mathbf{v}}$  the geodesic with initial tangent vector  $\mathbf{v}$ . The restriction  $\partial_{\tilde{\theta}} \circ \gamma_{\mathbf{v}}$  of the polar coordinate vector field  $\partial_{\tilde{\theta}}$  on  $M_{\mathbf{p}}$  is  $t \mapsto t\mathcal{I}_{t\mathbf{v}}\mathbf{w}$ , where  $\mathbf{w}$  is a unit vector orthogonal to  $\mathbf{v}$ . Since the polar coordinate fields on  $M_{\mathbf{p}}$  and on  $M$  are  $\exp_{\mathbf{p}}$ -related,

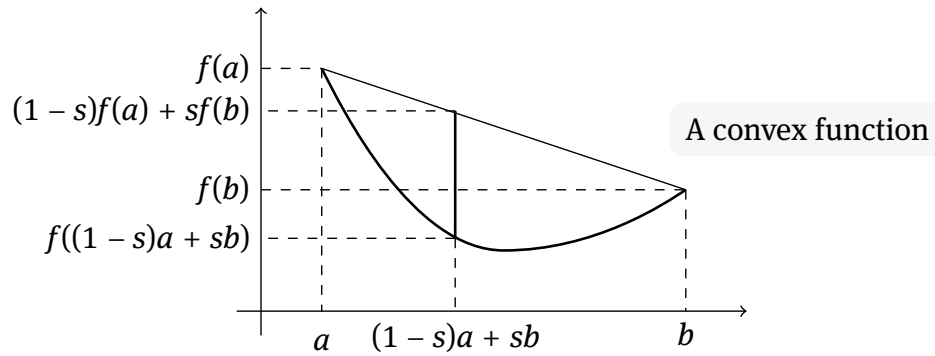
$$\partial_{\tilde{\theta}} \circ \mathbf{c}_{\mathbf{v}}(t) = \exp_{\mathbf{p}*} t\mathcal{I}_{t\mathbf{v}}\mathbf{w},$$

cf. also (6.2.1). In other words,  $\partial_{\tilde{\theta}} \circ \mathbf{c}_{\mathbf{v}}$  is one of the two Jacobi fields  $J$  orthogonal to  $\mathbf{c}_{\mathbf{v}}$  satisfying  $J(0) = \mathbf{0}$ ,  $|J'(0)| = 1$ . Thus, (7.9.3) says that in nonnegative (resp. nonpositive) curvature, this Jacobi field has concave (resp. convex) norm, see also Exercise 7.12.

**7.12.** Recall that a function  $f : I \rightarrow \mathbb{R}$  defined on an interval  $I$  is said to be *convex* if its graph lies below any secant line; explicitly, for  $a, b \in I$  with  $a < b$ ,

$$f((1 - s)a + sb) \leq (1 - s)f(a) + sf(b), \quad s \in [0, 1].$$

When the inequality above is strict for all  $s \in (0, 1)$  (and all  $a < b$ ),  $f$  is said to be *strictly convex*.  $f$  is said to be *concave* if  $-f$  is convex.



- (a) Show that a convex function is necessarily continuous. Prove that if  $f$  is  $C^2$ , then  $f$  is convex if and only if  $f'' \geq 0$ , and that if  $f'' > 0$ , then  $f$  is strictly convex.
- (b) A function  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  is said to be convex or concave if the restriction  $f \circ \mathbf{c}$  of  $f$  to any geodesic  $\mathbf{c}$  of  $\mathbb{R}^{n+1}$  has that property. Prove that if this is the case, then  $M = f^{-1}(a)$  has nonnegative curvature for any regular value  $a$  of  $f$ .
- (c) If  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  is convex, show that the set  $f^{-1}(-\infty, a)$  is convex for any  $a \in \mathbb{R}$ .
- (d) Prove that any convex, bounded function  $f : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$  is constant.

**7.13.** Prove that the graph of  $f : U \subset \mathbb{R}^2 \rightarrow \mathbb{R}$  has nonnegative curvature if  $f$  is either convex or concave (see Exercise 7.12).

**7.14.** Let  $f : (a, b) \rightarrow \mathbb{R}$ . Show that the cylinder over the graph of  $f$  (i.e., the hypersurface parametrized by  $(s, t) \mapsto (s, f(s), t)$ ) is flat.

**7.15.** Show that if the hypersurface  $M$  is positively curved at  $\mathbf{p}$ , then it is strictly convex at that point: i.e., there is a neighborhood  $U$  of  $\mathbf{p}$  in  $M$  such that  $U \cap \{\mathbf{p} + \pi_2(M_{\mathbf{p}})\} = \{\mathbf{p}\}$ .

**7.16.** Use Exercise 7.15 to show that any hypersurface with positive curvature is orientable. Is this still true for nonnegative curvature?

**7.17.** Suppose  $b$  is a semi-definite symmetric bilinear form on an inner product space  $(V, \langle \cdot, \cdot \rangle)$ . The eigenvalues of the associated self-adjoint operator  $L$ , where  $\langle L\mathbf{v}, \mathbf{w} \rangle = b(\mathbf{v}, \mathbf{w})$ , are then either all nonnegative or all nonpositive.

- (a) Write  $V = V_0 \oplus V_1$ , where  $V_0$  is the kernel of  $L$ , and  $V_1$  its orthogonal complement, and let  $\pi : V \rightarrow V_1$  denote the projection. Show that  $b(\mathbf{v}, \mathbf{w}) = b(\pi\mathbf{v}, \pi\mathbf{w})$  for all  $\mathbf{v}$  and  $\mathbf{w}$  in  $V$ . Conclude that

$$b(\mathbf{v}, \mathbf{v}) \cdot b(\mathbf{w}, \mathbf{w}) \geq b^2(\mathbf{v}, \mathbf{w}), \quad \mathbf{v}, \mathbf{w} \in V.$$

(b) Show that  $M$  has nonnegative curvature at a point if and only if the second fundamental form is semi-definite at that point.

**7.18.** A *quadric* in  $\mathbb{R}^3$  is a set of the form  $p^{-1}(0)$ , where  $p$  is a polynomial of degree  $\leq 2$ ; i.e.,

$$p(x_1, x_2, x_3) = \sum_{i,j=1}^3 a_{ij}x_i x_j + \sum_{i=1}^3 b_i x_i + c.$$

(a) Show that there exists a linear isometry  $L$  of  $\mathbb{R}^3$  and  $A_i, B_i, C \in \mathbb{R}$  such that

$$L(M) = \{(x_1, x_2, x_3) \mid \sum_{i=1}^3 A_i x_i^2 + B_i x_i + C = 0\}.$$

(b) Prove that there exists a translation  $T$  of  $\mathbb{R}^3$  such that

$$TL(M) = \{(x_1, x_2, x_3) \mid \sum_{i=1}^3 p_i(x_i) + \gamma = 0\},$$

where each  $p_i$  has the form  $p_i(x_i) = a_i x_i^2$  or  $p_i(x_i) = \beta_i x_i$ ; i.e., each  $B_i$  in (a) may be assumed to be zero if the corresponding  $A_i$  is nonzero.

(c) Show that a nonempty quadric is isometric to

– an ellipsoid

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} + \frac{z^2}{c^2} = 1,$$

a hyperboloid of one sheet

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} - \frac{z^2}{c^2} = 1,$$

or a hyperboloid of two sheets

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} - \frac{z^2}{c^2} = -1,$$

if none of the  $A_i$  vanish;

– an elliptic paraboloid

$$z = \frac{x^2}{a^2} + \frac{y^2}{b^2},$$

or a hyperbolic paraboloid

$$z = \frac{x^2}{a^2} - \frac{y^2}{b^2},$$

if exactly one  $A_i = 0$ , but the corresponding  $B_i \neq 0$ ;

– a line, or a cylinder over a line, parabola, ellipse, or hyperbola in the  $x$ - $y$  plane otherwise.

**7.19.** Prove that a surface  $M^2 \subset \mathbb{R}^3$  with negative curvature has no umbilical points.

*Hint:* There is a one sentence proof.

**7.20.** Let  $L$  be a self-adjoint operator on  $\mathbb{R}^3$ , and suppose that

$$M = \{\mathbf{x} \in \mathbb{R}^3 \mid \langle L\mathbf{x}, \mathbf{x} \rangle = 1\}$$

is nonempty. Assume also that 1 is a regular value of  $\mathbf{x} \mapsto \langle L\mathbf{x}, \mathbf{x} \rangle$ , so that  $M$  is a manifold (and a quadric).

(a) Show that  $\mathbf{x} \in M$  is umbilical if and only if

$$\det \begin{bmatrix} L\mathbf{x} & L\mathbf{u} & \mathbf{u} \end{bmatrix} = 0 \text{ for all } \mathbf{u} \perp L\mathbf{x}.$$

(b) Use (a) to prove that a 2-dimensional sphere (which corresponds to  $L$  being a multiple of the identity) is totally umbilic.

(c) Suppose  $L$  has 2 distinct eigenvalues  $\lambda_1 > 0$  and  $\lambda_2$ , where the  $\lambda_2$ -eigenspace is 2-dimensional. Show that there exist at least two points  $\pm\mathbf{x}$  of  $M$  that are umbilical. *Hint:*  $\mathbf{x}$  will be an eigenvector.

**7.21.** Let  $M^2$  be an oriented surface in  $\mathbb{R}^3$ , with unit normal  $\mathbf{n}$  and corresponding second fundamental tensor  $S$ . A curve  $\mathbf{c}$  in  $M$  is called a *line of curvature* if  $\dot{\mathbf{c}}(t)$  is an eigenvector of  $S$  at  $\mathbf{c}(t)$  for all  $t$ .

(a) Suppose  $P$  is a plane that intersects  $M$  orthogonally (in the sense that if  $\mathbf{c}$  parametrizes  $P \cap M$ , then  $\mathbf{n} \circ \mathbf{c}$  is tangent to the plane). Show that  $\mathbf{c}$  is a line of curvature, and that the corresponding principal curvature is the curvature of  $\mathbf{c}$ . Notice that this again proves that a sphere is totally umbilic.

(b) Part (a) shows that the meridians in a surface of revolution are lines of curvature. Prove that parallels are also lines of curvature. *Hint:* the restriction of  $\mathbf{n}$  to the parallel makes a constant angle with the axis of revolution.

(c) Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a smooth, even function with  $f(0) = 0$ , and  $M$  the surface of revolution obtained by revolving the curve  $t \mapsto (t, 0, f(t))$ ,  $t \geq 0$  about the  $z$ -axis. Show that the origin is umbilical.

**7.22.** Both Exercise 7.20 and 7.21 imply that if  $M$  is the ellipsoid of revolution

$$\frac{x^2 + y^2}{a^2} + \frac{z^2}{c^2} = 1,$$

then the points  $\pm(0, 0, c)$  on the axis of revolution are umbilical. In order to find any other possible umbilics, it is convenient to work with the description we gave of a surface of revolution.

(a) Suppose the profile curve  $\gamma = (\gamma_1, 0, \gamma_2)$  can be expressed as the graph of a function  $x = f(z)$  in the  $x$ - $z$  plane; i.e.,  $\gamma_1 = f$ ,  $\gamma_2(s) = s$ . Show that the principal curvatures from (7.8.3) become

$$\lambda_1 = \frac{-f''}{(1 + (f')^2)^{3/2}}, \quad \lambda_2 = \frac{1}{f(1 + (f')^2)^{1/2}},$$

whenever  $f \neq 0$ . Thus, umbilical points correspond to those values of  $s$  for which  $(ff'' + (f')^2)(s) + 1 = 0$ , and occur in (possibly degenerate) parallels.



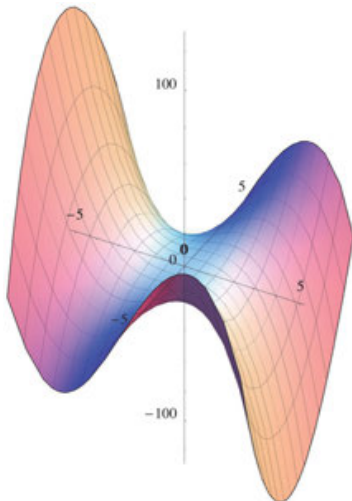
- (b) Prove that there are no other umbilics on the ellipsoid of revolution. This contrasts with the generic ellipsoid  $x^2/a^2 + y^2/b^2 + z^2/c^2 = 1$  which admits exactly four umbilics if  $a$ ,  $b$ , and  $c$  are all distinct.

**7.23.** Show that the paraboloid

$$z = \frac{x^2}{a^2} + \frac{y^2}{b^2}$$

has two umbilical points if  $a \neq b$  and only one (the vertex) if  $a = b$ .

**7.24.** A *parabolic umbilic* on a surface  $M$  in  $\mathbb{R}^3$  is a point where both principal curvatures are zero (in [4] such a point is said to be planar). In particular, the point is umbilical, and the sectional curvature and second fundamental form vanish there. Show that on the surface given by the graph of  $f$ , where  $f(x, y) = xy^2 - x^3$ , the origin is a parabolic umbilic. It is often called a *monkey saddle* because it has three downward slopes emanating from it, two for the legs and one for the tail.



**Fig. 7.8:** A “monkey saddle”

**7.25.** Recall that  $\mathbf{h}_1$ , with  $\mathbf{h}_1(s, t) = (t \cos s, t \sin s, s)$ , parametrizes a helicoid, and that  $\mathbf{h}_2$ , where  $\mathbf{h}_2(u, v) = (\cosh v \cos u, \cosh v \sin u, v)$ , yields a catenoid. Prove that both surfaces are locally isometric. *Hint:* reparametrize the catenoid by  $\tilde{\mathbf{h}}_2$ , where

$$\tilde{\mathbf{h}}_2(s, t) = \mathbf{h}_2(s, \operatorname{arccosh} \sqrt{1 + t^2}),$$

and consider  $\tilde{\mathbf{h}}_2 \circ \mathbf{h}_1^{-1}$ .

**7.26.** Let  $M = \{(x, y, z) \in \mathbb{R}^3 \mid e^z \cos x = \cos y\}$ .

- (a) Show that  $M$  is an orientable 2-dimensional manifold.  
 (b) Prove that  $M$  may be described as follows: given  $k, l \in \mathbb{Z}$ , consider the open squares

$$S_{k,l}^1 = \left(-\frac{\pi}{2} + 2k\pi, \frac{\pi}{2} + 2k\pi\right) \times \left(-\frac{\pi}{2} + 2l\pi, \frac{\pi}{2} + 2l\pi\right),$$

and

$$S_{k,l}^2 = \{(x + \pi, y + \pi) \mid (x, y) \in S_{k,l}^1\}.$$

Then  $M$  consists of the graph of  $f$ , where  $f(x, y) = \ln(\cos y / \cos x)$ , over each of these squares, together with the vertical lines that pass through the vertices of the squares.  $M$  is called *Scherk's surface*. It is named after the 19th century German mathematician who discovered it.

(c) Show that  $M$  is a minimal surface.

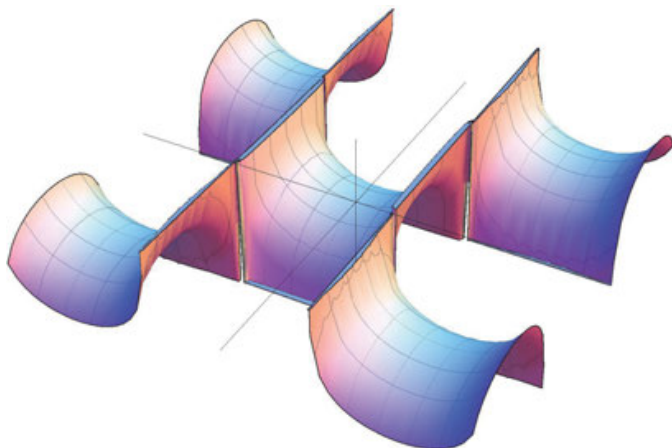


Fig. 7.9: Scherk's surface

**7.27.** The image of  $\mathbf{h}$ , where

$$\mathbf{h}(u, v) = \frac{1}{2}\left(u\left(1 - \frac{u^2}{3} + v^2\right), v\left(1 + \frac{v^2}{3} - u^2\right), u^2 - v^2\right),$$

is called *Enneper's surface*.  $\mathbf{h}$  is not strictly speaking a parametrization (and its image  $M$  is not a manifold) because it is not one-to-one. Nevertheless,  $M$  is locally a manifold.

- (a) Show that  $M$  is a minimal surface.  
 (b) Prove that  $M$  has strictly negative curvature

$$(K \circ \mathbf{h})(u, v) = -\frac{4}{9^3}(1 + u^2 + v^2)^4.$$

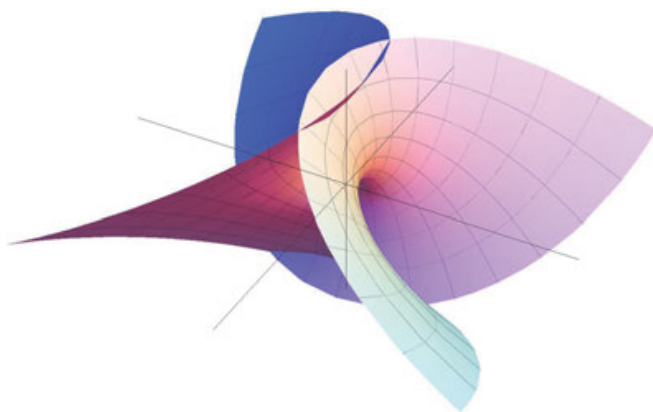


Fig. 7.10: Enneper's surface

## Appendix A

In this appendix, we review some basic properties of real numbers that are used throughout the text but are usually not discussed in detail in Calculus or other lower level math courses. This is by no means a comprehensive overview; rather, it emphasizes those features – such as the rationals being countable and dense – that provide interesting examples and counterexamples, particularly in the theory of limits and integration.

Recall that there is a nested sequence  $\mathbb{N} \subset \mathbb{Z} \subset \mathbb{Q} \subset \mathbb{R}$ , with  $\mathbb{N} = \{1, 2, \dots\}$  denoting the *natural numbers*,  $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$  the *integers*, and  $\mathbb{Q} = \{p/q \mid p \in \mathbb{Z}, q \in \mathbb{N}\}$  the *rationals*. A real number that is not rational is said to be *irrational*.

**Definition.**  $\alpha \in \mathbb{R}$  is said to be an *upper bound* (resp. a *lower bound*) of  $A \subset \mathbb{R}$  if for any  $a \in A$ ,  $\alpha \geq a$  (resp.  $\leq a$ ). An upper bound is said to be the *least upper bound* or *supremum* of  $A$  if it is less than or equal to any other upper bound of  $A$ . It is denoted  $\sup A$ . Similarly, a lower bound of  $A$  is called the *greatest lower bound* or *infimum* if it is greater or equal to any other lower bound. The infimum of  $A$ , if it exists, is denoted  $\inf A$ .

It is clear from the definition that infimum and supremum, if they exist, are unique. Notice that if we set  $-A = \{-a \mid a \in A\}$ , then  $A$  is bounded above by  $\alpha$  if and only if  $-A$  is bounded below by  $-\alpha$ . Furthermore,  $\alpha$  is less than or equal to any other upper bound of  $A$  if and only if  $-\alpha$  is greater than or equal to any other lower bound of  $-A$ ; i.e.,  $-\sup A = \inf(-A)$ , provided one of them exists. The supremum of a set may, or may not, belong to it; for example, the supremum of the interval  $(0, 1)$  is 1, which does not lie in the open interval. Indeed, 1 is an upper bound, and by definition, if  $\alpha$  is any upper bound of  $(0, 1)$ , then  $\alpha \geq 1$ ; i.e.,  $1 = \sup(0, 1)$ . When the supremum does belong to the set, it is called the *maximum*. Similar considerations apply to the infimum, which is called the *minimum* when it belongs to the set.

A useful criterion for determining the least upper bound of a set is the following:  $\alpha \in \mathbb{R}$  is the least upper bound of  $A$  if and only if

- (1)  $\alpha$  is an upper bound of  $A$ , and
- (2) for any  $\varepsilon > 0$ , there exists some  $a \in A$  such that  $\alpha - \varepsilon < a$ .

To see this, suppose first that  $\alpha = \sup A$ . Condition (1) is then verified by definition. If (2) were not, then there would exist some  $\varepsilon > 0$  such that  $\alpha - \varepsilon$  is greater than or equal to any element of  $A$ . But this would mean that  $\alpha - \varepsilon$  is an upper bound of  $A$  which is smaller than the least upper bound. Thus, (2) is also verified. Conversely, suppose  $\alpha$  satisfies both conditions (1) and (2). To establish that  $\alpha = \sup A$ , we only need to show that if  $\beta$  is an upper bound of  $A$ , then  $\alpha \leq \beta$ . So suppose that  $\beta$  is an upper bound of  $A$  but  $\beta < \alpha$ . Take  $\varepsilon = \alpha - \beta > 0$  in (2) to deduce that there exists some  $a \in A$  with  $\alpha - \varepsilon = \beta < a$ . This contradicts the fact that  $\beta$  is an upper bound of  $A$ .

A similar characterization with inequalities reversed applies to the infimum: a lower bound  $\alpha$  of  $A$  is the greatest lower bound if and only if for any  $\varepsilon > 0$ , there is some  $a \in A$  satisfying  $a < \alpha + \varepsilon$ .

A fundamental property of real numbers is the following:

**Completeness Axiom.** *Every nonempty subset of  $\mathbb{R}$  that is bounded above has a least upper bound.*

Notice that by the above discussion, the completeness axiom can also be phrased as: every nonempty subset of real numbers that is bounded below has a greatest lower bound.

Let us examine some consequences of the completeness axiom:

**Theorem (Well-ordering Principle).** *Every nonempty subset of the natural numbers has a smallest element (i.e., a minimum).*

*Proof.* Let  $A$  be a nonempty set of natural numbers. Being bounded below (by 0),  $A$  has an infimum  $\alpha$ . We claim  $\alpha \in \mathbb{N}$ : suppose  $\alpha$  is not a natural number. Using the fact that  $\alpha$  is the greatest lower bound of  $A$ , there exists  $a \in A$  such that  $\alpha < a < \alpha + 1/2$ . Since  $\alpha \neq a$ , the number  $\varepsilon = a - \alpha$  is positive, and there exists  $\tilde{a} \in A$  such that  $\alpha < \tilde{a} < \alpha + \varepsilon = a < \alpha + 1/2$ . This is impossible, because both  $a$  and  $\tilde{a}$  are natural numbers, so their distance cannot be less than  $1/2$ . Thus,  $\alpha \in \mathbb{N}$ , and must therefore also belong to  $A$ : otherwise there would exist a natural number in  $A$  at distance less than  $1/2$  from  $\alpha$ , which again is impossible.  $\square$

The next result is an extremely useful tool for proving properties involving natural numbers:

**Theorem (Mathematical Induction).** *Suppose  $P(n)$  is a statement for each  $n \in \mathbb{N}$  such that*

- (1)  $P(1)$  is true, and
- (2) for any  $k \in \mathbb{N}$ ,  $P(k + 1)$  is true whenever  $P(k)$  is.

*Then  $P(n)$  is true for all  $n \in \mathbb{N}$ .*

*Proof.* We argue by contradiction: suppose  $P(n)$  is not true for all  $n$ , so that the set  $A$  of all  $n \in \mathbb{N}$  for which  $P(n)$  is not true is nonempty. By the well-ordering principle,  $A$  has a smallest element  $k \in A \subset \mathbb{N}$ . By (1),  $k > 1$ , so that  $k - 1 \in \mathbb{N}$ , and  $P(k - 1)$  is true. This contradicts (2), since  $P(k)$  is false.  $\square$

**Theorem (Archimedean Principle).** *If  $a$  and  $b$  are real numbers with  $a > 0$ , there exists  $n \in \mathbb{N}$  such that  $na > b$ .*

*Proof.* Let us first rephrase the statement: since  $a > 0$ ,  $b/a \in \mathbb{R}$ , and we must show that there exists some  $n \in \mathbb{N}$  that is larger than  $b/a$ . In other words, the Archimedean principle just asserts that the set of natural numbers is not bounded above. Once again, the argument will be by contradiction. Suppose  $\mathbb{N}$  is bounded above. By the completeness

axiom,  $\mathbb{N}$  has a least upper bound  $\alpha$ . Since  $\alpha$  is the smallest upper bound, there must exist some natural number  $n$  such that  $\alpha - 1 < n \leq \alpha$ . The first inequality says that  $\alpha$  is less than the natural number  $n + 1$ , which is a contradiction.  $\square$

Our next endeavor is to compare the sizes of different sets of numbers. For finite sets, this is straightforward: sets  $A$  and  $B$  are said to have the same *cardinality* if there exists a bijection (i.e., a one-to-one and onto) map  $f : A \rightarrow B$ . It is therefore natural to extend this definition to infinite sets, and we do so. Strange things happen in the process, however. For example, intuitively,  $\mathbb{Z}$  should have at least twice as many elements as  $\mathbb{N}$ , since  $\mathbb{Z} = \mathbb{N} \cup (-\mathbb{N}) \cup \{0\}$ . Nevertheless both sets have the same cardinality: it is easy to check that the map  $f : \mathbb{Z} \rightarrow \mathbb{N}$ , where  $f(n) = 2n + 1$  if  $n \geq 0$ , and  $f(n) = -2n$  if  $n < 0$ , is a bijection.

**Definition.** A set  $A$  is said to be *countable* if there exists a one-to-one map  $f : A \rightarrow \mathbb{N}$ .

Thus, any finite set is countable, but so is the set of all integers. A useful alternative characterization of countability is the following:  $A$  is countable if and only if there exists a surjective  $f : \mathbb{N} \rightarrow A$ . Indeed, if  $f : \mathbb{N} \rightarrow A$  is surjective, we may construct a one-to-one map  $g : A \rightarrow \mathbb{N}$  by defining  $g(a)$  to be any one element in  $f^{-1}(a) = \{n \in \mathbb{N} \mid f(n) = a\}$ . Conversely, if  $A$  is countable, and  $g : A \rightarrow \mathbb{N}$  is one-to-one, we obtain a surjective map  $f : \mathbb{N} \rightarrow A$  by setting  $f(n) = a$  if  $a$  is the (necessarily unique) element of  $A$  such that  $g(a) = n$ , and  $f(n) = a_0$  for some fixed  $a_0 \in A$  if there is no  $a \in A$  with  $g(a) = n$ .

Since by definition, a map  $f : \mathbb{N} \rightarrow A$  is a sequence, we may write any countable set  $A$  as  $A = \{a_1, a_2, \dots\}$ , where  $a_n = f(n)$ . More generally,  $\mathbb{N}$  may be replaced in the above discussion by any countably infinite set (i.e., any set with the same cardinality as  $\mathbb{N}$ ).

**Proposition.** A countable union of countable sets is countable.

*Proof.* If  $E_i$ ,  $i \in \mathbb{N}$ , is countable, we may write  $E_i = \{x_{i1}, x_{i2}, \dots\}$ . Define  $f : \cup_{i=1}^{\infty} E_i \rightarrow \mathbb{N}$  by  $f(x_{ij}) = 2^i 3^j$ .  $f$  is one-to-one, for if  $f(x_{ij}) = f(x_{kl})$ , then  $2^i 3^j = 2^k 3^l$ , and  $2^{i-k} = 3^{l-j}$ . This can only hold if  $i = k$  and  $j = l$ .  $\square$

It follows that the Cartesian product  $A \times B$  of countable sets  $A$  and  $B$  is again countable: if  $A = \{a_1, a_2, \dots\}$ , then  $A \times B = \cup_{i=1}^{\infty} \{a_i\} \times B$  is a countable union of countable sets. In particular,

**Theorem.** The set  $\mathbb{Q}$  of rational numbers is countable.

*Proof.* The map  $f : \mathbb{Z} \times \mathbb{N} \rightarrow \mathbb{Q}$ ,  $f(m, n) = m/n$ , is onto.  $\square$

A set which is not countable is said to be *uncountable*. There are many such:

**Theorem.** The interval  $(0, 1)$  is uncountable. In particular, the set of irrational numbers between 0 and 1 is uncountable.

*Proof.* Any number  $a \in (0, 1)$  admits a decimal expansion

$$a = .a_1a_2a_3\cdots = \sum_{i=1}^{\infty} \frac{a_i}{10^i}, \quad a_i \in \{0, 1, 2, \dots, 9\}.$$

In fact, for any  $b > 0$ , let  $[b]$  denote the unique  $n \in \mathbb{N} \cup \{0\}$  such that  $n \leq b < n + 1$ ; i.e.,  $[b]$  is the largest integer smaller than or equal to  $b$  (specifically,  $[b] = \sup\{n \in \mathbb{N} \mid n \leq b\}$  – notice that the least upper bound of a set of integers, if it exists, is an integer, so the supremum is in fact a maximum). It is not hard to check that in the formula above,  $a_n$  may be taken to be  $\lceil 10(10^{n-1}a - [10^{n-1}a]) \rceil$ . The reason we write “may be taken” is that this decimal expansion need not be unique; e.g.,  $0.5 = 0.499999\dots$ . This, however, is the only way in which the expansion can fail to be unique: more precisely, suppose  $a = \sum(a_k/10^k) = \sum(b_k/10^k) = b$ , and let  $i$  be the smallest integer such that  $a_i \neq b_i$ ; assuming without loss of generality that  $a_i > b_i$ , we must have  $a_i = b_i + 1$ ,  $a_{k+i} = 0$  and  $b_{k+i} = 9$  for all  $k \geq 1$ . To see this, observe that

$$0 = 10^i \sum_{k=1}^{\infty} \frac{a_k - b_k}{10^k} = \sum_{k=i}^{\infty} \frac{a_k - b_k}{10^{k-i}} = a_i - b_i + \sum_{k=1}^{\infty} \frac{a_{k+i} - b_{k+i}}{10^k}. \quad (\text{A.1})$$

Now,  $a_i - b_i$  is an integer between 1 and 9, whereas the last series in the above expression satisfies in absolute value

$$\left| \sum_{k=1}^{\infty} \frac{a_{k+i} - b_{k+i}}{10^k} \right| \leq \sum_{k=1}^{\infty} \frac{|a_{k+i} - b_{k+i}|}{10^k} \leq \sum_{k=1}^{\infty} \frac{9}{10^k} = 1.$$

Thus, if the expression in (A.1) is to equal zero, then  $a_i - b_i$  must equal 1 and  $a_{k+i} - b_{k+i} = -9$ , as claimed.

Now, suppose that  $(0, 1)$  is countable, so that  $(0, 1) = \{a_1, a_2, \dots\}$ , with  $a_i = .a_{i1}a_{i2}\dots$ . The contradiction will arise once we exhibit some  $b \in (0, 1)$  that is not equal to any  $a_i$ . So set  $b = .b_1b_2\dots$ , where  $b_i = 1$  if  $a_{ii} = 2$  and  $b_i = 2$  otherwise. Then  $b$  differs in its  $i$ -th digit from that of  $a_i$ . Since the only numbers that appear in the decimal expansion of  $b$  are 1 and 2, this expansion is unique, so that  $b$  is not equal to any  $a_i$ . This shows that  $(0, 1)$  is uncountable.

The assertion that the set  $A$  of irrationals between 0 and 1 is uncountable is now clear: if it were countable, then  $(0, 1)$  would be the union of two countable sets, namely  $A$  and the set of rationals between 0 and 1, and would then also be countable.  $\square$

**Theorem.** For any real numbers  $a < b$ , the interval  $(a, b)$  contains infinitely many rational and infinitely many irrational numbers.

*Proof.* We first show that the interval contains at least one rational and one irrational, beginning with rational: it may be assumed that  $a > 0$ , for  $k + a > 0$  if  $k$  is a sufficiently large natural number, and if  $r$  is a rational between  $k + a$  and  $k + b$ , then  $r - k$  is one between  $a$  and  $b$ . We must find natural numbers  $m, n$  such that  $a < m/n < b$ , or equivalently, such that  $na < m < nb$ . This in turn requires that  $nb - na > 1$ . So choose

any  $n \in \mathbb{N}$  with  $n > 1/(b - a)$ ,  $1/a$  (which is possible by the Archimedean principle). If  $na \in \mathbb{N}$ , then  $r = (na + 1)/n$  is a rational between  $a$  and  $b$ . If  $na \notin \mathbb{N}$ , consider the set  $E$  of all natural numbers no larger than  $na$ . Since  $na > 1$ ,  $E$  is nonempty, and being bounded above admits a supremum  $k_0$ . Furthermore,  $k_0 \in \mathbb{N}$ , so it is strictly smaller than  $na$ , and

$$na < k_0 + 1 < na + 1 < nb.$$

We may then take  $r$  to equal  $(k_0 + 1)/n$ .

To exhibit an irrational, let  $\alpha \in (a, b) \cap \mathbb{Q}$  (the existence of such an  $\alpha$  has just been established), and use the Archimedean Principle to find  $N \in \mathbb{N}$  such that  $\sqrt{2}/N < b - \alpha$ . Then  $\alpha + \sqrt{2}/N$  is an irrational in  $(a, b)$ .

To see that there are infinitely many rationals, let  $\alpha_1 \in (a, b) \cap \mathbb{Q}$ , and construct inductively an infinite sequence of distinct rationals by choosing some  $\alpha_{n+1} \in (a, \alpha_n) \cap \mathbb{Q}$ . A similar argument shows that there are infinitely many irrationals.  $\square$

The Theorem says that the closure of  $\mathbb{Q}$  is the whole real line. This is often expressed by saying that  $\mathbb{Q}$  is *dense* in  $\mathbb{R}$ .

**Examples and Remarks.** (i) Given any positive number  $a$ , there exists an increasing sequence  $\{b_n\}$  of rationals with  $\lim_{n \rightarrow \infty} b_n = a$ . In fact, writing  $a$  as an infinite decimal  $a = a_1.a_2 \dots$  with  $a_1 \in \mathbb{N} \cup \{0\}$ , and  $a_i \in \{0, 1, \dots, 9\}$  for  $i > 1$ , define  $b_n = a_1.a_2 \dots a_n$ .  $\{b_n\}$  is clearly an increasing sequence, and it converges to  $a$  because  $a - b_n < 9/10^{n+1}$ . This shows, in particular, that the completeness axiom does not hold for the rational numbers: if  $a$  as above is irrational, then the set  $\{x \in \mathbb{Q} \mid x < a\}$  is bounded above, but its supremum  $a$  does not lie in  $\mathbb{Q}$ .

(ii) Let  $f : (0, 1) \rightarrow \mathbb{R}$  be given by  $f(x) = 0$  if  $x$  is irrational,  $f(x) = 1/q$  if  $x = p/q$  is rational, and  $p, q$  have no common factors. Then  $f$  is continuous at every irrational point and discontinuous elsewhere; i.e.,  $f$  has a countable number of discontinuities. In fact, if  $x_0$  is rational, there exists a sequence  $\{x_n\}$  of irrationals converging to  $x_0$  (for example, let  $x_n$  be an irrational number in  $(x_0 - 1/n, x_0) \cap (0, x_0)$ ), so that  $f(x_n) = 0$  does not converge to  $f(x_0)$ . Next, suppose  $x_0$  is irrational. It is enough to show that if  $x_n \rightarrow x_0$ ,  $x_n \in \mathbb{Q} \cap (0, 1)$ , then  $f(x_n) \rightarrow 0$ ; i.e., if  $x_n = p_n/q_n$ , then  $q_n \rightarrow \infty$ . Suppose not. Then there exists  $N \in \mathbb{N}$  such that for any  $M \in \mathbb{N}$ ,  $q_n \leq N$  for some  $n \geq M$ . In other words, there exists a subsequence  $\{q_{n_k}\}$  whose terms are all less than or equal to  $N$ . This means that the subsequence  $\{p_{n_k}/q_{n_k}\}$  can only take on the values  $1/N, 2/N, \dots, (N - 1)/N$ . Since  $x_0$  does not equal any of these values, this subsequence does not converge to  $x_0$ , and hence neither does the original one.





## Appendix B

The basic aim of this section is to show that any invertible linear transformation is a finite composition of certain elementary ones. It has been relegated to an appendix because it is only used in the proof of the change of variables theorem for integrals. Nevertheless, it has several useful applications to rank and systems of linear equations.

**Definition B.1.** A linear map  $L : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is said to be an *elementary transformation* if it belongs to one of the following groups:

- type 1: there exist indices  $1 \leq i < j \leq n$  such that  $Le_i = e_j$ ,  $Le_j = e_i$ , and  $Le_k = e_k$  for  $k \neq i, j$ ;
- type 2: there exists  $1 \leq i \leq n$ ,  $a \neq 0$ , such that  $Le_i = ae_i$ , and  $Le_k = e_k$  for  $k \neq i$ ;
- type 3: there exist  $1 \leq i, j \leq n$ ,  $i \neq j$ ,  $a \in \mathbb{R}$ , such that  $Le_i = e_i + ae_j$ , and  $Le_k = e_k$  for  $k \neq i$ .

The matrix of an elementary transformation with respect to the standard basis is called an *elementary matrix*.

It is clear from the definition that an elementary transformation is an isomorphism, and that its inverse is an elementary transformation of the same type. It follows that an elementary matrix is invertible, and its inverse is an elementary matrix of the same type. Properties of elementary transformations yield corresponding properties for elementary matrices, and vice-versa. The reader may wish to keep this in mind since we will usually only outline properties of one rather than of both.

Thus, an elementary matrix is of

- type 1 if it is obtained by interchanging two columns of the  $n \times n$  identity matrix  $I_n$ ,
- type 2 if it is obtained by multiplying a column of  $I_n$  by a nonzero scalar, and
- type 3 if it is obtained by adding a multiple of one column of  $I_n$  to another column.

These three operations are called *elementary column operations*. Replacing columns by rows, one obtains three types of *elementary row operations*. It is straightforward to check that an elementary matrix of a given type is also obtained by performing an elementary row operation of that type on  $I_n$ .

Given an  $m \times n$  matrix  $A$ , we will denote by  $A_i$  its  $i$ -th row, and by  $A^j$  its  $j$ -th column. Notice that if  $e_i$  is the  $i$ -th coordinate vector in the standard basis written as a column, then  $A_i = e_i^T A$ , and  $A^j = Ae_j$ .

One can also perform elementary row or column operations on any matrix. This is actually equivalent to multiplying that matrix by an elementary one:

**Proposition B.1.** Let  $B$  denote the matrix obtained by performing an elementary row (resp. column) operation on  $A \in M_{m,n}$ , and  $E$  the elementary matrix obtained by performing that same operation on  $I_m$  (resp.  $I_n$ ). Then  $B = EA$  (resp.  $AE$ ).

*Proof.* One has to check the claim for each type of operation. We illustrate this for a row operation of type 2, and leave the others as an exercise: let  $E$  denote the elementary matrix obtained by interchanging rows  $i$  and  $j$  of  $I_n$ . Then the  $i$ -th row of  $EA$  is

$$(EA)_i = E_i A = \mathbf{e}_j^T A = A_j = B_i,$$

as claimed, and by symmetry,  $(EA)_j = B_j$ . If, on the other hand,  $k \neq i, j$ , then

$$(EA)_k = E_k A = \mathbf{e}_k^T A = A_k = B_k,$$

so that  $EA = B$ , as claimed.  $\square$

**Proposition B.2.** *Multiplying a matrix by an elementary one does not affect its rank.*

*Proof.* Recall that the rank of  $A$  is the rank of the linear transformation  $L_A$ . If  $E$  is an elementary matrix, then

$$\begin{aligned} \text{rank}(EA) &= \text{rank}(L_E \circ L_A) = \dim L_E(L_A(\mathbb{R}^n)) = \dim L_A(\mathbb{R}^n) = \text{rank}(L_A) \\ &= \text{rank}(A) \end{aligned}$$

since  $L_E$  is an isomorphism. Similarly, the nullity of  $L_{AE} = L_A \circ L_E$  equals the nullity of  $L_A$ , so that by Theorem 1.2.2,  $AE$  and  $A$  have the same rank.  $\square$

Next, we show that any matrix  $A$  can be transformed by means of finitely many elementary row and column operations into a diagonal matrix (i.e.  $a_{ij} = 0$  if  $i \neq j$ ) with diagonal entries 1 or zero. By the Proposition, the number of 1's equals the rank of the matrix. The proof is best illustrated by means of an example:

**Example B.1.** Let

$$A = \begin{bmatrix} 1 & 2 & 3 & 1 \\ 2 & 1 & 1 & 1 \\ 1 & -1 & 1 & 0 \end{bmatrix}.$$

Subtract 2 times row 1 from row 2, and subtract row 1 from 3 to obtain

$$\begin{bmatrix} 1 & 2 & 3 & 1 \\ 0 & -3 & -5 & -1 \\ 0 & -3 & -2 & -1 \end{bmatrix}.$$

Subtract row 2 from row 3 to get

$$\begin{bmatrix} 1 & 2 & 3 & 1 \\ 0 & -3 & -5 & -1 \\ 0 & 0 & 3 & 0 \end{bmatrix}.$$

Next, divide row 2 by -3 and row 3 by 3:

$$\begin{bmatrix} 1 & 2 & 3 & 1 \\ 0 & 1 & \frac{5}{3} & \frac{1}{3} \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

Subtract 2 times the first column from the second, 3 times the first column from the third, and one time the first column from the fourth:

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & \frac{5}{3} & \frac{1}{3} \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

The final step consists in subtracting  $5/3$  times the second column from the third and  $1/3$  times the second column from the fourth. The result is

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

**Theorem B.1.** *Let  $A$  be an  $m \times n$  matrix of rank  $r$ . Then a finite number of elementary row and column operations transforms  $A$  into the  $m \times n$  matrix  $B$ , where  $b_{ij} = 0$  if  $i \neq j$ ,  $b_{ii} = 1$  if  $i \leq r$ , and  $b_{ii} = 0$  if  $i > r$ .*

*Proof.* If  $A$  has rank zero, then it is the zero matrix, and it is already in the desired form. So assume  $A \neq 0$ . The proof will be by induction on the number  $m$  of rows of  $A$ . Suppose  $m = 1$ . By assumption, there exists a nonzero element in  $A$ . By interchanging columns if necessary, we may assume it appears in the first column, so that  $A = [a_1 \ a_2 \ \dots \ a_n]$ , with  $a_1 \neq 0$ . Multiplying the first column by  $1/a_1$ , and subtracting  $a_i$  times the first column from the  $i$ -th one for  $i = 2, \dots, n$  establishes the claim. Next, assume the theorem holds for any matrix with  $k$  rows, and let  $A$  be a  $(k + 1) \times n$  matrix. Since  $A \neq 0$ , some  $a_{ij} \neq 0$ . Interchange rows 1 and  $i$ , and columns 1 and  $j$  to move  $a_{ij}$  into the first row and first column. Divide the first row (or column) by  $a_{ij}$  to obtain a matrix with 1 in the first row and column. Now subtract  $a_{1j}$  times the first column from the  $j$ -th one for  $j = 2, \dots, n$ , and subtract  $a_{i1}$  times the first row from the  $i$ -th one for  $i = 2, \dots, k + 1$ . The net result is now a matrix of the form

$$\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & B & \\ 0 & & & \end{bmatrix},$$

where  $B$  is  $k \times (n - 1)$ . Finally, applying the induction hypothesis to  $B$ , one can transform it by means of finitely many elementary operations into a diagonal matrix of the specified form.  $\square$

**Corollary B.1.** *Every isomorphism  $L : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a finite composition of elementary transformations.*

*Proof.* Let  $A$  denote the matrix of  $L$  with respect to the standard basis, so that  $L = L_A$  (see Example 1.2.3). Then  $A$  is invertible and can be transformed into the identity

matrix  $I_n$  by means of a finite sequence of elementary row and column operations. Thus, by Proposition B.1, there exist elementary matrices  $E_1, \dots, E_k$  and  $E'_1, \dots, E'_l$  such that

$$E_1 E_2 \cdots E_k A E'_1 E'_2 \cdots E'_l = I_n.$$

Multiply both sides of the above identity on the left by  $E_k^{-1} \cdots E_1^{-1}$  and on the right by  $E'_l{}^{-1} \cdots E'_1{}^{-1}$  to obtain

$$A = E_k^{-1} \cdots E_1^{-1} E'_l{}^{-1} \cdots E'_1{}^{-1}.$$

Since the inverse of an elementary matrix is elementary,  $A$  equals a product of elementary matrices. But then

$$L_A = L_{E_k^{-1}} \circ \cdots \circ L_{E'_1{}^{-1}}$$

equals a composition of elementary transformations, as claimed.  $\square$

# Bibliography

- [1] N. Aders, “Amsler-Polarplanimeter-2”, Own Work. Licensed under Creative Commons Attribution-Share Alike 3.0 via Wikimedia Commons.
- [2] T. Apostol, *Mathematical Analysis*, second edition, Addison-Wesley (1974).
- [3] G. Ascoli, “Vedute sintetiche sugli strumenti integratori”, *Rend. Sem. Mat. Fis. Milano* **18** (1947).
- [4] L. Auslander, *Differential Geometry*, Harper & Row (1967).
- [5] M. Berger, B. Gostiaux, *Géométrie Différentielle*, Armand Colin (1972).
- [6] D. Burago, Y. Burago, S. Ivanov, *A course in Metric Geometry*, Graduate Studies in Mathematics **33**, American Mathematical Society (2001).
- [7] M. do Carmo, *Differential geometry of curves and surfaces* (translated from the Portuguese), Prentice-Hall (1976).
- [8] M. do Carmo, *Differential forms and applications* (translated from the 1971 Portuguese original), Springer-Verlag (1994).
- [9] M. Hirsch and S. Smale, *Differential Equations, Dynamical Systems, and Linear Algebra*, Academic Press (1974).
- [10] O. Knill, “The planimeter and the theorem of Green”, URL: <http://www.math.harvard.edu/~knill/teaching/math21a2000/planimeter>
- [11] J. Milnor, *Topology from the Differentiable Viewpoint*, University Press of Virginia (1965).
- [12] M. Spivak, *Calculus on Manifolds*, Addison-Wesley (1965).
- [13] M. Spivak, *A Comprehensive Introduction to Differential Geometry*, Volumes I–V, Publish or Perish (1975–1979).
- [14] S. Sternberg, *Lectures on Differential Geometry*, Prentice-Hall (1964).
- [15] D. Struik, *Differential Geometry*, Addison-Wesley (1950).
- [16] J. Thorpe, *Elementary Topics in Differential Geometry*, Springer-Verlag (1979).



# Index

- absolute convergence 41
- acceleration 101
- adjoint representation 172
- admissible cover 185
- affine hyperplane 26
- alternating map 13
- alternating tensor 224
- Ampère's law 248
- associated linear map 19
- atlas 117
  
- basis 4
  - positively oriented 228
- binormal 331
- boundary 31, 241
- boundary point 31
- bounded set 31
- box 32
- bundle projection 142, 169
  
- cardioid 217
- Cartan's lemma 231
- Cartesian product 121
- cartesian product 1
- catenoid 174
- Cauchy sequence 40
- Cauchy-Schwarz inequality 20
- Cauchy-Schwarz inequality 23, 169
- centroid 213
- characteristic function 183
- characteristic polynomial 86
- chart 117
- circle of curvature 315
- class  $C^\infty$  65
- class  $C^k$  65
- closed form 245, 252
- closed set 30
- closure 31
- cohomology space 253
- compact set 32
- complete metric space 285
- complete set 53
- complete vector field 135
- cone 168, 176
- conjugate point 274
- conjugation 172
- connected 53, 283
- conservative force 252
- constant curvature 159
- constant curvature space 159
- content zero 215
- continuous 43
- contractible 254
- contraction 75
- convergence 36
- convex 49, 73, 296
  - strongly 296
- convex function 334
- convex hypersurface 318
- coordinate 1
- coordinate plane 191
- coordinate vector 6
- coordinate vector fields 130
- countable 341
- covariant derivative 99, 143
- cover 108
  - locally finite 108
  - subordinate 109
- critical point 89, 195
- critical value 195
- cross product 27, 258, 303
- curl 246, 259
- curvature
  - mean 174
  - negative 308
  - nonnegative 308
  - nonpositive 308
  - positive 308
  - sectional 159
- curvature of a curve in the plane 315
- curvature tensor 156
- curve 58, 69
  - length of 70
  - normal 73
  - parametrized by arc length 73
  - regular 73
- cut locus 295
  - tangential 295
- cylindrical coordinates 203

- decomposable 229
- dense 343
- density 211
- derivative 57, 94, 127
- determinant 15
  - of linear transformation 18
- diameter 53, 290, 295
- diffeomorphism 75, 127
- differentiable 57, 124
- differential 127
- differential form 232
  - integral of 237
- dimension 5
- direct sum 23
- directional derivative 101
- directrix 321
- distance 30, 283
- divergence 247, 258
- double tangent space
  - canonical inner product on 268
- dual basis 25
- dual space 24
  
- eigenspace 84
- eigenvalue 84
- eigenvector 84
- elementary matrix 345
- elementary transformation 345
- ellipsoid 313, 335
- Enneper's surface 338
- equivalence relation 72
- Euclidean motion 163
- exact form 245, 252
- exponential 47
- exponential map 151
- exterior derivative operator 234
- exterior product 225
- extremum 89
  
- first fundamental tensor 154
- flat 157
- flow 97, 135
  
- gamma function 217
- Gauss map 304
- Gaussian curvature 308
- geodesic 149
- geodesic spray 150
- glide rotation 322
- gradient 101, 258
  
- graph 41
- greatest lower bound 339
- Green's identities 264
- group 137
  - Lie 137
  
- harmonic 264
- helicoid 322
- Hessian 91
- Hodge star operator 258
- homeomorphism 52, 184
- homomorphism 139
- homotopic 254
- homotopy
  - equivalent 254
- hyperboloid 313, 322, 335
- hyperplane 26, 304
- hypersurface 117
  
- imbedding 113
- immersed submanifold 117
- immersion 113
- implicit function theorem 79
- improper integral 209
- incompressible 248
- infimum 35, 339
- injectivity radius 285
- inner product 20, 21
- integrable function 179
- integral 179
  - of a form 237
  - lower 178
- integral curve 134
- interior (of a set) 31
- interior point 30
- intermediate value theorem 54
- inverse 2, 8
- irrotational 265
- isometric 83
- isometric map 163
- isometry 163, 261
- isomorphism 6
- iterated integral 188
  
- Jacobi identity 105
- Jacobian matrix 59
- Jordan-measurable 183
  
- k-form 224
  - decomposable 229



- Killing field 173
- kinetic energy 239
- Kronecker delta 21, 83
  
- Lagrange multipliers 128
- Laplacian 264, 309
- least upper bound 32, 339
- left multiplication 8
- left translation 138
- left-invariant 138
- level set of a function 101
- Lie algebra 105
  - isomorphic 115
- Lie bracket 103
- Lie derivative 259
- Lie group 137
- limit (of a map) 41
- limit point 39, 41
- line 294
- line of curvature 336
- linear combination 4
- linear dependence 4
- linear independence 4
- linear isometry 50, 82
- linear transformation 6
  - nullity 11
  - rank 11
  - associated 19
- local flow 97
- locally finite 108
- locally symmetric space 298
- lower bound 339
- lower integral 178
- lower sum 177
  
- Möbius strip 302, 322
- manifold 117
  - with boundary 240
  - connected 283
- matrix 2
  - change of basis 10
  - identity 8
  - invertible 8
  - orthogonal 55, 83
  - product 7
  - skew-symmetric 3, 55
  - symmetric 3
- maximum 89, 339
- maximum principle 264
- Maxwell’s equations 248
  
- mean curvature 174, 308
- mean value theorem 111, 264
- measure zero 181
- meridians 323
- metric space 30
  - complete 41
- minimal geodesic 270
- minimal surface 174
- minimum 339
- moment of inertia 214
- monkey saddle 337
- moving frame 260
- multilinear 153
- multilinear map 13
- musical isomorphisms 24
  
- negative curvature 308
- nonnegative curvature 308
- nonpositive curvature 308
- norm 30
- normal geodesic 149
- normal space 143
- normal vector 25
  
- one-form 24
- one-parameter group 116
- open cover 32
- open map 34
- open set 30
- operator 47, 82
  - self-adjoint 82
  - skew-adjoint 82
- operator norm 21
- orientable 301
  - manifold 236
- orientation 228, 236, 302
  - induced on  $\partial M$  242
  - standard on  $\mathbb{R}^n$  228
- orientation-preserving 237
- orthogonal 21
- orthogonal complement 23
- orthogonal group 112
- orthogonal projection 24
- orthogonal transformation 50
- orthonormal basis 21
- outward unit normal 242
  
- parabolic umbilic 337
- paraboloid 335
- parallel translation 145

- parallel vector field 145
- parallelizable 170
- parallels 323
- parametric equations of line 26
- parametrization 117
- partial derivative 58
- partial sum 41
- partition 177
- partition of unity 109
- path connected 54, 73
- permutation 12
  - even 13
  - sign of 13
- piecewise smooth
  - curve 69
  - vector field 277
  - variation 277
- planimeter 249
- polar coordinates 77, 202
- pole 294
- polynomial 42, 112
- position vector field 96
- positive curvature 308
- positive definite (bilinear form) 87
- positively oriented basis 228
- pre-image 44
- principal curvatures 308
- principal normal 331
- projectable set 191
- pseudosphere 325
- pullback 223
  
- quadratic form 112
- quadric 312, 335
- quotient space 253
  
- rank 79
- rational function 44
- ray 294, 327
- refinement 177
- regular value 121
- related vector fields 105, 130
- remainder 90
- reparametrization 72
- restriction 121
- Riemann sum 210
- Riemannian homogeneous space 294
- Riemannian manifold 261
- Riemannian metric 223
- right circular cone 204
  
- rigid motion 313
- ruled surface 320
- rulings 320
  
- saddle point 92
- scalar product 87
- Scherk's surface 338
- second fundamental form 156
- second fundamental tensor 155
- second partial derivatives 65
- sectional curvature 159
- self-adjoint operator 82
- separation 53
- sequence 36
  - geometric 38
  - bounded 36
- Serret-Frénet formulas 331
- similar matrices 10
- sink 248
- skew-adjoint operator 82
- skew-symmetric tensor 224
- smooth function 65
- source 248
- span 4
- special linear group 112
- speed 58
- spherical coordinates 78, 130, 207
- stereographic projection 119
- strongly convex 296
- subordinate 109
- subsequence 39
- subspace 3
- support 67
- supremum 35, 339
- surface 117
- symmetric bilinear form 87
  - semi-definite 309
  - positive definite 87
- symmetric tensor 224
  
- tangent bundle 141
- tangent space 94, 122
- tangent vector 94
- tangential cut locus 295
- tensor 154, 221
  - first fundamental 154
  - second fundamental 155
- tensor field 154, 222
- tensor product 221
- topologist's sine curve 52, 73

- torque 211
- torus 174
  - flat 174
- totally geodesic 167, 307
- totally umbilic 307
- trace 50
- tractrix 325
- transposition 12
- triangle inequality 30
  
- umbilical 307
- uniformly continuous 45
- unit tangent sphere bundle 295
- upper bound 32, 339
- upper half-space 240
- upper integral 178
- upper sum 177
  
- variation of a curve 267
- variational vector field 272
- vector 2
  - coordinate 6
  
- vector field 95
  - along a map 96, 142
  - complete 98, 135
  - coordinate 130
  - divergence of 258
  - geodesic 173
  - Jacobi 271
  - Killing 173
  - left-invariant 138
  - parallel 145
  - related 130
- vector space 2
  - isomorphic 6
- velocity field 96
- velocity vector 58
- volume 29, 177, 183, 187, 239
- volume form 229, 239
  
- wedge product 225
- work 239
  
- zero section 151



This book offers an introduction to differential geometry for the non-specialist. It includes most of the required material from multivariable calculus, linear algebra, and basic analysis. An intuitive approach and a minimum of prerequisites make it a valuable companion for students of mathematics and physics.

The main focus is on manifolds in Euclidean space and the metric properties they inherit from it. Among the topics discussed are curvature and how it affects the shape of space, and the generalization of the fundamental theorem of calculus known as Stokes' theorem.

- ▶ A thorough introduction to differential geometry
- ▶ Covers all important concepts and theorems
- ▶ Many examples including applications to physics



**Gerard Walschap**

is a professor of mathematics at the University of Oklahoma, and specializes in differential geometry.



[www.degruyter.com](http://www.degruyter.com)

ISBN 978-3-11-036949-6