

Mathematics 1b

Functions of several variables

Text by Ole Christensen & Jakob Lemvig

Found a typo? Write me an email: jakle@dtu.dk

All figures by Nikoline Mai Bøgely Rehn

LaTeX layout: MatNat Compendium

Preface

Calculus of several variables is a fundamental tool in mathematics that helps us understand and quantify how things change and accumulate in a multi-dimensional context. This branch of calculus extends the concepts of differentiation and integration beyond the familiar single-variable scenarios, allowing us to explore more complicated systems. Differentiation in several variables, for instance, helps us determine the rate at which quantities change in multiple directions, providing critical insights into phenomena such as the changing rate of a growing plant in different environments. Integration, on the other hand, offers a way to aggregate or sum up effects over a range, like calculating the total area under a curve but in multiple dimensions. These mathematical concepts are instrumental in a variety of practical applications, from optimizing the performance of neural networks in AI to understanding the dynamics of force fields in physics.

Contents

Preface	ii
Contents	iii
List of Figures	vi
List of Tables	viii
0 Preliminaries	1
0.1 Scalars	1
0.2 Sets	2
0.3 Functions $f : A \rightarrow B$	2
0.4 Matrices and vectors	4
1 Functions of Several Variables	8
1.1 Scalar functions of one variable	9
1.2 Scalar functions of several variables	10
1.3 Vector functions of several variables	12
1.4 Visualizing functions	15
1.5 What is multivariate calculus useful for?	22
2 Inner Product Spaces and the Spectral Theorem	23
2.1 Inner product spaces	24
2.2 Open and closed sets	36
2.3 Projections onto a line	40
2.4 Orthonormal basis	42
2.5 The Gram-Schmidt process	45
2.6 Unitary and orthogonal matrices	49
2.7 Diagonalizable matrices	53

2.8	The Spectral Theorem	54
2.9	Positive definite and semi-definite matrices	62
3	Continuity and Differentiability	65
3.1	Analysis of functions of one variable	65
3.2	Continuity of vector functions of several variables	71
3.3	Partial derivatives of first order and the gradient vector	73
3.4	Directional derivatives	78
3.5	Partial derivatives of second order and the Hessian matrix	80
3.6	Differentiability of scalar functions of several variables	83
3.7	The chain rule for scalar functions	88
3.8	Differentiability of vector functions of several variables	90
4	Taylor Approximation	94
4.1	The tangent for a function of one variable	95
4.2	Taylor polynomials for functions of one variable	98
4.3	Taylor's formula for functions of one variable	101
4.4	The tangent plane for functions of several variables	105
4.5	Taylor polynomials for functions of several variables	107
4.6	Taylor's formula for functions of several variables	111
4.7	Taylor polynomials for vector functions of several variables	113
5	Local and Global Extrema of Functions	115
5.1	The range of functions of one variable	116
5.2	The range of functions of several variables	120
6	Integration	128
6.1	The Riemann integral of functions of one variable	128
6.2	Anti-derivatives of functions of one variable	132
6.3	The Riemann integral of functions of two variables	137
6.4	Change of variables in \mathbb{R}^2	144
6.5	Polar coordinates	149
6.6	The Riemann integral of functions of several variables	153
6.7	The Riemann integral of vector functions	159
7	Vector Fields	161
7.1	Parametric curves and surfaces	162
7.2	Line and surface integrals	167
7.3	Vector fields and gradient fields	172
7.4	Line integrals of vector fields and computing anti-derivatives	178
7.5	Surface integrals of vector fields	184

Appendices	187
A Additional Proofs	188
A.1 Proof of Taylor's theorem	188
B List of symbols	192
Bibliography	194
Index	195

List of Figures

1.1	Domain of the function in Example 1.2.2 is everything strictly <i>outside</i> the circle of radius one centered at $(1, 0)$	12
1.2	Input-output view of a vector function of n variables. The function \mathbf{f} takes $\mathbf{x} \in \mathbb{R}^n$ as input and outputs $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^k$. We have plotted both \mathbb{R}^n and \mathbb{R}^k as \mathbb{R}^2 , but we should allow ourselves to mentally think of n and k taking on larger values.	15
1.3	The graph of the function $f(x_1, x_2) = x_1^2 + x_2^2 - 5$ with domain $\text{dom}(f) = [-3, 3] \times [-3, 3]$	16
1.4	Illustration of the function $\mathbf{V}(x_1, x_2) = (-x_2/3, x_1/3)$ on $[-4, 4] \times [-4, 4]$	17
1.5	The function $f(x_1, x_2) = x_1^2 + x_2^2 + 3$ and some of its level sets (only shown for $(x, y) \in [-3, 3] \times [-3, 3]$.)	18
2.1	Illustration of the triangle inequality (2.13) in \mathbb{R}^2	28
2.2	Illustration of the Pythagorean theorem in \mathbb{R}^2	34
2.3	Orthogonal projection in \mathbb{R}^2	41
2.4	Illustration of the Gram-Schmidt process in \mathbb{R}^2	48
3.1	The graph of the function $f(x) = \text{ReLU}(x)$	68
4.1	The function $f(x) = \sin(x)$ and its first-degree Taylor polynomial at $x_0 = 0$ (dashed).	96
4.2	The function f in (4.4) and its first-degree Taylor polynomial at $x = 0.07$, shown on the interval $[0.03, 0.1]$ (dashed).	97
4.3	The functions in (4.5).	98
4.4	The function $f(x) = e^x$ and its fifth-degree Taylor polynomial at $x_0 = 1$ (dashed).	100
4.5	The function f in (4.9) and its fourth-degree Taylor polynomial at $x = 0.07$	101

4.6	The function $f(x, y) = \sin(x^2 + y^2)$ and its Taylor polynomials of first degree and second degree at the point $(0, 0)$	112
6.1	Riemann sums of the function $f(x) = x^{-1}$ on $x \in [1/2, 1]$	130
6.2	The domain B in (6.48).	143
6.3	The domain B in (6.52).	145
6.4	The subset $B \subset \mathbb{R}^2$ given by Equation (6.62) on page 151 with $\alpha = \frac{\pi}{6}$, $\beta = \frac{4\pi}{3}$, $\phi_1(\theta) = 5$ and $\phi_2(\theta) = \sqrt{\frac{180\theta}{\pi}}$	151
6.5	Illustration of cylindrical and spherical coordinates of a point in \mathbb{R}^3	158
7.1	The cylinder as in Equation (7.13) with $r = 1$ and $h = 4$	170
7.2	The graph surface associated with the function h in (7.17).	171

List of Tables

CHAPTER 0

Preliminaries

In this initial chapter, we recall some concepts from calculus and linear algebra that we be used throughout the book. The chapter is designed as a reference section to acquaint you with the essential notations and terminology, mainly from linear algebra, used throughout this book. It is meant to be consulted as needed, rather than read sequentially. Use it to clarify concepts and to familiarize yourself with the language of vector spaces and matrices.

0.1 Scalars

We denote the real numbers by \mathbb{R} and the complex numbers by \mathbb{C} . We use \mathbb{F} as a common symbol for either \mathbb{R} or \mathbb{C} . It is possible to consider other fields \mathbb{F} , e.g., the rational numbers \mathbb{Q} , but in this book \mathbb{F} always denote \mathbb{R} or \mathbb{C} . The numbers in \mathbb{F} are called scalar. Hence, a *scalar* is always either a real or complex number.

Let $c = c_1 + ic_2$ be a complex scalar, where $c_1 = \operatorname{Re} c$ and $c_2 = \operatorname{Im} c$. Let us recall some rules (where d is another complex scalar):

- (i) $c\bar{c} = c_1^2 + c_2^2 = |c|^2$
- (ii) $\overline{c+d} = \bar{c} + \bar{d}$ and $\overline{cd} = \bar{c}\bar{d}$
- (iii) $c + \bar{c} = 2c_1 = 2\operatorname{Re} c$
- (iv) $(c-d)^2 = c^2 + d^2 - 2cd$ and $(c-d)(c+d) = c^2 - d^2$

0.2 Sets

Recall that a set A is a way to “bundle” elements together in one object. If we for example want to write down a set consisting of the numbers 0 and 1, we simply write $\{0, 1\}$ (or $\{1, 0\}$ or more complicated $\{1, 0, 1, 1\}$).

We often construct sets using set-builder notation. The set of all even integers, expressed in set-builder notation, is given by:

$$\{ n \in \mathbb{Z} \mid \exists k \in \mathbb{Z}, n = 2k \}, \quad (1)$$

where \exists reads “there exists”. In set-builder notation the vertical bar \mid is a separator and is read “such that”. Hence, Equation (1) is read “integers n such that there exists an integer k for which $n = 2k$ ”. One often uses more complicated expressions left of the vertical bar \mid to obtain simpler set-builder notation, e.g., the even integers can be written $\{ 2k \mid k \in \mathbb{Z} \}$.

Sometimes the vertical \mid is replaced with a colon $:$, but the meaning is the same and we will only use it if it increases the readability, e.g., in case the variable or the predicate contains other vertical bars such as:

$$\{ \mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| \leq 1 \}$$

where $\|\mathbf{x}\|$ denotes the norm of a vector (it will be introduced in Equation (2.1) on page 24).

0.3 Functions $f : A \rightarrow B$

Definition 0.3.1 Function

Let A and B denote two sets. A *function* is a correspondence, which to each element $x \in A$ in a unique way associates an element $y \in B$. Denoting the function by f , this will be written as

$$f : A \rightarrow B, y = f(x). \quad (2)$$

The set A is called the *domain* of the function f and will be denoted by $\text{dom}(f)$. The set B is called the *codomain* set and will be denoted by $\text{co-dom}(f)$.

A function can be given in terms of an explicit expression, e.g., we can define the quadratic function by

$$f : \mathbb{R} \rightarrow \mathbb{R}, f(x) = x^2. \quad (3)$$

However, a function can also be given in terms of, e.g., a table, which to each element $x \in A$ simply lists the corresponding element $y = f(x) \in B$.

Note that f is the *function* while $f(x)$ is the *function value*, i.e., a real number (scalar) if $B = \mathbb{R}$. A function f consists of a set of inputs (the domain $x \in A$), a set within all outputs occur (the codomain $y \in B$), and a rule for assigning each input to *exactly one* output, which we write as $x \mapsto y$. Hence, we can, alternatively, define the quadratic function as

$$f : \mathbb{R} \rightarrow \mathbb{R}, f = x \mapsto x^2,$$

where $x \mapsto x^2$ reads “ x (the input) maps to x^2 (the output)”.

Two functions $f : A \rightarrow B$ and $g : C \rightarrow D$ are equal precisely if they have the same domain, the same codomain, *and* they assign the same values to each of the elements of their domain. Hence, if we define $g : [0, 1] \rightarrow \mathbb{R}$, $g = x \mapsto x^2$, then quadratic functions f and g are not the same function even though the expressions $f(x) = x^2 = g(x)$ agree. This is because $A = \mathbb{R}$ while $C = [0, 1]$ in this situation, in particular, $f(2) = 4$ while $g(2)$ is not defined.

The definition of a function does not require that A and B are sets of scalars — they can be arbitrary sets. For example, the correspondence which to each item in a shop associates its brand, is a function. Here the set A is the collection of all items in the shop, and B is the collection of all brands. Note that we indeed for each item in the shop can identify a unique brand. On the other hand, different items can very well be associated with the same brand.

Note that the concept of the codomain B should be clearly distinguished from the range of the function f . Indeed, the set B in Definition 0.3.1 is just a set, within which all function values occur! We do *not* claim that each of the elements in B occur as a function value. On the other hand, the *range* or *image* of a function f is defined as precisely the set of all function values $f(x)$, where x runs through the set A :

$$\text{im}(f) = \{ f(x) \mid x \in A \}. \quad (4)$$

The $\text{im}(f)$ is sometimes written as $f(A)$. For example, for the function f in (3), the codomain is \mathbb{R} , while the range is the set of nonnegative real scalars, i.e., the interval $[0, \infty[$.

Definition 0.3.2 Surjective, Injective, Bijective

Consider a function $f : A \rightarrow B$.

- (i) The function f is said to be *surjective* if the range of f equals the set

B ; that is, for each $y \in B$ there exists some $x \in A$ such that $f(x) = y$.

- (ii) The function f is said to be *injective* if different elements in the set A get mapped to different elements in B ; that is, if

$$x_1, x_2 \in A, x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2).$$

- (iii) The function f is said to be *bijective* if f is as well surjective as injective.

Example 0.3.1

Let us illustrate the concepts in Definition 0.3.2 with a number of examples.

- (a) The function $f : \mathbb{R} \rightarrow \mathbb{R}, f(x) = x^2$ is neither injective nor surjective.
- (b) The function $f : \mathbb{R} \rightarrow \mathbb{R}, f(x) = x^3$ is injective and surjective, hence bijective.
- (c) The function $f : \mathbb{R} \rightarrow \mathbb{R}, f(x) = e^x$ is injective, but not surjective.

0.4 Matrices and vectors

An $m \times n$ *matrix* is an array of real or complex numbers of the form

$$A = [a_{ij}] = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix},$$

where $a_{i,j} \in \mathbb{R}$ or \mathbb{C} is the (i, j) entry of the matrix A . We sometimes write a_{ij} (without a comma) instead of $a_{i,j}$. The set of all $m \times n$ matrices with entries in \mathbb{F} is denoted by $\mathbf{M}_{m \times n}(\mathbb{F})$. Real $m \times n$ matrices are then denoted $\mathbf{M}_{m \times n}(\mathbb{R})$ and complex $m \times n$ matrices are denoted $\mathbf{M}_{m \times n}(\mathbb{C})$. The set of square matrices, i.e., $m = n$, are denoted $\mathbf{M}_n(\mathbb{F}) := \mathbf{M}_{n \times n}(\mathbb{F})$.

The $m \times n$ zero matrix is denoted by $\mathbf{0}_{m \times n}$. The zero (column) vector $\mathbf{0}_{n \times 1}$ is denoted by $\mathbf{0}_n$ or simply $\mathbf{0}$. The $n \times n$ identity matrix is denoted I_n or just I .

We recall some properties and terminology of matrix theory.

Inverse. We say that $A \in M_n$ is an *invertible matrix* if there exists a $B \in M_n$ such that

$$AB = I_n \quad \text{and} \quad BA = I_n. \quad (5)$$

Such a matrix B is an inverse of A . If A has no inverse, then A is non-invertible or singular. Either of the equalities in (5) implies the other, that is, for *square*¹ matrices $A, B \in M_n$, we have $AB = I$ if and only if $BA = I$.

As we know not every square matrix has an inverse, but a matrix has at most one inverse. If A is invertible, we can therefore speak of *the* inverse of A . If A is invertible, then the inverse of A is denoted by A^{-1} . A matrix and its inverse satisfy:

$$AA^{-1} = I = A^{-1}A.$$

If A is invertible, we define $A^{-k} = (A^{-1})^k$ for $k = 1, 2, \dots$. Let A and B be matrices so that AB is well-defined, let j, k be integers, and let c be a scalar.

- (i) $(A^k)^j = A^{jk} = (A^j)^k$.
- (ii) $(A^{-1})^{-1} = A$.
- (iii) If $c \neq 0$, then $(cA)^{-1} = c^{-1}A^{-1}$.
- (iv) If $A \in M_{m \times n}(\mathbb{F})$ and $B \in M_{n \times m}(\mathbb{F})$, then $(AB)^{-1} = B^{-1}A^{-1}$ if AB is invertible (this requires $m \leq n$).

Transpose. The transpose of $A = [a_{ij}] \in M_{m \times n}$ is the matrix $A^T \in M_{n \times m}$ whose (i, j) entry is a_{ji} . Let A and B be matrices of appropriate sizes and let c be a scalar.

- (i) $(A^T)^T = A$.
- (ii) $(A + B)^T = A^T + B^T$.
- (iii) $(cA)^T = cA^T$.
- (iv) $(AB)^T = B^T A^T$.
- (v) If A is invertible, then $(A^T)^{-1} = (A^{-1})^T$. We write $(A^{-1})^T = A^{-T}$.

Conjugate. The conjugate of $A \in M_{m \times n}$ is the matrix $\bar{A} \in M_{m \times n}$ whose (i, j) entry is \bar{a}_{ij} , the complex conjugate of a_{ij} . Thus,

$$\overline{(\bar{A})} = A, \quad \overline{A + B} = \bar{A} + \bar{B}, \quad \overline{AB} = \bar{A}\bar{B}.$$

If A is a real matrix, then $\bar{A} = A$.

¹Be careful! This statement is false for non-square matrices.

Adjoint (conjugate transpose). The adjoint of A is defined by $A^* = \overline{A^T}$, meaning that we take the transpose of the matrix, and then take the complex conjugate of each entry (or the other way around; the order does not matter which you can easily check). Hence, the (i, j) entry of A^* is $\overline{a_{j,i}}$. Note, that for a real matrix A taking transpose and adjoint is the same operation, i.e., $A^* = A^T$. The adjoint of a matrix is also called the conjugate transpose since this is what the operation does. Let A and B be matrices of appropriate sizes and let c be a scalar.

- (i) $I^* = I_n$.
- (ii) $(0_{m \times n})^* = 0_{n \times m}$.
- (iii) $(A^*)^* = A$.
- (iv) $(A + B)^* = A^* + B^*$.
- (v) $(cA)^* = \bar{c}A^*$.
- (vi) $(AB)^* = B^*A^*$.
- (vii) If A is invertible, then $(A^*)^{-1} = (A^{-1})^*$. We write $(A^{-1})^* = A^{-*}$.

Special Types of Square Matrices. Let $A \in M_n(\mathbb{C})$.

- (i) If $A^* = A$, then A is *Hermitian*.
- (ii) If $A^T = A$, then A is *symmetric*.
- (iii) If $A^*A = I$, then A is *unitary*; if A is real and $A^T A = I$, then A is *real orthogonal*.
- (iv) If $A^*A = AA^*$, then A is *normal*.
- (v) If $A^2 = A$, then A is *idempotent*. . An projection matrix is an idempotent matrix.
- (vi) If A is Hermitian and has (strictly) positive eigenvalues, then A is *positive definite*.
- (vii) If A is Hermitian and has non-negative eigenvalues, then A is *positive semi-definite*.
- (viii) If A is idempotent and Hermitian, then A is an *orthogonal projection matrix*.

Trace The *trace* of $A = [a_{ij}] \in \mathbf{M}_n$ is the sum of the diagonal entries of A :

$$\operatorname{tr} A = \sum_{i=1}^n a_{ii}.$$

Let A and B be matrices of appropriate sizes and let c be a scalar.

- (i) $\operatorname{tr}(cA + B) = c \operatorname{tr} A + \operatorname{tr} B$.
- (ii) $\operatorname{tr}(A^T) = \operatorname{tr} A$.
- (iii) $\operatorname{tr}(cA) = c \operatorname{tr} A$.
- (iv) $\operatorname{tr}(A^*) = \operatorname{tr} A$.

CHAPTER 1

Functions of Several Variables

The concept of a function is probably the most fundamental topic in mathematical analysis. It provides the language for formulating basic laws in physics and engineering in a mathematical language and hereby calculate how we can obtain a desired effect in a given physical/mechanical system. Questions like

What is the necessary temperature of a gas in a tank that is necessary in order for the pressure in the gas to reach a certain value?

can easily be phrased and solved using functions.

We expect the basic concepts of functions $f : A \rightarrow B$ in [Section 0.3](#) on page 2 to be known to the reader. Standard calculus is concerned with the case where the domain of the function [Equation \(2\)](#) on page 2 is the set of real scalars, like in [\(3\)](#), or a subset hereof. Our goal in the current chapter is to generalize our knowledge of such functions to the case where A and B are sets of (column) vectors in \mathbb{R}^n and \mathbb{R}^k , respectively. For example, this will allow us to consider real-valued functions of the form

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = x + y^2 \tag{1.1}$$

and vector-valued functions of the form

$$f : \mathbb{R} \rightarrow \mathbb{R}^3, \quad f(t) = \begin{bmatrix} \cos(t) \\ \sin(t) \\ t \end{bmatrix}.$$

In [Section 1.2](#) we introduce scalar (also called real-valued) functions depending on several real variables, and in [Section 1.3](#) we turn to the general case of vector-valued functions of several variable. In order to get a good visual understanding of a function of several variables, we need the so-called level sets which is the topic of [Section 1.4](#).

1.1 Scalar functions of one variable

Functions depending on a single real variable are well-known from elementary calculus. In this section we will remind the reader about some of the most important concepts and results. We consider functions $f : I \rightarrow \mathbb{R}$, where either $I = \mathbb{R}$ or I is an open interval of the form $] - \infty, a[$, $]a, \infty[$ or $]a, b[$ for some $a, b \in \mathbb{R}$.

We already saw the first example of such a function in (3). In general, when dealing with a function, it is important to specify as well the expression for the function as its domain. For example, the natural logarithm $f(x) = \ln(x)$ is only defined for positive values of the variable x ; thus, the function must be described as

$$f : \text{dom}(f) = \{x \in \mathbb{R} \mid x > 0\} \rightarrow \mathbb{R}, \quad f(x) = \ln(x).$$

From elementary calculus we are already familiar with several classes of functions, e.g., polynomials, exponential functions, and trigonometric functions.

Example 1.1.1

Functions also appear naturally within all branches of science and engineering, typically denoted by different letters than in the mathematical literature:

- (a) Consider an (ideal) gas that is contained in a tank with volume V . Assuming that the gas contains n particles (measured in moles), the pressure P in the tank and the temperature T of the gas are related by $PV = nRT$, where R is the ideal gas constant (approximately 0.0821.) This can also be written as

$$P = \frac{nRT}{V}. \quad (1.2)$$

If the volume V of the tank is constant, (1.2) shows how the pressure in the tank varies with changing temperature. This means precisely that the pressure P is a function of the temperature T , i.e., we can write

$$P(T) = \frac{nRT}{V}.$$

- (b) Consider an electric circuit, where a battery is connected with a resistant R . By Ohm's law, the voltage V delivered by the battery and the current I in the circuit are related by $V = RI$. Thus $I = V/R$. Considering the voltage V to be constant, this relationship shows how

1.2. Scalar functions of several variables

the current I depends on the resistant R . In other words, the current is a function of R ,

$$I(R) = \frac{V}{R}.$$

Remark 1.1.1

Note that when introducing the functions in [Example 1.1.1](#), we have deliberately been sloppy. The problem is that we have only specified the function expression of, e.g., the pressure $P(T) = nRT/V$ in the tank, but not the domain and codomain of the function. Both the domain and codomain is part of the definition of the function so here we should specify $P : [0, \infty[\rightarrow \mathbb{R}$ (or perhaps we should even replace the domain $[0, \infty[$ with $[0, T_P]$ where $T_P \approx 1.416 \cdot 10^{32}$ is the Planck temperature measured in Kelvin).

1.2 Scalar functions of several variables

In this section we will take the first steps to generalize the theory for functions of one variable to functions of several variables. The desire to do this comes from science and engineering, where we often face the need of considering several changing parameters simultaneously. Let us illustrate this with an example.

Example 1.2.1

We consider again the two examples in [Example 1.1.1](#) on the previous page:

- (a) The pressure P in the tank is expressed by the formula (1.2). If we imagine that as well the volume V of the tank as the temperature T can vary, the pressure becomes a function of two variables:

$$P(T, V) = \frac{nRT}{V}.$$

- (b) If we replace the battery by a power supply with an adjustable output, the current I in the circuit becomes a function of two variables, namely, the output V of the power supply and the resistant R ,

$$I(R, V) = \frac{V}{R}.$$

We will now consider functions that are defined either on \mathbb{R}^n for some $n \in \mathbb{N}$, or a subset hereof. We saw the first example of such a function in

1.2. Scalar functions of several variables

(1.1). We will typically denote the variables by x_1, x_2, \dots, x_n ; however, for functions of just two variables we will also apply the shorter names x and y for the variables.

The following class of scalar functions of n variables called *quadratic forms* will play an important role throughout the text.

Definition 1.2.1 Quadratic form

Let $A \in M_{n \times n}(\mathbb{R})$ be a square, non-zero matrix, let $\mathbf{b} \in \mathbb{R}^n$ be a column vector, and $c \in \mathbb{R}$ a scalar. The scalar function $q : \mathbb{R}^n \rightarrow \mathbb{R}$ given by

$$q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + \mathbf{x}^T \mathbf{b} + c \quad (1.3)$$

where \mathbf{x} is considered as a column vector in \mathbb{R}^n , is called a *quadratic form*.¹

Note that if $n = 1$, the quadratic form is a polynomial function of degree two since (1.3) reduces to $q(x) = Ax^2 + bx + c$. Indeed, a quadratic form is a multivariate polynomial of degree at most two.

Determining a domain. As for functions of a single variable, it is important to be able to describe the domain for a function of n variables. Typically this is more involved than for functions of one variable as illustrated by the following example.

Example 1.2.2

We would like to describe the domain for the function expression:

$$f(x_1, x_2) = \ln(x_1^2 - 2x_1 + x_2^2). \quad (1.4)$$

The natural logarithm is only defined on positive scalars so the function f is defined on the set of $(x_1, x_2) \in \mathbb{R}^2$ for which $x_1^2 - 2x_1 + x_2^2 > 0$. Now, by a direct calculation,

$$x_1^2 - 2x_1 + x_2^2 = (x_1 - 1)^2 - 1 + x_2^2,$$

so the condition $x_1^2 - 2x_1 + x_2^2 > 0$ means precisely that

$$(x_1 - 1)^2 + x_2^2 > 1.$$

Thus the domain for the function f is

$$\text{dom}(f) = \{(x_1, x_2) \in \mathbb{R}^2 \mid (x_1 - 1)^2 + x_2^2 > 1\}.$$

¹The term *quadratic form* is often used to denote the quadratic term $\mathbf{x}^T A \mathbf{x}$ (without the linear $\mathbf{x}^T \mathbf{b}$ and constant terms c).

1.3. Vector functions of several variables

Note that $\text{dom}(f)$ has a geometrical description, see Figure 1.1. It is the set of points $(x_1, x_2) \in \mathbb{R}^2$ that are located *outside* the circle in \mathbb{R}^2 centered at the point $(1, 0)$ and with radius 1.

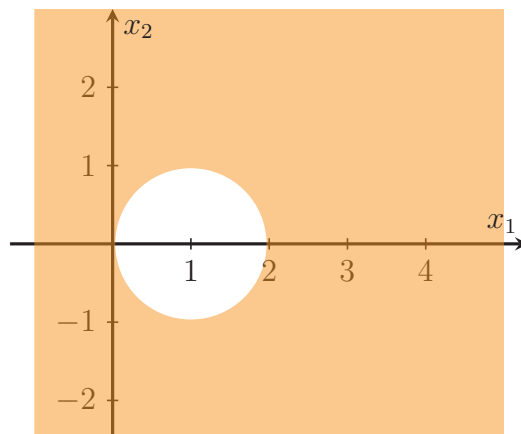


Figure 1.1: Domain of the function in Example 1.2.2 is everything strictly *outside* the circle of radius one centered at $(1, 0)$.

Note that the domain $\text{dom}(f)$ found in Example 1.2.2 is the largest possible domain of the real-valued function expression (1.4). It is perfectly fine to introduced *another* function with a smaller domain, e.g.,

$$g :]2, \infty[\times \mathbb{R} \rightarrow \mathbb{R}, \quad g(x_1, x_2) = \ln(x_1^2 - 2x_1 + x_2^2).$$

1.3 Vector functions of several variables

For functions of n variables the vector notation from linear algebra is very useful. For example, instead of writing a function value as $f(x_1, x_2, \dots, x_n)$, we can define $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and then write the function value as $f(\mathbf{x})$. Visually, this is similar to our standard writing $f(x)$ for a function of one variable, except that now the boldface \mathbf{x} indicates that the input is a vector. Note that strictly speaking (x_1, x_2, \dots, x_n) is a n -tuple real numbers. We will, however, most often consider \mathbf{x} as a column vector in \mathbb{R}^n :

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

1.3. Vector functions of several variables

or (less often) a row vector $\mathbf{x} = [x_1, \dots, x_n]$ in \mathbb{R}^n . Even when we consider \mathbf{x} as a vector, we will also, when there is no danger of confusion, use the tuple notation $\mathbf{x} = (x_1, x_2, \dots, x_n)$. Unless stated otherwise, we consider vectors in \mathbb{R}^n as column vectors.

The vector notation allow us to formulate several definitions and key results for functions of n variable in a very similar way as for functions of one variable (compare, e.g., the definition of continuity in Definition 3.1.1 on page 66 with the subsequent Definition 3.2.1 on page 71).

A vector function is simply a function that takes column vectors in \mathbb{R}^n as input and outputs column vectors in \mathbb{R}^k :

Definition 1.3.1 Vector function

Let $n, k \in \mathbb{N}$. A *vector function of several variables* is a function of the form

$$\mathbf{f}: \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k, \text{ where the domain } \text{dom}(\mathbf{f}) \text{ is a subset of } \mathbb{R}^n.$$

Since $\mathbf{f}(\mathbf{x})$ is a vector in \mathbb{R}^k for each $\mathbf{x} \in \text{dom}(\mathbf{f})$, we will write $\mathbf{f} = (f_1, \dots, f_k)$ or

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_k(\mathbf{x}) \end{bmatrix}, \quad (1.5)$$

where $f_i: \text{dom } f \rightarrow \mathbb{R}$ is the i th *coordinate function* of \mathbf{f} . In other words, the i th coordinate function is defined *from* \mathbf{f} by setting $f_i(\mathbf{x}) := (\mathbf{f}(\mathbf{x}))_i$, where $(\mathbf{f}(\mathbf{x}))_i$ is the i th coordinate of the vector $\mathbf{f}(\mathbf{x})$. The coordinate functions are *scalar* functions of n variables, and they all have domain $\text{dom}(\mathbf{f})$. A vector function $\mathbf{f} = (f_1, \dots, f_m): \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k$ can often be analyzed by considering the scalar-valued coordinate functions $f_i: \text{dom}(f) \rightarrow \mathbb{R}$ separately.

We already know examples of vector functions from basic linear algebra, namely, *linear mapping* from \mathbb{R}^n to \mathbb{R}^k .

Example 1.3.1

Let $A \in M_{k \times n}(\mathbb{R})$ be a rectangular matrix and define the linear map $L_A: \mathbb{R}^n \rightarrow \mathbb{R}^k$ by $L_A \mathbf{x} = A\mathbf{x}$, i.e., $L_A \mathbf{x} = A\mathbf{x}$. Then L_A is a vector function of n variables, and its image is the column space of the matrix A , i.e., $\text{im}(L_A) = \text{col}(A)$. The i th coordinate function of L_A is given by $(L_A)_i(\mathbf{x}) = a_{i,1}x_1 + a_{i,2}x_2 + \dots + a_{i,n}x_n$ for $i = 1, \dots, k$.

Linear mappings L_A is a subclass of *affine mappings* which are vector

1.3. Vector functions of several variables

functions of the form:

$$\mathbf{x} \mapsto A\mathbf{x} + \mathbf{b} \quad \text{for } \mathbf{x} \in \mathbb{R}^n \quad (1.6)$$

where $\mathbf{b} \in \mathbb{R}^n$.

Let us recap some terminology:

- (a) $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a *scalar* (or real) function² of n variables. This is just a vector function $f: \mathbb{R}^n \rightarrow \mathbb{R}^k$ where $k = 1$.
- (b) $f: \mathbb{R}^n \rightarrow \mathbb{R}^k$ is a *vector* function of n variables.

Let us also introduce some new terminology:

- (c) If $k = n$, the vector function $f: \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^n$ where $\text{dom}(\mathbf{f}) \subseteq \mathbb{R}^n$ is called a *vector field*.
- (d) If $k = 1$, the scalar function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is also called a scalar field (often the scalar will be a physical quantity with units).
- (e) If $n = 1$ and the (vector) function is continuous, we obtain *curves* in \mathbb{R}^k . Strictly speaking, the curve is the image $\text{im}(\mathbf{f})$ of the vector function, not the function itself. If $n = 1$ and $k \geq 3$, we obtain *space curves*, and one could think of a particle path in \mathbb{R}^k , e.g., a helix in \mathbb{R}^3 :

$$\mathbf{r}(t) = [\cos(t), \sin(t), t]^T \quad \text{for } t \in I \quad (1.7)$$

where I is usually a finite interval $[a, b]$, but could be infinite, e.g., \mathbb{R} . The vector function $\mathbf{r}: I \rightarrow \mathbb{R}^n$ is called a *parametrization* of the curve $\text{im}(\mathbf{r})$.

- (f) If $n = 2$ and the (vector) function is continuous, we obtain *surfaces* in \mathbb{R}^k . A vector function of the form $\mathbf{r}: I \times J \rightarrow \mathbb{R}^n$, where I, J are intervals, is called a *parametrization* of the surface $\text{im}(\mathbf{r})$. An example is the parametrization of the (unit) sphere in \mathbb{R}^3 by

$$\mathbf{r}(u, v) = [\cos(u) \cos(v), \sin(u) \cos(v), \sin(v)]^T \quad (1.8)$$

where $u \in [0, 2\pi[$ is the longitude and $v \in [0, \pi]$ is the latitude.

We will see that parametrizations of geometrical objects become crucial for the development of integration in higher dimensions, e.g., the computation of a flux through a surface.

²also called scalar-valued or real-valued functions

Remark 1.3.1 Complex functions

In this book we only consider functions with domain $\text{dom}(\mathbf{f}) \subseteq \mathbb{R}^n$ and with codomain $\text{co-dom}(\mathbf{f}) \subseteq \mathbb{R}^k$. We could without too much effort also allow for complex vector-valued functions, e.g., $\text{co-dom}(\mathbf{f}) \subseteq \mathbb{C}^k$. For $k = 1$ such functions are known from [Mathematics 1a] as *complex-valued functions*. Given a complex vector-valued function $\mathbf{f} : \text{dom}(\mathbf{f}) \rightarrow \mathbb{C}^k$, one can for any $\mathbf{x} \in \mathbb{R}^n$, write $\mathbf{f}(\mathbf{x}) = \mathbf{f}_1(\mathbf{x}) + i\mathbf{f}_2(\mathbf{x})$, where $\mathbf{f}_1(\mathbf{x}) = \text{Re}(\mathbf{f}(\mathbf{x}))$ is the coordinate-wise real part of the vector $\mathbf{f}(\mathbf{x})$ and $\mathbf{f}_2(\mathbf{x}) = \text{Im}(\mathbf{f}(\mathbf{x}))$ is the coordinate-wise imaginary part of vector $\mathbf{f}(\mathbf{x})$. In this way, any complex vector-valued function $\mathbf{f} : \text{dom}(\mathbf{f}) \rightarrow \mathbb{C}^k$ gives rise to two real valued-functions $\text{Re}(\mathbf{f}) : \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k$ defined as $\mathbf{x} \mapsto \text{Re}(\mathbf{f}(\mathbf{x}))$ and $\text{Im}(\mathbf{f}) : \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k$ defined as $\mathbf{x} \mapsto \text{Im}(\mathbf{f}(\mathbf{x}))$. Hence, we can reduce the analysis of a complex vector-valued function to two functions taking values in \mathbb{R}^k .

1.4 Visualizing functions

It is often helpful to have a certain “geometric understanding” of a function $\mathbf{f} : \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k$ where $\text{dom}(\mathbf{f}) \subseteq \mathbb{R}^n$. In Figure 1.2 we see the “input vectors” of the function \mathbf{f} on the left and the “output vectors” on the right.

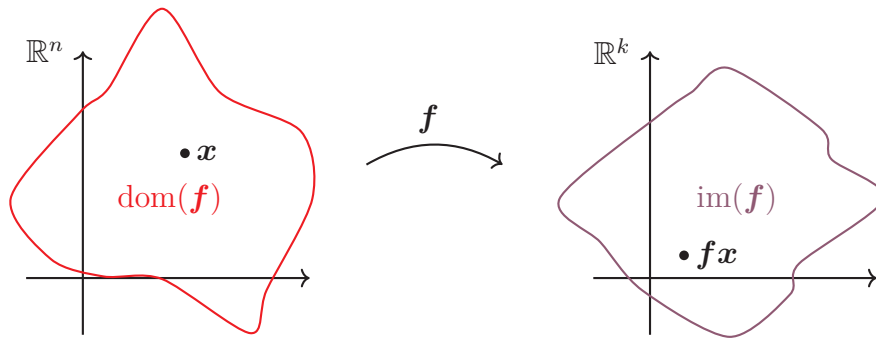


Figure 1.2: Input-output view of a vector function of n variables. The function \mathbf{f} takes $\mathbf{x} \in \mathbb{R}^n$ as input and outputs $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^k$. We have plotted both \mathbb{R}^n and \mathbb{R}^k as \mathbb{R}^2 , but we should allow ourselves to mentally think of n and k taking on larger values.

In general, it is not possible to visualize functions when k and n are large. However, for small values of k and n it is useful to study the graph and the level sets of a function.

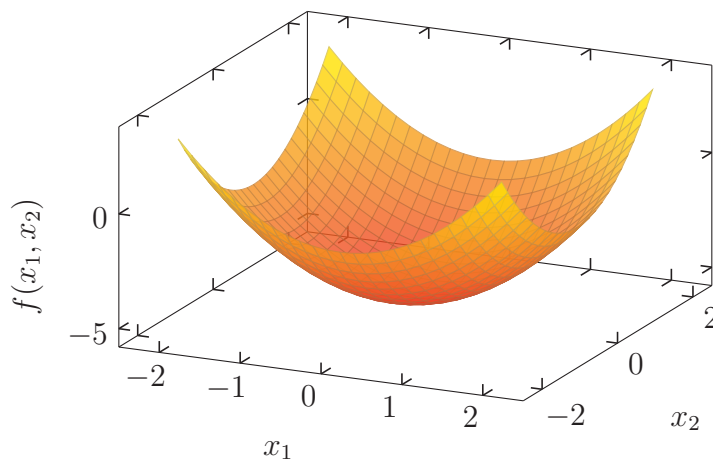


Figure 1.3: The graph of the function $f(x_1, x_2) = x_1^2 + x_2^2 - 5$ with domain $\text{dom}(f) = [-3, 3] \times [-3, 3]$.

Graph of a function. From standard calculus we are used to extract information on a function $f : \mathbb{R} \rightarrow \mathbb{R}$ by drawing or plotting its graph. Formally, the *graph* of a scalar function of one variable consists of all points in \mathbb{R}^2 of the form $(x, f(x))$, where x belongs to the domain $\text{dom}(f)$ of f .

Similarly, for a vector function $\mathbf{f} : \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k$, where $\text{dom}(\mathbf{f})$ is a subset of \mathbb{R}^n , the *graph* consists of all points in $\mathbb{R}^n \times \mathbb{R}^k = \mathbb{R}^{n+k}$ of the form $(\mathbf{x}, \mathbf{f}(\mathbf{x}))$, where $\mathbf{x} \in \text{dom}(\mathbf{f})$. We can only illustrate such a set graphically for the cases where $n + k \leq 3$. Examples of such functions are *scalar* functions of *two* variables, i.e., $n = 2$ and $k = 1$ in Definition 1.3.1.

Example 1.4.1

We consider the scalar function of two variables:

$$f : [-3, 3] \times [-3, 3] \rightarrow \mathbb{R}, \quad f(x_1, x_2) = x_1^2 + x_2^2 - 5.$$

The graph of this function is given by

$$\{(x_1, x_2, x_1^2 + x_2^2 - 5) \mid (x_1, x_2) \in [-3, 3] \times [-3, 3]\} \subset \mathbb{R}^3$$

See Figure 1.3 for an illustration of the graph.

Visualizing vector fields Vector fields $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ($n = 2, 3$) can be visualized by plotting the “output” vectors on a regular grid corresponding to the coordinates of the “input” vectors. Similar to the graph of a function,

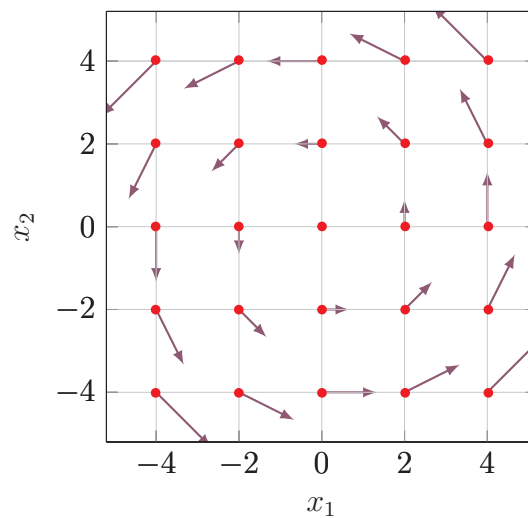


Figure 1.4: Illustration of the function $\mathbf{V}(x_1, x_2) = (-x_2/3, x_1/3)$ on $[-4, 4] \times [-4, 4]$.

this idea mixes the domain and codomain of a function into one object. Let us show the idea by an example.

Example 1.4.2

Consider the function

$$\mathbf{V} : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \quad \mathbf{V}(x_1, x_2) = (-x_2/3, x_1/3).$$

This is an example of a *vector-field*, which maps a point $\mathbf{x} \in \mathbb{R}^2$ to a point in $\mathbf{V}(\mathbf{x}) \in \mathbb{R}^2$. The new point $\mathbf{V}(\mathbf{x})$ is then interpreted as a vector originating from the point \mathbf{x} . This is illustrated in Figure 1.4. E.g., the vector-field \mathbf{V} evaluated at $(x_1, x_2) = (0, 2)$ is the vector $(0, 2/3)$, whereas at $(x_1, x_2) = (1, 4)$, it is $(-4/3, 1/3)$.

Level sets of a function. It can be difficult to get a good impression of a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ from its graph. In such cases it is often useful to consider the so-called *level sets*. They consist of points in the domain that are mapped to a pre-defined fixed vector:

Definition 1.4.1 Level set

Consider a vector function $\mathbf{f} : \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k$, where $\text{dom}(\mathbf{f}) \subseteq \mathbb{R}^n$. Given

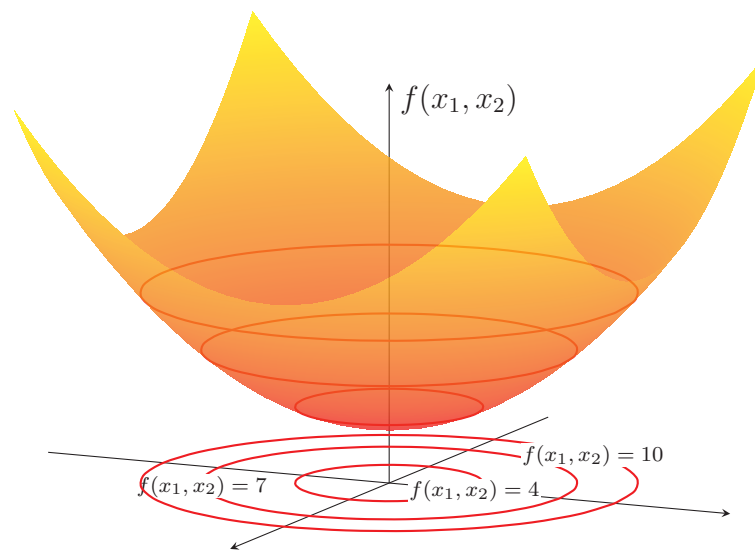


Figure 1.5: The function $f(x_1, x_2) = x_1^2 + x_2^2 + 3$ and some of its level sets (only shown for $(x, y) \in [-3, 3] \times [-3, 3]$.)

a vector $\mathbf{c} \in \mathbb{R}^k$, the corresponding level set is the subset of \mathbb{R}^n given by

$$\{\mathbf{x} \in \text{dom}(\mathbf{f}) \mid \mathbf{f}(\mathbf{x}) = \mathbf{c}\}.$$

Example 1.4.3

Consider the scalar function f

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = x_1^2 + x_2^2 + 3.$$

For $c \in \mathbb{R}$ the level sets are given by

$$\{(x_1, x_2) \in \text{dom}(f) \mid f(x_1, x_2) = c\}.$$

Since $x_1^2 + x_2^2 \geq 0$ for all (x, y) , the level set associated with any $c < 3$ is empty. For $c = 3$, the level set consists of the point $(x_1, x_2) = (0, 0)$, and for $c > 3$, it consists of the $(x_1, x_2) \in \mathbb{R}^2$ for which $x_1^2 + x_2^2 = c - 3$. Thus, for $c > 3$, the level set is a circle centered at $(0, 0)$ with radius $\sqrt{c - 3}$.

Level sets are particularly useful if they correspond to “recognizable” geometric shapes as in [Example 1.4.3](#). In the following example we remind the reader of the mathematical equations that characterize some of the most well-known geometric shapes.

Example 1.4.4

Fix a point $(c_1, c_2) \in \mathbb{R}^2$.

- (a) (*Circle*) A circle in \mathbb{R}^2 centered at the point (c_1, c_2) and with radius r consists of precisely the $(x, y) \in \mathbb{R}^2$ for which

$$(x - c_1)^2 + (y - c_2)^2 = r^2. \quad (1.9)$$

- (b) (*Ellipsoid*) Given $a, b > 0$, the set of points $(x, y) \in \mathbb{R}^2$ for which

$$\frac{(x - c_1)^2}{a^2} + \frac{(y - c_2)^2}{b^2} = 1, \quad (1.10)$$

forms an ellipsoid centered at (c_1, c_2) . The ellipsoid has symmetry axes $x = c_1, y = c_2$ and half-axes $a > 0, b > 0$.

- (c) (*Hyperbola*) Given $a, b > 0$, the set of points $(x, y) \in \mathbb{R}^2$ for which

$$\frac{(x - c_1)^2}{a^2} - \frac{(y - c_2)^2}{b^2} = 1, \quad (1.11)$$

forms a hyperbola centered at (c_1, c_2) . The hyperbola has symmetry axes $x = c_1, y = c_2$. Similarly, a hyperbola with symmetry axes $y = c_1, x = c_2$, is given by the equation

$$\frac{(y - c_2)^2}{b^2} - \frac{(x - c_1)^2}{a^2} = 1. \quad (1.12)$$

Example 1.4.5 Parabola

The graph of a second-order polynomial function of one variable is of the form

$$y = ax^2 + bx + c.$$

Hence, the graph forms a parabola with symmetry axis $x = -\frac{b}{2a}$.

Note that we can also consider higher-dimensional versions of the sets in Example 1.4.4.

Example 1.4.6 Sphere

Fixing a point $(c_1, c_2, c_3) \in \mathbb{R}^3$, the sphere in \mathbb{R}^3 centered at the point

(c_1, c_2, c_3) and with radius $r > 0$ is given by

$$(x - c_1)^2 + (y - c_2)^2 + (z - c_3)^2 = r^2. \quad (1.13)$$

In the rest of this section we will consider the level sets for quadratic forms, see [Definition 1.2.1](#), of a particular simple form, namely,

$$q : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad q(x, y) = \lambda_1 x^2 + \lambda_2 y^2 + b_1 x + b_2 y + c, \quad (1.14)$$

where $\lambda_1, \lambda_2, b_1, b_2, c \in \mathbb{R}$. This function is indeed a quadratic form according to [\(1.3\)](#) with the choices

$$A = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}, \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}.$$

For such functions we are able to give explicit descriptions of the level sets in terms of the geometric shapes discussed in [Example 1.4.4](#) and [Example 1.4.5](#). Note that the expression for the analyzed function q in [\(1.14\)](#) does not involve the product xy since A was assumed to be a diagonal matrix. The case where the product xy appears is significantly more complicated and will be handled in [Example 2.8.3](#) on page 60.

Example 1.4.7

For the function q in [\(1.14\)](#), the level set associated with a given $k \in \mathbb{R}$ consists precisely of the $(x, y) \in \mathbb{R}^2$ for which

$$\lambda_1 x^2 + \lambda_2 y^2 + b_1 x + b_2 y + c = k,$$

or,

$$\lambda_1 x^2 + \lambda_2 y^2 + b_1 x + b_2 y = k - c. \quad (1.15)$$

For some (not all) values of the parameters $\lambda_1, \lambda_2, b_1, b_2, c$ we will now describe the geometric shape of the level set.

(a) If $\lambda_1 = \lambda_2 = 0$ and $b_2 \neq 0$, the equation [\(1.15\)](#) takes the form

$$b_1 x + b_2 y = k - c,$$

or

$$y = -\frac{b_1}{b_2} x + \frac{k - c}{b_2}.$$

Thus, the level set is a straight line with slope $-b_1/b_2$.

(b) If $\lambda_2 = 0$ and $\lambda_1 \neq 0, b_2 \neq 0$, (1.15) can be written as

$$y = -\frac{\lambda_1}{b_2}x^2 - \frac{b_1}{b_2}x + \frac{k-c}{b_2}.$$

Thus, the level set is a parabola. Similarly, the case $\lambda_1 = 0, \lambda_2 \neq 0, b_1 \neq 0$ yields a parabola.

(c) If $\lambda_1 \neq 0, \lambda_2 \neq 0$, the level set can be identified as one of the sets in [Example 1.4.4](#). The technique to do so is always the same; let us illustrate it on the equation

$$4x^2 + y^2 + 8x - 6y - 3 = 0. \quad (1.16)$$

Let us rewrite (1.16) as

$$4(x^2 + 2x) + y^2 - 6y = 3. \quad (1.17)$$

We will now consider the terms involving x and y separately. Thus, we first look at the expression $x^2 + 2x$. Notice that all the definitions in [Example 1.4.4](#) involve terms of the form

$$(x - c_1)^2 = x^2 - 2c_1x + c_1^2$$

for some $c_1 \in \mathbb{R}$. In order to rewrite $x^2 + 2x$ such that a term of this form occur, we see that it is necessary that we choose c_1 such that $-2c_1 = 2$, i.e., $c_1 = -1$. With this choice of c_1 , we have that

$$(x - c_1)^2 = (x + 1)^2 = x^2 + 2x + 1 = (x^2 + 2x) + 1;$$

therefore,

$$x^2 + 2x = (x + 1)^2 - 1. \quad (1.18)$$

Thus, we see that we do not obtain that the term $x^2 + 2x$ can be written exactly as $(x - c_1)^2$, but we will see what to do with the occurring number -1 soon. Before doing so, we will now run the same procedure on the term $y^2 - 6y$. Exactly as before, we see that in order to introduce a term

$$(y - c_2)^2 = y^2 - 2c_2y + c_2^2,$$

it is necessary to take c_2 such that $-2c_2 = -6$, i.e., $c_2 = 3$. With this choice,

$$(y - c_2)^2 = (y - 3)^2 = y^2 - 6y + 9 = (y^2 - 6y) + 9,$$

1.5. What is multivariate calculus useful for?

i.e.,

$$y^2 - 6y = (y - 3)^2 - 9, \quad (1.19)$$

Inserting the expressions (1.18) and (1.19) in (1.17) yields that

$$4((x + 1)^2 - 1) + (y - 3)^2 - 9 = 3,$$

or,

$$4(x + 1)^2 + (y - 3)^2 = 16.$$

Finally, this equation can be rewritten as

$$\frac{(x + 1)^2}{4} + \frac{(y - 3)^2}{16} = 1;$$

comparing with the characterizations in [Example 1.4.4](#) shows that the set is an ellipsoid centered at the point $(-1, 3)$, with symmetry axes $x = -1, y = 3$ and half-axes $a = 2, b = 4$.

1.5 What is multivariate calculus useful for?

The derivative quantifies the sensitivity of change of a function's output with respect to its input. From standard calculus we know that the derivative of a scalar function $f : \mathbb{R} \rightarrow \mathbb{R}$ at a point $x = x_0$ is, whenever it exists, the slope of the tangent line to the graph of the function at $(x_0, f(x_0))$. It is not at all clear what the derivative of a vector function of several variables should be. In fact, we will see in [Chapter 3](#) that a proper understanding of differentiability of functions of several variables truly requires a multi-dimensional analysis of the problem. We here simply mention that the derivative of the linear mapping L_A in [Example 1.3.1](#) is A , that is, a $k \times n$ matrix, and that the derivative of the quadratic form considered in [Definition 1.2.1](#) is $(A + A^T)\mathbf{x} + \mathbf{b}$ which is a vector in \mathbb{R}^n and, consider as a function of \mathbf{x} , an affine function from \mathbb{R}^n to \mathbb{R}^n .

CHAPTER 2

Inner Product Spaces and the Spectral Theorem

In this chapter we will develop the theory of linear algebra that will be needed in later chapters. Continuing from the lecture notes [Mathematics 1a], we will introduce new central aspects of linear algebra that has its origins from plane geometry, e.g., lengths and angles in the two-dimensional real plane. This leads to the development of axioms for inner products and norms, extending the concepts of length and perpendicularity (or orthogonality as we will say) to abstract vector spaces. From the concept of a norm, we will be able to introduce open and closed sets in a vector space, Section 2.2. That section provides the necessary background for studying continuity and differentiability of functions of several real variables.

The main result of the chapter is the *spectral theorem* that shows us that symmetric matrices have a basis of eigenvectors that are orthogonal to each other. The name *spectral theory* was introduced by David Hilbert, and it concerns theories extending the eigenvalue problem of a single square matrix to classes of matrices and linear transformations. The set of eigenvalues of a square matrix is called the *spectrum* and the matrix classes we will study are real symmetric, Hermitian and normal matrices.

Abstract vectors spaces can be considered over arbitrary scalar fields \mathbb{F} , however, we will restrict our attention to either real or complex vector spaces. That is, we will always assume that $\mathbb{F} = \mathbb{R}$ or $\mathbb{F} = \mathbb{C}$. We are mainly interested in the real case, however, as even real matrices can have a complex spectrum (i.e., complex eigenvalues), it is natural to study spectral theory in the setting of both real and complex vector spaces. We will sometimes write formulas both for real and complex case, but we will in some cases only state the complex formulas as they also apply to the real case.

While we consider abstract vector spaces, we will focus on the space of

real column vectors of size $n \times 1$, that is, the vector space \mathbb{R}^n over the reals \mathbb{R} , as well as the space of complex column vectors of size $n \times 1$, that is, the vector space \mathbb{C}^n over the field \mathbb{C} of complex scalars.

2.1 Inner product spaces

Inner product and norm in \mathbb{R}^n . In dimensions two and three, the length of a vector, i.e., the distance from the origin to its endpoint, can be computed by the Pythagorean rule. For example in \mathbb{R}^2 , the length of the vector $\mathbf{x} = [x_1, x_2]^T$ is given by

$$\|\mathbf{x}\| = \sqrt{x_1^2 + x_2^2}.$$

It is natural to generalize this formula for all $n \in \mathbb{N}$, and we define the length of the vector $\mathbf{x} \in \mathbb{R}^n$ as

$$\|\mathbf{x}\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}. \quad (2.1)$$

Instead of the word “length” we will say that the non-negative real number $\|\mathbf{x}\|$ is the *norm* of \mathbf{x} .

The dot product of two vectors $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$, $\mathbf{y} = [y_1, y_2, \dots, y_n]^T$ in \mathbb{R}^n is defined as

$$\mathbf{x} \cdot \mathbf{y} = x_1 y_1 + x_2 y_2 + \dots + x_n y_n. \quad (2.2)$$

In dimensions two and three, this definition reduces to the well-known formula, e.g., the dot product in \mathbb{R}^2 is $\mathbf{x} \cdot \mathbf{y} = x_1 y_1 + x_2 y_2$, where $\mathbf{x} = [x_1, x_2]^T$ and $\mathbf{y} = [y_1, y_2]^T$. While the notation $\mathbf{x} \cdot \mathbf{y}$ and the term “dot product” is often used for the inner product in (2.2), we often prefer the notation $\langle \mathbf{x}, \mathbf{y} \rangle$ and call this the *standard inner product* on \mathbb{R}^n .

The norm of \mathbf{x} can be computed from the dot product of \mathbf{x} with itself:

$$\|\mathbf{x}\| = \sqrt{\mathbf{x} \cdot \mathbf{x}}. \quad (2.3)$$

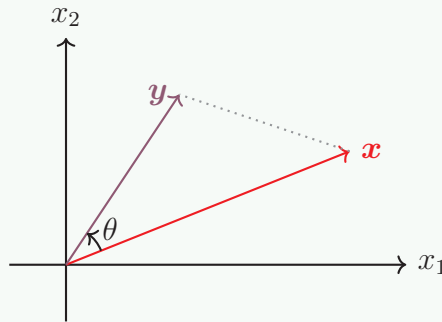
With the dot product and the norm the vector space \mathbb{R}^n is sometimes said to have euclidean structure or geometry.

Using the notion of matrix transpose, one can write the standard inner product in \mathbb{R}^n by $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^T \mathbf{x}$. Note that this is a slight abuse of notation because $\langle \mathbf{x}, \mathbf{y} \rangle$ is a real scalar while $\mathbf{x}^T \mathbf{y}$ is a 1×1 matrix as it is the matrix-matrix multiplication of a $1 \times n$ and $n \times 1$ matrix, respectively. Note also that $\mathbf{y}^T \mathbf{x} = \mathbf{x}^T \mathbf{y}$, and we use the notation $\mathbf{y}^T \mathbf{x}$ only to be consistent with the formula in the complex case.

The angle between two vectors is closely related to the norm and inner product as the following example shows.

Example 2.1.1

Consider the triangle below, whose vertices are given by the origin $\mathbf{0} = [0, 0]^T$ and the real, non-zero vectors $\mathbf{x} = [x_1, x_2]^T$ and $\mathbf{y} = [y_1, y_2]^T$.



Let θ denote the angle between the vectors \mathbf{x} and \mathbf{y} . One can show (e.g., using the law of cosines) that

$$\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x}\| \|\mathbf{y}\| \cos \theta.$$

If $\theta = \pm 90^\circ$, we say that \mathbf{x} and \mathbf{y} are orthogonal to each other. Note that \mathbf{x} and \mathbf{y} are *orthogonal* if and only if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

From the formula

$$\cos \theta = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

we see that the angle θ between two vectors can be found from the values of inner product $\langle \mathbf{x}, \mathbf{y} \rangle$ and the norms $\|\mathbf{x}\|$ and $\|\mathbf{y}\|$.

Inner product and norm in \mathbb{C}^n . Let us now define norm and inner product for the vector space \mathbb{C}^n over the scalar field \mathbb{C} . We cannot directly use (2.3) as for \mathbb{R}^n since the expression inside the square root might be complex if \mathbf{x} is a complex vector; recall that x^2 is rarely a real number if x is complex (in fact, x^2 is real only if x is real or if x is purely imaginary). However, for a complex number $x = a + ib$, $a, b \in \mathbb{R}$, we have $|x|^2 = x\bar{x} = a^2 + b^2 \geq 0$. Hence, if $\mathbf{x} \in \mathbb{C}^n$ is given by

$$\mathbf{x} = \begin{bmatrix} a_1 + ib_1 \\ a_2 + ib_2 \\ \vdots \\ a_n + ib_n \end{bmatrix},$$

it is natural to define its norm $\|\mathbf{x}\|$ by

$$\|\mathbf{x}\| = \sqrt{\sum_{k=1}^n (a_k^2 + b_k^2)} = \sqrt{\sum_{k=1}^n |x_k|^2}.$$

This definition guarantees that $\|\mathbf{x}\|$ is a *non-negative* real number which is what we would expect from a quantity that is the “length” of a vector $\mathbf{x} \in \mathbb{C}^n$.

Since we want to keep the relation $\|\mathbf{x}\|^2 = \langle \mathbf{x}, \mathbf{x} \rangle$ as for \mathbb{R}^n , we can define the *standard inner product on \mathbb{C}^n* by

$$g\langle \mathbf{x}, \mathbf{y} \rangle = x_1\bar{y}_1 + x_2\bar{y}_2 + \dots + x_n\bar{y}_n = \sum_{k=1}^n x_k\bar{y}_k;$$

it follows that $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ since $x_k\bar{x}_k = |x_k|^2 \geq 0$ for all $k = 1, \dots, n$.

Using the notion of matrix adjoints, one can write the standard inner product in \mathbb{C}^n by $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^*\mathbf{x}$. Again, this is an abuse of notation (that we will ignore) because $\langle \mathbf{x}, \mathbf{y} \rangle$ is a complex scalar while $\mathbf{y}^*\mathbf{x}$ is a 1×1 matrix.

Inner products and norms in abstract spaces. Let V be a vector space over \mathbb{F} . If $\mathbb{F} = \mathbb{R}$, we say that V is a real vector space, and if $\mathbb{F} = \mathbb{C}$, we say that V is a complex vector space. An inner product is a function that assigns to each *pair* of vectors (\mathbf{x}, \mathbf{y}) a scalar in \mathbb{F} , denoted by $\langle \mathbf{x}, \mathbf{y} \rangle$. The inner product function $(\mathbf{x}, \mathbf{y}) \mapsto \langle \mathbf{x}, \mathbf{y} \rangle$ has to satisfy the following properties:

Definition 2.1.1 Inner product

Let $\mathbb{F} = \mathbb{R}$ or \mathbb{C} . An *inner product* on V is a function

$$\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{F}$$

that satisfies for any $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$ and $c, d \in \mathbb{F}$ the following properties:

- (i) Non-negativity: $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$
- (ii) Non-degeneracy: $\langle \mathbf{x}, \mathbf{x} \rangle = 0$ if and only if $\mathbf{x} = 0$,
- (iii) (Conjugate) symmetry: $\langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle}$; note, that for a real vector space, this property is just symmetry, $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$;
- (iv) Linearity: $\langle c\mathbf{x} + d\mathbf{y}, \mathbf{z} \rangle = c\langle \mathbf{x}, \mathbf{z} \rangle + d\langle \mathbf{y}, \mathbf{z} \rangle$

A vector space V with an inner product is called an *inner product space*.

It is important to remember that the quantity $\langle \mathbf{x}, \mathbf{y} \rangle$ is always just a scalar: it is a real number for real vector spaces and a (possible) complex number for complex vector spaces. For complex vector space the inner product is conjugate linear in the second entry (and linear in the first as required in Item (iv)). Conjugate linearity means that

$$\langle \mathbf{x}, c\mathbf{y} + d\mathbf{z} \rangle = \bar{c}\langle \mathbf{x}, \mathbf{y} \rangle + \bar{d}\langle \mathbf{x}, \mathbf{z} \rangle$$

which follows from (iii) and (iv) in Definition 2.1.1. Loosely speaking, we are allowed to move scalars “out” from inside the inner product, but we have to remember to do a complex conjugation if the scalar comes from the second entry of the inner product. Finally, we mention that the inner products with the zero vector are always zero, i.e.,

$$\langle \mathbf{0}, \mathbf{y} \rangle = 0 \text{ and } \langle \mathbf{x}, \mathbf{0} \rangle = 0.$$

To see this, we use linearity of the inner product to compute $\langle \mathbf{0}, \mathbf{y} \rangle = \langle 0\mathbf{x}, \mathbf{y} \rangle = 0 \cdot \langle \mathbf{x}, \mathbf{y} \rangle = 0$.

A space V together with an inner product on it is called an *inner product space*. Given an inner product space V , we define the *norm* on V by

$$\|\cdot\| : V \rightarrow [0, \infty[, \quad \|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} \tag{2.4}$$

Note that by the non-negativity property of the inner product $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$, the expression inside the square root in (2.4) satisfies $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ hence $\|\mathbf{x}\| \geq 0$. In fact, the norm is a function that behaves like the length of vectors in the plane \mathbb{R}^2 , e.g., the length of (-5) times the vector is 5 times the length of the original vector. This last property is a special case of Item (iii) in the following general result.

Theorem 2.1.1 Norm

Let V be an inner product space over \mathbb{F} with inner product $\langle \cdot, \cdot \rangle$ and derived norm $\|\cdot\|$ given by (2.4). Then the norm satisfies:

- (i) Non-negativity: $\|\mathbf{x}\|$ is real and non-negative
- (ii) Non-degeneracy: $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$
- (iii) Scaling: $\|c\mathbf{x}\| = |c| \|\mathbf{x}\|$
- (iv) Triangle inequality: $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$,

for any $\mathbf{x}, \mathbf{y} \in V$ and $c \in \mathbb{F}$.

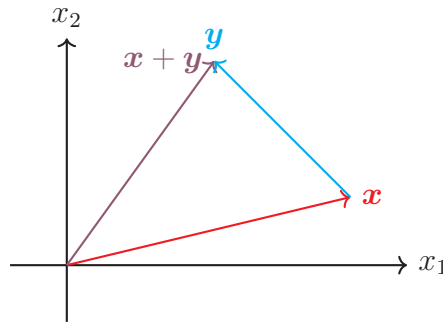


Figure 2.1: Illustration of the triangle inequality (2.13) in \mathbb{R}^2 .

Proof. (i): Non-negativity follows from non-negativity of the inner product in Definition 2.1.1.

(ii): The equation $\|\mathbf{x}\| = 0$ holds if and only if $\langle \mathbf{x}, \mathbf{x} \rangle = 0$ holds which, by Definition 2.1.1, holds if and only if $\mathbf{x} = 0$.

(iii): The scaling law follows from (conjugate) linearity and conjugate symmetry of the inner product:

$$\|c\mathbf{x}\| = \langle c\mathbf{x}, c\mathbf{x} \rangle^{1/2} = (c\bar{c} \langle \mathbf{x}, \mathbf{x} \rangle)^{1/2} = (|c|^2 \langle \mathbf{x}, \mathbf{x} \rangle)^{1/2} = |c| \|\mathbf{x}\|.$$

(iv): The inequality is illustrated in Figure 2.1. We postpone the proof of the triangle inequality to Theorem 2.1.7 on page 35. ■

The standard inner products on \mathbb{R}^n and on \mathbb{C}^n introduced above are in fact inner products, that is, they satisfy the requirements in Definition 2.1.1, see Exercise 2.1.2. Let \mathbb{F}^n denote either \mathbb{R}^n and on \mathbb{C}^n . Then the standard inner products can be written as:

$$\langle \mathbf{x}, \mathbf{y} \rangle = x_1 \bar{y}_1 + x_2 \bar{y}_2 + \dots + x_n \bar{y}_n = \sum_{k=1}^n x_k \bar{y}_k, \quad (2.5)$$

since this formula reduces to the dot product $\mathbf{x} \cdot \mathbf{y}$ if \mathbf{y} is a real vector. The standard inner product on \mathbb{F}^n has the following useful property when taking adjoints (or transposes) of a matrix $A \in M_n(\mathbb{F})$:

$$\langle A\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, A^*\mathbf{y} \rangle \quad \text{and} \quad \langle \mathbf{x}, A\mathbf{y} \rangle = \langle A^*\mathbf{x}, \mathbf{y} \rangle. \quad (2.6)$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{F}^n$. If $\mathbb{F} = \mathbb{R}$ one can write A^T in place of A^* in (2.6), e.g., $A\mathbf{x} \cdot \mathbf{y} = \mathbf{x} \cdot A^T\mathbf{y}$.

Exercise 2.1.2

Show that the standard inner product on \mathbb{F}^n in (2.5) satisfies the assumptions of Definition 2.1.1.

There are many different inner products on an inner product space, thus, also on \mathbb{F}^n . Hence, one sometimes add subscripts to distinguish them, e.g., one writes $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbb{F}^n}$ for the standard inner product in (2.5). The *standard* inner product is also called the *canonical* inner product. However, this particular inner product is nothing special, it is just the one we will always use *unless* stated otherwise. As an example of another inner product on \mathbb{F}^n , we can define $\langle \mathbf{x}, \mathbf{y} \rangle_A := \langle A\mathbf{x}, A\mathbf{y} \rangle_{\mathbb{F}^n}$ where $A \in M_n(\mathbb{F})$ is assumed to be an invertible matrix.

Exercise 2.1.3

Show that $\langle \cdot, \cdot \rangle_A$ is not an inner product on \mathbb{F}^n if A is singular. (Hint: The non-degeneracy condition of Definition 2.1.1 fails.)

Inner products and norms in other spaces. Let us equip some of the vector spaces from Section 9.1 in Mathematics 1a with an inner product.

Example 2.1.2

Recall the vector space of matrices $\mathbb{F}^{m \times n}$ from Example 9.12 in Mathematics 1a. In this book we also use the notation $M_{m \times n}(\mathbb{F})$ as an alternative to $\mathbb{F}^{m \times n}$. This vector space becomes an inner product space under the *Frobenius inner product*

$$\langle A, B \rangle_F = \text{trace}(B^*A).$$

The Frobenius norm is then

$$\|A\|_F = \sqrt{\langle A, A \rangle_F} = \sqrt{\text{trace}(A^*A)} = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{i,j}|^2 \right)^{1/2}.$$

The Frobenius norm is often used in machine learning and many other areas of applied mathematics.

We leave it as an exercise for the reader to prove that Frobenius inner product indeed is an inner product.

Example 2.1.3

Let P_n be the complex vector space of polynomials of degree at most n .

Given a finite, nonempty real interval $[a, b]$, we consider the vector space $P_n([a, b])$ of degree n complex polynomial functions restricted to $[a, b]$. For polynomials $p, q \in V$ we define

$$\langle p, q \rangle_{L^2} = \int_a^b p(x) \overline{q(x)} dx \quad (2.7)$$

This is known as the L^2 inner product on P_n over the interval $[a, b]$, and the inner product plays an important role in many areas of mathematics and physics, e.g., in signal processing and quantum mechanics.

Unit and orthogonal vectors. Suppose V is an inner product space over \mathbb{F} with inner product $\langle \cdot, \cdot \rangle$ and derived norm $\| \cdot \|$. A vector $\mathbf{x} \in V$ with the property that

$$\| \mathbf{x} \| = 1 \quad (2.8)$$

is called a *unit norm vector* or simply a *unit vector*. We can always create a unit norm vector from a non-zero vector $\mathbf{y} \neq \mathbf{0}$ that points in the same direction as \mathbf{y} . In fact, we just define \mathbf{u} by:

$$\mathbf{u} = \frac{\mathbf{y}}{\| \mathbf{y} \|} = \frac{1}{\| \mathbf{y} \|} \mathbf{y}$$

and see that

$$\| \mathbf{u} \| = \left\| \frac{1}{\| \mathbf{y} \|} \mathbf{y} \right\| = \frac{1}{\| \mathbf{y} \|} \| \mathbf{y} \| = 1,$$

where we have used the scaling law in [Theorem 2.1.1](#). Since \mathbf{u} is a (very simple) linear combination of \mathbf{y} , it follows that $\mathbf{u} \in \text{span}\{\mathbf{y}\}$ and thus \mathbf{u} and \mathbf{y} point in the “same direction”. We have added quotation marks here since V is an abstract vector space, e.g., \mathbf{u} and \mathbf{y} could be polynomials, so we take “same direction” to mean exactly $\mathbf{u} = c\mathbf{y}$ where c is a non-negative real number. If $c < 0$ in the relation $\mathbf{x} = c\mathbf{y}$, one can say that \mathbf{x} and \mathbf{y} has “opposite” direction, but if c is complex and non-real this picture does not work. (In fact, in complex vector space it is usually sufficient to know if two vectors span the same subspace, i.e., whether $\text{span}\{\mathbf{u}\} = \text{span}\{\mathbf{y}\}$).

With [Example 2.1.1](#) on page 25 in mind, the following definition of orthogonality is rather natural.

Definition 2.1.2 Orthogonality

Let V be an inner product space over \mathbb{F} with inner product $\langle \cdot, \cdot \rangle$.

- (i) Two vectors $\mathbf{x}, \mathbf{y} \in V$ are *orthogonal* if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$. A list/set of vectors

is orthogonal if the vectors in the list/set are pairwise orthogonal.

- (ii) Two vectors $\mathbf{u}, \mathbf{v} \in V$ are *orthonormal* if they are both unit vectors and orthogonal. A list/set of vectors is orthonormal if the vectors in the list/set are unit vectors and pairwise orthogonal.
- (iii) Let S denote a (nonempty) list/set of vectors in V . Then the orthogonal complement of S in V consists of vectors that are orthogonal to all vectors in S . It is denoted S^\perp and given by

$$S^\perp = \{\mathbf{x} \in V : \langle \mathbf{x}, \mathbf{s} \rangle = 0 \text{ for all } \mathbf{s} \in S\}$$

If $S = \emptyset$, then $S^\perp := V$

Note that also $\langle \mathbf{y}, \mathbf{x} \rangle = 0$ if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ by conjugate symmetry of the inner product. We write $\mathbf{x} \perp \mathbf{y}$ if \mathbf{x} and \mathbf{y} are orthogonal. It follows from Definition 2.1.2 that a list of vectors $\mathbf{u}_1, \dots, \mathbf{u}_n$ is orthonormal if and only if

$$\langle \mathbf{u}_j, \mathbf{u}_k \rangle = \begin{cases} 1 & \text{if } j = k, \\ 0 & \text{if } j \neq k. \end{cases} \quad (2.9)$$

Example 2.1.4

We consider \mathbb{R}^3 as an inner product space over the reals \mathbb{R} with the standard inner product (2.5) (which is the same as the dot product (2.2)). Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^3$ be given by

$$\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} -2/3 \\ 2/3 \\ -1/3 \end{bmatrix}$$

The vectors \mathbf{x} and \mathbf{y} are *orthogonal*, written $\mathbf{x} \perp \mathbf{y}$, since

$$\langle \mathbf{x}, \mathbf{y} \rangle = 1(-2/3) + 2(2/3) + 2(-1/3) = 0$$

However, the two vectors are *not orthonormal* as only \mathbf{y} is a unit vector. Let us check:

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{1^2 + 2^2 + 2^2} = \sqrt{5} \neq 1$$

and

$$\begin{aligned} \|\mathbf{y}\| &= \sqrt{\langle \mathbf{y}, \mathbf{y} \rangle} = \sqrt{(-2/3)^2 + (2/3)^2 + (-1/3)^2} \\ &= \sqrt{4/9 + 4/9 + 1/9} = \sqrt{1} = 1 \end{aligned}$$

Example 2.1.5

We consider the vector space $M_{m \times n}(\mathbb{F})$ of $n \times m$ matrices equipped with the Frobenius inner product introduced in Example 2.1.2. In Example 9.12 in Mathematics 1a the standard basis of $E^{(i,j)}$ matrices was introduced. The matrices $E^{(i,j)}$ are defined as $n \times m$ matrices of zeros except at index (i, j) , where the matrix has a one. The list containing the nm matrices $E^{(i,j)}$ is orthonormal, i.e.,

$$\langle E^{(i,j)}, E^{(i',j')} \rangle_F = \text{trace}((E^{(i',j')})^* E^{(i,j)}) = \begin{cases} 1 & \text{if } j = j' \text{ and } i = i', \\ 0 & \text{if } j \neq j' \text{ or } i \neq i'. \end{cases}$$

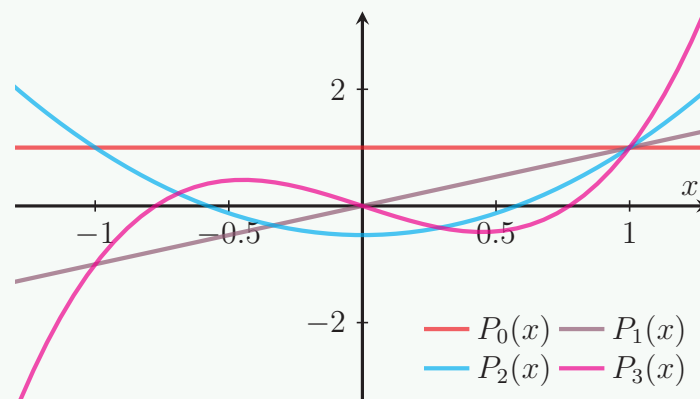
This can be verified by computing the diagonal of $(E^{(i',j')})^* E^{(i,j)}$. Note that we are only interested in the diagonal of this matrix product as the Frobenius inner product only takes the trace of the product.

Example 2.1.6

The Legendre polynomials are a family of orthogonal polynomials with respect to the L^2 -inner product introduced in (2.7) on the interval $[-1, 1]$. The first four Legendre polynomials are given by:

$$P_0(x) = 1, \quad P_1(x) = x, \quad P_2(x) = \frac{1}{2}(3x^2 - 1), \quad P_3(x) = \frac{1}{2}(5x^3 - 3x)$$

and plotted here (note we are only interested in the polynomial inside the interval $[-1, 1]$):



It is straightforward to check that P_0, P_1, P_2, P_3 is an orthogonal list in $P_3([-1, 1])$. Let us check one of these orthogonalities:

$$\begin{aligned}\langle P_0, P_2 \rangle &= \int_{-1}^1 P_0(x)P_2(x) \, dx = \int_{-1}^1 \frac{1}{2} (3x^2 - 1) \, dx \\ &= \frac{1}{2} [x^3 - x]_{x=-1}^1 = \frac{1}{2} ((1^3 - 1) - ((-1)^3 - (-1))) = \frac{1}{2}(0 - 0) = 0\end{aligned}$$

hence $P_0 \perp P_2$. Note that the list is not *orthonormal* as the polynomials are not normalized, e.g.,

$$\|P_1\|^2 = \int_{-1}^1 |P_1(x)|^2 \, dx = \int_{-1}^1 x^2 \, dx = \left[\frac{1}{3}x^3 \right]_{x=-1}^1 = \frac{2}{3} \neq 1$$

The only vector that is orthogonal to all other vectors is the zero vector. Let us state and prove this simple but very useful fact; it is an essential property in the so-called weak formulation of partial differential equations used in, e.g., the finite element method.

Lemma 2.1.4

Let V be an inner product space over \mathbb{F} with inner product $\langle \cdot, \cdot \rangle$. If $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle$ for all $\mathbf{x} \in V$, then $\mathbf{y} = \mathbf{z}$. In particular, if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ for all $\mathbf{x} \in V$, then $\mathbf{y} = \mathbf{0}$.

Proof. Suppose $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle$ for all $\mathbf{x} \in V$. Then $\langle \mathbf{x}, \mathbf{y} - \mathbf{z} \rangle = 0$ for all $\mathbf{x} \in V$. Taking $\mathbf{x} = \mathbf{y} - \mathbf{z}$ yields $\langle \mathbf{y} - \mathbf{z}, \mathbf{y} - \mathbf{z} \rangle = 0$. By non-degeneracy of the inner product it follows that $\mathbf{y} - \mathbf{z} = \mathbf{0}$, i.e., $\mathbf{y} = \mathbf{z}$. ■

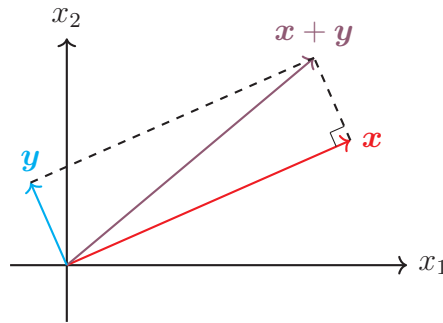
The Pythagorean theorem, known from two-dimensional and three-dimensional euclidean spaces¹, just think of \mathbb{R}^2 and \mathbb{R}^3 , can also be generalized to inner product spaces.

Theorem 2.1.5 Pythagoras

Let V be an inner product space over \mathbb{F} with inner product $\langle \cdot, \cdot \rangle$ and derived norm $\|\cdot\|$. If \mathbf{x} and \mathbf{y} are orthogonal, then

$$\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 \tag{2.10}$$

¹Strictly speaking, a euclidean vector is a geometric vector, often represented by a directed line segment (an “arrow”), that has a length and a direction. In physics this representation is often used, e.g., as the velocity or force vector. Here, we will think of euclidean vectors as elements of $\mathbf{x} \in \mathbb{R}^n$ (an “arrow” from $\mathbf{0}$ to \mathbf{x}) equipped with the standard inner product.

Figure 2.2: Illustration of the Pythagorean theorem in \mathbb{R}^2 .

Proof. Assume $\langle \mathbf{x}, \mathbf{y} \rangle = 0$. Then $\langle \mathbf{y}, \mathbf{x} \rangle = \overline{\langle \mathbf{x}, \mathbf{y} \rangle} = \bar{0} = 0$. Now just compute:

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|^2 &= \langle \mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y} \rangle \\ &= \langle \mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{y}, \mathbf{x} \rangle + \langle \mathbf{y}, \mathbf{y} \rangle \\ &= \|\mathbf{x}\|^2 + 0 + 0 + \|\mathbf{y}\|^2 \\ &= \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2. \end{aligned}$$

■

Two important inequalities for the inner product and norm. The inner product and the derived norm satisfies the following useful inequality:

Theorem 2.1.6 Cauchy/Schwarz' inequality

Let V be an inner product space over \mathbb{F} with inner product $\langle \cdot, \cdot \rangle$ and derived norm $\|\cdot\|$. Then

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\| \quad (2.11)$$

Proof. Let $\mathbf{x}, \mathbf{y} \in V$. If \mathbf{y} is the zero vector, then both sides of the inequality are equal to zero, and the inequality is trivially satisfied. Thus, we can assume that $\mathbf{y} \neq \mathbf{0}$.

Suppose first that $\langle \mathbf{x}, \mathbf{y} \rangle$ is real. We know that, for all $t \in \mathbb{R}$,

$$\|\mathbf{x} + t\mathbf{y}\|^2 \geq 0. \quad (2.12)$$

since the norm is always non-negative. Expanding the left-hand side using the properties of the norm, we have

$$\|\mathbf{x} + t\mathbf{y}\|^2 = \langle \mathbf{x} + t\mathbf{y}, \mathbf{x} + t\mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{x}, t\mathbf{y} \rangle + \langle t\mathbf{y}, \mathbf{x} \rangle + \langle t\mathbf{y}, t\mathbf{y} \rangle$$

$$= \|\mathbf{x}\|^2 + 2t\langle \mathbf{x}, \mathbf{y} \rangle + t^2\|\mathbf{y}\|^2,$$

where we have used that $\langle \mathbf{x}, \mathbf{y} \rangle$ is real. The expression above is a quadratic polynomial in t with the leading coefficient $\|\mathbf{y}\|^2 > 0$, which means it opens upwards. Therefore, since the polynomial is non-negative for all $t \in \mathbb{R}$, the discriminant of the polynomial must be zero or negative. Indeed, if the discriminant is positive, the polynomial has two real, distinct roots and therefore taking negative values for t between the two roots, contradicting (2.12). Hence, the discriminant is non-positive, that is,

$$D = (2\langle \mathbf{x}, \mathbf{y} \rangle)^2 - 4\|\mathbf{x}\|^2\|\mathbf{y}\|^2 \leq 0.$$

Rearranging this inequality and dividing through by 4, we get

$$\langle \mathbf{x}, \mathbf{y} \rangle^2 \leq \|\mathbf{x}\|^2\|\mathbf{y}\|^2$$

and by taking the square root of both sides, we obtain

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\|\|\mathbf{y}\|$$

Now, if $c := \langle \mathbf{x}, \mathbf{y} \rangle$ is not real, we replace \mathbf{x} by $\mathbf{x}' = \frac{\bar{c}}{|c|}\mathbf{x}$. Then $\|\mathbf{x}\| = \|\mathbf{x}'\|$ and $|\langle \mathbf{x}, \mathbf{y} \rangle| = |\langle \mathbf{x}', \mathbf{y} \rangle|$. Since $\langle \mathbf{x}', \mathbf{y} \rangle = \frac{\bar{c}}{|c|}c = |c|$ is real, the general case follows from the special case we considered above. ■

For $V = \mathbb{R}^n$ the Cauchy/Schwarz' inequality, when squaring both sides, reads

$$(x_1y_1 + \cdots + x_ny_n)^2 \leq (x_1^2 + \cdots + x_n^2)(y_1^2 + \cdots + y_n^2),$$

where $\mathbf{x} = [x_1, \dots, x_n]^T$ and $\mathbf{y} = [y_1, \dots, y_n]^T$ are real vectors. It is not obvious that such a relation should hold. In fact, it is rather difficult to show this Cauchy/Schwarz' inequality by means of algebraic methods alone.

The second inequality we want to prove here is the triangle inequality for inner product space. Intuitively, the length of the sum of two vectors $\mathbf{x} + \mathbf{y}$ should be maximized if the two vectors \mathbf{x} and \mathbf{y} point in the same direction. Moreover, the maximal length of $\mathbf{x} + \mathbf{y}$ should not be able to be bigger than the sum of the individual lengths of the two vectors. The precise statement behind these ideas can be formulated as follows:

Theorem 2.1.7 Triangle inequality

Let V be an inner product space over \mathbb{F} with inner product $\langle \cdot, \cdot \rangle$ and derived norm $\|\cdot\|$, and let $\mathbf{x}, \mathbf{y} \in V$. Then

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\| \tag{2.13}$$

with equality if and only if \mathbf{x} and \mathbf{y} points in the same direction, i.e., one vector is a non-negative scalar multiple of the other.

Proof. By properties of the inner product and Cauchy/Schwarz' inequality, we compute:

$$\begin{aligned}
 \|\mathbf{x} + \mathbf{y}\|^2 &= \langle \mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y} \rangle \\
 &= \langle \mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{y}, \mathbf{x} \rangle + \langle \mathbf{y}, \mathbf{y} \rangle \\
 &= \langle \mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle + \overline{\langle \mathbf{x}, \mathbf{y} \rangle} + \langle \mathbf{y}, \mathbf{y} \rangle \\
 &= \|\mathbf{x}\|^2 + 2 \operatorname{Re} \langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|^2 \\
 &\leq \|\mathbf{x}\|^2 + 2|\langle \mathbf{x}, \mathbf{y} \rangle| + \|\mathbf{y}\|^2 \\
 &\leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 \\
 &= (\|\mathbf{x}\| + \|\mathbf{y}\|)^2
 \end{aligned}$$

where we have used the following well-known properties: $z + \bar{z} = 2 \operatorname{Re} z$ and $\operatorname{Re} z \leq |z|$ for $z \in \mathbb{C}$ and $a^2 + 2ab + b^2 = (a + b)^2$ for $a, b \in \mathbb{R}$.

Finally, to prove the “with equality”-statement in the theorem, we need to prove that the two inequalities in the above computations become equalities precisely when \mathbf{x} and \mathbf{y} points in the same direction. We leave the details to the reader. ■

2.2 Open and closed sets

For functions of one variable the concept of an *open interval* is crucial in order to study continuity and differentiability. The goal of this section we will introduce the corresponding concept in an inner product space V and therefore, in particular, for \mathbb{R}^n .

From the norm $\|\cdot\|$ on an inner product space V we will define open and closed sets. This is a so-called *topology* on V , and it is precisely the topology of \mathbb{R}^n that is crucial when we want to consider continuity and differentiability of functions of several variables.

Throughout this section, we let V be an inner product space over \mathbb{F} with inner product $\langle \cdot, \cdot \rangle$ and derived norm $\|\mathbf{x}\| = \langle \mathbf{x}, \mathbf{x} \rangle^{1/2}$. However, you can safely think of \mathbb{R}^n with the euclidean norm whenever you see V . To simplify notation, we will also write vectors $\mathbf{x} = [x_1, \dots, x_n]^T$ in \mathbb{R}^n as (x_1, \dots, x_n) when this does not lead to confusion.

Definition 2.2.1 Open ball

Fix a vector $\mathbf{x}_0 \in V$. The (open) ball at \mathbf{x}_0 with radius $r > 0$ is the set

$$B(\mathbf{x}_0, r) := \{ \mathbf{x} \in V \mid \|\mathbf{x} - \mathbf{x}_0\| < r \}. \quad (2.14)$$

In \mathbb{R}^2 , the ball at a point \mathbf{x}_0 corresponds precisely to the set of points inside a circle with center at the point.

Example 2.2.1

Let $\mathbf{x}_0 = (1, 2) \in \mathbb{R}^2$ and $r = 2$. Then

$$\begin{aligned} B(\mathbf{x}_0, r) &= B((1, 2), r) \\ &= \{ \mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x} - (1, 2)\| < 2 \} \\ &= \{ (x_1, x_2) \in \mathbb{R}^2 \mid \sqrt{(x_1 - 1)^2 + (x_2 - 2)^2} < 2 \} \\ &= \{ (x_1, x_2) \in \mathbb{R}^2 \mid (x_1 - 1)^2 + (x_2 - 2)^2 < 2^2 \}. \end{aligned}$$

This is precisely the points in \mathbb{R}^2 that are located within a circle of radius 2 centered at the point $\mathbf{x}_0 = (1, 2)$.

Definition 2.2.2 Open set

Consider a subset U of an inner product space V . Assume that for each $\mathbf{x}_0 \in U$, there exists a ball $B(\mathbf{x}_0, \epsilon)$ at \mathbf{x}_0 that is entirely contained in U , that is, there exists $\epsilon > 0$ such that

$$B(\mathbf{x}_0, \epsilon) \subseteq U.$$

Then the set U is said to be *open*.

The concept of an open set in \mathbb{R}^n , in fact, generalize the concept of open intervals in \mathbb{R} as the following example shows.

Example 2.2.2 (a) An interval of the form $]a, b[$ in \mathbb{R} is an open set.

Indeed, for each $x_0 \in]a, b[$, letting ϵ denote the minimal distance to either a or b , the set $B(x_0, \epsilon/2)$ in \mathbb{R} is completely contained in $]a, b[$.

(b) A ball $B(\mathbf{x}_0, r)$ in \mathbb{R}^n is itself an open set.

(c) Letting $a > 0$, a set in \mathbb{R}^2 of the form $[-a, a]^2 = [-a, a] \times [-a, a]$ is not open. Indeed, letting, e.g., $\mathbf{x}_0 = (a, a)$, every ball $B(\mathbf{x}_0, r)$ will contain points outside the set $[-a, a]^2$.

We now define the concept of a *closed set* U in V . Interestingly, this is defined via a condition on the *complement* of the set U , i.e., the set $V \setminus U$.

Definition 2.2.3 Closed set

A set U in V is said to be *closed* if its complement $V \setminus U$ is an open set.

Example 2.2.3 (a) An interval of the form $[a, b]$ in \mathbb{R} is a closed set. Indeed,

$$\mathbb{R} \setminus [a, b] =]-\infty, a[\cup]b, \infty[,$$

which is a union of two open sets and hence open.

(b) Let $\mathbf{x}_0 \in \mathbb{R}^2$ and $r > 0$. Then the set

$$\{\mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x} - \mathbf{x}_0\| \leq r\}$$

is closed. Compare the set with the definition of the ball in Definition 2.2.1.

(c) Letting $a > 0$, a set in \mathbb{R}^2 of the form $[-a, a]^2 = [-a, a] \times [-a, a]$ is closed.

Note that many sets in V and \mathbb{R}^n are neither open nor closed (e.g., intervals of the form $]a, b[$ in \mathbb{R}). On the other hand, there are two special sets with the particular property of being both open and closed, namely, the entire space V and the empty set \emptyset .

We need a few more concepts related to sets in \mathbb{R}^n .

Definition 2.2.4 Boundary of a set

Consider a subset M of V . The boundary of M , to be denoted by ∂M , consists of precisely the points $\mathbf{x}_0 \in V$ for which *each* ball $B(\mathbf{x}_0, \epsilon)$ contains points from the set M and from the complement $V \setminus M$.

Note that the points in the boundary ∂M do not necessarily belong to the set M .

Example 2.2.4

Consider an interval $]a, b[$ in \mathbb{R} . Then the boundary consists of precisely the two points a and b , i.e., $\partial]a, b[= \{a, b\}$.

Similarly, for the set $\{\mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x} - \mathbf{x}_0\| \leq r\}$ considered in Example 2.2.3(b), the boundary is

$$\partial\{\mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x} - \mathbf{x}_0\| \leq r\} = \{\mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x} - \mathbf{x}_0\| = r\}.$$

The reader is encouraged to make a graphical illustration of this result. In

words, it says that the boundary of a disc with radius simply consists of the corresponding circle. If we, e.g., consider the case where $r = 1$ and $\mathbf{x}_0 = \mathbf{0} = (0, 0)$, this is precisely the unit circle. It can be described as the set of points

$$\mathbf{r}(t) = \begin{bmatrix} \cos(t) \\ \sin(t) \end{bmatrix} \quad (2.15)$$

where $t \in [0, 2\pi]$. We say that the vector function \mathbf{r} provides a *parametrization* of the boundary curve; cf. Equation (1.7) on page 14 and the discussion following the equation.

Example 2.2.5

Let $a > 0$ and consider the rectangle

$$[-a, a]^2 = [-a, a] \times [-a, a]$$

in \mathbb{R}^2 . Then the boundary consists of exactly what our intuitive understanding of the word says, namely, the four line pieces that bound the rectangle. It is a little bit more complicated to give an exact mathematical description than in Example 2.2.4, due to the fact that the boundary consists of four pieces. One way of describing the boundary is as the union

$$(\{a\} \times [-a, a]) \cup (\{-a\} \times [-a, a]) \cup ([-a, a] \times \{a\}) \cup ([-a, a] \times \{-a\}).$$

Due to this more complicated structure of the boundary, the parametrization is more complicated. Indeed, we have to choose separate parametrizations for each of the four curves. The piece of the boundary consisting of the points $\{a\} \times [-a, a]$ can be parametrized as

$$\mathbf{r}(t) = \begin{bmatrix} a \\ t \end{bmatrix}, \quad (2.16)$$

where $t \in [-a, a]$.

Definition 2.2.5 Closure and interior of a set

Consider a subset M of V . The *closure* of the set M , to be denoted by \overline{M} , consists of the union of the set M and its boundary ∂M , i.e.,

$$\overline{M} = M \cup \partial M.$$

The *interior* of the set M , to be denoted by M° , consists of the points in

M that are *not* boundary points, i.e.,

$$M^\circ = M \setminus \partial M.$$

Using the concept of the norm of a vector in V , we will now define what it means that a set in V is *bounded*.

Definition 2.2.6 Bounded set

A set U in V is said to be *bounded* if there exists $r > 0$ such that $\|\mathbf{x}\| \leq r$ for all $\mathbf{x} \in U$.

Geometrically, a set U is bounded if there exists a sufficiently big ball in V that contains U .

Example 2.2.6 (a) A ball $B(\mathbf{x}_0, r)$ in \mathbb{R}^n is itself a bounded set.

(b) Letting $a > 0$, a set in \mathbb{R}^n of the form $[-a, a]^n$ is bounded. Indeed, if $\mathbf{x} = (x_1, x_2, \dots, x_n) \in [-a, a]^n$, then $|x_k| \leq a$ for all $k = 1, \dots, n$. Hence

$$\|\mathbf{x}\|^2 = x_1^2 + x_2^2 + \dots + x_n^2 \leq a^2 + a^2 + \dots + a^2 = na^2;$$

this implies that $\|\mathbf{x}\| \leq \sqrt{na}$ for all $\mathbf{x} \in [-a, a]^n$.

(c) The graph of the function $y = x^2$, i.e., the set

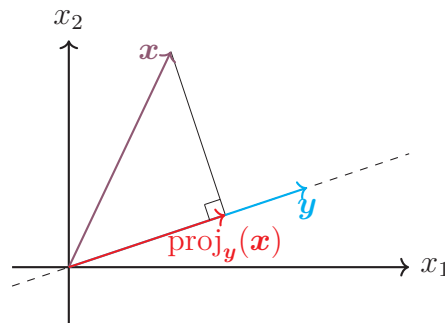
$$\{(x, y) \in \mathbb{R}^2 \mid y = x^2\}$$

is not a bounded set in \mathbb{R}^2 .

2.3 Projections onto a line

In \mathbb{R}^2 and \mathbb{R}^3 we can project a vector (orthogonally) onto another vector as shown in Figure 2.3. The aim of this section is to generalize this construction to arbitrary inner product spaces.

Let \mathbf{y} be a non-zero vector in an inner product space V . We consider $Y := \text{span}\{\mathbf{y}\}$ as the line spanned by \mathbf{y} ; Y is of course a one-dimensional subspace of V . The *orthogonal projection* of a vector $\mathbf{x} \in V$ onto the

Figure 2.3: Orthogonal projection in \mathbb{R}^2 .

subspace $Y = \text{span}\{\mathbf{y}\}$ is defined as:

$$\text{proj}_Y : V \rightarrow V, \quad \text{proj}_Y(\mathbf{x}) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{y}, \mathbf{y} \rangle} \mathbf{y}$$

Note that the image of proj_Y indeed is Y , i.e., $\text{im}(\text{proj}_Y) = Y$.

We can simplify the expression for proj_Y by normalizing the spanning vector \mathbf{y} . So let $\mathbf{u} = \mathbf{y}/\|\mathbf{y}\|$. Then $Y = \text{span}\{\mathbf{y}\} = \text{span}\{\mathbf{u}\}$ and

$$\text{proj}_Y(\mathbf{x}) = \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{u}. \quad (2.17)$$

With abuse of notation we will also write the projection proj_Y onto Y as $\text{proj}_{\mathbf{y}}$, where \mathbf{y} is any vector that spans Y .

Exercise 2.3.1

(a) Verify the formula in (2.17). (b) Show that $\text{proj}_Y : V \rightarrow V$ is a linear mapping. (c) Show that $\text{proj}_Y \circ \text{proj}_Y = \text{proj}_Y$. Thus, it does not matter if we apply the mapping once or twice (or n times), the output is the same. In other words, nothing happens to the input of the mapping the second time we apply the projections. Why is this a reasonable property of a projection?

Example 2.3.1

We consider the vector space \mathbb{F}^n with the standard inner product $\langle \cdot, \cdot \rangle$. Recall that the standard inner product is given by

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^* \mathbf{x} = \sum_{k=1}^n x_k \overline{y_k}$$

for $\mathbf{x}, \mathbf{y} \in \mathbb{F}^n$. Let $\mathbf{u} \in \mathbb{F}^n$ be a unit vector and let $Y = \text{span}\{\mathbf{u}\}$ be the one-dimensional subspace spanned by \mathbf{u} . Then, for any $\mathbf{x} \in \mathbb{F}^n$:

$\text{proj}_Y(\mathbf{x}) = \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{u} = \mathbf{u} \langle \mathbf{x}, \mathbf{u} \rangle = \mathbf{u} (\mathbf{u}^* \mathbf{x}) = (\mathbf{u} \mathbf{u}^*) \mathbf{x} = P \mathbf{x}$. The $n \times n$ matrix $P = \mathbf{u} \mathbf{u}^*$ is called a projection matrix and $\mathbf{u} \mathbf{u}^*$ is called an outer product (as a matrix product of a column vector and a row vector). The projection matrix is Hermitian (i.e., $P^* = P$) and idempotent (i.e., $P^2 = P$) as is readily verified:

$$P^* = (\mathbf{u} \mathbf{u}^*)^* = (\mathbf{u}^*)^* \mathbf{u}^* = \mathbf{u} \mathbf{u}^* = P$$

and

$$P^2 = (\mathbf{u} \mathbf{u}^*)^2 = \mathbf{u} \mathbf{u}^* \mathbf{u} \mathbf{u}^* = \mathbf{u} (\mathbf{u}^* \mathbf{u}) \mathbf{u}^* = \mathbf{u} \mathbf{u}^* = P$$

since $\mathbf{u}^* \mathbf{u} = \langle \mathbf{u}, \mathbf{u} \rangle = \|\mathbf{u}\|^2 = 1$. We finally note that $\text{col } P = Y$.

Example 2.3.2

Consider $\mathbf{y} = [1, 2, 2]^T \in \mathbb{R}^3$, and let $Y = \text{span}\{\mathbf{y}\}$. We compute $\langle \mathbf{y}, \mathbf{y} \rangle = 1^2 + 2^2 + 2^2 = 9$. The orthogonal projection onto Y is therefore given by

$$\text{proj}_Y(\mathbf{x}) = \frac{1}{9} \left\langle \mathbf{x}, \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \right\rangle \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$$

for $\mathbf{x} \in \mathbb{R}^3$. Alternatively, we normalize the vector \mathbf{y} to a unit vector $\mathbf{u} = [1/3, 2/3, 2/3]^T \in Y$ and compute

$$P = \mathbf{u} \mathbf{u}^* = \begin{bmatrix} 1/3 \\ 2/3 \\ 2/3 \end{bmatrix} [1/32/32/3] = \frac{1}{9} \begin{bmatrix} 1 & 2 & 2 \\ 2 & 4 & 4 \\ 2 & 4 & 4 \end{bmatrix} \quad (2.18)$$

Hence, the projection $\text{proj}_Y(\mathbf{x})$ can also be computed by $P \mathbf{x}$ for any $\mathbf{x} \in \mathbb{R}^3$.

2.4 Orthonormal basis

Recall from [Definition 9.14 in Mathematics 1a](#) that a finite ordered basis in a vector space V over \mathbb{F} is a list of linearly independent vectors $\alpha := \mathbf{v}_1, \dots, \mathbf{v}_n$ with the property that any vector \mathbf{x} in V can be written as a linear combination of the vectors from the list α , that is,

$$\mathbf{x} = \sum_{k=1}^n c_k \mathbf{v}_k \quad \text{for all } \mathbf{x} \in V. \quad (2.19)$$

In this text we will simply say that α is a *basis* as all bases will be finite and ordered. Recall also that if V is an n -dimensional vector space, then

any list of n linearly independent vectors is a basis.

One can show that for a given vector $\mathbf{x} \in V$ the scalars (also called coefficients) $c_k \in \mathbb{F}$, $k = 1, \dots, n$, are unique. The column vector of coefficients is the coordinate vector, and it is denoted

$${}_{\alpha}\mathbf{x} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix}.$$

That is, there is a one-to-one correspondence between vectors and coordinate vectors. Hence, instead of working with (potentially complicated) vectors \mathbf{x} (e.g., a polynomial function), we can always choose to work with column (coordinate) vectors in \mathbb{F}^n . However, to actually find ${}_{\alpha}\mathbf{x}$ we usually have to solve a linear system of equations which is not always feasible in applications. To address this problem, we introduce a particular nice class of bases called orthonormal bases.

Definition 2.4.1 Orthonormal basis

Let V be an inner product vector space over \mathbb{F} . A list β of vectors in V is called an *orthonormal basis* of V if the list is orthonormal and a basis.

Suppose $\beta := \mathbf{u}_1, \dots, \mathbf{u}_n$ is an orthonormal basis in an inner product vector space V over \mathbb{F} . We fix $j = 1, \dots, n$ and compute

$$\langle \mathbf{x}, \mathbf{u}_j \rangle = \left\langle \sum_{k=1}^n c_k \mathbf{u}_k, \mathbf{u}_j \right\rangle = \sum_{k=1}^n c_k \langle \mathbf{u}_j, \mathbf{u}_k \rangle$$

However, all the terms in the sum $\sum_{k=1}^n c_k \langle \mathbf{u}_j, \mathbf{u}_k \rangle$ are zero *unless* k is equal to the fixed number j since, by orthonormality of β ,

$$\langle \mathbf{u}_j, \mathbf{u}_k \rangle = \begin{cases} 1 & \text{if } j = k, \\ 0 & \text{if } j \neq k. \end{cases}$$

Therefore,

$$\langle \mathbf{x}, \mathbf{u}_j \rangle = 0 + \dots + 0 + c_j + 0 + \dots + 0 = c_j. \quad (2.20)$$

Hence, to find the coefficients c_k , $k = 1, \dots, n$, we just have to compute the inner products $\langle \mathbf{x}, \mathbf{u}_k \rangle$.

So, for an orthonormal basis, it is straightforward to compute the coordinate vector of a vector \mathbf{x} using the formula:

$${}_{\beta}\mathbf{x} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} \langle \mathbf{x}, \mathbf{u}_1 \rangle \\ \langle \mathbf{x}, \mathbf{u}_2 \rangle \\ \vdots \\ \langle \mathbf{x}, \mathbf{u}_n \rangle \end{bmatrix}$$

and the expansion formula (2.19) becomes:

$$\mathbf{x} = \sum_{k=1}^n \langle \mathbf{x}, \mathbf{u}_k \rangle \mathbf{u}_k \quad \text{for all } \mathbf{x} \in V.$$

Lemma 2.4.1

Suppose $\mathbf{u}_1, \dots, \mathbf{u}_n$ is a list of orthonormal vectors in an inner product space. Then $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linear independent.

Proof. Suppose $\mathbf{u}_1, \dots, \mathbf{u}_n$ is orthonormal. Consider the expression:

$$\mathbf{0} = \sum_{k=1}^n c_k \mathbf{u}_k$$

To show linear independence of $\mathbf{u}_1, \dots, \mathbf{u}_n$, we have to show that all the coefficients c_k , $k = 1, \dots, n$, are zero. From the computation (2.20) we know that $c_k = \langle \mathbf{0}, \mathbf{u}_k \rangle$ for $k = 1, \dots, n$. However, for any $k = 1, \dots, n$, it then follows that $c_k = \langle \mathbf{0}, \mathbf{u}_k \rangle = \langle 0\mathbf{x}, \mathbf{u}_k \rangle = 0\langle \mathbf{x}, \mathbf{u}_k \rangle = 0$ where $\mathbf{x} \in V$. ■

Corollary 2.4.2

Let V be an n -dimensional inner product vector space over \mathbb{F} , and let $\beta := \mathbf{u}_1, \dots, \mathbf{u}_n$ be a list of n vectors. Then β is an orthonormal basis if and only if β is orthonormal, that is, if β satisfies (2.9).

Proof. Suppose $\mathbf{u}_1, \dots, \mathbf{u}_n$ is orthonormal. By Lemma 2.4.1 the list is linear independent. Since the list consists of n vectors in an n -dimensional vector space, it must span the vector space. Hence, $\mathbf{u}_1, \dots, \mathbf{u}_n$ is a basis. It is assumed to be orthonormal so we conclude that $\mathbf{u}_1, \dots, \mathbf{u}_n$ is an orthonormal basis. ■

Example 2.4.1

We consider the standard basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ for \mathbb{F}^n , where \mathbf{e}_k is a size n column vector with all zeros except a one at index k , that is $\mathbf{e}_k = [0, \dots, 0, 1, 0, \dots, 0]^T$. Since $\mathbf{e}_1, \dots, \mathbf{e}_n$ is a list of n orthonormal vectors, it follows from Corollary 2.4.2 that the standard basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ is an orthonormal basis.

The standard basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ for \mathbb{F}^n is obviously an orthonormal basis for both \mathbb{R}^n and \mathbb{C}^n .

Remark 2.4.1

Recall that we consider \mathbb{R}^n as a vector space over the field \mathbb{R} and \mathbb{C}^n as a vector space over the field \mathbb{C} . Suppose $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathbb{R}^n$. If $\mathbf{u}_1, \dots, \mathbf{u}_n$ is an orthonormal basis for \mathbb{R}^n , then $\mathbf{u}_1, \dots, \mathbf{u}_n$ is also an orthonormal basis for \mathbb{C}^n . (Can you give a proof?)

Exercise 2.4.3

Show that the basis considered in Example 2.1.5 on page 32 is, in fact, an orthonormal basis.

2.5 The Gram-Schmidt process

In the previous section we saw that orthonormal bases are particularly nice bases to work with. Orthonormal bases are also important in numerical linear algebra, where orthogonality is frequently employed to stabilize computations. Given any list of vectors in an inner product space, the Gram-Schmidt process is used to find an orthonormal basis for the subspace spanned by those vectors. The Gram-Schmidt process is closely related to the QR decomposition of matrices used in many different applications in mathematics and engineering.

Suppose we have a list of linearly independent vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_\ell$ in an inner product space V . The goal of the Gram-Schmidt procedure is now to generate a list of orthonormal vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_\ell$ that span the same subspace. The process is as follows:

- (i) (Initialization) Set $\mathbf{w}_1 = \mathbf{v}_1$ and $\mathbf{u}_1 = \frac{1}{\|\mathbf{w}_1\|} \mathbf{w}_1$.
- (ii) (Orthogonalization & Normalization) For $k = 2, \dots, \ell$ do as follows: construct \mathbf{w}_k by removing the component of \mathbf{v}_k along each of $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k-1}$:

$$\mathbf{w}_k := \mathbf{v}_k - \sum_{j=1}^{k-1} \text{proj}_{\mathbf{u}_j}(\mathbf{v}_k)$$

where $\text{proj}_{\mathbf{u}_j}(\mathbf{v}_k)$ is the orthogonal projection of \mathbf{v}_k onto the subspace $\text{span}\{\mathbf{u}_j\}$. Then normalize \mathbf{w}_k to make sure it is a unit vector:

$$\mathbf{u}_k := \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|}. \quad (2.21)$$

Example 2.5.1

Let $V \in M_{3 \times 2}(\mathbb{R})$ be given by

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

Denote the two column vectors of V by \mathbf{v}_1 and \mathbf{v}_2 . The Gram-Schmidt process goes as follows. In the initialization, we compute the norm of \mathbf{v}_1

$$\|\mathbf{v}_1\| = (1^2 + 2^2 + 2^2)^{1/2} = \sqrt{9} = 3$$

and normalize \mathbf{v}_1 to $\mathbf{u}_1 = (1/3)\mathbf{v}_1 = [1/3, 2/3, 2/3]^T$. The orthogonal projection $\text{proj}_{\mathbf{u}_1}$ onto the subspace spanned by \mathbf{u}_1 was already found in Equation (2.18) on page 42, so here we can directly compute the projection of \mathbf{v}_2 :

$$\text{proj}_{\mathbf{u}_1}(\mathbf{v}_2) = P\mathbf{v}_2 = \frac{1}{9} \begin{bmatrix} 1 & 2 & 2 \\ 2 & 4 & 4 \\ 2 & 4 & 4 \end{bmatrix} \begin{bmatrix} -2 \\ 4 \\ 0 \end{bmatrix} = \begin{bmatrix} 2/3 \\ 4/3 \\ 4/3 \end{bmatrix}$$

We then construct \mathbf{w}_2 by

$$\mathbf{w}_2 = \mathbf{v}_2 - \text{proj}_{\mathbf{u}_1}(\mathbf{v}_2) = \begin{bmatrix} -2 \\ 4 \\ 0 \end{bmatrix} - \begin{bmatrix} 2/3 \\ 4/3 \\ 4/3 \end{bmatrix} = \begin{bmatrix} -8/3 \\ 8/3 \\ -4/3 \end{bmatrix}$$

Finally, we need to normalize \mathbf{w}_2 as

$$\mathbf{u}_2 = \frac{1}{\|\mathbf{w}_2\|} \mathbf{w}_2 = \begin{bmatrix} -2/3 \\ 2/3 \\ -1/3 \end{bmatrix}$$

Let us summarize. The list of vectors $\mathbf{v}_1, \mathbf{v}_2$ is a basis for the column space $\text{col } V = \text{span}\{\mathbf{v}_1, \mathbf{v}_2\}$. From the Gram-Schmidt process we found vectors $\mathbf{u}_1, \mathbf{u}_2$ that constitute an *orthonormal* basis for the *same* subspace $\text{col } V = \text{span}\{\mathbf{u}_1, \mathbf{u}_2\}$.

Exercise 2.5.1

This exercise is a continuation of Example 2.5.1. Suppose

$$V = \begin{bmatrix} 1 & -2 & -1 \\ 2 & 4 & -1 \\ 2 & 0 & 3 \end{bmatrix}$$

Perform the Gram-Schmidt process on the three column vectors of V . Show that the outcome is the three column vectors of U given by

$$U = \begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & -\frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & -\frac{1}{3} \\ \frac{2}{3} & -\frac{1}{3} & \frac{2}{3} \end{bmatrix}$$

Note that, if a list $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$, as in [Exercise 2.5.1](#), spans the entire space \mathbb{R}^3 , and if we are only asked to find an orthonormal basis for the span of the list (the entire space), there is no need to do the Gram-Schmidt process. We can simply pick the standard basis $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$. What makes the orthonormal basis found via the Gram-Schmidt process special is the property (2.22) stated below. Let us prove this property and at the same time show that the Gram-Schmidt process always works as claimed:

Theorem 2.5.2 Gram-Schmidt

Suppose $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_\ell$ is a list of linearly independent vectors in an inner product space V . Then the list $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_\ell$ given by (2.21) is an orthonormal list with the property

$$\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\} = \text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\} \tag{2.22}$$

for $k = 1, \dots, \ell$.

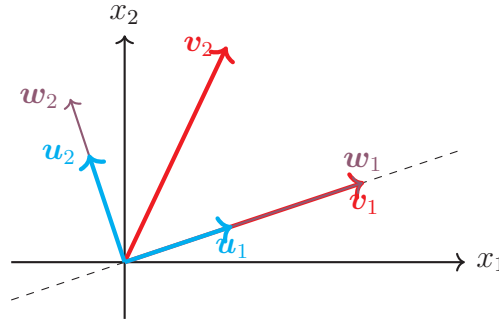
Proof. The proof goes by induction on ℓ . The base case $\ell = 1$ is easy. Since \mathbf{v}_1 is assumed to be linearly independent, it cannot be the zero vector. Hence, $\mathbf{u}_1 = \frac{\mathbf{v}_1}{\|\mathbf{v}_1\|}$ is well-defined. A list of only one vector is orthonormal precisely when the vector is of norm one. We compute $\|\mathbf{u}_1\| = \|\frac{\mathbf{v}_1}{\|\mathbf{v}_1\|}\| = \frac{1}{\|\mathbf{v}_1\|}\|\mathbf{v}_1\| = 1$.

For the induction step, let $k \leq \ell$. Suppose that, given $k - 1$ linearly independent vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}$, the vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k-1}$ are orthonormal and

$$\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}\} = \text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k-1}\}$$

Since $\mathbf{v}_1, \dots, \mathbf{v}_k$ are linearly independent, it is impossible to write \mathbf{v}_k as a linear combination of $\mathbf{v}_1, \dots, \mathbf{v}_{k-1}$, that is,

$$\mathbf{v}_k \notin \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}\} = \text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k-1}\}$$


 Figure 2.4: Illustration of the Gram-Schmidt process in \mathbb{R}^2 .

In particular, $\mathbf{v}_k \neq \sum_{j=1}^{k-1} \text{proj}_{\mathbf{u}_j}(\mathbf{v}_k)$ hence $\mathbf{w}_k \neq \mathbf{0}$, and the normalization $\mathbf{u}_k = \mathbf{w}_k / \|\mathbf{w}_k\|$ in (2.21) is well-defined.

We need to prove that \mathbf{u}_k is orthogonal to $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k-1}$. For $1 \leq m \leq k-1$, we first compute

$$\begin{aligned} \langle \mathbf{w}_k, \mathbf{u}_m \rangle &= \left\langle \mathbf{v}_k - \sum_{j=1}^{k-1} \text{proj}_{\mathbf{u}_j}(\mathbf{v}_k), \mathbf{u}_m \right\rangle = \langle \mathbf{v}_k, \mathbf{u}_m \rangle - \left\langle \sum_{j=1}^{k-1} \langle \mathbf{v}_k, \mathbf{u}_j \rangle \mathbf{u}_j, \mathbf{u}_m \right\rangle \\ &= \langle \mathbf{v}_k, \mathbf{u}_m \rangle - \sum_{j=1}^{k-1} \langle \mathbf{v}_k, \mathbf{u}_j \rangle \langle \mathbf{u}_j, \mathbf{u}_m \rangle = \langle \mathbf{v}_k, \mathbf{u}_m \rangle - \sum_{j=1}^{k-1} \langle \mathbf{v}_k, \mathbf{u}_j \rangle \delta_{j,m} \\ &= \langle \mathbf{v}_k, \mathbf{u}_m \rangle - \langle \mathbf{v}_k, \mathbf{u}_m \rangle = 0 \end{aligned}$$

and therefore

$$\langle \mathbf{u}_k, \mathbf{u}_m \rangle = \left\langle \frac{\mathbf{w}_k}{\|\mathbf{w}_k\|}, \mathbf{u}_m \right\rangle = \frac{1}{\|\mathbf{w}_k\|} \langle \mathbf{w}_k, \mathbf{u}_m \rangle = \frac{1}{\|\mathbf{w}_k\|} \cdot 0 = 0.$$

Since \mathbf{u}_k is defined as a linear combinations of $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$, it follows by the induction hypothesis that

$$\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\} \subseteq \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$$

However, the dimension of both these subspaces are k as both lists of vectors are linearly independent by Lemma 2.4.1. We conclude that $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\}$ is a k -dimensional subspace of a k -dimensional space $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$, and hence these subspaces must agree². ■

²Strictly speaking, we here assume that the reader is familiar with the fact that the only k -dimensional subspace of a k -dimensional vector space is the vector space itself. We refer the reader to Theorem 2.2.9 in [1] for a proof.

Exercise 2.5.3

Consider the list $\alpha = 1, x, x^2, x^3$ of polynomials in $P_3([-1, 1])$ equipped with the L^2 -inner product. Argue that α is a list of linearly independent vectors. Apply the Gram-Schmidt process to α and show that it outputs a *normalized*³ version of the Legendre polynomials from [Example 2.1.6](#) on page 32.

Note that we have formulated the Gram-Schmidt process only for a list of linearly independent vectors. Mathematically, it is not difficult to extend the Gram-Schmidt process to vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_\ell$ that are not linearly independent. The idea is simply to skip a vector \mathbf{v}_k in the process whenever the list $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ becomes linearly dependent. Algorithmically, we then update $k \rightarrow k + 1$ without appending a vector to the output list. Effectively, we remove vectors from the list $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_\ell$ to obtain a linearly independent list of vectors. Numerically, however, the “trick” often leads to instability issues, and it is usually recommended to use an implementation of the QR factorization.

2.6 Unitary and orthogonal matrices

Definition 2.6.1 Unitary and real orthogonal matrix

A square matrix $U \in M_n(\mathbb{C})$ is unitary if $U^*U = I$. A unitary matrix $Q \in M_n(\mathbb{R})$ that is real is called *real orthogonal*.

Since $U^*U = I$ if and only if U is invertible with U^* being the inverse, it is easy to invert unitary matrices. If Q is real orthogonal, the adjoint and the transpose are equal, hence it satisfies $Q^TQ = I$. So, for unitary U and real orthogonal Q , we have:

$$U^{-1} = U^*(= \overline{U}^T) \quad Q^{-1} = Q^T = Q^*.$$

Example 2.6.1

The following 3×3 matrix is real orthogonal and therefore also unitary.

³There is no standard rule of how to “normalize” a polynomial. Considered as a vector in normed vector space, it is natural to normalize polynomials to have norm one, but there are other choices, e.g., one can require that $P(1) = 1$ or that the coefficient of the highest order term is 1.

2.6. Unitary and orthogonal matrices

$$U = \begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & -\frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & -\frac{1}{3} \\ \frac{2}{3} & -\frac{1}{3} & \frac{2}{3} \end{bmatrix} \quad (2.23)$$

This is verified by computing: $U^T U = U^* U = I$.

Consider the vector space \mathbb{F}^n with the standard inner product $\langle \cdot, \cdot \rangle$, i.e., $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^* \mathbf{x}$. Suppose $U \in M_n(\mathbb{C})$ is unitary and let \mathbf{u}_k denote the k th column of U . Then $U^* U = I$ can be written:

$$\begin{bmatrix} \mathbf{u}_1^* \\ \mathbf{u}_2^* \\ \vdots \\ \mathbf{u}_n^* \end{bmatrix} \begin{bmatrix} \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_n \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1^* \mathbf{u}_1 & \mathbf{u}_1^* \mathbf{u}_2 & \cdots & \mathbf{u}_1^* \mathbf{u}_n \\ \mathbf{u}_2^* \mathbf{u}_1 & \mathbf{u}_2^* \mathbf{u}_2 & \cdots & \mathbf{u}_2^* \mathbf{u}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{u}_n^* \mathbf{u}_1 & \mathbf{u}_n^* \mathbf{u}_2 & \cdots & \mathbf{u}_n^* \mathbf{u}_n \end{bmatrix} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

which, in turn, is just

$$\langle \mathbf{u}_j, \mathbf{u}_k \rangle = \begin{cases} 1 & \text{if } j = k, \\ 0 & \text{if } j \neq k. \end{cases}$$

Hence, a matrix U is unitary, i.e., $U^* U = I$, if and only if the list of its column vectors is an orthonormal list. By [Corollary 2.4.2](#) on page 44 we conclude that the column vectors of a square matrix constitute an orthonormal basis of \mathbb{C}^n if and only if the matrix is unitary. This is one of the key properties of unitary matrices, and therefore also of real orthogonal matrices. We collect a few more characterizing properties in the following result.

Theorem 2.6.1

Let $U \in M_n(\mathbb{C})$. The following assertions are equivalent:

- (i) U is unitary, that is, $U^* U = I$.
- (ii) The columns of U are an orthonormal basis of \mathbb{C}^n .
- (iii) $U U^* = I$.
- (iv) U^* is unitary.
- (v) The rows of U are an orthonormal basis of \mathbb{C}^n .
- (vi) For all $\mathbf{x}, \mathbf{y} \in \mathbb{F}^n$: $\langle U \mathbf{x}, U \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$.

2.6. Unitary and orthogonal matrices

Proof. We have already argued that Item (i) \Leftrightarrow Item (ii). Recall that for two square matrices the identity $AB = I$ is equivalent to $BA = I$ (which, in turn, is equivalent to B being invertible with $B = A^{-1}$). This shows Item (i) \Leftrightarrow Item (iii) with $A = U$ and $B = U^*$. To see Item (iii) \Leftrightarrow Item (iv), we simply write Item (iii) as $(U^*)^*U^* = I$ which by definition is unitarily of U^* .

A list of vectors is an orthonormal basis for \mathbb{F}^n if and only if the list of vectors *complex conjugated* is an orthonormal basis for \mathbb{F}^n . The rows of U are the complex conjugates of the columns of U^* . Hence, from Item (i) \Leftrightarrow Item (ii) we see that Item (iv) \Leftrightarrow Item (v).

Assume now that Item (i) holds, that is, assume $U^*U = I$. Then we can compute, for $\mathbf{x}, \mathbf{y} \in \mathbb{F}^n$,

$$\langle U\mathbf{x}, U\mathbf{y} \rangle = \langle U^*U\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle,$$

which is Item (vi). Assume then Item (vi) holds. By the same computation as just before, we then have $\langle U^*U\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$. By Lemma 2.1.4 on page 33 we conclude that $U^*U\mathbf{x} = \mathbf{x}$ for all \mathbf{x} which shows that Item (i) holds. This completes the proof. ■

Using the so-called polarization identities, one can show that the conditions in Theorem 2.6.1 are also equivalent to:

(iii) For all $\mathbf{x} \in \mathbb{F}^n$: $\|U\mathbf{x}\| = \|\mathbf{x}\|$.

From Items (iii) and (vi) we say that linear maps associated with unitary matrices *preserves* the length of vectors and angle between them.

Exercise 2.6.2

Consider the vector space \mathbb{F}^n with the standard inner product $\langle \cdot, \cdot \rangle$, i.e., $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^*\mathbf{x}$. Suppose $U \in M_n(\mathbb{C})$ is unitary and let \mathbf{u}_k denote the k th column of U . As shown above $U^*U = I$ is just another way of stating that $\mathbf{u}_1, \dots, \mathbf{u}_n$ is orthonormal. Show that $UU^* = I$ is just another way of writing

$$\mathbf{x} = \sum_{k=1}^n \langle \mathbf{x}, \mathbf{u}_k \rangle \mathbf{u}_k \quad \text{for all } \mathbf{x} \in \mathbb{F}^n.$$

(Hint: $UU^* = I$ states that $UU^*\mathbf{x} = \mathbf{x}$ for all $\mathbf{x} \in \mathbb{F}^n$. Start by relating $U^*\mathbf{x}$ to $\langle \mathbf{x}, \mathbf{u}_k \rangle$.)

The equivalences in Theorem 2.6.1 can be formulated for real orthogonal matrices. Here, we focus on the most important equivalence.

Corollary 2.6.3

Let $Q \in M_n(\mathbb{R})$. The following assertions are equivalent:

- (i) Q is (real) orthogonal, that is, $Q^T Q = I$.
- (ii) The columns of Q are an orthonormal basis of \mathbb{R}^n .

Proof. Since $Q \in M_n(\mathbb{R})$ is real, it is unitary if and only if it is orthogonal. So, by Theorem 2.6.1, Q is orthogonal if and only if the list of the columns of Q is an orthonormal basis of \mathbb{C}^n . However, the column vectors are real vectors by assumption, so they also constitute an orthonormal basis of \mathbb{R}^n as the standard inner product in \mathbb{C}^n is identical to the dot product on \mathbb{R}^n for real vectors. ■

Example 2.6.2

The $n \times n$ Fourier matrix F_n is a complex matrix $F_n = [f_{i,j}]_{i,j=0}^{n-1}$ whose (i, j) entry is given by

$$f_{i,j} = \frac{1}{\sqrt{n}} \omega_n^{ij} \quad \text{where } \omega_n = e^{2\pi i/n}$$

The complex number ω_n is called the n th root of unity (try to sketch it in the complex plane). It can be shown that F_n is unitary for every $n \in \mathbb{N}$. The column vectors of F_n are called the discrete Fourier basis of \mathbb{C}^n . The linear mapping

$$\text{DFT} : \mathbb{C}^n \rightarrow \mathbb{C}^n, \text{DFT} = \mathbf{x} \mapsto F_n^* \mathbf{x}$$

is the so-called *Discrete Fourier Transform* (DFT). If n is a power of 2, i.e., $n = 2^m$ for some m , it can be computed very fast⁴ using an algorithm called the Fast Fourier Transform (FFT). This transform is extremely useful and has applications in all of engineering.

Note that if $W \in M_{k \times n}(\mathbb{F})$ has orthonormal columns, it is *only* unitary if it is square $k = n$. However, if $k < n$, that is, we have fewer orthonormal column vectors than needed to span \mathbb{F}^n , we can extend the matrix with additional columns so that the extended matrix becomes unitary. This fact is a consequence of the following simple result.

⁴To be precise: The speed of the FFT depends on the number of factors in the prime factorization of n . For $n = 2^m$ the algorithm uses a constant multiple of $n \log(n)$ floating-point operations, while standard matrix-vector multiplication uses Cn^2 operations.

Lemma 2.6.4 extension to orthonormal basis

Suppose $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_\ell$ is an orthonormal list of vectors in an inner product space V of dimension n . Then there are $n - \ell$ vectors $\mathbf{u}_{\ell+1}, \dots, \mathbf{u}_n$ so that the combined list $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ is an orthonormal basis for V .

Proof. Let $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_\ell$ is an orthonormal list. Pick $n - r$ vectors⁵ $\mathbf{v}_1, \dots, \mathbf{v}_{n-r}$ so that the combined list $\alpha := \mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_\ell, \mathbf{v}_1, \dots, \mathbf{v}_{n-r}$ is linearly independent, that is, so that $\text{span}\{\alpha\} = V$. Now, if we perform the Gram-Schmidt process on the list α , we obtain the desired orthonormal basis. ■

2.7 Diagonalizable matrices

Recall that two square matrices A and B in $M_n(\mathbb{F})$ are called *similar* if there is an invertible matrix $S \in M_n(\mathbb{F})$ such that $B = S^{-1}AS$. A complex square matrix is said to be diagonalizable (over \mathbb{C}) if and only if it is *similar* with a diagonal matrix. We are interested in the cases where S is not just invertible, but unitary. This leads to the following definitions.

Definition 2.7.1 Diagonalizable matrix

A matrix $A \in M_n(\mathbb{C})$ is:

- (i) *diagonalizable* if there is an invertible matrix $S \in M_n(\mathbb{C})$ such that $S^{-1}AS$ is diagonal.
- (ii) *unitarily diagonalizable* if there is a unitary matrix $U \in M_n(\mathbb{C})$ such that U^*AU is diagonal.
- (iii) *real orthogonally diagonalizable* if there is a real orthogonal matrix $Q \in M_n(\mathbb{R})$ such that $Q^T A Q$ is diagonal.

In [Chapter 11 in Mathematics 1a](#) it is shown that $A \in M_n(\mathbb{C})$ is diagonalizable if and only if A has n linearly independent eigenvectors which in turn happens if and only if the algebraic and geometric multiplicities agree for all eigenvalues⁶.

However, it may not be apparent directly from [Definition 2.7.1](#) that the notion of diagonalibility of a matrix is closely related to eigenvectors and,

⁵It is almost obvious that this can be done. However, it does require a proof. The proof is “easy” and can be done by induction of the definition of linear independence. We leave the details to the reader.

⁶See [Lemma 11.25 in Mathematics 1a](#) and [Corollary 11.29 in Mathematics 1a](#).

in fact, the column vectors of S , U and Q in Definition 2.7.1 constitute a basis of eigenvectors. To see this, suppose that A is diagonalizable. Then $S^{-1}AS$ is diagonal for some invertible matrix S . Let us denote this diagonal matrix by $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. From $S^{-1}AS = \Lambda$ we get $AS = S\Lambda$ by multiplying with S from the left. If we denote the i th column of S by \mathbf{s}_i , we get

$$AS = A[\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n] = [A\mathbf{s}_1, A\mathbf{s}_2, \dots, A\mathbf{s}_n] \quad (2.24)$$

and

$$S\Lambda = [\lambda_1\mathbf{s}_1, \lambda_2\mathbf{s}_2, \dots, \lambda_n\mathbf{s}_n]. \quad (2.25)$$

Since the matrix-matrix products in (2.24) and (2.25) are equal, each column must be equal, thus:

$$A\mathbf{s}_i = \lambda_i\mathbf{s}_i, \quad \text{for } i = 1, \dots, n. \quad (2.26)$$

Hence, we see that $(\lambda_i, \mathbf{s}_i)$ is actually an eigenvalue-eigenvector pair of A . Since S is assumed to be invertible, its column vectors are linearly independent, hence they constitute a basis.

If A is either unitarily diagonalizable or real orthogonal diagonalizable, the same computations carries through. We just have to use that U^* and Q^T are inverse of a unitary matrix U and a real orthogonal matrix Q , respectively. Moreover, in these cases the eigenbasis is an *orthonormal basis* for \mathbb{R}^n and \mathbb{C}^n , respectively, by Theorem 2.6.1.

2.8 The Spectral Theorem

The spectral theorem can be formulated for three classes of square matrices:

- (i) real, symmetric matrices A , that is, matrices satisfying $A = A^T$ with real entries.
- (ii) Hermitian matrices A , that is, matrices satisfying $A = A^*$ with complex entries.
- (iii) normal matrices A , that is, matrices satisfying $AA^* = A^*A$ with complex entries.

Real, symmetric matrices are Hermitian matrices, and Hermitian matrices are normal matrices. These statements are easy to prove. E.g., if A is Hermitian, it satisfies $A = A^*$ and therefore also $AA^* = AA = A^*A$ which is the condition for being normal.

Example 2.8.1

Consider the 3×3 matrices A, B, C defined by

$$A = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{bmatrix}, B = \begin{bmatrix} 2 & 3+i & 1 \\ 3-i & 2 & -i \\ 1 & i & 2 \end{bmatrix}, C = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 2 & 1 \\ 1 & 0 & 2 \end{bmatrix}.$$

The matrix A is real and symmetric; hence, A is also Hermitian and normal. The matrix B is Hermitian (and normal), but not symmetric nor real. The matrix C is neither symmetric nor Hermitian, but it is normal as can be seen by computing CC^* and C^*C and verifying that they agree.

The different versions of the spectral theorem are very similar so we will for brevity focus on the case of real and symmetric matrices. Let $A \in M_n(\mathbb{R})$ be a symmetric matrix. The spectral theorem, [Theorem 2.8.5](#), we will prove below, implies that A can be decomposed as the matrix product:

$$A = Q\Lambda Q^T \tag{2.27}$$

where $Q \in M_n(\mathbb{R})$ is an orthogonal matrix and Λ is a diagonal matrix. The formula (2.27) is called the *spectral decomposition* of A . From [Section 2.7](#) we then know that Λ contains the eigenvalues of A and that the columns of Q contain the corresponding eigenvectors.

Before we prove the spectral decomposition, let us discuss how we can compute (2.27) given that the formula is actually holds. So, let A be a real and symmetric matrix. We first compute the eigenvalues of A . Counted with algebraic multiplicity there are exactly n eigenvalues. Hence, we let $\lambda_1, \dots, \lambda_n$ be the eigenvalues of A . Note that the value of λ_i in this list has been repeated according to the algebraic multiplicity of the eigenvalue. We define $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. As we will see a real, symmetric matrix A has n linearly independent eigenvectors (recall that this is not the case for all $n \times n$ matrices). The list of eigenvectors should be ordered to match the eigenvalues so that the i th eigenvector has eigenvalue λ_i . After applying the Gram-Schmidt procedure to this list of n linearly independent eigenvectors, we end up with an orthonormal basis of eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$. Setting $Q = [\mathbf{q}_1, \dots, \mathbf{q}_n]$ provides us with the spectral decomposition $A = Q\Lambda Q^T$.

Lemma 2.8.1

If $A \in M_n(\mathbb{C})$ is Hermitian, i.e., $A = A^*$, then all eigenvalues of A are real.

Proof. Let $\lambda \in \mathbb{C}$ be an eigenvalue of a Hermitian matrix A , and let $\mathbf{v} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$ be an associated eigenvector. Hence, $A\mathbf{v} = \lambda\mathbf{v}$. Then,

since $A = A^*$,

$$\lambda \langle \mathbf{v}, \mathbf{v} \rangle = \langle \lambda \mathbf{v}, \mathbf{v} \rangle = \langle A\mathbf{v}, \mathbf{v} \rangle = \langle \mathbf{v}, A^*\mathbf{v} \rangle = \langle \mathbf{v}, A\mathbf{v} \rangle = \langle \mathbf{v}, \lambda \mathbf{v} \rangle = \bar{\lambda} \langle \mathbf{v}, \mathbf{v} \rangle.$$

Since $\|\mathbf{x}\|^2 = \langle \mathbf{x}, \mathbf{x} \rangle$ for all vectors \mathbf{x} , we have shown that $\lambda \|\mathbf{v}\|^2 = \bar{\lambda} \|\mathbf{v}\|^2$. An eigenvector \mathbf{v} cannot be the zero vector so $\|\mathbf{v}\| \neq 0$. Hence, we can divide through by $\|\mathbf{v}\|^2$, and we arrive at $\lambda = \bar{\lambda}$, which shows that λ is real. ■

Eigenvectors of a Hermitian matrix have the following orthogonality property: two eigenvectors associated with two *different* eigenvalues are orthogonal.

Exercise 2.8.2

Suppose $A \in M_n(\mathbb{C})$ is Hermitian, i.e., $A = A^*$. Let λ_1 and λ_2 be eigenvalues of A with associated eigenvectors \mathbf{u}_1 and \mathbf{u}_2 , respectively. Show that if $\lambda_1 \neq \lambda_2$, then $\mathbf{u}_1 \perp \mathbf{u}_2$, i.e., $\langle \mathbf{u}_1, \mathbf{u}_2 \rangle = 0$

A real and symmetric matrix is automatically Hermitian since $A = A^T = A^*$. Hence, real, symmetric matrices inherits the same properties as Hermitian matrix, in particular, they have real eigenvalues by Lemma 2.8.1. However, not only will a real and symmetric matrix have real eigenvalues, it will also have a real eigenvector:

Lemma 2.8.3

A real, symmetric matrix A has a real eigenvector.

Proof. The characteristic polynomial of A is a degree n polynomial and by the fundamental theorem of algebra, we know that it has at least one root λ_1 . By Lemma 2.8.1 we know that λ_1 is real. The fact that λ_1 is an eigenvalue means that $A - \lambda_1 I$ is singular, hence there is a non-zero vector \mathbf{q}_1 such that $A\mathbf{q}_1 = \lambda_1 \mathbf{q}_1$. We need to prove that we can find a *real* eigenvector of A associated with λ_1 . For the (possibly complex) non-zero eigenvector \mathbf{q}_1 we write $\mathbf{q}_1 = \text{Re } \mathbf{q}_1 + i \text{Im } \mathbf{q}_1$, where we take the real and imaginary parts coordinate wise. The equation $A\mathbf{q}_1 = \lambda_1 \mathbf{q}_1$ holds if and only if the real and imaginary parts of the left and right hand side agrees. Since both A and λ_1 are real, we arrive at the two equations:

$$A \text{Re } \mathbf{q}_1 = \lambda_1 \text{Re } \mathbf{q}_1 \quad \text{and} \quad A \text{Im } \mathbf{q}_1 = \lambda_1 \text{Im } \mathbf{q}_1$$

Note that both $\text{Re } \mathbf{q}_1$ and $\text{Im } \mathbf{q}_1$ are real vectors and that they cannot both be the zero vector (as then \mathbf{q}_1 would be the zero vector). Hence, we have shown the existence of at least one real eigenvector. ■

Exercise 2.8.4

Find the eigenvalues and eigenvectors of the Hermitian matrix $A = \begin{bmatrix} 0 & i \\ -i & 0 \end{bmatrix}$. The eigenvalues of A are real by Lemma 2.8.1, but are there any real eigenvectors of A ?

We are now finally ready to state and prove the spectral theorem for real, symmetric matrices.

Theorem 2.8.5 Spectral Theorem (the real case)

Let $A \in M_n(\mathbb{R})$. The following assertions are equivalent:

- (i) A is a symmetric matrix.
- (ii) A is real orthogonally diagonalizable, that is, there is a real orthogonal matrix $Q \in M_n(\mathbb{R})$ such that $Q^T A Q$ is diagonal.
- (iii) \mathbb{R}^n has an orthonormal basis consisting of eigenvectors of A .

Proof. (i) \Rightarrow (ii): Suppose A is real and symmetric. We need to show that A is real orthogonally diagonalizable. The proof is by induction on the size n of the matrix A .

The result is trivial for $n = 1$: The matrix $A \in M_1(\mathbb{R})$ is just a scalar, i.e., $A = [a_{1,1}]$ for $a_{1,1} \in \mathbb{R}$. Hence, $\lambda_1 = a_{1,1}$ is the eigenvalue of A and with $Q = [1]$ or $Q = [-1]$, we have $A = Q \Lambda Q^T$. Note that both $Q = [1]$ or $Q = [-1]$ are real orthogonal 1×1 matrices.

Now, we turn to the induction step. So, let $n > 1$ and assume the implication (i) \Rightarrow (ii) is true for any matrix of size $n - 1$.

Let λ_1 be an eigenvalue of A . By Lemma 2.8.1 it is real, and by Lemma 2.8.3, we can find a real eigenvector \mathbf{u}_1 associated with λ_1 , i.e., $A\mathbf{u}_1 = \lambda_1\mathbf{u}_1$. If \mathbf{u}_1 is not normalized, i.e., of length one, we can replace it by $\mathbf{u}_1/\|\mathbf{u}_1\|$ which is an eigenvector of norm one. Hence, we can assume that $\|\mathbf{u}_1\| = 1$. Note that there is no risk of division by zero since eigenvectors are always non-zero.

By a result in Lemma 2.6.4 on page 53 we can extend \mathbf{u}_1 to an orthonormal basis for \mathbb{R}^n . The first element in this list is still \mathbf{u}_1 since it was a normalized vector, but the remaining $n - 1$ vectors are not necessarily eigenvectors of A . We let the last $n - 1$ vectors in the orthonormal basis be columns in an $n \times (n - 1)$ matrix V_1 . Then $[\mathbf{u}_1, V_1]$ is an orthogonal matrix since the columns are an orthonormal basis of \mathbb{R}^n .

The matrix $V_1^T AV_1$ is symmetric since it is equal to its transpose:

$$(V_1^T AV_1)^T = V_1^T A^T (V_1^T)^T = V_1^T AV_1 \quad (2.28)$$

So, by our induction hypothesis, the $(n-1) \times (n-1)$ symmetric matrix $V_1^T AV_1$ is real orthogonally diagonalizable, that is, there is a real orthogonal $(n-1) \times (n-1)$ -matrix Q_1 such that $Q_1^T V_1^T AV_1 Q_1$ is diagonal. We denote this diagonal matrix by Λ_1 , i.e.,

$$\Lambda_1 = Q_1^T V_1^T AV_1 Q_1. \quad (2.29)$$

We define the $n \times (n-1)$ matrix $U_1 := V_1 Q_1$. By construction the columns of U_1 are linear combinations of the columns of V_1 ; we can express this as $\text{col } U_1 \subseteq \text{col } V_1$ ⁷. Since \mathbf{u}_1 is orthogonal to all columns in V_1 , it is therefore also orthogonal to all columns of U_1 . The conclusion that \mathbf{u}_1 is orthogonal to all columns in U_1 can be summarized as

$$U_1^T \mathbf{u}_1 = \mathbf{0}_{n-1}, \quad (2.30)$$

where $\mathbf{0}_{n-1}$ is the column vector of $n-1$ zeros. We claim that the matrix $U := [\mathbf{u}_1, U_1]$ is real orthogonal. To see this claim, we have to show that $U^T U = I$, but this follows from the following computation:

$$\begin{aligned} U^T U &= \begin{bmatrix} \mathbf{u}_1^T \\ U_1^T \end{bmatrix} \begin{bmatrix} \mathbf{u}_1 & U_1 \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{u}_1^T \mathbf{u}_1 & \mathbf{u}_1^T U_1 \\ U_1^T \mathbf{u}_1 & U_1^T U_1 \end{bmatrix} = \begin{bmatrix} 1 & \mathbf{0}_{n-1}^T \\ \mathbf{0}_{n-1} & I_{n-1} \end{bmatrix} = I_n, \end{aligned}$$

since $U_1^T U_1 = Q_1^T V_1^T V_1 Q_1 = Q_1^T I_{n-1} Q_1 = I_{n-1}$.

From equation (2.30), we see that $U_1^T A \mathbf{u}_1 = \lambda_1 U_1^T \mathbf{u}_1 = \lambda_1 \mathbf{0}_{n-1} = \mathbf{0}_{n-1}$. Note also that Equation (2.29) reads $\Lambda_1 = U_1^T A U_1$. Using these relations, we can compute:

$$\begin{aligned} U^T A U &= U^T A \begin{bmatrix} \mathbf{u}_1 & U_1 \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1^T \\ U_1^T \end{bmatrix} \begin{bmatrix} A \mathbf{u}_1 & A U_1 \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{u}_1^T A \mathbf{u}_1 & \mathbf{u}_1^T A U_1 \\ U_1^T A \mathbf{u}_1 & U_1^T A U_1 \end{bmatrix} = \begin{bmatrix} \lambda_1 & \mathbf{0}_{n-1}^T \\ \mathbf{0}_{n-1} & \Lambda_1 \end{bmatrix}, \end{aligned}$$

where in the last step we have used that $\mathbf{u}_1^T A \mathbf{u}_1 = \langle A \mathbf{u}_1, \mathbf{u}_1 \rangle = \langle \lambda_1 \mathbf{u}_1, \mathbf{u}_1 \rangle = \lambda_1 \|\mathbf{u}_1\|^2 = \lambda_1$, and

$$\mathbf{u}_1^T A U_1 = (U_1^T A \mathbf{u}_1)^T = \mathbf{0}_{n-1}^T$$

⁷In fact, $\text{col } U_1 = \text{col } V_1$, but we will not need this property.

2.8. The Spectral Theorem

Hence, we have exhibited an orthogonal matrix $U \in \mathbf{M}_n(\mathbb{R})$ such that $U^T A U$ is diagonal. This completes the induction proof.

(ii) \Rightarrow (iii): We assume that there is a real orthogonal matrix $Q \in \mathbf{M}_n(\mathbb{R})$ such that $Q^T A Q$ is diagonal $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Write $Q = [\mathbf{q}_1, \dots, \mathbf{q}_n]$. As in (2.26), we then have:

$$A\mathbf{q}_i = \lambda_i \mathbf{q}_i, \quad \text{for } i = 1, \dots, n.$$

So $\mathbf{q}_1, \dots, \mathbf{q}_n$ are eigenvectors of A . Since Q is real orthogonal, the list $\mathbf{q}_1, \dots, \mathbf{q}_n$ is an orthonormal basis (of eigenvectors).

(iii) \Rightarrow (i): Suppose that $\mathbf{q}_1, \dots, \mathbf{q}_n$ are an orthonormal basis of \mathbb{R}^n and let $\lambda_1, \dots, \lambda_n$ be the associated eigenvalues, i.e.,

$$A\mathbf{q}_i = \lambda_i \mathbf{q}_i \text{ for } i = 1, \dots, n. \tag{2.31}$$

Then $Q = [\mathbf{q}_1, \dots, \mathbf{q}_n]$ is unitary, and we can write Equation (2.31) as matrix-matrix products as in Section 2.7:

$$A Q = Q \Lambda.$$

Since $Q^{-1} = Q^T$, we can express A as $A = Q \Lambda Q^T$. From this expression of A , it is easy to see that it is symmetric:

$$A^T = (Q \Lambda Q^T)^T = (Q^T)^T \Lambda^T Q^T = Q \Lambda Q^T = A.$$

■

Corollary 2.8.6

Let $A \in \mathbf{M}_n(\mathbb{R})$ be a symmetric matrix. Then *the spectral decomposition*

$$A = Q \Lambda Q^T$$

holds, where Λ is a diagonal matrix whose diagonal consists of the eigenvalues of A (repeated according to their algebraic multiplicity) and where $Q \in \mathbf{M}_n(\mathbb{R})$ is a real orthogonal matrix, whose columns are eigenvectors of A .

Proof. From 2.8.5 (ii) we have that $Q^T A Q$ is diagonal. Denote this diagonal matrix by $\Lambda := Q^T A Q$. Then

$$Q \Lambda Q^T = Q Q^T A Q Q^T = I_n A I_n = A$$

since Q is real orthogonal, i.e., $Q Q^T = I$. ■

The spectral decomposition shows how we pick the simplest possible matrix representations of linear mappings defined through real, symmetric matrices by changing basis to the orthonormal eigenvector basis.

Example 2.8.2

Let $A \in M_n(\mathbb{R})$ be a symmetric matrix and define the linear map $L_A : \mathbb{R}^n \rightarrow \mathbb{R}^n$ by $L_A = \mathbf{x} \mapsto A\mathbf{x}$, i.e., $L_A\mathbf{x} = A\mathbf{x}$.

If we let e denote the standard basis for \mathbb{R}^n , i.e., $e = \mathbf{e}_1, \dots, \mathbf{e}_n$, then the matrix representation of L_A is simply ${}_e[L_A]_e = A$. However, the spectral decomposition provide us with a much simpler matrix representation of L_A as follows:

Let $A = Q\Lambda Q^T$ be the spectral decomposition of A , where $Q = [\mathbf{q}_1, \dots, \mathbf{q}_n]$ contains the orthonormal basis of eigenvectors. Let $q := \mathbf{q}_1, \dots, \mathbf{q}_n$. Recall from [Lemma 11.10 in Mathematics 1a](#):

$${}_e[L_A]_e = {}_eM_q {}_q[L_A]_q {}_qM_e,$$

where ${}_eM_q$ is the change-of-basis matrix from q -basis to e -basis. This is in fact just the spectral decomposition of A in disguise. We first note that $Q = {}_eM_q$ and $Q^T = Q^{-1} = {}_qM_e$. Hence, the spectral decomposition of A can be written

$${}_e[L_A]_e = A = Q\Lambda Q^T = {}_eM_q \Lambda {}_qM_e.$$

We conclude that ${}_q[L_A]_q = \Lambda$.

The spectral theorem can also be used to simplify the study of functions, in particular, quadratic forms. Consider a quadratic form $q : \mathbb{R}^n \rightarrow \mathbb{R}$. The following simple result says that we may always assume that the matrix $A \in M_n(\mathbb{R})$ in [Equation \(1.3\)](#) on page 11 is *symmetric*.

Lemma 2.8.7

Let $q : \mathbb{R}^n \rightarrow \mathbb{R}$ be a quadratic form:

$$q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + \mathbf{x}^T \mathbf{b} + c,$$

where $A \in M_n(\mathbb{R})$, $\mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$. Let $B := (A + A^T)/2$. Then B is symmetric and

$$q(\mathbf{x}) = \mathbf{x}^T B \mathbf{x} + \mathbf{x}^T \mathbf{b} + c$$

Exercise 2.8.8

Prove [Lemma 2.8.7](#).

Example 2.8.3 Reduction of quadratic forms

Let $q : \mathbb{R}^n \rightarrow \mathbb{R}$ be a quadratic form. By [Lemma 2.8.7](#) we can write q as

$$q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + \mathbf{x}^T \mathbf{b} + c$$

where $A \in M_n(\mathbb{R})$ is symmetric, $\mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$.

By Corollary 2.8.6, we have the spectral decomposition

$$A = Q \Lambda Q^T,$$

where Λ is a diagonal matrix whose diagonal consists of the eigenvalues of A (repeated according to their algebraic multiplicity) and where $Q \in M_n(\mathbb{R})$ is a real orthogonal matrix. The column vectors \mathbf{q}_i of Q constitute an orthonormal basis of \mathbb{R}^n , i.e., $q = \mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$ is an orthonormal basis and Q^T is the change-of-basis matrix from e -basis to standard q -basis.

We consider the column vector \mathbf{x} as the coordinates of the vector \mathbf{x} with respect to the e -basis. However, there is of course no difference on \mathbf{x} and ${}_e\mathbf{x}$, but if we now define $\mathbf{y} = Q^T \mathbf{x}$ we see that $\mathbf{y} = {}_q\mathbf{x}$ represents the coordinates of \mathbf{x} with respect to the eigenvector basis q . That is, given $\mathbf{x} \in \mathbb{R}^n$, we write as $\mathbf{x} = y_1 \mathbf{q}_1 + y_2 \mathbf{q}_2 + \dots + y_n \mathbf{q}_n$ instead of $\mathbf{x} = x_1 \mathbf{e}_1 + x_2 \mathbf{e}_2 + \dots + x_n \mathbf{e}_n$.

Let us see why this change-of-basis is useful. Consider the second order terms:

$$\mathbf{x}^T A \mathbf{x} = \mathbf{x}^T Q \Lambda Q^T \mathbf{x} = \mathbf{y}^T \Lambda \mathbf{y} = \lambda_1 y_1^2 + \lambda_2 y_2^2 + \dots + \lambda_n y_n^2,$$

hence there are no cross-terms of the form $y_i y_j$ ($i \neq j$) in the q -basis.

We can therefore write q as

$$q(\mathbf{x}) = \lambda_1 y_1^2 + \lambda_2 y_2^2 + \dots + \lambda_n y_n^2 + \mathbf{y}^T Q^T \mathbf{b} + c$$

where $\mathbf{y} = {}_q\mathbf{x}$. Hence, we can analyze quadratic forms in a simplified form:

$$\lambda_1 y_1^2 + \lambda_2 y_2^2 + \dots + \lambda_n y_n^2 + \mathbf{y}^T \mathbf{b}' + c$$

where $\mathbf{b}' := Q^T \mathbf{b}$. This reduced expression is exactly the quadratic forms considered in (1.14), for which we described its levelsets in Example 1.4.7 on page 20.

Finally, with a little more work, we can prove the complex version of the spectral theorem. We leave the proof to the reader.

Theorem 2.8.9 Spectral Theorem (the complex case)

Let $A \in M_n(\mathbb{C})$. The following assertions are equivalent:

- (i) A is a normal matrix.

- (ii) A is unitarily diagonalizable, that is, there is a unitary matrix $U \in M_n(\mathbb{C})$ such that U^*AU is diagonal.
- (iii) \mathbb{C}^n has an orthonormal basis consisting of eigenvectors of A .

Proof. Most of the proof of Theorem 2.8.5 generalize to the complex case. The only real difficulty is proving that $V_1^*AV_1$, see (2.28), is normal which is needed in order to apply our induction hypothesis. It requires a bit of work, and we leave the details to the reader. ■

From Theorem 2.8.9 we obtain a spectral decomposition of normal matrices A . The proof of the following result is essentially identical to the proof of Corollary 2.8.6 on page 59 so we do not repeat it.

Corollary 2.8.10

Let $A \in M_n(\mathbb{C})$ be a normal matrix. Then *the spectral decomposition*

$$A = U\Lambda U^*$$

holds, where Λ is a diagonal matrix whose diagonal consists of the eigenvalues of A (repeated according to their algebraic multiplicity) and where $U \in M_n(\mathbb{C})$ is a unitary matrix, whose columns are eigenvectors of A .

2.9 Postive definite and semi-definite matrices

Hermitian matrices have, as we have seen, real eigenvalues. Hermitian matrices with only *nonnegative* or *positive* eigenvalues are of particular interest, and they are the subject of this chapter. They can be seen as the matrix counterparts to nonnegative or positive real numbers. These matrices frequently appear in various fields: in statistics, they emerge as correlation matrices and in the normal equations for least squares fitting problems; in Lagrangian mechanics, they represent the kinetic energy functional; in quantum mechanics, they are crucial as density matrices; and in mathematics, they characterize inner products.

Definition 2.9.1

Let $A \in M_n(\mathbb{F})$.

- (i) A is positive definite if A is Hermitian and $\langle A\mathbf{x}, \mathbf{x} \rangle > 0$ for all nonzero $\mathbf{x} \in \mathbb{F}^n$.

2.9. Postive definite and semi-definite matrices

- (ii) A is positive semi-definite if A is Hermitian and $\langle A\mathbf{x}, \mathbf{x} \rangle \geq 0$ for all $\mathbf{x} \in \mathbb{F}^n$.
- (iii) A is negative definite if $-A$ is positive definite.
- (iv) A is negative semi-definite if $-A$ is positive semi-definite.

Since $\langle A\mathbf{0}, \mathbf{0} \rangle = 0$, we see that a positive definite matrix A is also a positive semi-definite matrix. Note also that the inequalities in Definition 2.9.1 implicitly assume that $\langle A\mathbf{x}, \mathbf{x} \rangle$ is a real scalar. This potential issue is solved by assuming A is Hermitian since $\langle A\mathbf{x}, \mathbf{x} \rangle \in \mathbb{R}$ whenever A is Hermitian.

We already know several examples of positive semi-definite matrices. The zero $n \times n$ matrix and the projection matrix (from Example 2.3.1 on page 41) are positive semi-definite matrices, and the identity matrix is positive definite since $\langle \mathbf{x}, \mathbf{x} \rangle > 0$ for all nonzero $\mathbf{x} \in \mathbb{F}^n$. It is easy to construct other (and less trivial) positive semi-definite matrices as we will see in the next example.

Example 2.9.1

Let $B \in M_{m \times n}(\mathbb{F})$. We claim that both $BB^* \in M_m(\mathbb{F})$ and $B^*B \in M_n(\mathbb{F})$ are positive semi-definite matrices. Let us prove this claim for $A = B^*B$. The matrix A is Hermitian since

$$A^* = (B^*B)^* = B^*(B^*)^* = B^*B = A.$$

From the computation, using that the norm is always non-negative,

$$\langle A\mathbf{x}, \mathbf{x} \rangle = \langle B^*B\mathbf{x}, \mathbf{x} \rangle = \langle B\mathbf{x}, (B^*)^*\mathbf{x} \rangle = \langle B\mathbf{x}, B\mathbf{x} \rangle = \|B\mathbf{x}\|^2 \geq 0,$$

we see that A is a positive semi-definite matrix.

The following characterization of positive definite (positive semi-definite) is a consequence of the Spectral Theorem.

Theorem 2.9.1

Let $A \in M_n(\mathbb{F})$ be Hermitian. Then the following are equivalent⁸:

- (i) A is positive definite (positive semi-definite)
- (ii) The eigenvalues of A are positive (non-negative)

2.9. Postive definite and semi-definite matrices

- (iii) There exists a positive (non-negative) constant c such that $\langle A\mathbf{x}, \mathbf{x} \rangle \geq c\langle \mathbf{x}, \mathbf{x} \rangle$ for all $\mathbf{x} \in \mathbb{F}^n$.

Proof. We will only prove that A is positive semi-definite if and only if A is Hermitian and all its eigenvalues are non-negative and leave the rest of the proof to the reader. If A is positive semi-definite, then A is Hermitian, and for an eigenvalue-eigenvector pair (λ, \mathbf{u}) , where we have normalized \mathbf{u} , we have

$$\lambda = \lambda\langle \mathbf{u}, \mathbf{u} \rangle = \langle \lambda\mathbf{u}, \mathbf{u} \rangle = \langle A\mathbf{u}, \mathbf{u} \rangle \geq 0.$$

Suppose A is Hermitian with non-negative eigenvalues $\lambda_1, \dots, \lambda_n$. Let $A = U\Lambda U^*$ be the spectral decomposition of A , where U is a unitary matrix and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Define $B = U\Gamma U^*$, where $\Gamma = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$. Then

$$B^*B = (U\Gamma U^*)^*U\Gamma U^* = U\Gamma^*U^*U\Gamma U^* = U\Gamma^2U^* = U\Lambda U^* = A$$

From [Example 2.9.1](#) it follows that A is positive semi-definite. ■

If $A \in M_n(\mathbb{F})$ is Hermitian, then it is normal. From the Spectral Theorem it follows that $A = U\Lambda U^*$ for some unitary matrix U and diagonal matrix Λ . [Theorem 2.9.1](#) tells us that the matrix A is positive definite precisely when all the eigenvalues in the diagonal of Λ are positive.

⁸There are two theorems stated here. The equivalence holds either for all the statements in the parentheses or for all the statements without the parentheses.

CHAPTER 3

Continuity and Differentiability

The current chapter forms the core of the book, and we will see that the introduced topics play crucial roles in the later chapters.

We expect the reader is familiar with calculus of scalar functions of one variable, but we will briefly review some of the theory in [Section 3.1](#) for the reader's convenience. As for functions of a single real variable, the topics of continuity and differentiability are central in the analysis of functions of several variables. We start by considering continuity in [Section 3.2](#); the good news here is that having defined continuity of functions of one variable in a precise way, we can easily extend the definition to functions of several variables. [Section 3.3](#) introduces the concept of partial derivatives for functions of several variables, in much the same way as differentiability is defined for functions of one variable. This is generalized in [Section 3.5](#), where we consider partial derivatives of higher order and the so-called Hessian matrix.

The material in [Sections 3.3–3.5](#) unfortunately shows that some of the standard results for functions of one variable do not immediately generalize to functions of several variables. For this reason we introduce the topic of differentiability in [Section 3.6](#). [Section 3.7](#) discusses the so-called chain-rule, an important result that generalizes the well-known results for differentiation of a composed functions to the case of functions of several variables. Finally, in [Section 3.8](#), we generalize the previous discussions of differentiability to vector functions of several variables.

3.1 Analysis of functions of one variable

The mathematical analysis of functions is based on the concepts *continuity* and *differentiability*. We first consider functions of the form $f : I \rightarrow \mathbb{R}$,

3.1. Analysis of functions of one variable

where either $I = \mathbb{R}$ or I is an open interval of the form $] - \infty, a[$, $]a, \infty[$ or $]a, b[$ for some $a, b \in \mathbb{R}$.

Definition 3.1.1 Continuity

A function $f : I \rightarrow \mathbb{R}$ is continuous at $x_0 \in I$ if

$$f(x) \rightarrow f(x_0) \text{ whenever } x \rightarrow x_0. \quad (3.1)$$

For short, the function f is said to be *continuous* if it is continuous at all $x_0 \in I$.

The condition (3.1) can alternatively be stated as

$$|f(x) - f(x_0)| \rightarrow 0 \text{ whenever } |x - x_0| \rightarrow 0. \quad (3.2)$$

It is phrased as “ $f(x)$ tends to $f(x_0)$ if x tends to x_0 ”. Intuitively, this means that the function values $f(x)$ are “close” to the function value $f(x_0)$ for points x that are “close” to x_0 . The mathematically exact way of formulating the definition is:

For any given $\epsilon > 0$ there exists a $\delta > 0$ such that if $x \in I$ satisfies $|x - x_0| < \delta$ then $|f(x) - f(x_0)| < \epsilon$ holds.

We can write this much shorter using the so-called universal quantifier \forall (which is read “for all”) and the existential quantifier \exists (which is read “there exists”).

$$\forall \epsilon > 0 \exists \delta > 0 \forall x \in I : |x - x_0| < \delta \Rightarrow |f(x) - f(x_0)| < \epsilon. \quad (3.3)$$

It is well-known that, e.g., polynomials, exponential functions, and trigonometric functions are continuous within their domains. We also know that standard algebraic manipulations on continuous functions again produce continuous functions. To be more precise,

- If f and g are continuous functions defined on the same domain, the functions $f + g$ and fg are continuous;
- If f and g are continuous functions and the range of g is contained in the domain of f , the composed function $h(x) := f(g(x))$ is continuous. Recall that we denote the composition by $h = f \circ g$.

Thus, in most cases continuity of a given function can be checked without referring explicitly to Definition 3.1.1.

Example 3.1.1

The function

$$f(x) = x^2 + 3 \cos(x^4), \quad x \in \mathbb{R},$$

is continuous, because all the involved functions x^2, x^4 , and $\cos(x)$ are continuous.

Definition 3.1.2 Differentiability

A function $f : I \rightarrow \mathbb{R}$ is differentiable at $x_0 \in I$ if the expression

$$\frac{\Delta f}{h} := \frac{f(x_0 + h) - f(x_0)}{h} \quad (3.4)$$

has a limit as $h \rightarrow 0$; if this is the case, the limit is denoted by $f'(x_0)$, i.e.,

$$f'(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}. \quad (3.5)$$

For short, the function f is said to be differentiable if it is differentiable at all $x_0 \in I$. If f is differentiable and $f'(x_0)$ depends continuously on x_0 , the function f is said to be continuously differentiable. The value $f'(x_0)$ is called the derivative of the function f at the point x_0 .

As continuity, the concept of differentiability behaves in a convenient way regarding the standard mathematical operations: a sum, a product, or a composition of two differentiable functions $f, g : \mathbb{R} \rightarrow \mathbb{R}$ is again differentiable, and the derivatives are given by:

$$(f + g)'(x) = f'(x) + g'(x), \quad (3.6)$$

$$(fg)'(x) = f(x)g'(x) + f'(x)g(x), \quad (3.7)$$

and

$$(f \circ g)'(x) = f'(g(x))g'(x). \quad (3.8)$$

Example 3.1.2

Let $c \in \mathbb{R}$ denote any constant, and consider the function

$$h : \mathbb{R} \rightarrow \mathbb{R}, \quad h(x) := \sin(x^2 + c).$$

Then $h(x) = (f \circ g)(x)$ with $f(x) = \sin(x)$ and $g(x) = x^2 + c$. Thus h is differentiable, and

$$h'(x) = 2x \cos(x^2 + c).$$

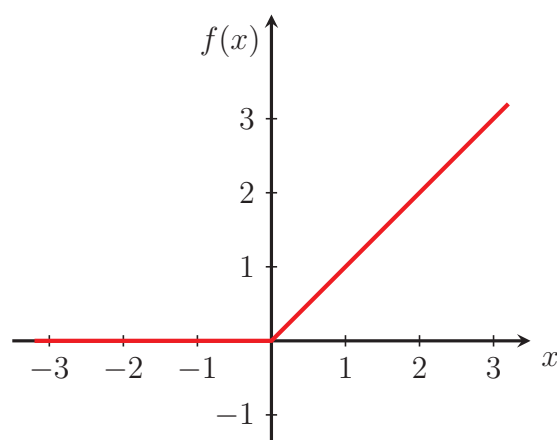


Figure 3.1: The graph of the function $f(x) = \text{ReLU}(x)$.

Note that in order to check that a function $f : I \rightarrow \mathbb{R}$ is differentiable, we need to examine differentiability at each single point $x_0 \in I$. There exist functions that are differentiable at certain points, but not in all points:

Example 3.1.3

The *rectified linear unit* function used as a so-called activation function in neural networks in machine learning is given by:

$$\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}, \quad \text{ReLU}(x) = \max(x, 0).$$

We can write the function on explicit form as

$$\text{ReLU}(x) = \begin{cases} x & \text{for } x \geq 0, \\ 0 & \text{for } x < 0. \end{cases}$$

The function ReLU is differentiable for all $x \neq 0$, but not for $x = 0$; see Figure 3.1.

Vector-valued functions of one variable. We now turn to *vector* functions of one variable. These are functions of the form $\mathbf{f} : I \rightarrow \mathbb{R}^k$, where I is an interval of \mathbb{R} . Such functions are used to parameterize curves in \mathbb{R}^k as we discussed in [Item \(e\)](#) on page 14. We have indeed already met curves, e.g., in our discussion of the boundary of a set in \mathbb{R}^2 , see [Example 2.2.4](#) and [Example 2.2.5](#). Here are a few more illustrative examples to keep in mind. As is typical when considering curves, we use the symbols \mathbf{r} and $t \in I$ instead of \mathbf{f} and $x \in I$.

Example 3.1.4 (a) Consider the function

$$\mathbf{r} : [0, 2] \rightarrow \mathbb{R}^2, \mathbf{r}(t) := \begin{bmatrix} t \\ t \end{bmatrix} = t \begin{bmatrix} 1 \\ 1 \end{bmatrix}. \quad (3.9)$$

Note that by letting the parameter t traverse through the interval $[0, 2]$, the vectors $\mathbf{r}(t)$ traverse through the straight line starting at the point $\mathbf{r}(0) = (0, 0)$ and ending at the point $\mathbf{r}(2) = (2, 2)$.

(b) Similarly, we can consider the function

$$\mathbf{r} : [0, \pi] \rightarrow \mathbb{R}^2, \mathbf{r}(t) := \begin{bmatrix} \cos t \\ \sin t \end{bmatrix}. \quad (3.10)$$

When we let the parameter t traverse through the interval $[0, \pi]$, the vectors $\mathbf{r}(t)$ traverse through “half the unit circle”, starting at the point $\mathbf{r}(0) = (1, 0)$ and ending at the point $\mathbf{r}(\pi) = (-1, 0)$.

More generally, the concept of a function $\mathbf{r} : I \rightarrow \mathbb{R}^k$ in a natural fashion leads to the concept of a *curve*:

Definition 3.1.3 Curve in \mathbb{R}^k

Let I denote an interval in \mathbb{R} , and consider a function $\mathbf{r} : I \rightarrow \mathbb{R}^k$. The image of \mathbf{r}

$$\text{im}(\mathbf{r}) = \{ \mathbf{r}(t) \mid t \in I \} \quad (3.11)$$

is called a curve in \mathbb{R}^k . The function \mathbf{r} is said to provide a parametrization of the curve, and the variable $t \in I$ is called the parameter of \mathbf{r} . Curves will often be denoted by \mathcal{C} .

Note that a curve \mathcal{C} can be parameterized in different ways; that is, a set of the form (3.11) will occur by different choices of the function \mathbf{r} . For example, the function

$$\mathbf{r} : [0, 1] \rightarrow \mathbb{R}^2, \mathbf{r}(t) := \begin{bmatrix} 2t \\ 2t \end{bmatrix} \quad (3.12)$$

yields exactly the same set of points in \mathbb{R}^2 as the function in (3.9) — we just “travel through the curve with double speed”. When we speak about a curve, we will always consider a *fixed* choice of a parametrization.

As in Equation (1.5) on page 13 we write $\mathbf{r}(t) = (r_1(t), r_2(t), \dots, r_k(t))$ using the coordinate functions $r_i : I \rightarrow \mathbb{R}$. The coordinate functions are

3.1. Analysis of functions of one variable

scalar-valued functions of one variable, which we can analyze using the result presented in the current section. E.g., if each coordinate function is differentiable, we can differentiate \mathbf{r} by $\mathbf{r}'(t) := (r'_1(t), r'_2(t), \dots, r'_k(t))$ for $t \in I$.

Definition 3.1.4 Tangent vector

Let $\mathbf{r} : I \rightarrow \mathbb{R}^k$. We say that \mathbf{r} is (continuously) differentiable if all the functions r_1, r_2, \dots, r_k are (continuously) differentiable^a. The vector $\mathbf{r}'(t) = (r'_1(t), r'_2(t), \dots, r'_k(t))$ is the *tangent vector* at $\mathbf{r}(t)$. If $\mathbf{r}'(t) \neq \mathbf{0}$ for all $t \in I$, the parametrization \mathbf{r} is *regular*.

^aIf I is not an open interval, we assume that r_1, r_2, \dots, r_k are continuously differentiable on an open interval containing I .

In Definition 3.1.3, the starting point is the function \mathbf{r} from which, we obtain a curve. In applications we often face the opposite scenario, where a curve is given and we subsequently want to find a parametrization. In the next example we describe a parametrization of the straight line connecting two given points in \mathbb{R}^k .

Example 3.1.5

Given arbitrary points $\mathbf{a}, \mathbf{b} \in \mathbb{R}^k$, the straight line starting at \mathbf{a} and ending at \mathbf{b} can be parameterized as

$$\mathbf{r} : [0, 1] \rightarrow \mathbb{R}^k, \quad \mathbf{r}(t) = \mathbf{a} + t(\mathbf{b} - \mathbf{a}) = (1 - t)\mathbf{a} + t\mathbf{b}.$$

Example 3.1.6

We consider the parametrization of the half circle in (3.10). Since $\mathbf{r}'(t) = (-\sin t, \cos t) \neq (0, 0)$ for all $t \in [0, \pi]$, the parametrization is continuously differentiable and regular. Fix $t_0 \in I$. The *tangent line* or just tangent at $\mathbf{r}(t_0)$ can then be written as

$$\mathbf{r}(t_0) + s\mathbf{r}'(t_0) = (\cos t_0, \sin t_0) + s(-\sin t_0, \cos t_0), \quad s \in \mathbb{R}.$$

E.g., with $t_0 = \pi/2$, we find the tangent at the mid-point $\mathbf{r}(\pi/2) = (0, 1)$ to be $(0, 1) + s(-1, 0) = (-s, 1)$ for $s \in \mathbb{R}$.

3.2 Continuity of vector functions of several variables

We are now ready to start the approach of generalizing the key concepts in mathematical analysis from functions of one variable to functions of n variables. We begin with continuity.

Recall the way we described continuity for a function of one variable, $f : I \rightarrow \mathbb{R}$, in (3.2). Here, the term $|x - x_0|$ measures the distance from the point $x \in I$ to the fixed point $x_0 \in I$. In \mathbb{R}^n we measure the distance between two points \mathbf{x} and \mathbf{x}_0 by the norm $\|\mathbf{x} - \mathbf{x}_0\|$ introduced in Equation (2.1) on page 24. In fact, the norm $\|\cdot\|$ on \mathbb{R}^n for $n = 1$ is just the absolute value. Therefore, continuity for a function of n variables can be defined precisely as we did for functions of one variable, simply by replacing $|x - x_0|$ by $\|\mathbf{x} - \mathbf{x}_0\|$, where \mathbf{x} and \mathbf{x}_0 now are vectors. Similarly, we replace the $|f(x) - f(x_0)|$ by $\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)\|$. The precise definition reads as follows:

Definition 3.2.1 Continuity of vector functions

Let A denote a set in \mathbb{R}^n . A vector function $\mathbf{f} : A \rightarrow \mathbb{R}^k$ is continuous at $\mathbf{x}_0 \in A$ if

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)\| \rightarrow 0 \quad \text{whenever} \quad \|\mathbf{x} - \mathbf{x}_0\| \rightarrow 0. \quad (3.13)$$

We will also write the condition (3.13) as

$$\mathbf{f}(\mathbf{x}) \rightarrow \mathbf{f}(\mathbf{x}_0) \quad \text{whenever} \quad \mathbf{x} \rightarrow \mathbf{x}_0. \quad (3.14)$$

The function \mathbf{f} is continuous if it is continuous at all $\mathbf{x}_0 \in A$.

Definition 3.2.1 is a proper generalization of the definition for a function of one variable. Indeed, if we write $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{x}_0 = (x_{1,0}, \dots, x_{n,0})$, then

$$\|\mathbf{x} - \mathbf{x}_0\| = \sqrt{(x_1 - x_{1,0})^2 + \dots + (x_n - x_{n,0})^2};$$

in the case $n = 1$ this corresponds to $\|\mathbf{x} - \mathbf{x}_0\| = \sqrt{(x_1 - x_{1,0})^2} = |x_1 - x_{1,0}|$, which is the expression appearing in (3.2).

As for a function of one variable, the condition (3.14) means that the function values $\mathbf{f}(\mathbf{x})$ are “close” to the function value $\mathbf{f}(\mathbf{x}_0)$ for vectors \mathbf{x} that are “close” to \mathbf{x}_0 . The mathematically exact way of phrasing this is:

For any given $\epsilon > 0$ there exists a $\delta > 0$ such that

$$\forall \mathbf{x} \in A : \|\mathbf{x} - \mathbf{x}_0\| < \delta \Rightarrow \|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)\| < \epsilon. \quad (3.15)$$

3.2. Continuity of vector functions of several variables

As for functions of one variable, sums, multiples, and compositions of continuous functions are again continuous. Often we can therefore argue for the continuity of a function without explicitly to refer to Definition 3.2.1.

Example 3.2.1 (a) The functions

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(x, y) = x \text{ and } g : \mathbb{R}^2 \rightarrow \mathbb{R}, g(x, y) = y,$$

are continuous. These functions are called coordinate projections.

(b) The functions

$$\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^2, \mathbf{f}(x) = [x, 0]^T \text{ and } \mathbf{g} : \mathbb{R} \rightarrow \mathbb{R}^2, \mathbf{g}(x) = [0, x]^T$$

are also continuous.

(c) The function

$$\mathbf{h} : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \mathbf{h}(x, y) = [x^2 + xy + \sin(x + y), e^{xy}]^T$$

is a composition of continuous functions, and hence continuous.

Special attention is needed for functions that are defined by multiple expressions as illustrated in the next example.

Example 3.2.2

Consider the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(x, y) = \begin{cases} \frac{x^2 y}{x^4 + y^2} & \text{whenever } (x, y) \neq (0, 0), \\ 0 & \text{whenever } (x, y) = (0, 0). \end{cases}$$

Then f is continuous for all $(x, y) \neq (0, 0)$ since it is a composition of continuous functions. The point $(x, y) = (0, 0)$ needs special attention – and the function is indeed not continuous at this point! In order to see this, consider points in \mathbb{R}^2 on the curve $y = x^2$, $x > 0$, i.e., points of the form $(x, y) = (x, x^2)$. For any such points we see that

$$f(x, y) = f(x, x^2) = \frac{x^2 x^2}{x^4 + (x^2)^2} = \frac{x^4}{2x^4} = \frac{1}{2};$$

that is, along the curve $y = x^2$, the function f takes the constant value $\frac{1}{2}$. Since there are points on the curve $y = x^2$ arbitrarily close to the point $(0, 0)$ and $f(0, 0) = 0$, this implies that the function f cannot be continuous at $(x, y) = (0, 0)$.

3.3. Partial derivatives of first order and the gradient vector

Going from continuity of scalar functions to vector functions is not a real complication. In fact, we can check continuity of a vector function by looking at each coordinate separately:

Theorem 3.2.1 Continuity of coordinate functions

A vector function $\mathbf{f} = (f_1, \dots, f_k) : A \rightarrow \mathbb{R}^k$ is continuous at \mathbf{x}_0 if and only if all the functions $f_i : A \rightarrow \mathbb{R}$, $i = 1, \dots, k$ are continuous at \mathbf{x}_0 .

Proof. For $i = 1, \dots, k$ we compute:

$$\begin{aligned} |f_i(\mathbf{x}) - f_i(\mathbf{x}_0)| &= \sqrt{|f_i(\mathbf{x}) - f_i(\mathbf{x}_0)|^2} \\ &\leq \sqrt{|f_1(\mathbf{x}) - f_1(\mathbf{x}_0)|^2 + \dots + |f_k(\mathbf{x}) - f_k(\mathbf{x}_0)|^2} = \|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)\|. \end{aligned}$$

Thus, if $\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)\| < \epsilon$, then $|f_i(\mathbf{x}) - f_i(\mathbf{x}_0)| < \epsilon$. It follows from Equation (3.15) on page 71 that if $\mathbf{f} = (f_1, \dots, f_k) : A \rightarrow \mathbb{R}^k$ is continuous at \mathbf{x}_0 , then so is f_i for each i .

Conversely, assume all the coordinate functions f_i are continuous at \mathbf{x}_0 and we are given an $\epsilon > 0$. For each functions f_i we can find a $\delta_i > 0$ such that $\|\mathbf{x} - \mathbf{x}_0\| < \delta_i \Rightarrow |f_i(\mathbf{x}) - f_i(\mathbf{x}_0)| < \epsilon/k$. We now put $\delta = \min\{\delta_1, \dots, \delta_k\}$. If $\|\mathbf{x} - \mathbf{x}_0\| < \delta$ then $|f_i(\mathbf{x}) - f_i(\mathbf{x}_0)| < \epsilon/k$ for all $i = 1, \dots, k$ and hence

$$\begin{aligned} \|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}_0)\| &= \sqrt{|f_1(\mathbf{x}) - f_1(\mathbf{x}_0)|^2 + \dots + |f_k(\mathbf{x}) - f_k(\mathbf{x}_0)|^2} \\ &< \sqrt{k \left(\frac{\epsilon}{k}\right)^2} = \frac{\epsilon}{\sqrt{k}} < \epsilon. \end{aligned}$$

■

We are usually interested in continuity in more than in a single point. This leads to the following definition.

Definition 3.2.2 C^0 vector function

A vector function $\mathbf{f} : A \rightarrow \mathbb{R}^k$ is called *continuous* if it is continuous at all points $\mathbf{x} \in A$. Continuous vector functions are called C^0 vector function.

3.3 Partial derivatives of first order and the gradient vector

In Section 3.2 we saw that continuity for functions of n variables can be defined in a similar fashion as for functions of one variable. In contrast,

3.3. Partial derivatives of first order and the gradient vector

the theory for differentiability is significantly more involved for functions of several variables than for functions of one variable.

In the next three sections we will develop the theory of differentiability of *scalar* functions of several variables. In [Section 3.8](#) we will return to vector functions.

Let us begin with an introductory example.

Example 3.3.1

Consider the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(x, y) := xy^2 + x^2 + y^3.$$

The function f is obviously continuous as it is composed of well-known continuous functions, but we have not yet defined how we should interpret differentiability for a function of two variables.

However, note that if we *fix* an arbitrary value of the variable y and consider the function f only as a function of the remaining variable x , we obtain the function h_y (the subscript is chosen to indicate that the function depends on the chosen y) given by

$$h_y : \mathbb{R} \rightarrow \mathbb{R}, h_y(x) = f(x, y) = xy^2 + x^2 + y^3.$$

For all (fixed) values of y , the function h_y is a differentiable function of the variable x , and

$$h'_y(x) = y^2 + 2x. \quad (3.16)$$

Note that the expression (3.16) is obtained simply by considering y as a constant and applying standard rules for differentiation of a function of one variable.

We will now introduce a new terminology to describe the situation. Being a little ahead of the formal definition (to be given in [Definition 3.3.1](#)) we will say that *the partial derivative of the function f with respect to the variable x , to be denoted by $\frac{\partial f}{\partial x}(x, y)$, exists, and is given by*

$$\frac{\partial f}{\partial x}(x, y) = h'_y(x) = y^2 + 2x. \quad (3.17)$$

Similarly, if we *fix* an arbitrary value of the variable x and consider the function f only as a function of the remaining variable y , we obtain the function

$$k_x : \mathbb{R} \rightarrow \mathbb{R}, k_x(y) = f(x, y) = xy^2 + x^2 + y^3.$$

3.3. Partial derivatives of first order and the gradient vector

Regardless of the chosen x , the function k_x is a differentiable function of the variable y , and

$$k'_x(y) = 2xy + 3y^2. \quad (3.18)$$

We say that *the partial derivative of the function f with respect to the variable y , to be denoted by $\frac{\partial f}{\partial y}(x, y)$, exists, and is given by*

$$\frac{\partial f}{\partial y}(x, y) = k'_x(y) = 2xy + 3y^2. \quad (3.19)$$

Note that in order to obtain the expressions (3.17) and (3.19), it is not necessary to introduce the functions h_y and k_x . Since, in order to calculate the partial derivative $\frac{\partial f}{\partial x}(x, y)$, we simply consider y as a *constant* and differentiate the function f with respect to the variable x . Similarly, we calculate the partial derivative $\frac{\partial f}{\partial y}(x, y)$ by differentiating the expression $f(x, y)$ with respect to the variable y when considering x as a *constant*. Let us illustrate this by one more example.

Example 3.3.2

Consider the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = \sin(x^2 + y).$$

Then (compare with the calculations in Example 3.1.2)

$$\frac{\partial f}{\partial x}(x, y) = 2x \cos(x^2 + y)$$

and

$$\frac{\partial f}{\partial y}(x, y) = \cos(x^2 + y).$$

Already in Equation (3.5) we saw that a function of one variable can be differentiable at some points and non-differentiable at other points. The same complication occurs for functions of several variables. Therefore we would like to stress the fact that the existence of the partial derivatives for a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ has to be checked *individually* for each point $(x, y) \in \mathbb{R}^2$.

Example 3.3.3

Consider the function

3.3. Partial derivatives of first order and the gradient vector

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(x, y) = \text{ReLU}(x) y^2.$$

For each fixed $x \in \mathbb{R}$, the function $k_x(y) = \text{ReLU}(x) y^2$ is differentiable for all choices of $y \in \mathbb{R}$; thus, the partial derivative $\frac{\partial f}{\partial y}(x, y)$ exists for all $(x, y) \in \mathbb{R}^2$, and

$$\frac{\partial f}{\partial y}(x, y) = 2 \text{ReLU}(x) y.$$

The situation is slightly more complicated when we fix $y \in \mathbb{R}$ and want to consider f as a function of x . The reason is that as we saw in [Example 3.1.3](#), the function ReLU is not differentiable at $x = 0$. In order to handle this, let us write the function f as

$$f(x, y) = \begin{cases} xy^2 & \text{whenever } x \geq 0, \\ 0 & \text{whenever } x < 0. \end{cases}$$

From here we see that for any fixed $y \in \mathbb{R}$, the function $h_y(x) = \text{ReLU}(x) y^2$ is differentiable whenever $x \neq 0$; thus, for $(x, y) \in \mathbb{R}^2$ with $x \neq 0$ the partial derivative $\frac{\partial f}{\partial x}(x, y)$ exists and

$$\frac{\partial f}{\partial x}(x, y) = \begin{cases} y^2 & \text{whenever } x > 0 \\ 0 & \text{whenever } x < 0 \end{cases}$$

However, when we fix any $y \neq 0$, the function h_y is not differentiable at $x = 0$, as we saw in [Example 3.1.3](#); that is, the partial derivatives $\frac{\partial f}{\partial x}(0, y)$ do not exist.

Complications like in [Example 3.3.3](#) makes it evident that we need a formal definition of the partial derivatives. We will now state such a definition for a function of n variables (x_1, x_2, \dots, x_n) .

Definition 3.3.1 Partial derivative

Let U be an open subset of \mathbb{R}^n , and let $f : U \rightarrow \mathbb{R}$ be a scalar function. Consider a point $(x_1, x_2, \dots, x_n) \in U$. Fix some $j = 1, \dots, n$, and assume that the expression

$$\frac{f(x_1, x_2, \dots, x_j + h, \dots, x_n) - f(x_1, x_2, \dots, x_j, \dots, x_n)}{h}$$

has a limit as $h \rightarrow 0$. Then we say that the *partial derivative* $\frac{\partial f}{\partial x_j}$ exists at the point $\mathbf{x} = (x_1, x_2, \dots, x_n)$, and we define

3.3. Partial derivatives of first order and the gradient vector

$$\frac{\partial f}{\partial x_j}(\mathbf{x}) := \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_j + h, \dots, x_n) - f(x_1, x_2, \dots, x_j, \dots, x_n)}{h}.$$

Let us explicitly formulate the definition of the partial derivative with respect to the first variable x_1 . If the expression

$$\frac{f(x_1 + h, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h}$$

has a limit as $h \rightarrow 0$, the partial derivative $\frac{\partial f}{\partial x_1}$ exists at the point (x_1, x_2, \dots, x_n) and it is given by

$$\frac{\partial f}{\partial x_1}(x_1, x_2, \dots, x_n) = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h}.$$

If all the partial derivatives $\frac{\partial f}{\partial x_k}$, $k = 1, \dots, n$ exist at a given point (x_1, \dots, x_n) , we will collect the information in a single vector:

Definition 3.3.2 Gradient vector

If all the partial derivatives $\frac{\partial f}{\partial x_k}$, $k = 1, \dots, n$, exist at a given point (x_1, \dots, x_n) , the *gradient vector* of the function f at the point $\mathbf{x} = (x_1, \dots, x_n)$ is defined as the vector

$$\nabla f(\mathbf{x}) := \left(\frac{\partial f}{\partial x_1}(\mathbf{x}), \frac{\partial f}{\partial x_2}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}) \right). \quad (3.20)$$

The symbol ∇ is pronounced “nabla”. The gradient vector will be considered as a *column* vector unless stated otherwise.

Example 3.3.4

For the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = xy^2 + x^2 + y^3,$$

we calculated the partial derivatives in [Example 3.3.1](#). Formulated in terms of the gradient vector, the result says that

$$\nabla f(x, y) = (y^2 + 2x, 2xy + 3y^2). \quad (3.21)$$

Example 3.3.5

Consider the function

$$f : \mathbb{R}^3 \rightarrow \mathbb{R}, f(x_1, x_2, x_3) = x_1x_2 + x_2^2x_3^3.$$

Then all the partial derivatives exist at all points $(x_1, x_2, x_3) \in \mathbb{R}$, and

$$\begin{aligned} \nabla f(x_1, x_2, x_3) &= \left(\frac{\partial f}{\partial x_1}(x_1, x_2, x_3), \frac{\partial f}{\partial x_2}(x_1, x_2, x_3), \frac{\partial f}{\partial x_3}(x_1, x_2, x_3) \right) \\ &= (x_2, x_1 + 2x_2x_3^3, 3x_2^2x_3^2). \end{aligned}$$

Recall that if a function $f : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable, then f is also continuous. For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ it is therefore natural to ask whether existence of all the partial derivatives implies that the function is continuous. The answer turns out to be no:

Example 3.3.6

Consider again the function in [Example 3.2.2](#),

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(x, y) = \begin{cases} \frac{x^2y}{x^4+y^2} & \text{whenever } (x, y) \neq (0, 0), \\ 0 & \text{whenever } (x, y) = (0, 0). \end{cases}$$

For $(x, y) \neq (0, 0)$, the function f is composed of functions that are differentiable with respect to as well x as y , i.e., the partial derivatives $\frac{\partial f}{\partial x}(x, y)$ and $\frac{\partial f}{\partial y}(x, y)$ exist. In order to check whether the partial derivative $\frac{\partial f}{\partial x}(x, y)$ exists at the point $(x, y) = (0, 0)$, we apply [Definition 3.3.1](#) and calculate

$$\frac{f(h, 0) - f(0, 0)}{h} = \frac{0 - 0}{h} = 0;$$

since this expression indeed has a limit (namely, 0) as $h \rightarrow 0$, we conclude that also the partial derivative $\frac{\partial f}{\partial x}(0, 0)$ exists. By a similar argument, also the partial derivative $\frac{\partial f}{\partial y}(0, 0)$ exists, i.e., both partial derivatives exist in all points. However, as we saw in [Example 3.2.2](#) the function f is not continuous.

In order to avoid issues like in [Example 3.3.6](#) we need to consider a stronger differentiability condition than just existence of the partial derivatives. We will do so in [Section 3.6](#).

3.4 Directional derivatives

When computing the partial derivative $\frac{\partial f}{\partial x_i}$ in $\mathbf{x}_0 \in \mathbb{R}^n$, one moves along the i th coordinate axes, in the sense that $\frac{\partial f}{\partial x_i}(\mathbf{x}_0)$ is equal to the usual

derivative as in 3.1.2 of the function $h \mapsto f(\mathbf{x}_0 + h\mathbf{e}_i), \mathbb{R} \rightarrow \mathbb{R}$, evaluated at $h = 0$, where $\mathbf{e}_i = [0, \dots, 0, 1, 0, \dots, 0]^T$ is the i th standard basis vector of \mathbb{R}^n . However, there is nothing special about the coordinate axes, and it is just as natural to move along other directions. This is the idea of directional derivatives.

Definition 3.4.1 Directional derivative

Let \mathbf{v} be a non-zero vector in \mathbb{R}^n , let U be an open subset of \mathbb{R}^n , and let $f : U \rightarrow \mathbb{R}$ be a scalar function. Define $\mathbf{u} = \mathbf{v}/\|\mathbf{v}\|$. Consider a point $(x_1, x_2, \dots, x_n) \in U$. Assume that the expression

$$\frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h}$$

has a limit as $h \rightarrow 0$. Then we say that the *directional derivative* $\nabla_{\mathbf{v}}f$ exists at $\mathbf{x} = (x_1, x_2, \dots, x_n)$, and we define

$$\nabla_{\mathbf{v}}f(\mathbf{x}) := \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{u}) - f(\mathbf{x})}{h}.$$

Note that the directional derivative $\nabla_{\mathbf{v}}f$ evaluated at a point $\mathbf{x} \in U$ is just a scalar in \mathbb{R} . Moreover, $\nabla_{\mathbf{e}_i}f(\mathbf{x}) = \frac{\partial f}{\partial x_i}(\mathbf{x})$.

The length of the vector \mathbf{v} in the definition of $\nabla_{\mathbf{v}}f$ is irrelevant, and we see that only the direction matters. In particular, if $\mathbf{w} \in \text{span}\{\mathbf{v}\}$ is a non-zero vector, then $\nabla_{\mathbf{w}}f = \pm \nabla_{\mathbf{v}}f$. (The minus sign is needed if \mathbf{w} and \mathbf{v} points in opposite directions, i.e., $\mathbf{w} = c\mathbf{v}$ for $c < 0$.)

The directional derivative can easily be computed from the gradient vector using the standard inner product in \mathbb{R}^n .

Lemma 3.4.1

If \mathbf{u} is a unit vector, then

$$\nabla_{\mathbf{u}}f(\mathbf{x}) = \langle \mathbf{u}, \nabla f(\mathbf{x}) \rangle \tag{3.22}$$

for $\mathbf{x} \in \mathbb{R}^n$.

We leave the proof to the reader. In case \mathbf{v} is not necessarily of norm one, the formula in Equation (3.22) should be changed to $\nabla_{\mathbf{v}}f(\mathbf{x}) = \frac{1}{\|\mathbf{v}\|} \langle \mathbf{v}, \nabla f(\mathbf{x}) \rangle$.

The directional derivative $|\nabla_{\mathbf{v}}f(\mathbf{x})|$ is maximized if \mathbf{v} is a unit vector in the span of the gradient vector as shown in Exercise 3.4.2. Loosely speaking, we say that the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ increases and decreases the most along the gradient vector.

Exercise 3.4.2

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function for which all directional derivatives exist at $\mathbf{x} \in \mathbb{R}^n$. Show that $\mathbf{u} = \nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|$ is a unit vector. Show that the scalar $|\nabla_{\mathbf{v}} f(\mathbf{x})|$ is maximized for $\mathbf{v} = \pm \mathbf{u}$. *Hint: Equation (3.22) and Theorem 2.1.6 on page 34*

3.5 Partial derivatives of second order and the Hessian matrix

We will now give a short introduction to higher-order partial derivatives. In Chapter 4 and Chapter 5 we will see that they play a central role in connection with Taylor polynomials and optimization problems for functions of n variables (like $f''(x)$ does for a function of a single variable!).

Consider again a function $f : U \rightarrow \mathbb{R}$, where U is an open set in \mathbb{R}^n . If the partial derivative $\frac{\partial f}{\partial x_k}$ exists at all points $(x_1, \dots, x_n) \in U$ for some index $k = 1, \dots, n$, we can consider the function $\frac{\partial f}{\partial x_k}$ as a new function of n variables,

$$\frac{\partial f}{\partial x_k} : U \rightarrow \mathbb{R}.$$

Thus, we can ask whether the function $\frac{\partial f}{\partial x_k}$ has partial derivatives, say, with respect to x_j for some $j = 1, \dots, n$, at a point $(x_1, \dots, x_n) \in U$; if this is the case, the partial derivative will be denoted by

$$\frac{\partial^2 f}{\partial x_j \partial x_k}(x_1, \dots, x_n).$$

In the special case where $j = k$ we will use the shorter notation

$$\frac{\partial^2 f}{\partial x_k^2}(x_1, \dots, x_n) := \frac{\partial^2 f}{\partial x_k \partial x_k}(x_1, \dots, x_n).$$

Example 3.5.1

For the function

$$f : \mathbb{R}^3 \rightarrow \mathbb{R}, \quad f(x_1, x_2, x_3) = x_1 x_2 + x_2^2 x_3^3,$$

we saw in Example 3.3.5 on page 77 that, for $\mathbf{x} = (x_1, x_2, x_3)$,

$$\frac{\partial f}{\partial x_1}(\mathbf{x}) = x_2, \quad \frac{\partial f}{\partial x_2}(\mathbf{x}) = x_1 + 2x_2 x_3^3, \quad \frac{\partial f}{\partial x_3}(\mathbf{x}) = 3x_2^2 x_3^2.$$

3.5. Partial derivatives of second order and the Hessian matrix

All these partial derivatives have partial derivatives with respect to x_1, x_2 , and x_3 ; by direct calculation we get that

$$\begin{aligned} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) &= 0, & \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}) &= 1, & \frac{\partial^2 f}{\partial x_3 \partial x_1}(\mathbf{x}) &= 0, \\ \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}) &= 1, & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{x}) &= 2x_3^3, & \frac{\partial^2 f}{\partial x_3 \partial x_2}(\mathbf{x}) &= 6x_2x_3^2, \\ \frac{\partial^2 f}{\partial x_1 \partial x_3}(\mathbf{x}) &= 0, & \frac{\partial^2 f}{\partial x_2 \partial x_3}(\mathbf{x}) &= 6x_2x_3^2, & \frac{\partial^2 f}{\partial x_3^2}(\mathbf{x}) &= 6x_2^2x_3, \end{aligned}$$

where $\mathbf{x} = (x_1, x_2, x_3)$.

By inspection of the second-order partial derivatives in [Example 3.5.1](#) we notice that for all values of $j, k = 1, 2, 3$ we have that

$$\frac{\partial^2 f}{\partial x_j \partial x_k}(x_1, x_2, x_3) = \frac{\partial^2 f}{\partial x_k \partial x_j}(x_1, x_2, x_3);$$

that is, the order in which the partial derivatives is taken does not matter. This turns out to hold in general, provided that an extra assumption introduced in [Section 3.6](#) is satisfied – see [Theorem 3.6.5](#).

All information about the second-order partial derivatives are collected in the so-called *Hessian matrix* defined as follows.

Definition 3.5.1 Hessian matrix

Assuming that all the second-order partial derivatives exist for a scalar function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the associated *Hessian matrix* is the $n \times n$ matrix given by:

$$\mathbf{H}_f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{bmatrix} \quad (3.23)$$

For a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ the Hessian matrix takes the form

$$\mathbf{H}_f(x, y) = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2}(x, y) & \frac{\partial^2 f}{\partial x \partial y}(x, y) \\ \frac{\partial^2 f}{\partial y \partial x}(x, y) & \frac{\partial^2 f}{\partial y^2}(x, y) \end{bmatrix}. \quad (3.24)$$

3.5. Partial derivatives of second order and the Hessian matrix

Example 3.5.2

For the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) := \sin(x^2 + y),$$

we saw in Example 3.3.2 that

$$\frac{\partial f}{\partial x}(x, y) = 2x \cos(x^2 + y), \quad \text{and} \quad \frac{\partial f}{\partial y}(x, y) = \cos(x^2 + y).$$

Thus

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2}(x, y) &= 2 \cos(x^2 + y) - 4x^2 \sin(x^2 + y) \\ \frac{\partial^2 f}{\partial x \partial y}(x, y) &= -2x \sin(x^2 + y) \\ \frac{\partial^2 f}{\partial y \partial x}(x, y) &= -2x \sin(x^2 + y) \\ \frac{\partial^2 f}{\partial y^2}(x, y) &= -\sin(x^2 + y). \end{aligned}$$

The Hessian matrix is therefore

$$\begin{aligned} \mathbf{H}_f(x, y) &= \begin{bmatrix} \frac{\partial^2 f}{\partial x^2}(x, y) & \frac{\partial^2 f}{\partial x \partial y}(x, y) \\ \frac{\partial^2 f}{\partial y \partial x}(x, y) & \frac{\partial^2 f}{\partial y^2}(x, y) \end{bmatrix} \\ &= \begin{bmatrix} 2 \cos(x^2 + y) - 4x^2 \sin(x^2 + y) & -2x \sin(x^2 + y) \\ -2x \sin(x^2 + y) & -\sin(x^2 + y) \end{bmatrix} \end{aligned}$$

Example 3.5.3

For the function f in Example 3.5.1, the Hessian matrix is

$$\begin{aligned} \mathbf{H}_f(x_1, x_2, x_3) &= \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(x_1, x_2, x_3) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(x_1, x_2, x_3) & \frac{\partial^2 f}{\partial x_1 \partial x_3}(x_1, x_2, x_3) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(x_1, x_2, x_3) & \frac{\partial^2 f}{\partial x_2^2}(x_1, x_2, x_3) & \frac{\partial^2 f}{\partial x_2 \partial x_3}(x_1, x_2, x_3) \\ \frac{\partial^2 f}{\partial x_3 \partial x_1}(x_1, x_2, x_3) & \frac{\partial^2 f}{\partial x_3 \partial x_2}(x_1, x_2, x_3) & \frac{\partial^2 f}{\partial x_3^2}(x_1, x_2, x_3) \end{bmatrix} \\ &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2x_3^2 & 6x_2x_3^2 \\ 0 & 6x_2x_3^2 & 6x_2^2x_3 \end{bmatrix}. \end{aligned}$$

Example 3.5.4 Quadratic forms

Let $q : \mathbb{R}^n \rightarrow \mathbb{R}$ be a quadratic form. Recall from Lemma 2.8.7 on page 60 that we may assume that

$$q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + \mathbf{x}^T \mathbf{b} + c,$$

where $A \in M_n(\mathbb{R})$ is symmetric, $\mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$. It can be shown using Example 2.8.3 on page 60 that the Hessian \mathbf{H}_q is constant and equals

$$\mathbf{H}_q = 2A.$$

Exercise 3.5.1

Prove the claims made in Example 3.5.4.

3.6 Differentiability of scalar functions of several variables

We have already mentioned that differentiability is a much more complicated issue for functions of several variables than for functions of one variable. One reason is that Definition 3.1.2 is restricted to the one-dimensional case: in contrast to our approach to continuity, the condition can not be changed to a condition on \mathbb{R}^n simply by exchanging the point x_0 by a vector \mathbf{x}_0 . The argument for this is that if we want to consider the numerator in (3.4) with x_0 replaced by a vector \mathbf{x}_0 , i.e., an expression of the form

$$f(\mathbf{x}_0 + h) - f(\mathbf{x}_0),$$

then h is forced to be a vector in \mathbb{R}^n , and therefore the division with h in (3.4) becomes meaningless as we have no concept of division by a vector.

In order to generalize the concept of differentiability to higher dimensions, we therefore need to reformulate the condition in Definition 3.1.2 in such a way that we obtain a condition that can be transferred into \mathbb{R}^n . This is done next.

Lemma 3.6.1

Consider a function $f : I \rightarrow \mathbb{R}$, where I is an open interval in \mathbb{R} , and let $x_0 \in I$. Fix $c \in \mathbb{R}$. Then the following are equivalent:

- (i) The function f is differentiable at x_0 and $f'(x_0) = c$.

3.6. Differentiability of scalar functions of several variables

(ii) There exists a function $\varepsilon : \mathbb{R} \rightarrow \mathbb{R}$ such that $\varepsilon(h) \rightarrow 0$ as $h \rightarrow 0$ and

$$f(x_0 + h) - f(x_0) - ch - \varepsilon(h)|h| = 0 \quad \text{for all } h \in \mathbb{R}. \quad (3.25)$$

Proof. We first note that we in (ii) can replace (3.25) with the equation

$$f(x_0 + h) - f(x_0) - ch - \varepsilon(h)h = 0 \quad \text{for all } h \in \mathbb{R}.$$

Indeed, the two ε functions are related by a changing the sign for $h < 0$. We will in the rest of the prove use this version of (3.25).

Now assume that (i) holds. Then, by (3.5)

$$\frac{f(x_0 + h) - f(x_0)}{h} \rightarrow f'(x_0) = c \quad \text{as } h \rightarrow 0.$$

Now define the function ε by

$$\varepsilon(h) := \begin{cases} \frac{f(x_0+h)-f(x_0)}{h} - c & \text{whenever } h \neq 0, \\ 0 & \text{whenever } h = 0. \end{cases}$$

Then $\varepsilon(h) \rightarrow 0$ as $h \rightarrow 0$, and by multiplying the expression for $\varepsilon(h)$ with h and collecting all terms on one side of the equality sign, we see that (3.25) holds.

On the other hand, if (ii) holds, then

$$\varepsilon(h)h = f(x_0 + h) - f(x_0) - ch.$$

Dividing by h now yields that

$$\varepsilon(h) = \frac{f(x_0 + h) - f(x_0)}{h} - c;$$

applying now that $\varepsilon(h) \rightarrow 0$ as $h \rightarrow 0$, we conclude that

$$\frac{f(x_0 + h) - f(x_0)}{h} \rightarrow c \quad \text{as } h \rightarrow 0.$$

This means precisely that the function f is differentiable at x_0 and that $f'(x_0) = c$, which proves (i). \blacksquare

Let us discuss condition (ii) in Lemma 3.6.1. First, we remark that the constant c appearing in (ii) is not arbitrary, indeed, it is equal to $f'(x_0)$. Second, we note that ch for a given constant c defines a linear map $h \mapsto ch$,

3.6. Differentiability of scalar functions of several variables

$\mathbb{R} \rightarrow \mathbb{R}$. Hence, the derivative of a function is closely related to linear maps in h . This linear map also “appears” in the tangent line of the graph at $(x_0, f(x_0))$ ¹.

In contrast to our original definition of differentiability, the condition (ii) in Lemma 3.6.1 has a natural generalization to functions of several variables. Let us explain this in detail.

Consider a function $f : U \rightarrow \mathbb{R}$, where U is an open set in \mathbb{R}^n , and let $\mathbf{x}_0 \in U$. In order to generalize the condition (3.25) we will do the following:

- 1) We will replace the scalar h by a (column) vector $\mathbf{h} \in \mathbb{R}^n$.
- 2) We will replace the scalar c by a (column) vector $\mathbf{c} \in \mathbb{R}^n$ and the quantity ch by the scalar $\mathbf{c}^T \mathbf{h}$.
- 3) We will replace the function $\varepsilon : \mathbb{R} \rightarrow \mathbb{R}$ by a function $\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$ of n variables, and we replace the expression $\varepsilon(h)|h|$ by $\varepsilon(\mathbf{h})\|\mathbf{h}\|$.

Note that in 2) $\mathbf{h} \mapsto \mathbf{c}^T \mathbf{h}$ is a linear map from \mathbb{R}^n to \mathbb{R} defined by the $1 \times n$ “matrix” \mathbf{c}^T . Note also that the $\mathbf{c}^T \mathbf{h}$ is just the standard inner product on \mathbb{R}^n , i.e., $\mathbf{c}^T \mathbf{h} = \mathbf{h} \cdot \mathbf{c} = \langle \mathbf{h}, \mathbf{c} \rangle$. In 3), the norm $\|\mathbf{h}\|$ reduces to $|h_1|$ if $n = 1$ and $\mathbf{h} = [h_1]$.

These modifications lead to the following definition of differentiability:

Definition 3.6.1 Differentiability

Consider a function $f : U \rightarrow \mathbb{R}$, where U is an open set in \mathbb{R}^n , and let $\mathbf{x}_0 \in U$. We say that the function f is differentiable at the point \mathbf{x}_0 if there exist a vector $\mathbf{c} \in \mathbb{R}^n$ and a function $\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$ such that $\varepsilon(\mathbf{h}) \rightarrow 0$ as $\mathbf{h} \rightarrow 0$ and

$$f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) - \mathbf{c}^T \mathbf{h} - \varepsilon(\mathbf{h})\|\mathbf{h}\| = 0, \quad \forall \mathbf{h} \in \mathbb{R}^n. \quad (3.26)$$

We will now show that there is a strong link between Definition 3.6.1 and the concept of partial derivatives in Section 3.3. In particular, if the function f is differentiable at the point \mathbf{x}_0 , all the partial derivatives $\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n}$ exist at the point \mathbf{x}_0 and the vector \mathbf{c} in (3.26) equals the gradient vector $\nabla f(\mathbf{x}_0)$. Note that this is similar to the situation in Lemma 3.6.1 on page 83, where the constant c is equal to $f'(x_0)$. The linear map $\mathbf{h} \mapsto \nabla f(\mathbf{x}_0)^T \mathbf{h}$ is called the *differential* of f at \mathbf{x}_0 .

¹The graph of this linear map is a line in \mathbb{R}^2 through $(0, 0)$. If we add $(x_0, f(x_0))$ to each point in this set, we obtain the tangent line.

Theorem 3.6.2

Assume that the function $f : U \rightarrow \mathbb{R}$ is differentiable at the point \mathbf{x}_0 . Then all the partial derivatives $\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n}$ exist at the point \mathbf{x}_0 and the vector \mathbf{c} in Definition 3.6.1 equals the gradient vector $\nabla f(\mathbf{x}_0)$.

Proof. Assume that f is differentiable at the point \mathbf{x}_0 and choose the vector \mathbf{c} and the function $\varepsilon(\mathbf{h})$ such that (3.26) is satisfied. We will now consider certain special choices of the vector \mathbf{h} , namely $\mathbf{h} = (h, 0, \dots, 0)$, where $h \in \mathbb{R}$. The matrix product $\mathbf{c}^T \mathbf{h}$ becomes:

$$\mathbf{c}^T \mathbf{h} = \begin{bmatrix} c_1 & c_2 & \cdots & c_n \end{bmatrix} \begin{bmatrix} h \\ 0 \\ \vdots \\ 0 \end{bmatrix} = c_1 h.$$

Since $\|\mathbf{h}\| = \sqrt{h^2 + 0 + \cdots + 0} = |h|$, we also see that

$$\varepsilon(\mathbf{h})\|\mathbf{h}\| = \varepsilon(\mathbf{h})|h|$$

Inserting into (3.26), we obtain

$$f(\mathbf{x}_0 + (h, 0, \dots, 0)) - f(\mathbf{x}_0) - c_1 h - \varepsilon(h, 0, \dots, 0)|h| = 0,$$

where

$$\varepsilon(h, 0, \dots, 0) \rightarrow 0 \quad \text{whenever } h \rightarrow 0.$$

Applying now Lemma 3.6.1, we see that the partial derivative $\frac{\partial f}{\partial x_1}$ exists at the point \mathbf{x}_0 and that $\frac{\partial f}{\partial x_1}(\mathbf{x}_0) = c_1$. A similar argument applies to the other partial derivatives, and the conclusion is that

$$\mathbf{c} = (c_1, c_2, \dots, c_n) = \left(\frac{\partial f}{\partial x_1}(\mathbf{x}_0), \frac{\partial f}{\partial x_2}(\mathbf{x}_0), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}_0) \right) = \nabla f(\mathbf{x}_0)$$

as claimed. ■

Theorem 3.6.2 shows that differentiability of a function $f : U \rightarrow \mathbb{R}$ implies existence of the partial derivatives. On the other hand, existence of the partial derivatives is not enough to guarantee differentiability! Indeed, in Example 3.3.6 we saw an example of a function, for which all the partial derivatives exist; but the function is not continuous, and hence, by Theorem 3.6.4 on the next page, not differentiable. However, a slightly weaker result hold: existence *and* continuity of all partial derivatives implies differentiability.

3.6. Differentiability of scalar functions of several variables

Theorem 3.6.3

Let $U \subseteq \mathbb{R}^n$ be an open set and let $f : U \rightarrow \mathbb{R}$ be a scalar function. If the partial derivatives $\frac{\partial f}{\partial x_k}$ exist at all points and are continuous, then f is differentiable for all $\mathbf{x} \in U$.

Proof. For a proof we refer the reader to [Theorem 3.22 in the additional notes](#). ■

Example 3.6.1

Let $A \in \mathbb{R}^{n \times n}$, $\mathbf{b} \in \mathbb{R}^n$, and $c \in \mathbb{R}$. A quadratic form $q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$ is differentiable at all points $\mathbf{x} \in \mathbb{R}^n$, and the gradient vector at \mathbf{x}_0 is given by $\nabla q(\mathbf{x}) = (A + A^T)\mathbf{x} + \mathbf{b}$. Now, if A is symmetric, then $A + A^T = 2A$, and so the gradient vector becomes $\nabla q(\mathbf{x}) = 2A\mathbf{x} + \mathbf{b}$. Recall that the derivative of the second degree polynomial function $p : \mathbb{R} \rightarrow \mathbb{R}$, $p(x) = ax^2 + bx + c$ is $p'(x) = 2ax + b$. Hence, $\nabla q(\mathbf{x}) = (A + A^T)\mathbf{x} + \mathbf{b}$ is a proper generalization of this well-known fact.

Similarly to what we know for functions of one variable, differentiable functions of n variables are continuous:

Theorem 3.6.4

Assume that the function $f : U \rightarrow \mathbb{R}$ is differentiable at the point $\mathbf{x}_0 \in U$. Then f is also continuous at \mathbf{x}_0 .

Proof. Assume that f is differentiable at the point $\mathbf{x}_0 \in U$. Then (3.26) holds for some function $\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$ such that $\varepsilon(\mathbf{h}) \rightarrow 0$ as $\mathbf{h} \rightarrow 0$. It then follows that

$$f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) = \mathbf{c}^T \mathbf{h} + \varepsilon(\mathbf{h})\|\mathbf{h}\| \rightarrow 0 \text{ as } \mathbf{h} \rightarrow 0,$$

i.e., the function f is indeed continuous at the point \mathbf{x}_0 . ■

The concept of differentiability also yields the condition that implies that the order in which we take partial derivatives of second order is irrelevant:

Theorem 3.6.5

Let U be an open set in \mathbb{R}^n and consider a differentiable function $f : U \rightarrow \mathbb{R}$ for which all the second-order partial derivatives $\frac{\partial^2 f}{\partial x_j \partial x_k}$, $j, k = 1, \dots, n$, exist

on U and are continuous. Then

$$\frac{\partial^2 f}{\partial x_j \partial x_k}(x_1, x_2, \dots, x_n) = \frac{\partial^2 f}{\partial x_k \partial x_j}(x_1, x_2, \dots, x_n)$$

for all $j, k = 1, \dots, n$ and all $(x_1, x_2, \dots, x_n) \in U$.

The proof of [Theorem 3.6.5](#) is technical and will not be given here. Again, we refer to [Theorem 3.34 in the additional notes](#) for a proof.

3.7 The chain rule for scalar functions

The rule (3.8) concerning differentiation of a composed function $f \circ g$ has an important generalization to the case where f is a function of several variables. For the sake of convenient notation we will only state it for a function f that is defined and differentiable on all of \mathbb{R}^n . We will state the general case in [Theorem 3.8.4](#) on page 93.

Theorem 3.7.1

Consider a differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, and let $g_1, g_2, \dots, g_n : \mathbb{R} \rightarrow \mathbb{R}$ denote differentiable functions. Then the composed function

$$h : \mathbb{R} \rightarrow \mathbb{R}, \quad h(x) := f(g_1(x), g_2(x), \dots, g_n(x)) \quad (3.27)$$

is differentiable, and

$$\begin{aligned} h'(x) &= \frac{\partial f}{\partial x_1}(g_1(x), g_2(x), \dots, g_n(x))g_1'(x) \\ &\quad + \frac{\partial f}{\partial x_2}(g_1(x), g_2(x), \dots, g_n(x))g_2'(x) \\ &\quad + \dots + \frac{\partial f}{\partial x_n}(g_1(x), g_2(x), \dots, g_n(x))g_n'(x). \end{aligned} \quad (3.28)$$

The proof of [Theorem 3.7.1](#) follows closely the classical proof of (3.8). The expression (3.28) can be written in a more compressed form using the inner product between the gradient of f at the point $(g_1(x), g_2(x), \dots, g_n(x))$ and the vector $(g_1'(x), g_2'(x), \dots, g_n'(x))$. Combining $\mathbf{g}(x) = (g_1(x), g_2(x), \dots, g_n(x))$ into a vector function of one variable, we can express the derivative of $h(x) := f(\mathbf{g}(x))$ as follows.

Corollary 3.7.2

The expression (3.28) means precisely that

$$\begin{aligned} h'(x) &= \langle (g'_1(x), g'_2(x), \dots, g'_n(x)), \nabla f(g_1(x), g_2(x), \dots, g_n(x))) \rangle \\ &= \langle \mathbf{g}'(x), \nabla f(\mathbf{g}(x)) \rangle = \nabla f(\mathbf{g}(x))^T \mathbf{g}'(x). \end{aligned} \quad (3.29)$$

The result in [Theorem 3.7.1](#) and [Corollary 3.7.2](#) is known as the *chain rule*. The generalized chain rule for vector functions can be found in [Theorem 3.8.4](#) on page 93. The chain rule has many applications, e.g., it is fundamental in the concept of backpropagation used to estimate the gradient in machine learning under the training of neural networks and artificial intelligence models.

If the functions f and g_1, g_2, \dots, g_n in [Corollary 3.7.2](#) are *explicitly* given, we can directly use (3.27) to obtain a formula for the function h , which then can be differentiated “as usual” without reference to the chain rule. Hence, for us, the chain rule is mainly important for theoretical analysis as it rarely makes practical computations easier. Indeed, let us illustrate this with an example.

Example 3.7.1

Consider the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x_1, x_2) = x_1^2 + x_2$$

and let

$$g_1(x) = \sin(x), \quad g_2(x) = \cos(x).$$

Then

$$h(x) := f(g_1(x), g_2(x)) = \sin^2(x) + \cos(x),$$

so we directly see that h is differentiable, with

$$h'(x) = 2 \sin(x) \cos(x) - \sin(x). \quad (3.30)$$

Let us derive the same result by applying the chain rule as formulated in [Corollary 3.7.2](#). We see that the gradient of the function f is

$$\nabla f(x_1, x_2) = (2x_1, 1),$$

so

$$\nabla f(g_1(x), g_2(x)) = (2 \sin(x), 1).$$

Also,

$$(g'_1(x), g'_2(x)) = (\cos(x), -\sin(x)).$$

Thus, via (3.30)

$$\begin{aligned} h'(x) &= \langle \nabla f(g_1(x), g_2(x)), (g_1'(x), g_2'(x)) \rangle \\ &= \langle (2 \sin(x), 1), (\cos(x), -\sin(x)) \rangle \\ &= 2 \sin(x) \cos(x) - \sin(x), \end{aligned}$$

as we also saw in (3.30).

While the chain rule did not really make the calculations easier for us in Example 3.7.1, it is important to become familiar with the use of the gradient notation and the inner product in (3.29). Indeed, as we will see in Section 4.6 on page 111 this type of notation is indispensable in order to consider Taylor polynomials in higher dimensions.

3.8 Differentiability of vector functions of several variables

We will now turn to differentiability of a *vector* function $\mathbf{f} = (f_1, \dots, f_k) : U \rightarrow \mathbb{R}^k$, where U is an open set in \mathbb{R}^n . We again need to generalize Equation (3.26) on page 85, this time to an equality in \mathbb{R}^k since $\mathbf{f}(\mathbf{x}_0) \in \mathbb{R}^k$. The modifications are as follows.

- 1) We will replace the linear map $\mathbf{h} \mapsto \mathbf{c}^T \mathbf{h}$, $\mathbb{R}^n \rightarrow \mathbb{R}$ with a linear map $L : \mathbb{R}^n \rightarrow \mathbb{R}^k$, i.e., $\mathbf{h} \mapsto L(\mathbf{h})$
- 2) We will replace the function $\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$ by a vector function $\boldsymbol{\varepsilon} : \mathbb{R}^n \rightarrow \mathbb{R}^k$.

Definition 3.8.1 Differentiability of vector functions

Consider a vector function $\mathbf{f} : U \rightarrow \mathbb{R}^k$, where U is an open set in \mathbb{R}^n , and let $\mathbf{x}_0 \in U$. We say that the function \mathbf{f} is differentiable at the point \mathbf{x}_0 if there exist a linear map $L : \mathbb{R}^n \rightarrow \mathbb{R}^k$ and a function $\boldsymbol{\varepsilon} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ such that $\boldsymbol{\varepsilon}(\mathbf{h}) \rightarrow \mathbf{0}$ as $\mathbf{h} \rightarrow \mathbf{0}$ and

$$\mathbf{f}(\mathbf{x}_0 + \mathbf{h}) - \mathbf{f}(\mathbf{x}_0) - L(\mathbf{h}) - \boldsymbol{\varepsilon}(\mathbf{h}) \|\mathbf{h}\| = 0, \quad \text{for all } \mathbf{h} \in \mathbb{R}^n. \quad (3.31)$$

The linear map L is called the *differential* of \mathbf{f} at \mathbf{x}_0 and is denoted $d\mathbf{f}_{\mathbf{x}_0}$. From Lemma 10.28 in Mathematics 1a we know that L is of the form $L(\mathbf{h}) = \mathbf{J}_f \mathbf{h}$ where $\mathbf{J}_f \in M_{k \times n}(\mathbb{R})$ is a unique matrix called the *Jacobian matrix*.

3.8. Differentiability of vector functions of several variables

We already know examples of vector functions that can be differentiated. Indeed, a linear map is differentiable:

Lemma 3.8.1 Linear maps are differentiable

If $L : \mathbb{R}^n \rightarrow \mathbb{R}^k$ is linear then it is differentiable at all points $\mathbf{x}_0 \in \mathbb{R}^n$ and the differential is $dL_{\mathbf{x}_0} = L$. That is, if $A \in \mathbf{M}_{k \times n}(\mathbb{R})$, then the linear map $\mathbf{x} \mapsto A\mathbf{x}, \mathbb{R}^n \rightarrow \mathbb{R}^k$ is a differentiable vector function at all points $\mathbf{x}_0 \in \mathbb{R}^n$, and its Jacobian matrix is (the constant) matrix A .

Proof. Since $L : \mathbb{R}^n \rightarrow \mathbb{R}^k$ is linear it satisfies

$$L(\mathbf{x}_0 + \mathbf{h}) - L(\mathbf{x}_0) - L(\mathbf{h}) = \mathbf{0}$$

for all $\mathbf{x}_0, \mathbf{h} \in \mathbb{R}^n$. Thus, if we take $\boldsymbol{\varepsilon}$ to be the zero vector function $\boldsymbol{\varepsilon}(\mathbf{h}) = \mathbf{0}$ for all $\mathbf{h} \in \mathbb{R}^n$, then (3.31) is satisfied. ■

As for scalar functions, differentiability of vector functions implies continuity.

Theorem 3.8.2

Assume that the function $\mathbf{f} : U \rightarrow \mathbb{R}^k$ is differentiable at the point $\mathbf{x}_0 \in U$. Then \mathbf{f} is also continuous at \mathbf{x}_0 .

The proof of Theorem 3.6.4 on page 87 carries over verbatim to the setting of Theorem 3.8.2 so we will not repeat it here.

In terms of the Jacobian matrix (3.31) becomes

$$\mathbf{f}(\mathbf{x} + \mathbf{h}) = \mathbf{f}(\mathbf{x}) + \mathbf{J}_f \mathbf{h} + \boldsymbol{\varepsilon}(\mathbf{h}) \|\mathbf{h}\|. \quad (3.32)$$

We will return to this expressions in Section 4.7 on page 113.

As for scalar functions, where \mathbf{c} in (3.26) was equal to the gradient vector, the Jacobian matrix also has a very special form. If $\mathbf{f} : U \rightarrow \mathbb{R}^k$ is differentiable at \mathbf{x} , the Jacobian matrix \mathbf{J}_f is a $k \times n$ matrix given by:

$$\mathbf{J}_f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \frac{\partial f_1}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \frac{\partial f_2}{\partial x_1}(\mathbf{x}) & \frac{\partial f_2}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_2}{\partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_k}{\partial x_1}(\mathbf{x}) & \frac{\partial f_k}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_k}{\partial x_n}(\mathbf{x}) \end{bmatrix}. \quad (3.33)$$

The proof is similar to the proof of Theorem 3.6.2 on page 86 and will not be repeated here. Note that \mathbf{J}_f consists of gradient vectors of f_i as row

3.8. Differentiability of vector functions of several variables

vectors. This shows that if we take $k = 1$ in Definition 3.8.1, we recover Definition 3.6.1 on page 85.

As for continuity of vector functions, we can also check differentiability of a vector function by looking at each coordinate separately:

Theorem 3.8.3 Differentiability of coordinate functions

A vector function $\mathbf{f} = (f_1, \dots, f_k) : U \rightarrow \mathbb{R}^k$ is differentiable at \mathbf{x}_0 if and only if all the functions $f_i : U \rightarrow \mathbb{R}$, $i = 1, \dots, k$ are differentiable at \mathbf{x}_0 .

The proof is similar to the proof of Theorem 3.2.1 on page 73 and will not be given here.

When speaking of the Jacobian matrix, we should really say “the Jacobian matrix evaluated at \mathbf{x} ” as the matrix clearly depends on \mathbf{x} . If we want to stress this fact, we will write $\mathbf{J}_{\mathbf{f}}(\mathbf{x})$. Consider a vector function $\mathbf{f} : U \rightarrow \mathbb{R}^k$, where U is an open set in \mathbb{R}^n . The map $\mathbf{x} \mapsto \mathbf{J}_{\mathbf{f}}(\mathbf{x})$, $U \mapsto \mathbb{R}^{k \times m}$ is then called the *differential* of \mathbf{f} , and it is denoted $d\mathbf{f}_{\mathbf{x}} : U \mapsto \mathbb{R}^{k \times m}$. This map is clearly only well-defined for vector functions differentiable at *all* points $\mathbf{x} \in U$.

Definition 3.8.2 C^1 vector function

Let U be an open set in \mathbb{R}^n . A vector function $\mathbf{f} : U \rightarrow \mathbb{R}^k$ is called *differentiable* if it is differentiable at all points $\mathbf{x} \in U$. If all partial derivatives $\frac{\partial f_i}{\partial x_j} : U \rightarrow \mathbb{R}$ appearing in the Jacobian matrix are continuous functions, then the vector function \mathbf{f} is said to be *continuously differentiable* or, alternatively, \mathbf{f} is said to be a C^1 vector function.

A note on terminology: the class of C^0 vector functions consists of all continuous vector functions and the class C^1 consists of all differentiable functions whose derivative is continuous (such functions are also called continuously differentiable). In general, the class of C^m vector functions consists of functions for which the first m partial derivatives of *each* coordinate function exists and are continuous. The class C^∞ are functions so smooth that they belong to C^m for all values of m .

Just like for scalar functions, continuity and differentiability are preserved by the usual arithmetic operations, e.g., if $\mathbf{f} : U \rightarrow \mathbb{R}^k$ and $\mathbf{g} : U \rightarrow \mathbb{R}^k$ are differentiable at $\mathbf{x}_0 \in U$, then so is $\mathbf{f} + \mathbf{g}$. In particular, the composition of differentiable vector functions is differentiable. This result is the *generalized chain rule*.

Theorem 3.8.4 Generalized chain rule

Let $U \subseteq \mathbb{R}^n$ and $V \subseteq \mathbb{R}^k$ be open sets. Suppose that $\mathbf{f} : U \rightarrow \mathbb{R}^k$ and $\mathbf{g} : V \rightarrow \mathbb{R}^\ell$ are vector functions such that $\text{im}(\mathbf{f}) \subseteq V$, \mathbf{f} is differentiable at $\mathbf{x}_0 \in U$, and \mathbf{g} is differentiable at $\mathbf{y}_0 := \mathbf{f}(\mathbf{x}_0)$. Then $\mathbf{g} \circ \mathbf{f} : U \rightarrow \mathbb{R}^\ell$ is differentiable at \mathbf{x}_0 with differential $d(\mathbf{g} \circ \mathbf{f})_{\mathbf{x}_0} = d\mathbf{g}_{\mathbf{y}_0} \circ d\mathbf{f}_{\mathbf{x}_0}$ and Jacobian matrix

$$\mathbf{J}_{\mathbf{g} \circ \mathbf{f}}(\mathbf{x}_0) = \mathbf{J}_{\mathbf{g}}(\mathbf{y}_0)\mathbf{J}_{\mathbf{f}}(\mathbf{x}_0). \quad (3.34)$$

where $\mathbf{J}_{\mathbf{g} \circ \mathbf{f}} \in M_{\ell \times n}(\mathbb{R})$, $\mathbf{J}_{\mathbf{g}} \in M_{\ell \times k}(\mathbb{R})$, and $\mathbf{J}_{\mathbf{f}} \in M_{k \times n}(\mathbb{R})$.

The proof is not particularly difficult, however, for brevity, we will not give it here, but refer to [Theorem 3.28 in the additional notes](#). The generalized chain rule shows that the Jacobian matrix of composed vector functions is just a matrix-matrix product of other Jacobian matrices (evaluated at two different points). Note that, compared to the chain rule for scalar functions of one variable [Equation \(3.8\)](#) on page 67, the order of the Jacobian matrices in (3.34) is important. Indeed, $\mathbf{J}_{\mathbf{g}}(\mathbf{y}_0)$ is a matrix of size $\ell \times k$, while $\mathbf{J}_{\mathbf{f}}(\mathbf{x}_0)$ is of size $k \times n$. Hence, if $\ell \neq n$, only the matrix product in (3.34) is well-defined.

CHAPTER 4

Taylor Approximation

In mathematics, physics, and engineering we often face functions that are complicated to work with. The reason can be, e.g., that the function is given by a complicated expression, or that we only know the function in terms of its graph. In such cases it is helpful if we can approximate the function with a simpler function, for example, a polynomial. The precise meaning of this will be clear soon, but the idea is to search for a polynomial P such that the deviation between the function values $f(x)$ of the given “complicated function” f and $P(x)$ is small for all x in the considered interval. In many cases it is then possible to do the further calculations with the polynomial P instead of the function f , and still obtain useful and realistic results.

The starting point of Taylor approximations is the notion of differentiability from [Definition 3.8.1](#) on page 90. Suppose \mathbf{f} is differentiable at $\mathbf{x}_0 \in U$ as in [Definition 3.8.1](#). By a simple rearrangement of the terms in [\(3.31\)](#), we get, by setting $\mathbf{x} := \mathbf{x}_0 + \mathbf{h}$,

$$\begin{aligned}\mathbf{f}(\mathbf{x}) &= \mathbf{f}(\mathbf{x}_0 + \mathbf{h}) = \mathbf{f}(\mathbf{x}_0) + \mathbf{J}_f(\mathbf{x}_0)(\mathbf{h}) + \varepsilon(\mathbf{h})\|\mathbf{h}\| \\ &= \mathbf{f}(\mathbf{x}_0) + \mathbf{J}_f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \varepsilon(\mathbf{x} - \mathbf{x}_0)\|\mathbf{x} - \mathbf{x}_0\|\end{aligned}\quad (4.1)$$

for $\mathbf{h} \in \mathbb{R}^n$. The goal is to generalize the formula [\(4.1\)](#) to higher order derivatives and, in particular, to quantify the behavior of the ε -function. For now, we just know that $\varepsilon(\mathbf{x} - \mathbf{x}_0) \rightarrow \mathbf{0}$ as $\mathbf{x} - \mathbf{x}_0 \rightarrow \mathbf{0}$, i.e., as $\mathbf{x} \rightarrow \mathbf{x}_0$. The function $P_{1,\mathbf{f},\mathbf{x}_0} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ defined by $P_{1,\mathbf{f},\mathbf{x}_0}(\mathbf{x}) = \mathbf{f}(\mathbf{x}_0) + \mathbf{J}_f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)$ for $\mathbf{x} \in \mathbb{R}^n$ is the so-called first-degree Taylor polynomial for \mathbf{f} at the point \mathbf{x}_0 . Hence, from [Chapter 3](#), we already know that, for small $\mathbf{h} = \mathbf{x} - \mathbf{x}_0$, i.e., for \mathbf{x} close to \mathbf{x}_0 , we can approximate $\mathbf{f}(\mathbf{x})$ by a first degree polynomial $P_{1,\mathbf{f},\mathbf{x}_0}(\mathbf{x})$ (in n variables) and the error $\mathbf{f}(\mathbf{x}) - P_{1,\mathbf{f},\mathbf{x}_0}(\mathbf{x})$ goes to zero faster than $\|\mathbf{h}\|$. As we will see, we can obtain better approximations using higher degree polynomials.

4.1. The tangent for a function of one variable

However, let us take a step back and first consider the simple setup of approximating a scalar function f of *one* variable by a first-degree polynomial. This is just the well-known tangent line for a function of one variable and the topic of Section 4.1. This approximation is generalized and improved in Section 4.2, where we approximate the function f by higher-degree polynomials, the so-called Taylor polynomials. The associated error estimates, i.e., the estimates of the deviation between the function values $f(x)$ and the values $P(x)$ of the Taylor polynomials, are considered in Section 4.3. In the sections 4.4–4.6 we run precisely the same program for functions of several variables. Finally, in Section 4.7 we return to Taylor polynomials for vector functions.

The Taylor polynomials we will encounter in this chapter depend on the given function \mathbf{f} , the desired degree of the Taylor polynomial K , and a fixed “expansion” point \mathbf{x}_0 in the domain of \mathbf{f} . The polynomial will be denoted by $P_{K,\mathbf{f},\mathbf{x}_0}(\mathbf{x})$, $P_{K,\mathbf{f}}(\mathbf{x})$, or simply $P_K(\mathbf{x})$ if \mathbf{f} and \mathbf{x}_0 are clear from the context.

4.1 The tangent for a function of one variable

Definition 4.1.1

Let $I \subset \mathbb{R}$ denote an interval and consider a differentiable function $f : I \rightarrow \mathbb{R}$. Fix a point $x_0 \in I$. Then the first-degree approximating polynomial – or the *Taylor polynomial of first degree* – at the point x_0 is defined by

$$P_{1,f,x_0}(x) = P_1(x) = f(x_0) + f'(x_0)(x - x_0), \quad x \in \mathbb{R}. \quad (4.2)$$

The graph of the function P_1 , i.e., the straight line given by

$$y = f(x_0) + f'(x_0)(x - x_0), \quad x \in \mathbb{R}, \quad (4.3)$$

is called the *tangent* of the function f at the point x_0 .

Remark 4.1.1

The form of the polynomial $P_1 = P_{1,f,x_0}$ in (4.2) implies that two very particular properties hold, namely,

$$P_1(x_0) = f(x_0) \quad \text{and} \quad P_1'(x_0) = f'(x_0).$$

That is, the function value of P_1 at the point x_0 coincides with the function value $f(x_0)$, i.e., the graph of the approximating polynomial P_1 goes through

4.1. The tangent for a function of one variable

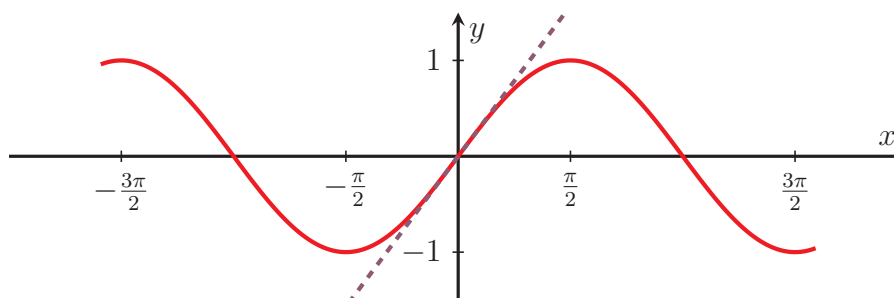


Figure 4.1: The function $f(x) = \sin(x)$ and its first-degree Taylor polynomial at $x_0 = 0$ (dashed).

the point $(x_0, f(x_0))$. The fact that also the first derivative of P_1 and f coincide at x_0 , implies that the polynomial P_1 provides an approximation of the function f in a (small) neighborhood of x_0 . The graph of P_1 it is a straight line with slope $f'(x_0)$; see the subsequent [Example 4.1.1](#) for an illustration of this. In particular, the slope $f'(x_0)$ tells how the function f evolves around the point x_0 : for example, a positive value of $f'(x_0)$ means that the function f is increasing at the point x_0 .

Example 4.1.1

Consider the function $f : \mathbb{R} \rightarrow \mathbb{R}$ given by $f(x) = \sin(x)$, $x \in \mathbb{R}$. Then $f'(x) = \cos(x)$. If we choose $x_0 = 0$, the Taylor polynomial of first degree at x_0 is given by

$$P_1(x) = \sin(0) + \cos(0)(x - 0) = x, \quad x \in \mathbb{R}.$$

We see from Figure 4.1 that P_1 indeed provides a quite good approximation of the function f for x close to $x_0 = 0$.

In practice it depends on a given application whether we can consider the Taylor polynomial of first degree as a satisfying approximation of a given function. Depending on the function f , the interval where P_1 is approximating the function f might be very small:

Example 4.1.2

Consider the function $f : \mathbb{R} \setminus \{0\} \rightarrow \mathbb{R}$ given by

$$f(x) = \frac{1}{x} \sin\left(\frac{1}{x}\right), \quad x \neq 0. \quad (4.4)$$

4.1. The tangent for a function of one variable

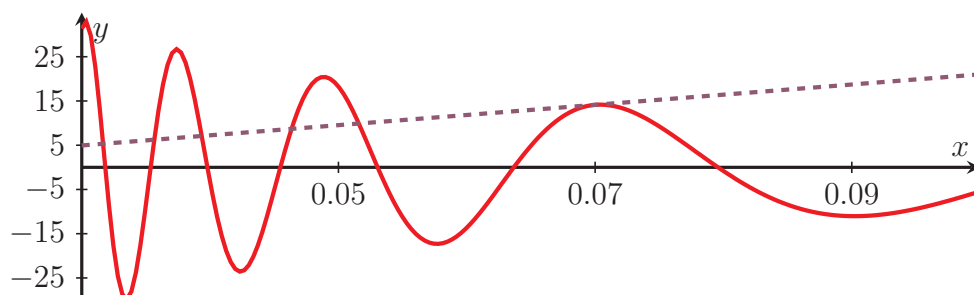


Figure 4.2: The function f in (4.4) and its first-degree Taylor polynomial at $x = 0.07$, shown on the interval $[0.03, 0.1]$ (dashed).

The graph of f is sketched on Figure 4.2 for $x \in [0.03, 0.1]$; note that the function has “infinitely” many oscillations around $x = 0$. Due to the strong oscillations of the function f it is clear that for values of x_0 close to 0, the Taylor polynomial of first degree at x_0 can only provide us with a reasonable approximation in a very small interval containing x_0 . In other words: overall, we can not claim that the first-degree Taylor polynomial provides us with a good approximation of the function f . See Figure 4.2 for an illustration of this.

We also note that the information provided by the first-degree Taylor polynomial is too limited to distinguish even very different functions:

Example 4.1.3

Figure 4.3 shows the graphs of the functions

$$f_1(x) = x^2, \quad f_2(x) = -x^2, \quad f_3(x) = x^3. \quad (4.5)$$

It is clear that these three functions behave very differently; nevertheless they all have the same first-degree Taylor polynomial at the point $x_0 = 0$, namely, $P_1(x) = 0$.

Motivated by our observations in Example 4.1.2 and Example 4.1.3 we will now turn our attention to a more general theory for approximation using higher-degree polynomials.

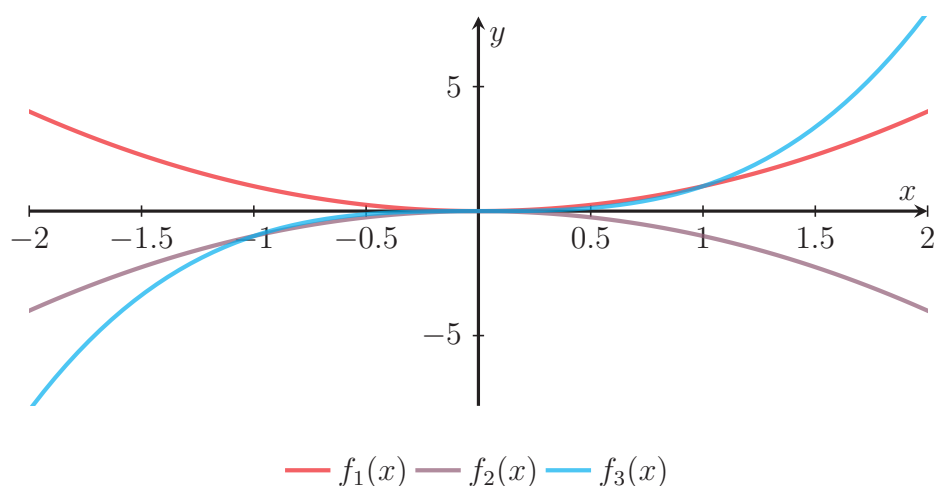


Figure 4.3: The functions in (4.5).

4.2 Taylor polynomials for functions of one variable

Example 4.1.2 and Example 4.1.3 demonstrate clearly that we need to develop a finer approximation theory than the one provided by first-degree Taylor polynomials. In order to do so, it is natural to use the expression (4.2) as the starting point. Indeed, as we saw in Remark 4.1.1 the graph of P_1 goes through the point $(x_0, f(x_0))$ and has the “correct slope”, so maybe we can obtain better approximation properties by adding “small correction terms” of the form $a(x - x_0)^2, b(x - x_0)^3, c(x - x_0)^4$, etc? This turns out indeed to be the case whenever the constants a, b, c etc are chosen correctly. We will now provide a general formula for an approximating polynomial of degree K for any positive integer K . In order to do so, we need to assume that the function f is K times differentiable; we will denote its derivatives at the point x_0 by $f'(x_0), f^{(2)}(x_0), f^{(3)}(x_0)$, etc.

Definition 4.2.1 Taylor polynomial of K th degree

Fix a positive integer K . Let $I \subset \mathbb{R}$ denote an interval and consider a K times differentiable function $f : I \rightarrow \mathbb{R}$. Fix a point $x_0 \in I$. Then the *Taylor polynomial of K th degree* $P_{K,f,x_0} = P_K$ at the point x_0 is defined by

$$P_K(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f^{(2)}(x_0)}{2!}(x - x_0)^2 + \cdots + \frac{f^{(K)}(x_0)}{K!}(x - x_0)^K. \quad (4.6)$$

4.2. Taylor polynomials for functions of one variable

Recall that for a positive integer K , the expression $K!$ (pronounced “ K factorial”) is defined by

$$K! = K(K - 1)(K - 2) \cdots 1.$$

For convenience we also define $0! = 1$. In particular, we thus have

$$0! = 1, \quad 1! = 1, \quad 2! = 2, \quad 3! = 6, \quad 4! = 24, \quad 5! = 120, \quad 6! = 720.$$

Note that for $K = 1$ the expression (4.6) indeed corresponds to our definition of the first degree Taylor polynomial. The K th Taylor polynomial can be written in a more compressed form as

$$P_K(x) = \sum_{k=0}^K \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k. \quad (4.7)$$

The compressed expression (4.7) allows to work with the K th degree Taylor polynomial on a computer, using only little extra effort compared with the first-degree Taylor polynomials.

Example 4.2.1

Consider the function

$$f(x) = e^x, \quad x \in \mathbb{R}.$$

The function is arbitrarily often differentiable, and

$$f'(x) = f^{(2)}(x) = \cdots = f^{(k)}(x) = \cdots = e^x.$$

For $K \in \mathbb{N}$ the K th Taylor polynomial at $x_0 = 0$ is given by

$$\begin{aligned} P_K(x) &= f(0) + \frac{f'(0)}{1!}x + \frac{f''(0)}{2!}x^2 + \cdots + \frac{f^{(K)}(0)}{K!}x^K \\ &= 1 + x + \frac{1}{2}x^2 + \frac{1}{3!}x^3 + \cdots + \frac{1}{K!}x^K \\ &= \sum_{k=0}^K \frac{1}{k!}x^k. \end{aligned} \quad (4.8)$$

Figure 4.4 shows that the fifth-degree Taylor polynomial at $x_0 = 1$ gives a very good approximation to the function in a quite large interval containing $x_0 = 1$.

4.2. Taylor polynomials for functions of one variable

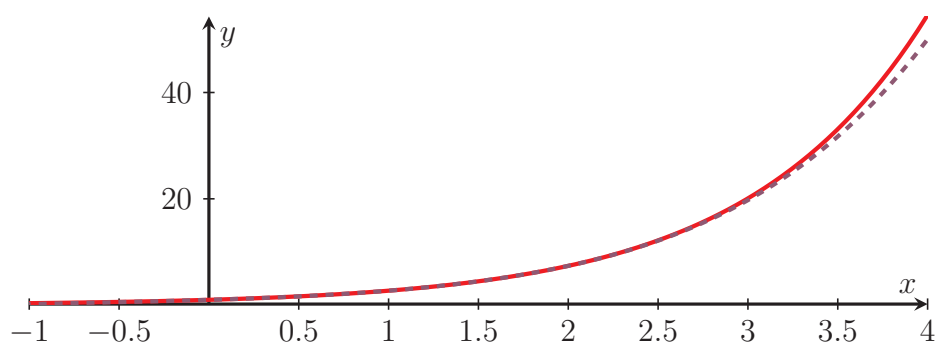


Figure 4.4: The function $f(x) = e^x$ and its fifth-degree Taylor polynomial at $x_0 = 1$ (dashed).

Example 4.2.2

For a function like in [Example 4.1.2](#),

$$f(x) = \frac{1}{x} \sin\left(\frac{1}{x}\right), \quad x \neq 0, \quad (4.9)$$

it is cumbersome to calculate higher order derivatives by hand, but for a mathematical computer program it is easy to do so and also to plot the graph. [Figure 4.5](#) shows the 4th degree Taylor polynomial at the point $x_0 = 0.07$. Compared with the first-degree Taylor polynomial shown on [Figure 4.2](#), we see that we get a good approximation of the function f over a significantly larger interval.

Remark 4.2.1

Note the very particular form of the “correction terms” in (4.6): each term of the form $(x - x_0)^k$, $k = 0, 1, \dots, K$, is multiplied with the corresponding factor $\frac{f^{(k)}(x_0)}{k!}$. This particular structure implies (compare with [Remark 4.1.1](#) on page 95) that the function value and the value of the first K derivatives of the polynomial P_K and the function f agree at the point x_0 , i.e.,

$$\begin{aligned} P_K(x_0) &= f(x_0), \\ P'_K(x_0) &= f'(x_0), \\ P_K^{(2)}(x_0) &= f^{(2)}(x_0) \\ &\dots, \\ P_K^{(K)}(x_0) &= f^{(K)}(x_0). \end{aligned}$$

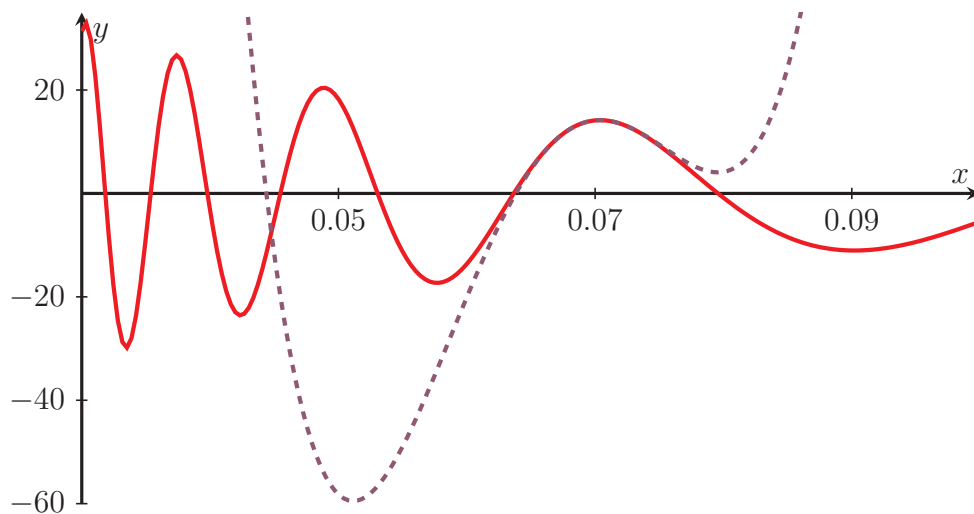


Figure 4.5: The function f in (4.9) and its fourth-degree Taylor polynomial at $x = 0.07$.

4.3 Taylor's formula for functions of one variable

Example 4.2.1 and Example 4.2.2 indicate that we indeed improve the approximation quality by considering higher degree Taylor polynomials. However, in order to apply the theory we need to have more exact knowledge about the approximation quality. In other words: we need to be able to estimate the difference between the function value $f(x)$ and the corresponding approximations $P_K(x)$ for x belonging to a suitable interval containing the point x_0 . In order to do so, we define the so-called *remainder term* by

$$R_K(x) := f(x) - P_K(x) = f(x) - \sum_{k=0}^K \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k. \quad (4.10)$$

Our goal is to develop methods to answer the following type of questions:

- Given a fixed value of $K \in \mathbb{N}$, how do we estimate the maximal deviation $|f(x) - P_K(x)|$ between $f(x)$ and $P_K(x)$ over a given interval I ?
- Given a fixed value of $K \in \mathbb{N}$, how do we specify an interval containing the point x_0 such that the maximal deviation between $f(x)$ and $P_K(x)$, measured over the interval, is below a certain prescribed threshold?

4.3. Taylor's formula for functions of one variable

- Given a fixed interval containing x_0 , how large do we have to choose the degree K of the Taylor polynomial P_K such that the maximal deviation between $f(x)$ and $P_K(x)$ is below a certain threshold, measured over the entire interval?

The technical tool to answer such questions is known under the name *Taylor's formula*.

Lemma 4.3.1 Taylor's formula

Let $I \subset \mathbb{R}$ denote an interval and assume that the function $f : I \rightarrow \mathbb{R}$ is arbitrarily often differentiable. Let $x_0 \in I$ and $K \in \mathbb{N}$. Then, for each $x \in I$ there exists a scalar ξ between x and x_0 such that

$$R_K(x) = f(x) - P_K(x) = \frac{f^{(K+1)}(\xi)}{(K+1)!} (x - x_0)^{K+1}. \quad (4.11)$$

The proof of Lemma 4.3.1 is quite lengthy and is given in Appendix A.1, see Lemma A.1.3. Note that Taylor's formula does not provide us with a completely explicit expression for the difference $f(x) - P_K(x)$, it simply says that when we fix $x \in I$, there exists a corresponding ξ such that (4.11) holds. Taylor's formula does not tell us how we can choose ξ ; it just tells us that such a ξ exists. In applications of Lemma 4.3.1 it is typically a technical challenge that different choices of x also lead to different values of ξ . Before we show how to apply Taylor's formula, let us restate the result:

Lemma 4.3.2 Taylor's formula

Let $I \subset \mathbb{R}$ denote an interval and assume that the function $f : I \rightarrow \mathbb{R}$ is arbitrarily often differentiable. Let $x_0 \in I$ and $K \in \mathbb{N}$. Then there exists a function $\varepsilon_K : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$R_K(x) = f(x) - P_K(x) = \varepsilon_K(x - x_0) (x - x_0)^K, \quad \text{for all } x \in I, \quad (4.12)$$

and $\varepsilon_K(x - x_0) \rightarrow 0$ as $x - x_0 \rightarrow 0$, i.e., $x \rightarrow x_0$.

Proof. Fix some $x_0 \in I$. The scalar ξ in the expression (4.11) actually depends on the considered $x \in I$ and the parameter K so for the sake of clarity of the proof we will denote it by $\xi_{x,K}$. Thus,

$$R_K(x) = \frac{f^{(K+1)}(\xi_{x,K})}{(K+1)!} (x - x_0)^{K+1},$$

4.3. Taylor's formula for functions of one variable

and hence

$$R_K(x) = (x - x_0)^K \left[\frac{f^{(K+1)}(\xi_{x,K})}{(K+1)!} (x - x_0) \right].$$

This corresponds precisely to (4.12) with

$$\varepsilon_K(x - x_0) := \frac{f^{(K+1)}(\xi_{x,K})}{(K+1)!} (x - x_0).$$

Also, when $x \rightarrow x_0$, it follows that $\xi_{x,K} \rightarrow x_0$. The function $f^{(K+1)}$ is differentiable by assumption, and hence continuous; therefore $f^{(K+1)}(\xi_x) \rightarrow f^{(K+1)}(x_0)$ as $x \rightarrow x_0$. This implies that

$$\varepsilon_K(x - x_0) = \frac{f^{(K+1)}(\xi_{x,K})}{(K+1)!} (x - x_0) \rightarrow \frac{f^{(K+1)}(x_0)}{(K+1)!} \cdot 0 = 0,$$

as $x \rightarrow x_0$, i.e., as $x - x_0 \rightarrow 0$, as claimed. ■

We are now ready to address the question of how to apply Taylor's formula to estimate the maximal deviation $|f(x) - P_K(x)|$ over a given interval. This is done via *Taylor's theorem*, stated next. It shows that under certain conditions, we can find a polynomial which is as close to f in a bounded interval as we wish:

Theorem 4.3.3 Taylor's theorem

Let $I \subset \mathbb{R}$ be an interval. Assume that the function $f : I \rightarrow \mathbb{R}$ is arbitrarily often differentiable and that there exists a constant $C > 0$ such that

$$|f^{(k)}(x)| \leq C, \quad \text{for all } k \in \mathbb{N} \text{ and all } x \in I. \quad (4.13)$$

Let $x_0 \in I$. Then for all $K \in \mathbb{N}$ and all $x \in I$,

$$\left| f(x) - \sum_{k=0}^K \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k \right| \leq \frac{C}{(K+1)!} |x - x_0|^{K+1}. \quad (4.14)$$

In particular, if I is a bounded interval, there exists for arbitrary $\epsilon > 0$ an $K_0 \in \mathbb{N}$ such that

$$\left| f(x) - \sum_{k=0}^K \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k \right| \leq \epsilon \quad \text{for all } x \in I \text{ and all } K \geq K_0. \quad (4.15)$$

The proof of Theorem 4.3.3 is given in Appendix A.1 on page 188.

4.3. Taylor's formula for functions of one variable

Note that in case we can choose $x_0 = 0$, Taylor's theorem takes a slightly simpler form. Under the assumptions in [Theorem 4.3.3](#) we can then approximate f near $x_0 = 0$ as well as we like with a polynomial

$$P_K(x) = \sum_{k=0}^K \frac{f^{(k)}(0)}{k!} x^k, \quad (4.16)$$

by choosing the degree K sufficiently high.

The inequality (4.14) can be used to decide how many terms we should include in the Taylor polynomial in order to reach a certain approximation of a given function on a fixed interval. The following example illustrates this for the exponential function, and also shows how the number of terms depend on the considered interval and on how well we want the function to be approximated.

Example 4.3.1

Consider the function

$$f(x) = e^x, \quad x \in \mathbb{R}.$$

We saw in [Example 4.2.1](#) that the K th Taylor polynomial at $x_0 = 0$ is given by

$$\begin{aligned} P_K(x) &= 1 + x + \frac{1}{2}x^2 + \frac{1}{3!}x^3 + \cdots + \frac{1}{K!}x^K \\ &= \sum_{k=0}^K \frac{1}{k!}x^k. \end{aligned} \quad (4.17)$$

We now aim at finding three values of $K \in \mathbb{N}$, namely:

(a) $K \in \mathbb{N}$ such that

$$\left| e^x - \sum_{k=0}^K \frac{x^k}{k!} \right| \leq 0.015 \quad \text{for all } x \in [-1, 2]. \quad (4.18)$$

(b) $K \in \mathbb{N}$ such that

$$\left| e^x - \sum_{k=0}^K \frac{x^k}{k!} \right| \leq 0.0015 \quad \text{for all } x \in [-1, 2]. \quad (4.19)$$

(c) $K \in \mathbb{N}$ such that

$$\left| e^x - \sum_{k=0}^K \frac{x^k}{k!} \right| \leq 0.015 \quad \text{for all } x \in [-5, 5]. \quad (4.20)$$

4.4. The tangent plane for functions of several variables

In order to do so, we first observe that the estimate (4.14) involves the constant C , which we can choose as the maximum of all derivatives of f on the considered interval. For the function f we have that $f^{(k)}(x) = e^x$, i.e., all derivatives equal the function itself.

For Item (a) we can thus take $C = e^2$. Now (4.14) shows that (4.18) is obtained if we choose K such that

$$\frac{e^2}{(K+1)!} 2^{K+1} \leq 0.015; \quad (4.21)$$

this is satisfied for all $K \geq 8$.

In Item (b), we obtain an appropriate value for $K \in \mathbb{N}$ by replacing the number 0.015 on the right-hand side of (4.21) by 0.0015; the obtained inequality is satisfied for all $K \geq 10$.

Note than in Item (c), we enlarge the interval on which we want to approximate f . In order to find an appropriate K -value in Item (a), we need to replace the previous value for C with $C = e^5$; we obtain the inequality

$$\frac{e^5}{(K+1)!} 5^{K+1} \leq 0.015,$$

which is satisfied for $K \geq 19$.

4.4 The tangent plane for functions of several variables

We will now take the first steps to generalize the theory for Taylor polynomials to higher dimensions. In this section we focus on first-degree Taylor polynomials; the case of second degree Taylor polynomials will be treated in Section 4.5.

Consider a function $f : U \rightarrow \mathbb{R}$, where U is an open set in \mathbb{R}^n , and fix a vector $\mathbf{x}_0 \in U$. Similarly to our generalizations of the concepts of continuity and differentiability, our task is to identify a general definition of the Taylor polynomial of first degree at the point \mathbf{x}_0 ; the definition must correspond to Definition 4.1.1 in the case $n = 1$.

Looking now at (4.2), the difference $x - x_0$ will then be replaced by the difference $\mathbf{x} - \mathbf{x}_0$, where $\mathbf{x} \in \mathbb{R}^n$. For a function of n variables it is natural to replace the derivative $f'(x_0)$ in (4.2) by the gradient vector at the point

4.4. The tangent plane for functions of several variables

\mathbf{x}_0 . The gradient vector was defined in Equation (3.20) on page 77 as

$$\nabla f(\mathbf{x}_0) = \left(\frac{\partial f}{\partial x_1}(\mathbf{x}_0), \frac{\partial f}{\partial x_2}(\mathbf{x}_0), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}_0) \right). \quad (4.22)$$

With these observations, the expression in (4.2) indeed has a generalization to \mathbb{R}^n if we consider the product between $f'(\mathbf{x}_0)$ and $(\mathbf{x} - \mathbf{x}_0)$ as the inner product in \mathbb{R} . We thus arrive at the following definition:

Definition 4.4.1 Taylor polynomial of first degree

Let U denote an open set in \mathbb{R}^n and consider a function $f : U \rightarrow \mathbb{R}$ for which all partial derivatives exist on the set U . Fix a point $\mathbf{x}_0 \in U$. Then the *Taylor polynomial of first degree* at the point \mathbf{x}_0 is defined by

$$P_{1,f,\mathbf{x}_0}(\mathbf{x}) = P_1(\mathbf{x}) = f(\mathbf{x}_0) + \langle (\mathbf{x} - \mathbf{x}_0), \nabla f(\mathbf{x}_0) \rangle, \quad \mathbf{x} \in \mathbb{R}^n. \quad (4.23)$$

The graph of the function P_1 , i.e., the set of points $(\mathbf{x}, P_1(\mathbf{x}))$, $\mathbf{x} \in \mathbb{R}^n$, is called the *tangent plane* of the function f at the point \mathbf{x}_0 .

We will often consider functions of two variables. In such cases we can apply the following explicit expression for the first-degree Taylor polynomial:

Lemma 4.4.1

Let U denote an open set in \mathbb{R}^2 and consider a differentiable function $f : U \rightarrow \mathbb{R}$. Fix a point $(x_0, y_0) \in U$. Then the Taylor polynomial of first degree at the point (x_0, y_0) is

$$P_1(x, y) = f(x_0, y_0) + \frac{\partial f}{\partial x}(x_0, y_0)(x - x_0) + \frac{\partial f}{\partial y}(x_0, y_0)(y - y_0), \quad (4.24)$$

where $(x, y) \in \mathbb{R}^2$.

Proof. Write $\mathbf{x} = (x, y)$ and $\mathbf{x}_0 = (x_0, y_0)$. Then (4.23) shows that

$$\begin{aligned} P_1(x, y) &= f(\mathbf{x}_0) + \langle (\mathbf{x} - \mathbf{x}_0), \nabla f(\mathbf{x}_0) \rangle \\ &= f(x_0, y_0) + \begin{bmatrix} \frac{\partial f}{\partial x}(x_0, y_0) & \frac{\partial f}{\partial y}(x_0, y_0) \end{bmatrix} \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix} \\ &= f(x_0, y_0) + \frac{\partial f}{\partial x}(x_0, y_0)(x - x_0) + \frac{\partial f}{\partial y}(x_0, y_0)(y - y_0), \end{aligned}$$

as claimed. ■

4.5. Taylor polynomials for functions of several variables

Remark 4.4.1

Lemma 4.4.1 shows that the first-degree Taylor polynomial at a point (x_0, y_0) for a function of two variables satisfies that

$$P_1(x_0, y_0) = f(x_0, y_0), \quad \frac{\partial P_1}{\partial x}(x_0, y_0) = \frac{\partial f}{\partial x}(x_0, y_0), \quad \frac{\partial P_1}{\partial y}(x_0, y_0) = \frac{\partial f}{\partial y}(x_0, y_0).$$

That is, the function f and its first-degree Taylor polynomial P_1 coincide at the point (x_0, y_0) , and the gradient of f and the gradient of P_1 also coincide at (x_0, y_0) . This property corresponds to what we have for the first-degree Taylor polynomial for a function of one variable; see Remark 4.1.1. The similar result (with a similar calculation) holds for a function of n variables.

Example 4.4.1

Let us return to the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = \sin(x^2 + y),$$

which we considered in Example 3.3.2. The calculations in Example 3.3.2 show that

$$\nabla f(x, y) = (2x \cos(x^2 + y), \cos(x^2 + y)).$$

Thus, at the point $(x_0, y_0) = (\sqrt{\pi}, \pi)$ we have

$$\nabla f(\sqrt{\pi}, \pi) = (2\sqrt{\pi} \cos(\pi + \pi), \cos(\pi + \pi)) = (2\sqrt{\pi}, 1).$$

Therefore the Taylor polynomial of first degree at the point $(x_0, y_0) = (\sqrt{\pi}, \pi)$ is

$$\begin{aligned} P_1(x, y) &= f(\sqrt{\pi}, \pi) + \left\langle \begin{bmatrix} x - \sqrt{\pi} \\ y - \pi \end{bmatrix}, \begin{bmatrix} 2\sqrt{\pi} \\ 1 \end{bmatrix} \right\rangle \\ &= \sin(\pi + \pi) + (x - \sqrt{\pi})2\sqrt{\pi} + (y - \pi) \cdot 1 \\ &= 2\sqrt{\pi}x + y - 3\pi, \quad (x, y) \in \mathbb{R}^2. \end{aligned}$$

4.5 Taylor polynomials for functions of several variables

As for functions of one variable, we can define higher-degree Taylor polynomials by adding higher-degree terms to the formula for the Taylor polynomial of first degree. We will here only consider the Taylor polynomial

4.5. Taylor polynomials for functions of several variables

of second degree as the notation for higher degree polynomials becomes rather complicated as one usually formulate such higher degree Taylor polynomials using tensors or multi-indices. We refer the interested reader to [the additional notes](#). To define second-degree Taylor polynomials of functions of several variables, we need the Hessian matrix $\mathbf{H}_f(\mathbf{x}_0)$ introduced in Equation (3.23) on page 81.

Definition 4.5.1 Taylor polynomial of second degree

Let U denote an open set in \mathbb{R}^n and consider a function $f : U \rightarrow \mathbb{R}$ for which all partial derivatives of order one and two exist on the set U . Fix a point $\mathbf{x}_0 \in U$. Then the *Taylor polynomial of second degree* $P_{2,f,\mathbf{x}_0} = P_2$ at the point \mathbf{x}_0 is defined by

$$P_2(\mathbf{x}) = f(\mathbf{x}_0) + \langle (\mathbf{x} - \mathbf{x}_0), \nabla f(\mathbf{x}_0) \rangle + \frac{1}{2} \langle (\mathbf{x} - \mathbf{x}_0), \mathbf{H}_f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \rangle \quad (4.25)$$

for $\mathbf{x} \in \mathbb{R}^n$.

If, in addition to the assumptions of Definition 4.5.1, all second-order partial derivatives are continuous in U , then the Hessian matrix is a real *symmetric* $n \times n$ matrix by Theorem 3.6.5 on page 87.

For functions of just two variables we can apply the following explicit expression for the second-degree Taylor polynomial:

Lemma 4.5.1

Let U denote an open set in \mathbb{R}^2 and consider a function $f : U \rightarrow \mathbb{R}$ for which all partial derivatives of second order exist *and are continuous* on the set U . Fix a point $(x_0, y_0) \in U$. Then the Taylor polynomial of second order at the point (x_0, y_0) is

$$P_2(x, y) = f(x_0, y_0) + \frac{\partial f}{\partial x}(x_0, y_0)(x - x_0) + \frac{\partial f}{\partial y}(x_0, y_0)(y - y_0) \quad (4.26)$$

$$\begin{aligned} &+ \frac{1}{2} \frac{\partial^2 f}{\partial x^2}(x_0, y_0)(x - x_0)^2 + \frac{1}{2} \frac{\partial^2 f}{\partial y^2}(x_0, y_0)(y - y_0)^2 \\ &+ \frac{\partial^2 f}{\partial x \partial y}(x_0, y_0)(x - x_0)(y - y_0), \quad (x, y) \in \mathbb{R}^2. \end{aligned} \quad (4.27)$$

Proof. The terms in the first line (4.26) corresponds exactly to the first-degree Taylor polynomial $f(\mathbf{x}_0) + \langle (\mathbf{x} - \mathbf{x}_0), \nabla f(\mathbf{x}_0) \rangle$, so we now focus on the new term

$$\frac{1}{2} \langle (\mathbf{x} - \mathbf{x}_0), \mathbf{H}_f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \rangle = \frac{1}{2} \langle \mathbf{H}_f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0), (\mathbf{x} - \mathbf{x}_0) \rangle,$$

4.5. Taylor polynomials for functions of several variables

where we have used that $\mathbf{H}_f(\mathbf{x}_0)$ is a symmetric matrix by [Theorem 3.6.5](#). Using the expression in [\(3.24\)](#) for the Hessian matrix at the point (x_0, y_0) , we see that in the two-dimensional case,

$$\begin{aligned} \mathbf{H}_f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) &= \mathbf{H}_f(x_0, y_0) \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial^2 f}{\partial x^2}(x_0, y_0) & \frac{\partial^2 f}{\partial x \partial y}(x_0, y_0) \\ \frac{\partial^2 f}{\partial y \partial x}(x_0, y_0) & \frac{\partial^2 f}{\partial y^2}(x_0, y_0) \end{bmatrix} \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial^2 f}{\partial x^2}(x_0, y_0)(x - x_0) + \frac{\partial^2 f}{\partial x \partial y}(x_0, y_0)(y - y_0) \\ \frac{\partial^2 f}{\partial y \partial x}(x_0, y_0)(x - x_0) + \frac{\partial^2 f}{\partial y^2}(x_0, y_0)(y - y_0) \end{bmatrix}. \end{aligned}$$

Therefore

$$\begin{aligned} \frac{1}{2} \langle \mathbf{H}_f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0), (\mathbf{x} - \mathbf{x}_0) \rangle &= \\ \frac{1}{2} \begin{bmatrix} x - x_0 & y - y_0 \end{bmatrix} &\begin{bmatrix} \frac{\partial^2 f}{\partial x^2}(x_0, y_0)(x - x_0) + \frac{\partial^2 f}{\partial x \partial y}(x_0, y_0)(y - y_0) \\ \frac{\partial^2 f}{\partial y \partial x}(x_0, y_0)(x - x_0) + \frac{\partial^2 f}{\partial y^2}(x_0, y_0)(y - y_0) \end{bmatrix} \\ &= \frac{1}{2} \frac{\partial^2 f}{\partial x^2}(x_0, y_0)(x - x_0)^2 + \frac{1}{2} \frac{\partial^2 f}{\partial y^2}(x_0, y_0)(y - y_0)^2 \\ &\quad + \frac{\partial^2 f}{\partial x \partial y}(x_0, y_0)(x - x_0)(y - y_0) \end{aligned}$$

as claimed. ■

The second-degree Taylor polynomial $P_{2,f,\mathbf{x}_0}(\mathbf{x})$ defined in [Equation \(4.25\)](#) on the preceding page is actually a class of functions we understand very well. In fact, P_{2,f,\mathbf{x}_0} is a *quadratic form*. We ask the reader to prove this fact in the next exercise.

Exercise 4.5.2

Show that the second-degree Taylor polynomial $P_{2,f,\mathbf{x}_0}(\mathbf{x})$ defined in [Equation \(4.25\)](#) on the previous page is a quadratic form as defined in [Definition 1.2.1](#) on page 11. You should express the matrix A , the column vector \mathbf{b} and the constant c in terms of f and its derivatives and \mathbf{x}_0 . *Hint:* Start by expanding all inner products in [\(4.25\)](#) using linearity of the inner product. Collect *all* constant terms (those not depending on \mathbf{x}) and set them equal to c . Proceed to the terms linear in \mathbf{x} .

4.5. Taylor polynomials for functions of several variables

Example 4.5.1

For the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = \sin(x^2 + y),$$

we saw in [Example 4.4.1](#) that, at the point $(x_0, y_0) = (\sqrt{\pi}, \pi)$, the Taylor polynomial of first degree is

$$P_1(x, y) = 2\sqrt{\pi}x + y - 3\pi, \quad (x, y) \in \mathbb{R}^2.$$

The second-order partial derivatives were calculated in [Example 3.5.2](#), where we saw that

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2}(x, y) &= 2 \cos(x^2 + y) - 4x^2 \sin(x^2 + y) \\ \frac{\partial^2 f}{\partial x \partial y}(x, y) &= -2x \sin(x^2 + y) \\ \frac{\partial^2 f}{\partial y \partial x}(x, y) &= -2x \sin(x^2 + y) \\ \frac{\partial^2 f}{\partial y^2}(x, y) &= -\sin(x^2 + y). \end{aligned}$$

Thus, at the point $(x_0, y_0) = (\sqrt{\pi}, \pi)$,

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2}(\sqrt{\pi}, \pi) &= 2 \cos(\pi + \pi) - 4\pi^2 \sin(\pi + \pi) = 2, \\ \frac{\partial^2 f}{\partial x \partial y}(\sqrt{\pi}, \pi) &= -2\sqrt{\pi} \sin(\pi + \pi) = 0, \\ \frac{\partial^2 f}{\partial y \partial x}(\sqrt{\pi}, \pi) &= -2\sqrt{\pi} \sin(\pi + \pi) = 0, \end{aligned}$$

and

$$\frac{\partial^2 f}{\partial y^2}(\sqrt{\pi}, \pi) = -\sin(\pi + \pi) = 0.$$

By [Lemma 4.5.1](#), the second-degree Taylor polynomial at the point $(x_0, y_0) = (\sqrt{\pi}, \pi)$ is

$$\begin{aligned} P_2(x, y) &= 2\sqrt{\pi}x + y - 3\pi + (x - \sqrt{\pi})^2 \\ &= x^2 + y - 2\pi. \end{aligned}$$

Example 4.5.2

For the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = \sin(x^2 + y^2),$$

direct calculation (try to do it) shows that the first-degree Taylor polynomial at the point $(x_0, y_0) = (0, 0)$ is

$$P_1(x, y) = 0,$$

and that the second-degree Taylor polynomial is

$$P_2(x, y) = x^2 + y^2.$$

Figure 4.6 shows the function f and the two Taylor polynomials.

4.6 Taylor's formula for functions of several variables

The Taylor polynomial P_2 of second degree is typically a better approximation of a given function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ than the first-degree Taylor approximation, but still it is just an approximation. That is, we can write

$$f(\mathbf{x}) = P_2(\mathbf{x}) + R_2(\mathbf{x}),$$

where the function $R_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ is the so called *remainder term*. Alternatively, we may write

$$R_2(\mathbf{x}) = f(\mathbf{x}) - P_2(\mathbf{x}).$$

In general we do not know exactly how the function R_2 look like, but as for the case of functions of one variable the important *Taylor's formula* quantifies its behavior:

Theorem 4.6.1 Taylor's formula

In the setup of Definition 4.5.1, the remainder-term R_2 has the form

$$R_2(\mathbf{x}) = f(\mathbf{x}) - P_2(\mathbf{x}) = \varepsilon(\mathbf{x} - \mathbf{x}_0) \|\mathbf{x} - \mathbf{x}_0\|^2 \quad (4.28)$$

for some function $\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$ for which $\varepsilon(\mathbf{x} - \mathbf{x}_0) \rightarrow 0$ whenever $\mathbf{x} \rightarrow \mathbf{x}_0$.

4.6. Taylor's formula for functions of several variables

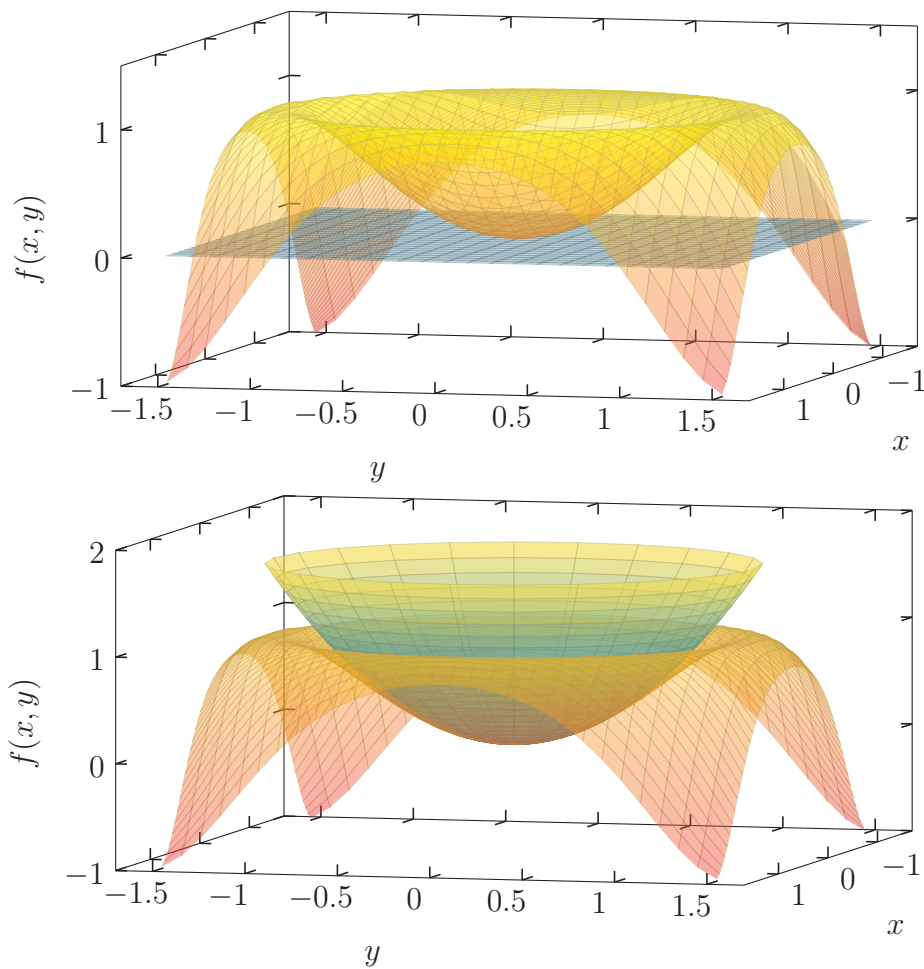


Figure 4.6: The function $f(x, y) = \sin(x^2 + y^2)$ and its Taylor polynomials of first degree and second degree at the point $(0, 0)$.

Theorem 4.6.1 can be proved by a number of tricks, reducing the problem to the one-dimensional case considered in Lemma 4.3.1; this leads to a result of a similar type (but of higher complexity) as in Lemma 4.3.1. Furthermore, there exist definitions of Taylor polynomials of arbitrary high degree, and corresponding error estimates. It is outside the scope of the current notes to go into an analysis of this.

4.7 Taylor polynomials for vector functions of several variables

We will now turn to Taylor polynomials of a *vector* function $\mathbf{f} = (f_1, \dots, f_k) : U \rightarrow \mathbb{R}^k$, where U , as usual, is an open set in \mathbb{R}^n . Luckily, the situation is not much more complicated than for scalar functions of several variables. Recall that each coordinate function $f_i : U \rightarrow \mathbb{R}$ is a scalar function of n variables. So to get the K th degree Taylor polynomial for \mathbf{f} at \mathbf{x}_0 , we “just” need to find the K th degree Taylor polynomial of each coordinate function f_i at the same point \mathbf{x}_0 . This leads to the following definition.

Definition 4.7.1 Taylor polynomial for vector functions

Let U denote an open set in \mathbb{R}^n and consider a function $\mathbf{f} = (f_1, \dots, f_k) : U \rightarrow \mathbb{R}^k$. Suppose all partial derivatives up to degree K of each coordinate function f_i , $i = 1, \dots, k$, exist on the set U . Fix a point $\mathbf{x}_0 \in U$. Then the *Taylor polynomial of K degree* $\mathbf{P}_{K,\mathbf{f},\mathbf{x}_0} = \mathbf{P}_K$ at the point \mathbf{x}_0 is defined as the vector function $\mathbf{P}_K : \mathbb{R}^n \rightarrow \mathbb{R}^k$ whose i th coordinate function is P_{K,f_i,\mathbf{x}_0} .

In Definition 4.7.1, one usually assumes that \mathbf{f} is a C^K vector function, see the discussion just below Definition 3.8.2 on page 92, as it guarantees that all the partial derivatives exist and are continuous. Note that Definition 4.7.1 is not very explicit, and we have only discussed how to write down P_{K,f_i,\mathbf{x}_0} for $K = 1$ and $K = 2$ in case $n > 1$. For $K = 1$ there is a simple formula for $\mathbf{P}_{1,\mathbf{f},\mathbf{x}_0}$ in terms of the Jacobian of \mathbf{f} . Suppose that $\mathbf{f} : U \rightarrow \mathbb{R}^k$ is a differentiable function at $\mathbf{x}_0 \in U$. Then the first-degree Taylor polynomial at the point \mathbf{x}_0 is given by:

$$\mathbf{P}_{1,\mathbf{f},\mathbf{x}_0}(\mathbf{x}) = \mathbf{f}(\mathbf{x}_0) + \mathbf{J}_{\mathbf{f}}(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0). \quad (4.29)$$

We actually derived this first-degree Taylor polynomial in Equation (4.1) on page 94 from Definition 3.8.1 on page 90. Moreover, we know the remainder term satisfies:

$$\mathbf{R}_{1,\mathbf{f},\mathbf{x}_0}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) - \mathbf{P}_{1,\mathbf{f},\mathbf{x}_0}(\mathbf{x}) = \boldsymbol{\varepsilon}(\mathbf{x} - \mathbf{x}_0) \|\mathbf{x} - \mathbf{x}_0\|,$$

where $\boldsymbol{\varepsilon}(\mathbf{x} - \mathbf{x}_0) \rightarrow \mathbf{0}$ as $\mathbf{x} \rightarrow \mathbf{x}_0$.

Exercise 4.7.1

Suppose that $\mathbf{f} = (f_1, \dots, f_k) : U \rightarrow \mathbb{R}^k$ is a differentiable function at $\mathbf{x}_0 \in U$, where $U \subseteq \mathbb{R}^n$ is an open set. Derive the formula in Equation (4.29) by applying the formula in Equation (4.23) on page 106 for each coordinate

4.7. Taylor polynomials for vector functions of several variables

function f_i , $i = 1, \dots, k$. *Hint:* Remember that the rows of the Jacobian matrix contain the gradient vector of each coordinate function.

It is rather technical to write down $\mathbf{P}_{K,f,x_0}(\mathbf{x})$ explicitly for $K > 1$, and it is beyond the scope of this text.

CHAPTER 5

Local and Global Extrema of Functions

Many important problems in engineering deals with *optimization*. For example, how can we – under given physical conditions – minimize the power consumption of an electric device? Or, how can we maximize the output from a wind mill or a solar panel? Typically, the value of the “item” to be optimized can be expressed as a function of a number of variables representing physical quantities. In such cases we are left with a purely mathematical question of how to maximize/minimize a function in a certain domain that is determined by the physical limitations for the involved variables. In this chapter we will discuss methods for optimization for *scalar* functions of several variables. We begin by reminding the reader of the one-dimensional case in [Section 5.1](#). The general case of a scalar function of several variables is then considered in [Section 5.2](#).

We will not consider vector functions in this chapter. In fact, one cannot optimize vector functions since such functions takes values in \mathbb{R}^k . Recall that it does not make sense to ask if a vector $\mathbf{y}_1 \in \mathbb{R}^k$ is smaller than or larger than $\mathbf{y}_2 \in \mathbb{R}^k$. Phrased mathematically, there is *no ordering* on \mathbb{R}^k so we *cannot* ask for the “smallest” or “largest” value of a function $\mathbf{f} : \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k$. We already see this issue for complex-valued scalar function, again, because there is no ordering on \mathbb{C} . However, the solution is evident: we optimize the norm of the function values instead. Hence, if $\mathbf{f} : \text{dom}(\mathbf{f}) \rightarrow \mathbb{F}^k$, we consider the function $\mathbf{x} \mapsto \|\mathbf{f}(\mathbf{x})\|, \text{dom}(\mathbf{f}) \rightarrow [0, \infty[$, which is a real-valued scalar function. The norm $\|\cdot\|$ is not necessarily the standard norm \mathbb{F}^k , and the choice often depends on the application in question.

For (real-valued) scalar functions, we will consider two types of optimization, a *global optimization* and a *local optimization*. Global

optimization aims at finding the maximum/minimum value of a given function over the entire domain. On the other hand, *local optimization* deals with an analysis of a given function in a small neighborhood of a given point. It provides methods to decide whether the given point at least yields the maximum/minimum value within a small neighborhood of the point. The main tool here is to approximate the function locally by a *quadratic form*, namely, a second-degree Taylor polynomial. Both local and global optimization are highly relevant in practice.

5.1 The range of functions of one variable

In this section we will remind the reader about standard results concerning the *image* or *range* of a function f that is defined on a closed and bounded interval $[a, b]$ in \mathbb{R} . The image set $\text{im}(f)$ of $f : [a, b] \rightarrow \mathbb{R}$ is also written as $f([a, b])$, and it is the set given by $\{f(x) \mid x \in [a, b]\}$. One of the fundamental theorems of analysis states that, if the function f defined on a closed and bounded interval is continuous, then the image set is itself a closed and bounded interval.

Theorem 5.1.1

Consider a continuous function $f : [a, b] \rightarrow \mathbb{R}$ defined on a bounded and closed interval $[a, b]$. Then the image of f is bounded and has a minimal value $m \in \mathbb{R}$ and maximal value $M \in \mathbb{R}$, that is,

$$\exists m, M \in \text{im}(f) \forall y \in \text{im}(f) : m \leq y \leq M.$$

Moreover, the image set has the form

$$f([a, b]) = [m, M]. \tag{5.1}$$

The proof of [Theorem 5.1.1](#) is beyond the scope of this text, and we refer the reader to [Theorem 2.9 in the additional notes](#) for a proof.

The value m in [Theorem 5.1.1](#) is called the *minimum* of the function f , and the value M is called the *maximum*. We use *extremum* to denote either the minimum or the maximum. [Theorem 5.1.1](#) tells us that there indeed exist $x_1 \in [a, b]$ such that $f(x_1) = m$, and we say that the function f *attains its minimum* in x_1 . Similarly, there exists $x_2 \in [a, b]$ such that $f(x_2) = M$, that is, f *attains its maximum* in x_2 . Furthermore, [Theorem 5.1.1](#) says that any y in the interval between m and M there exists an $x \in [a, b]$ so that $f(x) = y$.

5.1. The range of functions of one variable

For the validity of [Theorem 5.1.1](#) it is essential that all the assumptions are satisfied. The following examples illustrate this.

Example 5.1.1 (a) Consider the continuous function

$$f :]0, 1[\rightarrow \mathbb{R}, f(x) = \ln(x).$$

Then $\text{im}(f) = f(]0, 1[) =]-\infty, 0[$. That is, the image set is not bounded, and the function f neither has a minimum nor a maximum. Note that the domain $]0, 1[$ is non-closed, i.e., the conclusion does not contradict [Theorem 5.1.1](#). Note that $\ln(1) = 0$ is not a maximum of $\text{im}(f)$. It is, however, the supremum of $\text{im}(f)$ as the least upper bound of $]-\infty, 0[$. The function f does not attain a maximum since there is no $x \in \text{dom}(f)$ such that $f(x) = 1$.

(b) Consider the continuous function

$$f : [0, \infty[\rightarrow \mathbb{R}, f(x) = 3x^2.$$

Then $\text{im}(f) = f([0, \infty[) = [0, \infty[$. Hence, the image set is not a bounded interval, and the function does not have a maximum value. Note that the domain $[0, \infty[$ is not a bounded set so the conclusion does not contradict [Theorem 5.1.1](#).

(c) Consider the function

$$f : [-2, -1] \cup [1, 2] \rightarrow \mathbb{R}, f(x) = \begin{cases} -1 & \text{whenever } x \in [-2, -1], \\ 1 & \text{whenever } x \in [1, 2]. \end{cases}$$

Then the image set only consists of the two points $\{-1, 1\}$, and is thus not an interval. This does not contradict [Theorem 5.1.1](#); the function f is continuous (why?), and the domain $[-2, -1] \cup [1, 2]$ is closed and bounded, but it is not an interval.

(d) Consider the function

$$f : [-1, 1] \rightarrow \mathbb{R}, f(x) = \begin{cases} \frac{1}{x} & \text{whenever } x \neq 0, \\ 0 & \text{whenever } x = 0. \end{cases}$$

Then the image set is $f([-1, 1]) =]-\infty, -1] \cup \{0\} \cup [1, \infty[$. That is, even though the function f is defined on the closed and bounded interval $[-1, 1]$, the image set is not an interval; furthermore there is neither a minimum or a maximum. This does not contradict [Theorem 5.1.1](#) since the function f is not continuous.

5.1. The range of functions of one variable

Note that [Theorem 5.1.1](#) is a theoretical result: it tells us that a function f that satisfies the stated conditions has a minimum and a maximum, but not how to find these values. For this purpose we need the following result.

Theorem 5.1.2

Consider a continuous function $f : [a, b] \rightarrow \mathbb{R}$. Then, if f attains its minimum or maximum at the point x_0 , one of the following possibilities occur:

- (i) $x_0 = a$;
- (ii) $x_0 = b$;
- (iii) x_0 is a point where f is not differentiable;
- (iv) x_0 is a point where f is differentiable, and $f'(x_0) = 0$.

Proof. The result follows from a more general result proved in [Theorem 5.2.2](#). ■

In other words: in order to find the minimum and the maximum for the function f , it is enough to search among the *end points* a, b of the domain and the points x_0 where f is either not differentiable, or $f'(x_0) = 0$.

Example 5.1.2

Consider the function

$$f : [0, 2] \rightarrow \mathbb{R}, \quad f(x) = x^3 - x.$$

Then f satisfies the conditions in [Theorem 5.1.1](#). In order to find the minimum and the maximum, we now apply [Theorem 5.1.2](#). First, the function f is differentiable for all $x \in]0, 2[$, and

$$f'(x) = 3x^2 - 1.$$

Note that the solutions to the equation $3x^2 - 1 = 0$ are $x = \pm 1/\sqrt{3}$. However, $x = -1/\sqrt{3}$ does not belong to the considered domain $[0, 2]$ and will thus be discarded. According to [Theorem 5.1.2](#), the minimum and maximum of the function f are among the values

$$f(0) = 0, \quad f(2) = 6, \quad f\left(\frac{1}{\sqrt{3}}\right) = \frac{1}{3\sqrt{3}} - \frac{1}{\sqrt{3}} = -\frac{2}{3\sqrt{3}}.$$

We conclude that the minimum is $f\left(\frac{1}{\sqrt{3}}\right) = -\frac{2}{3\sqrt{3}}$, the maximum is $f(2) = 6$, and the image set is $[-\frac{2}{3\sqrt{3}}, 6]$.

5.1. The range of functions of one variable

Note that [Theorem 5.1.2](#) just tells us that in order to find the minimum and the maximum, we should examine the points $x \in [a, b]$ where $f'(x) = 0$; it can very well be that $f'(x) = 0$ and that the function f neither assumes a minimum nor a maximum at the point x . The following example illustrates this.

Example 5.1.3

Consider the function

$$f : [-2, 2] \rightarrow \mathbb{R}, f(x) = x^3.$$

The function f is increasing, so the image set is $f([-2, 2]) = [f(-2), f(2)] = [-8, 8]$. We see that $f'(x) = 3x^2$ and that $f'(0) = 0$; however, the function f neither has a minimum nor a maximum at the point $x = 0$. A point where the derivative is zero, but where f neither has a minimum nor a maximum, is called a *saddle point*.

Recall that the *second derivative* f'' can help us to decide whether a point x_0 with $f'(x_0) = 0$ qualifies as a candidate to give a minimum value or a maximum value.

Theorem 5.1.3 Second derivative test

Consider a function $f : [a, b] \rightarrow \mathbb{R}$ which is two times differentiable on the open interval $]a, b[$. Assume that $f'(x_0) = 0$ for some $x_0 \in]a, b[$. Then the following assertions hold:

- (i) If $f''(x_0) < 0$, there exists an open interval $I \subset]a, b[$ containing x_0 such that $f(x) < f(x_0)$ for all $x \in I \setminus \{x_0\}$.
- (ii) If $f''(x_0) > 0$, there exists an open interval $I \subset]a, b[$ containing x_0 such that $f(x) > f(x_0)$ for all $x \in I \setminus \{x_0\}$.

Proof. The result follows from a more general result proved in [Theorem 5.2.4](#). ■

Note that if the assumption in [Item \(i\)](#) of [Theorem 5.1.3](#) is satisfied, the conclusion only says that the function value $f(x_0)$ is larger than the function values in a small interval containing x_0 . That is, we can not conclude yet that $f(x_0)$ is a maximal value measured over the entire interval $[a, b]$. We say that $f(x_0)$ is a (strict) *local maximum value*. Similarly, in [Item \(ii\)](#) of [Theorem 5.1.3](#), the conclusion is phrased by saying that $f(x_0)$ is a (strict) *local minimum value*.

5.2. The range of functions of several variables

In case $f'(x_0) = 0$ and also $f''(x_0) = 0$, [Theorem 5.1.3](#) does not give us any information. In this case it is necessary with a closer examination to find out whether x_0 can lead to a minimum value or a maximum value. The following example illustrates this.

Example 5.1.4

Consider the following three functions, all of them on the interval $[-1, 1]$

$$f(x) = -x^4, \quad g(x) = x^4, \quad h(x) = x^3.$$

Then

$$f'(x) = -4x^3, \quad g'(x) = 4x^3, \quad h'(x) = 3x^2$$

and

$$f''(x) = -12x^2, \quad g''(x) = 12x^2, \quad h''(x) = 6x,$$

so

$$f'(0) = f''(0) = g'(0) = g''(0) = h'(0) = h''(0) = 0.$$

By inspection, we see that $x = 0$ leads to a maximum for the function f and a minimum for the function g . For the function h , the point neither yields a minimum value nor a maximum value.

5.2 The range of functions of several variables

We will now consider the image set for functions of n variables. Our first goal is to generalize [Theorem 5.1.1](#) to functions $f : B \rightarrow \mathbb{R}$, where B is a subset of \mathbb{R}^n . In order to do so, we need to find the correct higher-dimensional variants of the conditions in [Theorem 5.1.1](#), such that we avoid the issues encountered in [Example 5.1.1](#). Let us look at the conditions in [Theorem 5.1.1](#), one by one:

- (i) [Example 5.1.1](#) (i) shows that in order for the result in [Theorem 5.1.1](#) to hold, it is essential that the function f is defined on an interval $[a, b]$ and *not* on $]a, b[$; that is, the domain must contain the end points of the interval, or, in other words, the domain must be closed. For a function of n variables, $f : B \rightarrow \mathbb{R}$, we will therefore require that the domain B is *closed*.
- (ii) [Example 5.1.1](#) (ii) shows that in order for the result in [Theorem 5.1.1](#) to hold, it is essential that the function f is defined on a bounded set.

5.2. The range of functions of several variables

For a function of n variables, $f : B \rightarrow \mathbb{R}$, we will therefore require that B is a *bounded* set; see Definition 2.2.6.

- (iii) Example 5.1.1 (iii) shows that in order for the result in Theorem 5.1.1 to hold, it is essential that the function f is defined on an interval; that is, the result might not hold for functions that are defined on a union of two disjoint intervals. The next definition translates this into a condition on sets in \mathbb{R}^n .

Definition 5.2.1 Connected set

A set $B \subset \mathbb{R}^n$ is said to be connected if, for each choice of points $\mathbf{x}_1, \mathbf{x}_2 \in B$, there exists a continuous function $\mathbf{r} : [0, 1] \rightarrow B$ such that $\mathbf{r}(0) = \mathbf{x}_1$ and $\mathbf{r}(1) = \mathbf{x}_2$.

Since the function \mathbf{r} in Definition 5.2.1 is a curve, one sometimes speak of *curve connected sets*. The key point in Definition 5.2.1 is that the curve \mathbf{r} runs *entirely inside* the set B . In the one-dimensional setting, this condition excludes the possibility that the function is defined on the union of two disjoint intervals. It turns out that we now have identified the precise conditions that are necessary in order to generalize Theorem 5.1.1:

Theorem 5.2.1

Let $B \subset \mathbb{R}^n$ be a closed and bounded, and consider a continuous function $f : B \rightarrow \mathbb{R}$. Then f has a minimal and maximal value. If B is, in addition, also a connected set, then the image set has the form

$$\text{im}(f) = f(B) = [m, M], \quad (5.2)$$

for some $m, M \in \mathbb{R}$.

Thus, exactly like Theorem 5.1.1, the result guarantees the existence of a minimum value and a maximum value. In order to find these values, we will typically apply the following result, which generalizes Theorem 5.1.2.

Theorem 5.2.2

Let $A \subset \mathbb{R}^n$ be a set, and consider a function $f : A \rightarrow \mathbb{R}$. Then, if f attains its minimum or maximum at the point $\mathbf{x}_0 \in A$, one of the following possibilities occur:

- (i) $\mathbf{x}_0 \in A \cap \partial A$ (\mathbf{x}_0 is a boundary point in the domain);
- (ii) $\mathbf{x}_0 \in A^\circ$ is a point where f is not differentiable;

5.2. The range of functions of several variables

(iii) $\mathbf{x}_0 \in A^\circ$ is a point where f is differentiable, and $\nabla f(\mathbf{x}_0) = \mathbf{0}$.

Proof. Suppose that $\mathbf{x}_0 \in A$ is a point where f attains an extremum. It is obvious that \mathbf{x}_0 is either a boundary point $\mathbf{x}_0 \in \partial A$ or an interior point $\mathbf{x}_0 \in A^\circ$. If \mathbf{x}_0 is an interior point, then one of two things can happen: f is not differentiable at \mathbf{x}_0 or f is differentiable at \mathbf{x}_0 . Hence, the only thing we have to prove is that if f is differentiable at \mathbf{x}_0 , then $\nabla f(\mathbf{x}_0) = \mathbf{0}$. We postpone the proof of this fact to Lemma 5.2.3 on page 125. ■

Thus, under the assumptions in Theorem 5.2.2, if a function f is differentiable and we want to find its minimum value and maximum value, it is enough to examine the function values $f(\mathbf{x})$ for points \mathbf{x} belonging to the boundary of the domain of the function and for points \mathbf{x} in the interior of the domain for which $\nabla f(\mathbf{x}_0) = \mathbf{0}$. This motivates the following definition.

Definition 5.2.2 Stationary point

Let $A \subset \mathbb{R}^n$ be a set. A point $\mathbf{x}_0 \in A^\circ$ is called a *stationary point* for f if f is differentiable at \mathbf{x}_0 and $\nabla f(\mathbf{x}_0) = \mathbf{0}$.

Calculation of the stationary points require us to solve a number of equations:

Example 5.2.1

For the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, f(x, y) = xy^2 + x^2 + y^3,$$

we saw in Example 3.3.4 on page 77 that

$$\nabla f(x, y) = (y^2 + 2x, 2xy + 3y^2). \quad (5.3)$$

Thus, in order to find the stationary points, we need to solve the equations

$$y^2 + 2x = 0 \quad \wedge \quad 2xy + 3y^2 = 0. \quad (5.4)$$

Note that the first equation implies that $2x = -y^2$; inserting this in the second equation shows that $-y^3 + 3y^2 = 0$, or, $y^2(-y + 3) = 0$. Thus, we have $y = 0$ or $y = 3$. For $y = 0$, the first equation in (5.4) leads to $x = 0$, and for $y = 3$ we get $x = -9/2$. Thus, the stationary points are $(0, 0)$ and $(-9/2, 3)$.

Note that the set of stationary points might not consist of isolated points:

Example 5.2.2

For the function

$$f : \mathbb{R}^3 \rightarrow \mathbb{R}, \quad f(x_1, x_2, x_3) = x_1x_2 + x_2^2x_3^3,$$

we saw in [Example 3.3.5](#) that

$$\nabla f(x_1, x_2, x_3) = (x_2, x_1 + 2x_2x_3^3, 3x_2^2x_3^2).$$

Thus, (x_1, x_2, x_3) is a stationary point if and only if the equations

$$x_2 = 0, \quad x_1 + 2x_2x_3^3 = 0, \quad \text{and} \quad 3x_2^2x_3^2 = 0,$$

are satisfied. Using that $x_2 = 0$, the second equation shows that also $x_1 = 0$; also, the third equation is automatically satisfied, i.e., the equations do not put any restriction on x_3 . Thus, the stationary points are exactly the points of the form $(x_1, x_2, x_3) = (0, 0, t)$, where $t \in \mathbb{R}$.

Let us now consider an example, where we apply [Theorem 5.2.2](#) to find the minimum value and the maximum value for a given function.

Example 5.2.3

Let

$$B = \{(x, y) \in \mathbb{R}^2 \mid x^2 + y^2 \leq 4\},$$

and consider the function

$$f : B \rightarrow \mathbb{R}, \quad f(x, y) = xy + 1.$$

In order to apply [Theorem 5.2.2](#), we first calculate the stationary points in the interior of B . We have

$$\nabla f(x, y) = (y, x).$$

Thus, the only stationary point is $(x, y) = (0, 0)$, and $f(0, 0) = 1$.

We now need to examine the function values on the boundary of B . Note that B is a ball with radius 2. Thus, the boundary δB is the circle in \mathbb{R}^2 with radius $r = 2$. The points in this set are parameterized by $\mathbf{r}(t) = (2 \cos t, 2 \sin t)$, where $t \in [0, 2\pi]$. On the boundary, the function values for the function f are therefore precisely the values of the function $g : [0, 2\pi] \rightarrow \mathbb{R}, g = f \circ \mathbf{r}$, i.e., $g(t) = f(2 \cos t, 2 \sin t)$. We see that $g(t) = f(2 \cos t, 2 \sin t) = 4 \cos t \sin t + 1 = 2 \sin(2t) + 1$ for $t \in [0, 2\pi]$.

5.2. The range of functions of several variables

The sine function take all values in the interval $[-1, 1]$, so the image set of the function g is $g([0, 2\pi]) = [-1, 3]$; comparing this information with the function value of f in the stationary point $(0, 0)$, it follows that the minimum value for the function f is -1 and the maximum value is 3 , i.e., $\text{im}(f) = [-1, 3]$.

Frequently, in particular, in application with functions of many variables, it is impossible to find global extrema, and we have to settle with a local analysis of a given function f , that is, we search for extrema in a neighborhood of a certain point \mathbf{x}_0 . In such cases we apply the following definition, which is formalizing the discussion after [Theorem 5.1.3](#)

Definition 5.2.3 Local minimum and maximum

Let A be a set in \mathbb{R}^n and $f : A \rightarrow \mathbb{R}$ a function.

- (i) If $\mathbf{x}_0 \in A$ and $f(\mathbf{x}) \geq f(\mathbf{x}_0)$ for all $\mathbf{x} \neq \mathbf{x}_0$ in some ball $B(\mathbf{x}_0, \epsilon) \cap A$, that is,

$$\exists \epsilon > 0 \forall \mathbf{x} \in A \setminus \{\mathbf{x}_0\} : \|\mathbf{x} - \mathbf{x}_0\| < \epsilon \Rightarrow f(\mathbf{x}) \geq f(\mathbf{x}_0),$$

we say that the function f has a *local minimum* at the point \mathbf{x}_0 .

- (ii) If $\mathbf{x}_0 \in A$ and $f(\mathbf{x}) \leq f(\mathbf{x}_0)$ for all $\mathbf{x} \neq \mathbf{x}_0$ in some ball $B(\mathbf{x}_0, \epsilon) \cap A$, that is,

$$\exists \epsilon > 0 \forall \mathbf{x} \in A \setminus \{\mathbf{x}_0\} : \|\mathbf{x} - \mathbf{x}_0\| < \epsilon \Rightarrow f(\mathbf{x}) \leq f(\mathbf{x}_0),$$

we say that the function f has a *local maximum* at the point \mathbf{x}_0 .

- (iii) If the above conditions hold with strict inequalities, i.e., $f(\mathbf{x}) > f(\mathbf{x}_0)$ and $f(\mathbf{x}) < f(\mathbf{x}_0)$, we say that f has a *strict local minimum* and a *strict local maximum*, respectively.

Let us illustrate the notion of *local extrema* with an example of a function of one variable.

Example 5.2.4

Consider the function

$$f : [-2, 3] \rightarrow \mathbb{R}, \quad f(x) = \begin{cases} x^2 & x \in [-2, 1[, \\ 1 & x \in [1, 2], \\ 4(x - 5/2)^2 & x \in]2, 3]. \end{cases}$$

5.2. The range of functions of several variables

By plotting the graph of the function, it is not difficult to conclude the following. The function has local minima at $x = 0$, $x = 3/2$ and for any $x \in]1, 2[$. At $x = 0$ and $x = 3/2$ the two local minima are also strict global minima with value 0, while the local minima in $]1, 2[$ are neither strict nor global. The function has local maxima at $x = -2$, $x = 3$ and for any $x \in [1, 2]$. The local maxima at $x = -2$ is also a strict global maxima with value $f(-2) = 4$.

In order to find the points where a given function attains local maximum and local minimum values, we use the same approach as we saw in Theorem 5.2.2, that is, we search for points x_0 , where $\nabla f(x_0) = \mathbf{0}$. This gives us a necessary condition for a differentiable function having a local extremum in an interior point.

Lemma 5.2.3

Let $A \subseteq \mathbb{R}^n$. If $f : A \rightarrow \mathbb{R}$ has a local extremum in an interior point $\mathbf{x}_0 \in A^\circ$, where f is differentiable, then \mathbf{x}_0 is a stationary point, i.e., $\nabla f(\mathbf{x}_0) = \mathbf{0}$.

Proof. Assume f is differentiable at \mathbf{x}_0 , but $\nabla f(\mathbf{x}_0) \neq \mathbf{0}$. Let $\mathbf{e} = \nabla f(\mathbf{x}_0) / \|\nabla f(\mathbf{x}_0)\|$, and let $\mathbf{h} = t\mathbf{e}$ for some $t \in \mathbb{R}$. Then, there is a function $\varepsilon : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying $\varepsilon(\mathbf{h}) \rightarrow 0$ for $\mathbf{h} \rightarrow \mathbf{0}$ such that

$$\begin{aligned} f(\mathbf{x}_0 + \mathbf{h}) &= f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)\mathbf{h} + \varepsilon(\mathbf{h})\|\mathbf{h}\| \\ &= f(\mathbf{x}_0) + t\|\nabla f(\mathbf{x}_0)\| + \varepsilon(t\mathbf{e})|t| \\ &= f(\mathbf{x}_0) + t(\|\nabla f(\mathbf{x}_0)\| \pm \varepsilon(t\mathbf{e})). \end{aligned}$$

For sufficiently small values of $|t|$, the term $\|\nabla f(\mathbf{x}_0)\| \pm \varepsilon(t\mathbf{e})$ is positive, hence $f(\mathbf{x}_0 + \mathbf{e}t) = f(\mathbf{x}_0) + ct$ for some positive constant c . This shows that $f(\mathbf{x}_0)$ cannot be a local extremum. ■

Note that $\nabla f(\mathbf{x}_0) = \mathbf{0}$ in Lemma 5.2.3 is only a necessary condition. Indeed, the function $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = x^3$ has a stationary point at $x_0 = 0$, but $f(x_0)$ is neither a local minimum nor local maximum. This motivates the following definition.

Definition 5.2.4 Saddle point

A stationary point of $f : A \rightarrow \mathbb{R}$, $A \subseteq \mathbb{R}^n$, is called a *saddle point* if it is neither a local minimum nor a local maximum.

The name *saddle point* comes from the fact that the graph of $f(x, y) = x^2 - y^2$ near the stationary point $(0, 0)$ looks like a horse saddle.

5.2. The range of functions of several variables

Let us finally notice that under certain circumstances the Hessian matrix can be used to argue whether a stationary point is a candidate to yield a minimum value of a maximum value:

Theorem 5.2.4 Second partial derivative test

Let $U \subset \mathbb{R}^n$ be an open set, and assume that $\mathbf{x}_0 \in U$ is a stationary point for a function $f : U \rightarrow \mathbb{R}$. Assume that f is twice continuously differentiable, that is, assume that ∇f is a C^1 vector function (recall Definition 3.8.2 on page 92). Denote the Hessian matrix at \mathbf{x}_0 by $\mathbf{H}_f(\mathbf{x}_0)$. Then the following hold:

- (i) If $\mathbf{H}_f(\mathbf{x}_0)$ is positive definite (i.e., all eigenvalues are positive), then \mathbf{x}_0 is a strict local minimum.
- (ii) If $\mathbf{H}_f(\mathbf{x}_0)$ is negative definite (i.e., all eigenvalues are negative), then \mathbf{x}_0 is a strict local maximum.
- (iii) If $\mathbf{H}_f(\mathbf{x}_0)$ has both positive and negative eigenvalues, then \mathbf{x}_0 is a saddle point, that is, $f(\mathbf{x}_0)$ is neither a minimum nor maximum value (not even in a small neighborhood of \mathbf{x}_0 .)
- (iv) If $\mathbf{H}_f(\mathbf{x}_0)$ is singular (i.e., has $\lambda = 0$ as an eigenvalue) and all non-zero eigenvalues have the same sign, then a detailed examination is necessary in order to decide whether $f(\mathbf{x}_0)$ is a minimum value, a maximum value, or a saddle point.

Proof. Since $\mathbf{x}_0 \in U$ is a stationary point, i.e., $\nabla f(\mathbf{x}_0) = \mathbf{0}$, it follows from Taylor's formula for P_{2,f,\mathbf{x}_0} in Theorem 4.6.1 with $\mathbf{h} = \mathbf{x} - \mathbf{x}_0$ that

$$\begin{aligned} f(\mathbf{x}_0 + \mathbf{h}) &= f(\mathbf{x}_0) + \langle \mathbf{h}, \nabla f(\mathbf{x}_0) \rangle + \frac{1}{2} \langle \mathbf{h}, \mathbf{H}_f(\mathbf{x}_0) \mathbf{h} \rangle + \varepsilon(\mathbf{h}) \|\mathbf{h}\|^2 \\ &= f(\mathbf{x}_0) + \frac{1}{2} \langle \mathbf{h}, \mathbf{H}_f(\mathbf{x}_0) \mathbf{h} \rangle + \varepsilon(\mathbf{h}) \|\mathbf{h}\|^2, \end{aligned}$$

where $\varepsilon(\mathbf{h}) \rightarrow 0$ as $\mathbf{h} \rightarrow \mathbf{0}$. Since $\mathbf{H}_f(\mathbf{x}_0)$ is real and symmetric, it follows from the spectral theorem that $\mathbf{H}_f(\mathbf{x}_0) = Q\Lambda Q^T$ where Q is real orthogonal and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ contains the eigenvalues of the Hessian matrix. Let $\tilde{\mathbf{h}} = Q^T \mathbf{h}$, and let \tilde{h}_k be the k th entry of $\tilde{\mathbf{h}}$. Then

$$\langle \mathbf{h}, \mathbf{H}_f(\mathbf{x}_0) \mathbf{h} \rangle = \langle \tilde{\mathbf{h}}, \Lambda \tilde{\mathbf{h}} \rangle = \sum_{k=1}^n \lambda_k |\tilde{h}_k|^2$$

5.2. The range of functions of several variables

To prove [Item \(i\)](#), we assume that $\mathbf{H}_f(\mathbf{x}_0)$ is positive definite, i.e., $\lambda_k > 0$ for all k . Let $\lambda_{\min} > 0$ denote the smallest eigenvalue. Then

$$\langle \mathbf{h}, \mathbf{H}_f(\mathbf{x}_0)\mathbf{h} \rangle = \sum_{k=1}^n \lambda_k |\tilde{h}_k|^2 \geq \lambda_{\min} \sum_{k=1}^n |\tilde{h}_k|^2 = \lambda_{\min} \|\tilde{\mathbf{h}}\|^2 = \lambda_{\min} \|\mathbf{h}\|^2$$

Hence, for any $\mathbf{h} \in \mathbb{R}^n$ with $\|\mathbf{h}\| \neq 0$ sufficiently small, we have

$$\begin{aligned} f(\mathbf{x}_0 + \mathbf{h}) &\geq f(\mathbf{x}_0) + \frac{1}{2}\lambda_{\min}\|\mathbf{h}\|^2 + \varepsilon(\mathbf{h})\|\mathbf{h}\|^2 = \\ &f(\mathbf{x}_0) + \|\mathbf{h}\|^2(\lambda_{\min}/2 + \varepsilon(\mathbf{h})) > f(\mathbf{x}_0), \end{aligned}$$

which shows that $f(\mathbf{x}_0)$ is a strict local minimum.

The proofs of the other assertions are similar, and we leave the details to the reader. ■

Note that the complication in [Item \(iv\)](#) in case of a non-invertible Hessian matrix corresponds precisely to the information we got out of [Example 5.1.4](#) in the one-dimensional case. The complication is due to the fact the a second order Taylor approximation is not a sufficiently high order polynomial to examine the stationary point.

CHAPTER 6

Integration

Integration is one of the fundamental concepts in mathematics. Several definitions of integration exist in the mathematical literature, each tailored for specific purposes. Also, numerous problems in physics and engineering can be formulated and solved in terms of integrals. In this chapter we aim at a presentation of the classical Riemann integral, often just called *the integral*, as the other types of integrals only occur in more advanced literature. We begin in [Section 6.1](#) by considering integrals of scalar functions over bounded intervals $[a, b]$ in \mathbb{R} . In [Section 6.3](#) we first mimic the procedure to introduce integration over rectangles in \mathbb{R}^2 , and then show how to integrate scalar functions of two variables over more complicated domains in \mathbb{R}^2 . In [Section 6.4](#) we introduce a technique that can often be used to transform a “complicated integral” into a “simpler integral”; a special and particularly useful version of this technique is discussed in [Section 6.5](#). A natural generalization of integration in \mathbb{R}^2 to \mathbb{R}^n is presented in [Section 6.6](#). Finally, in [Section 6.7](#) we consider integration of vector functions of n variables.

6.1 The Riemann integral of functions of one variable

As prelude to the later and more advanced sections in this chapter, we will start by considering integration of functions $f : [a, b] \rightarrow \mathbb{R}$, where $[a, b]$ is a bounded and closed interval in \mathbb{R} . The integral can be defined in two equivalent ways, namely, either as a limit for certain sums (to be specified below) or as the anti-derivative of the function f . We will describe both ways. We expect the reader to have some basic knowledge about the topic, but maybe at a less technical level.

To get started, let us assume that the function $f : [a, b] \rightarrow \mathbb{R}$ is continuous

6.1. The Riemann integral of functions of one variable

and positive, that is, $f(x) \geq 0$ for all $x \in [a, b]$. Let us imagine that we want to find an approximate value for the area of the subset in \mathbb{R}^2 that is bounded by the lines $x = a$, $x = b$, the x -axis, and the graph of the function f . In order to do this, fix some $J \in \mathbb{N}$ and split the interval $[a, b]$ into J subintervals $Q_j, j = 1, \dots, J$, given by

$$Q_1 = [x_0, x_1], Q_2 = [x_1, x_2], \dots, Q_J = [x_{J-1}, x_J], \quad (6.1)$$

where $x_0, x_1, \dots, x_J \in [a, b]$ satisfy:

$$a = x_0 < x_1 < x_2 < \dots < x_{J-1} < x_J = b. \quad (6.2)$$

The length of each interval $Q_j = [x_{j-1}, x_j]$ is denoted by $\Delta x_j = x_j - x_{j-1}$. One such choice is to divide $[a, b]$ in J intervals of equal length $\Delta x = \Delta x_j = \frac{b-a}{J}$:

$$Q_1 = \left[a, a + \frac{b-a}{J} \right], Q_2 = \left[a + \frac{b-a}{J}, a + 2\frac{b-a}{J} \right], \dots, Q_J = \left[a + (J-1)\frac{b-a}{J}, b \right].$$

For each interval Q_j , we pick an arbitrary point $\xi_j \in Q_j$. Then the associated *Riemann sum*, to be denoted by S_J , is defined as

$$S_J = \sum_{j=1}^J f(\xi_j) \Delta x_j. \quad (6.3)$$

Intuitively, if J is “large”, the value S_J is a “good approximation” of the desired area; furthermore, typically we obtain “better approximations” by increasing J , i.e., by splitting the interval $[a, b]$ into finer intervals. See Figure 6.1 for an illustration of this.

The Riemann sums S_J indeed have a limit, denoted $I \in \mathbb{R}$, as $J \rightarrow \infty$ regardless of how we pick the subintervals Q_j and $\xi_j \in Q_j$; the limit I can be considered as the precise definition of the area of the considered domain. If the function f is not necessarily positive (but still continuous), one can prove that the Riemann sums S_J still converge to a real scalar I as $J \rightarrow \infty$, and the limit I is the “signed” area under the graph, see [Theorem 6.1.1](#) on page 131.

However, we often need to integrate functions that are *not* continuous on their entire domain. This leads to the following definition.

Definition 6.1.1 Riemann integral in \mathbb{R}

A function $f : [a, b] \rightarrow \mathbb{R}$ is said to be *Riemann integrable* if there exists a number $I \in \mathbb{R}$ with the following property: For any $\epsilon > 0$ there exists $\delta > 0$ such that for any subdivision of $[a, b]$ into J intervals $Q_j, j = 1, \dots, J$, with

6.1. The Riemann integral of functions of one variable

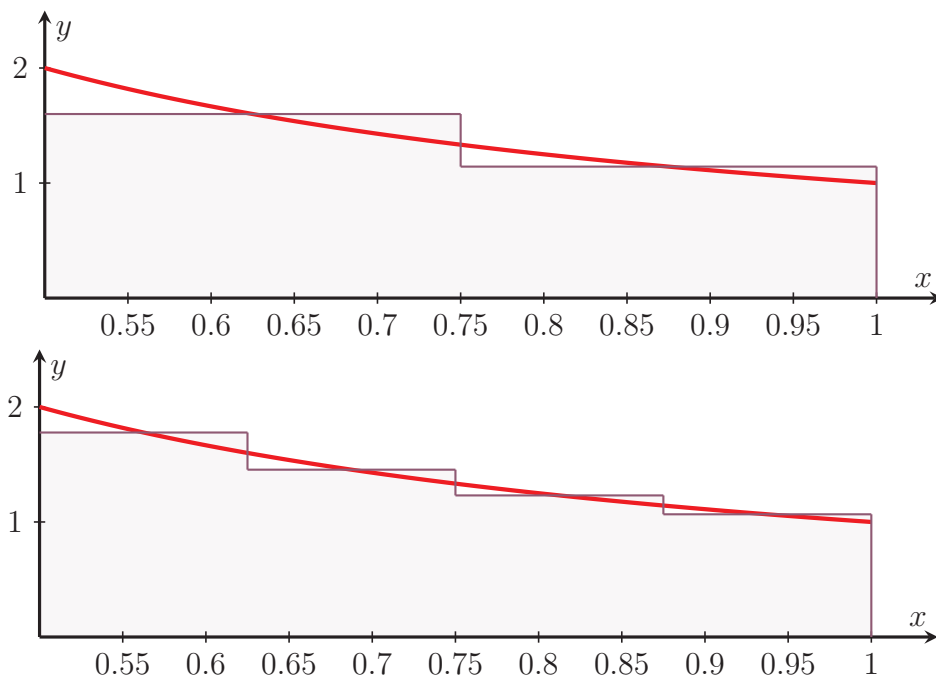


Figure 6.1: Riemann sums of the function $f(x) = x^{-1}$ on $x \in [1/2, 1]$. On the first figure, the interval $[1/2, 1]$ is split into two intervals (corresponding to $J = 2$.) The area “under the staircase” corresponds to the Riemann sum S_2 for the choice $\xi_1 = 5/8, \xi_2 = 7/8$. On the second figure, the interval $[1/2, 1]$ is split into four intervals (corresponding to $J = 4$.) The area “under the staircase” corresponds to the Riemann sum S_4 for the choice $\xi_1 = 9/16, \xi_2 = 11/16, \xi_3 = 13/16, \xi_4 = 15/16$. We see that already S_4 is a much better approximation to the area under the graph of the function f than S_2 .

$\max_j \Delta x_j < \delta$ and for any $\xi_j \in Q_j$, the Riemann sum $S_J = \sum_{j=1}^J f(\xi_j) \Delta x_j$ satisfies:

$$|I - S_J| < \epsilon.$$

In case f is Riemann integrable over the interval $[a, b]$, the real scalar I is called *the integral of f over $[a, b]$* , and it is written as

$$\int_a^b f(x) dx = I. \quad (6.4)$$

Often we will skip the name “Riemann” and simply speak about the integral of the function f over the interval $[a, b]$. It is important to remember

6.1. The Riemann integral of functions of one variable

that the symbol $\int_a^b f(x) dx$ has only one “meaning”: it is just the way we denote the limit for the Riemann sums S_J as $J \rightarrow \infty$, indeed,

$$\int_a^b f(x) dx = \lim_{J \rightarrow \infty} S_J = \lim_{J \rightarrow \infty} \sum_{j=1}^J f(\xi_j) \Delta x_j.$$

One can show that the number I from [Definition 6.1.1](#) is uniquely determined. It is also important to note that although different choices of the points x_j and ξ_j , $j = 1, \dots, J$, lead to different values for the Riemann sums S_J , the limit I is *independent* of the choice of these points.

As we already remarked, the Riemann sums always converge for continuous functions. Let us state this useful result as a theorem.

Theorem 6.1.1

Let $f : [a, b] \rightarrow \mathbb{R}$ be a function.

- (i) If f is Riemann integrable, then f is bounded, i.e., the image $\text{im}(f)$ is a bounded set in \mathbb{R} .
- (ii) If f is continuous, then f is Riemann integrable.

Let us state, again without proof, a number of useful properties of the Riemann integral.

Theorem 6.1.2 Rules of Riemann integration

Let $f, g : [a, b] \rightarrow \mathbb{R}$ be functions.

- (i) (*Linearity*) Let $c, d \in \mathbb{R}$. If f and g are Riemann integrable functions, then so is $cf + dg$ and

$$\int_a^b (cf(x) + dg(x)) dx = c \int_a^b f(x) dx + d \int_a^b g(x) dx. \quad (6.5)$$

- (ii) (*Monotonicity*) If f and g are Riemann integrable functions and $f(x) \leq g(x)$ for all $x \in [a, b]$, then

$$\int_a^b f(x) dx \leq \int_a^b g(x) dx. \quad (6.6)$$

- (iii) (*Triangle inequality*) If f and $|f|$ are Riemann integrable functions, then

$$\left| \int_a^b f(x) dx \right| \leq \int_a^b |f(x)| dx. \quad (6.7)$$

(iv) (*Insertion rule*) Let $x_0 \in]a, b[$. If f is Riemann integrable over $[a, x_0]$ and over $[x_0, b]$, then f is Riemann integrable over $[a, b]$ and

$$\int_a^b f(x) \, dx = \int_a^{x_0} f(x) \, dx + \int_{x_0}^b f(x) \, dx \quad (6.8)$$

It turns out that it is useful to define $\int_a^b f(x) \, dx$ not only for $b > a$ but also for $a < b$. We do this by setting, for $b \geq a$:

$$\int_b^a f(x) \, dx = \begin{cases} -\int_a^b f(x) \, dx & b > a \\ 0 & b = a \end{cases} \quad (6.9)$$

With this notation, we can phrase the insertion rule (6.8), given that each of the three integrals exist, as

$$\int_a^b f(x) \, dx = \int_a^c f(x) \, dx + \int_c^b f(x) \, dx$$

where $f : I \rightarrow \mathbb{R}$, $I \subset \mathbb{R}$ is any interval, and $a, b, c \in I$.

6.2 Anti-derivatives of functions of one variable

In general, it is cumbersome to calculate the integral $\int_a^b f(x) \, dx$ as in Equation (6.4) on page 130, even for simple functions like $\cos(x)$ and x^n (and one often has to approximate it on a computer using numerical integration). This is the point, where the second interpretation of the integral comes in handy. It is based on the following concept.

Definition 6.2.1 Anti-derivative

Let $I \subset \mathbb{R}$ be an arbitrary interval, and let $f : I \rightarrow \mathbb{R}$ be a function. A function $F : I \rightarrow \mathbb{R}$ is an *anti-derivative* of f on I if F is differentiable on I with

$$F'(x) = f(x), \quad \text{for all } x \in I. \quad (6.10)$$

Anti-derivatives are often denoted by a so-called *indefinite integral*:

$$F(x) = \int f(x) \, dx, \quad (6.11)$$

while the Riemann integral is called a *definite integral*.

Anti-derivatives are known for a large class of standard functions in the mathematical analysis:

Example 6.2.1

Let us list a number of functions from the classical mathematical analysis and their corresponding anti-derivatives. Note that in each case we only mention one anti-derivative.

$$\int x^n dx = \frac{1}{n+1} x^{n+1}, \quad \text{for } n \in \mathbb{R} \setminus \{-1\}, \quad (6.12)$$

$$\int \frac{1}{x} dx = \ln(x), \quad (6.13)$$

$$\int \cos(x) dx = \sin(x), \quad (6.14)$$

$$\int \sin(x) dx = -\cos(x), \quad (6.15)$$

$$\int e^x dx = e^x. \quad (6.16)$$

Be careful not to confuse $\int f(x) dx$ with the Riemann integral. The symbol $\int f(x) dx$ means a *function* whose derivative is f , while $\int_a^b f(x) dx$ is a real *scalar*. However, the anti-derivative is, in fact, closely related to the integral defined in [Definition 6.1.1](#) on page 129. More precisely, we have the following fundamental result.

Theorem 6.2.1 Fundamental Theorem of Calculus

Suppose $f : I \rightarrow \mathbb{R}$ is a continuous function. Then f has an anti-derivative on I given by

$$F(x) = \int_{x_0}^x f(y) dy \quad \text{for } x \in I, \quad (6.17)$$

where $x_0 \in I$ is a fixed, but arbitrary point in I .

It is easy to show using the mean value theorem that if F is an anti-derivative of the function f , then all anti-derivatives are given by the formula $F(x) + C$, where C is an arbitrary constant.

Hence, if a function has an anti-derivative, it actually has infinitely many anti-derivatives. On the other hand, there are many Riemann integrable functions that do not have an anti-derivative. It follows from [Theorem 6.2.1](#) that such functions cannot be continuous. A simple example of Riemann integrable function with no anti-derivative is considered in the next exercise.

Exercise 6.2.2

Consider the function $f : \mathbb{R} \rightarrow \mathbb{R}$ given by

6.2. Anti-derivatives of functions of one variable

$$f(x) = \begin{cases} 1 & x \in [c, d], \\ 0 & x \notin [c, d], \end{cases}$$

where $d > c$. Show first that f is Riemann integrable. Show then that f cannot have an anti-derivative.

Corollary 6.2.3

Let F denote an anti-derivative of a continuous function $f : I \rightarrow \mathbb{R}$. Then, for any $a, b \in I$,

$$\int_a^b f(x) \, dx = F(b) - F(a). \quad (6.18)$$

Often, the result in (6.18) is written in a slightly different form, namely as

$$\int_a^b f(x) \, dx = [F(x)]_a^b; \quad (6.19)$$

this form shows directly that we calculate the integral by *first* calculating an anti-derivative of the function f , and *then* insert the endpoints a and b of the considered interval.

Typically, application of Corollary 6.2.3 simplifies the calculation of the Riemann integral. Let us demonstrate this with a few examples.

Example 6.2.2

Using Corollary 6.2.3 and (6.14), we directly see that

$$\int_0^\pi \cos(x) \, dx = [\sin(x)]_0^\pi = \sin(\pi) - \sin(0) = 0.$$

Similarly, using (6.16),

$$\int_1^2 e^x \, dx = [e^x]_1^2 = e^2 - e.$$

Many complicated integrals can be simplified (basically to an application of some of the results in Example 6.2.1) by using certain rules, stated next. For the precise statement of the results, we need to assume that functions, whose derivative appears inside an integral, are continuously differentiable as defined in Definition 3.1.2 on page 67.

Theorem 6.2.4 Rules of indefinite integration

6.2. Anti-derivatives of functions of one variable

- (i) (*Partial integration*) If f is a continuous function, F an anti-derivative of f , and g is continuously differentiable, then

$$\int f(x)g(x) dx = F(x)g(x) - \int F(x)g'(x) dx. \quad (6.20)$$

- (ii) (*Integration by substitution*) If f is a continuous function, F an anti-derivative of f , and g a continuously differentiable function such that its range is contained in the domain of the function f , then

$$\int f(g(x))g'(x) dx = F(g(x)). \quad (6.21)$$

Obviously, [Theorem 6.2.4](#) can also be formulated for Riemann integration. In particular, [Equation \(6.20\)](#) becomes

$$\int_a^b f(x)g(x) dx = F(b)g(b) - F(a)g(a) - \int_a^b F(x)g'(x) dx.$$

Note the particular structure the right hand side: it does not give the “final result” for $\int_a^b f(x)g(x) dx$ since we still have to calculate $\int_a^b F(x)g'(x) dx$ subsequently. In other words, partial integration is only a useful technique if the particular structure of the functions f and g implies that it is easier to calculate the integral $\int F(x)g'(x) dx$ than the integral $\int f(x)g(x) dx$.

The difficult in integration is typically to identify an integration rule that can simplify the given integral. Worse than that, there exist integrals that can not be expressed in terms of the standard functions from mathematical analysis. Let us illustrate [Theorem 6.2.4](#) with a number of examples.

Example 6.2.3

We would like to calculate the integral

$$\int_1^2 2x \ln(x) dx. \quad (6.22)$$

In order to do so, we first aim at determining

$$\int 2x \ln(x) dx, x > 0. \quad (6.23)$$

We immediately see that the integrand in [\(6.23\)](#) can be considered as the product of two functions, namely $2x$ and $\ln(x)$. This makes it natural to try

6.2. Anti-derivatives of functions of one variable

to apply partial integration as in (6.20). First, let $f(x) = 2x$, $g(x) = \ln(x)$. Then $F(x) = x^2$ is an anti-derivative of f and $g'(x) = x^{-1}$. Thus, using Theorem 6.2.4(i),

$$\int 2x \ln(x) \, dx = x^2 \ln(x) - \int x^2 x^{-1} \, dx = x^2 \ln(x) - \int x \, dx = x^2 \ln(x) - \frac{1}{2}x^2.$$

In this particular case we could actually also take the “opposite choice” of the functions f and g ! Indeed, letting $f(x) = \ln(x)$, $g(x) = 2x$, the function $F(x) = x \ln(x) - x$ is an anti-derivative of f , and $g'(x) = 2$. Thus,

$$\int 2x \ln(x) \, dx = (x \ln(x) - x)2x - \int (x \ln(x) - x)2 \, dx \quad (6.24)$$

$$\begin{aligned} &= 2x^2 \ln(x) - 2x^2 - \int 2x \ln(x) \, dx + \int 2x \, dx \\ &= 2x^2 \ln(x) - x^2 - \int 2x \ln(x) \, dx. \end{aligned} \quad (6.25)$$

Note that at first sight, this calculation does not seem to be useful because it returns exactly the same integral in (6.25) as the one we start with in (6.24). However, exactly this observation allows us to *collect* the terms and write that

$$2 \int 2x \ln(x) \, dx = 2x^2 \ln(x) - x^2,$$

again showing that

$$\int 2x \ln(x) \, dx = x^2 \ln(x) - \frac{1}{2}x^2.$$

Finally, after having calculated the integral in (6.23), we can apply (6.19) to obtain that

$$\int_1^2 2x \ln(x) \, dx = \left[x^2 \ln(x) - \frac{1}{2}x^2 \right]_1^2 = (4 \ln(2) - 2) - \left(\ln(1) - \frac{1}{2} \right) = 4 \ln(2) - \frac{15}{2}.$$

Example 6.2.4

Fix some $a, b \in \mathbb{R}$ with $a \neq 0$ and consider the integral

$$\int \cos(ax + b) \, dx. \quad (6.26)$$

The function $\cos(ax + b)$ can naturally be considered as a composition of the functions $\cos(x)$ and $ax + b$, so it is natural to take $f(x) = \cos(x)$ and $g(x) = ax + b$, and aim at an application of the rule in Theorem 6.2.4(ii).

6.3. The Riemann integral of functions of two variables

With the mentioned choice of the functions $f(x)$ and $g(x)$ we indeed obtain that $f(g(x)) = \cos(ax + b)$. However, the integral does not get exactly the form in (6.21), as the term $g'(x)$ is missing. But in the current case we can get away with this problem by observing that $g'(x) = a$, i.e., the missing term is just a constant! Therefore we can write the integral (6.26) as

$$\begin{aligned}\int \cos(ax + b) dx &= \frac{1}{a} \int \cos(ax + b)a dx = \frac{1}{a} \int f(g(x))g'(x) dx \\ &= \frac{1}{a} F(g(x)) \\ &= \frac{1}{a} \sin(ax + b).\end{aligned}$$

Note that this calculation used the linearity rule in Theorem 6.2.4(i) in the first step!

Example 6.2.5

Using the results in Example 6.2.3 and Example 6.2.4,

$$\begin{aligned}\int \left(10x \ln(x) + 3 \cos\left(2x + \frac{\pi}{3}\right)\right) dx &= 5 \int 2x \ln(x) dx + 3 \int \cos\left(2x + \frac{\pi}{3}\right) dx \\ &= 5 \left(x^2 \ln(x) - \frac{1}{2}x^2\right) + \frac{3}{2} \sin\left(2x + \frac{\pi}{3}\right) \\ &= 5x^2 \ln(x) - \frac{5}{2}x^2 + \frac{3}{2} \sin\left(2x + \frac{\pi}{3}\right).\end{aligned}$$

6.3 The Riemann integral of functions of two variables

Integration in \mathbb{R}^2 (and \mathbb{R}^n) is in many regards completely similar to integration in \mathbb{R} . However, the notion of an anti-derivative does not generalize well to functions of several variables, and we will not return to this concept until Chapter 7. Moreover, even when defining the Riemann integral of a function of two variables, we will encounter technical issues arising from the fact that general sets in \mathbb{R}^2 can be significantly more complicated than just intervals in \mathbb{R} . In order to avoid these difficulties we will first consider integration over rectangles in \mathbb{R}^2 .

A *rectangle* in \mathbb{R}^2 is a set of the form

$$Q = \{(x, y) \in \mathbb{R}^2 \mid a_1 \leq x \leq b_1 \wedge a_2 \leq y \leq b_2\}. \quad (6.27)$$

6.3. The Riemann integral of functions of two variables

The rectangle will also be written as $Q = [a_1, b_1] \times [a_2, b_2]$. Consider now a continuous function $f : Q \rightarrow \mathbb{R}$, and assume first that all function values are positive, i.e., $f(x) \geq 0$ for all $x \in Q$. Our goal is to define the integral of f over Q ; parallel to our definition of the integral in \mathbb{R} , it should be interpreted as the *volume* of the region between the graph of the function $f : Q \rightarrow \mathbb{R}$ and the xy -plane.

We will follow the procedure for the Riemann integral in \mathbb{R} and first split the rectangle Q into smaller rectangles. Technically, this is done by fixing an $J \in \mathbb{N}$ and splitting the intervals $[a_1, b_1]$ and $[a_2, b_2]$ into J intervals of equal size. Denote the splitting points by x_1, x_2, \dots, x_{J-1} resp. y_1, y_2, \dots, y_{J-1} , and let, for notational convenience, $x_0 := a$, $x_J := b$, $y_0 := a_2$, and $y_J := b$; then

$$a_1 = x_0 < x_1 < x_2 < \dots < x_{J-1} < x_J = b_1, \quad (6.28)$$

and

$$a_2 = y_0 < y_1 < y_2 < \dots < y_{J-1} < y_J = b_2. \quad (6.29)$$

The partition of the rectangle Q into smaller rectangles is now defined via

$$Q_{i,j} = [x_{i-1}, x_i] \times [y_{j-1}, y_j], \quad i, j = 1, 2, \dots, J. \quad (6.30)$$

Note that the length of each subinterval $[x_{i-1}, x_i]$ is $\Delta x := \frac{b_1 - a_1}{J}$, and the length of each subinterval $[y_{j-1}, y_j]$ is $\Delta y := \frac{b_2 - a_2}{J}$. Now, for each $i, j = 1, 2, \dots, J$, pick an arbitrary point $\xi_{i,j} \in Q_{i,j}$ and define the associated *Riemann sum* to be denoted by S_J , by

$$S_J = \sum_{i=1}^J \sum_{j=1}^J f(\xi_{i,j}) \text{area}(Q_{i,j}) \quad (6.31)$$

where $\text{area}(Q_{i,j}) := \Delta x \Delta y$ is the area of each rectangle $Q_{i,j}$.

Intuitively, if J is “large”, the value S_J is a “good approximation” to the desired volume. Furthermore, we typically get better approximations by increasing the number of J . One can prove that the Riemann sums S_J indeed have a limit as $J \rightarrow \infty$; this limit can be considered as the definition of the *exact* value of the desired volume.

If the function f is not necessarily positive (but still continuous), one can prove that the Riemann sums S_J still converge. Hence, we can define the Riemann integral of the function f over the rectangle Q as

$$\int_Q f(x, y) \, d(x, y) := \lim_{J \rightarrow \infty} S_J = \lim_{J \rightarrow \infty} \sum_{i=1}^J \sum_{j=1}^J f(\xi_{i,j}) \text{area}(Q_{i,j}) \quad (6.32)$$

Although we now have a precise definition of integration over Q at hand, it is clear that it is tedious to apply the described procedure in practice.

6.3. The Riemann integral of functions of two variables

Fortunately, again the situation is similar to what we encountered in \mathbb{R} : we can calculate the integral using the concept of anti-derivatives. In order to explain this, consider again the rectangle Q in (6.27) and a continuous function $f : Q \rightarrow \mathbb{R}$. Fix an arbitrary value for $y \in [a_2, b_2]$; then the function $x \mapsto f(x, y)$ is a continuous function, and we can calculate the integral

$$\int_{a_1}^{b_1} f(x, y) \, dx, \quad (6.33)$$

e.g., as in Corollary 6.2.3. The expression in (6.33) is now a continuous function of the variable y ; we can therefore apply Corollary 6.2.3 again, and form the integral of this function, i.e.,

$$\int_{a_2}^{b_2} \left(\int_{a_1}^{b_1} f(x, y) \, dx \right) dy. \quad (6.34)$$

It can be shown that the outcome of this *double integration* is identical with the Riemann integral introduced in (6.32). In fact, we could also have integrated first with respect to y and then with respect to x , and we would have arrived at the same value of the integration. Let us formulate this as a theorem:

Theorem 6.3.1

Consider a rectangle Q in \mathbb{R}^2 of the form

$$Q = \left\{ (x, y) \in \mathbb{R}^2 \mid a_1 \leq x \leq b_1 \wedge a_2 \leq y \leq b_2 \right\}, \quad (6.35)$$

and let $f : Q \rightarrow \mathbb{R}$ denote a continuous function. Then

$$\int_Q f(x, y) \, d(x, y) = \int_{a_2}^{b_2} \left(\int_{a_1}^{b_1} f(x, y) \, dx \right) dy \quad (6.36)$$

$$= \int_{a_1}^{b_1} \left(\int_{a_2}^{b_2} f(x, y) \, dy \right) dx. \quad (6.37)$$

Theorem 6.3.1 says that integration over a rectangle in \mathbb{R}^2 can be reduced to calculation of *two* “standard integrals” over bounded intervals in \mathbb{R} . Let us illustrate formula (6.36) by an example.

Example 6.3.1

Let

$$Q = [2, 3] \times [0, 1] = \left\{ (x, y) \in \mathbb{R}^2 \mid x \in [2, 3] \wedge y \in [0, 1] \right\},$$

6.3. The Riemann integral of functions of two variables

and consider the function $f(x, y) = x^2y$. In order to calculate $\int_Q f(x, y) d(x, y)$, we apply (6.36), which yields that

$$\int_Q f(x, y) d(x, y) = \int_0^1 \left(\int_2^3 x^2y dx \right) dy. \quad (6.38)$$

According to the definition of the double integral, we first consider the variable y as a constant and calculate the *inner integral* with respect to x :

$$\int_2^3 x^2y dx = \left[\frac{1}{3}x^3y \right]_{x=2}^{x=3} = \frac{19}{3}y.$$

We now insert this result in (6.38) and perform the *outer integration* with respect to y :

$$\int_Q f(x, y) d(x, y) = \int_0^1 \frac{19}{3}y dy = \left[\frac{19}{3} \frac{1}{2}y^2 \right]_0^1 = \frac{19}{6}.$$

The two formulas (6.36) and (6.37) in Theorem 6.3.1 show that we can *switch the order of integration*. It is a very useful technique that can be used in case calculation of an integral simplifies by changing the order of integration.

We are now ready to address the technically much more complicated question of integration over a general subset B in \mathbb{R}^2 . We will assume that

(I) B is bounded, i.e.,

$$B \subset [a_1, b_1] \times [a_2, b_2] \quad (6.39)$$

for some $a_1, a_2, b_1, b_2 \in \mathbb{R}$.

(II) The boundary ∂B of B is formed by a finite number of continuously differentiable curves.

Recall that a curve is just (the image of) a vector function of one variable. We have already discussed boundary curves in Example 2.2.4 and Example 2.2.5. The technical condition Item (II) guarantees that the boundary ∂B of B has an area that is zero (more precisely, ∂B can be covered with a set of arbitrarily small area).

The subset B might have a geometrically complicated shape, so in general we can not split it into small rectangles as we did for the set Q in (6.27). However, we can perform such a splitting on the “bigger” rectangle $[a_1, b_1] \times [a_2, b_2]$ in (6.39) as we did in the beginning of this section. That is, we again fix an integer $J \in \mathbb{N}$, split the intervals $[a_1, b_1]$ and $[a_2, b_2]$ into

6.3. The Riemann integral of functions of two variables

J intervals (not necessarily of equal length), and form the corresponding subrectangles $Q_{i,j}$, $i, j = 1, 2, \dots, J$; see Equations (6.28) to (6.30).

Among the rectangles $Q_{i,j}$, some of them will be *outside* the subset B ; some of them will be *inside* B , and some of them will have parts inside B and parts outside B . We will exclusively consider the rectangles $Q_{i,j}$ that are completely contained in B , i.e., $Q_{i,j} \subset B$. For these sets $Q_{i,j}$, we again pick arbitrary points $\xi_{i,j} \in Q_{i,j}$, and form the corresponding *Riemann sum*, to be denoted by S_J , defined by

$$S_J = \sum_{\{(i,j)|Q_{i,j} \subset B\}} f(\xi_{i,j}) \text{area}(Q_{i,j}) \quad (6.40)$$

where $\text{area}(Q_{i,j}) := \Delta x_i \Delta y_j$, and Δx_i and Δy_j are the side lengths of the rectangle $Q_{i,j}$ and $1 \leq i, j \leq J$.

In words: the Riemann sum is formed precisely like we did in (6.31), except that we only consider the rectangles $Q_{i,j}$ that are completely contained in the given subset B . With this setup, one can again show that if the function f is continuous, the Riemann sum S_J tends to a limit as $J \rightarrow \infty$. The limit is denoted by

$$\int_B f(x, y) \, d(x, y), \quad \int_B f(x, y) \, d\mathbf{x}, \quad (6.41)$$

or by

$$\int_B f(x, y) \, dX, \quad (6.42)$$

and it is called the *Riemann integral of f over the subset B* .

For functions that are not continuous in all points in its domain, we need a slightly more technical limiting procedure, just as was the case in Definition 6.1.1 on page 129.

Definition 6.3.1 Riemann integral in \mathbb{R}^2

Consider a subset B in \mathbb{R}^2 , which satisfies the conditions **Items (I) and (II)** on the previous page. A function $f : B \rightarrow \mathbb{R}$ is said to be *Riemann integrable* if for any $\epsilon > 0$ there exists $\delta > 0$ such that

$$\left| \int_B f(x, y) \, d(x, y) - S_J \right| < \epsilon$$

for *any* Riemann sum S_J with $\max_i \Delta x_i < \delta$, $\max_j \Delta y_j < \delta$ and any $\xi_{i,j} \in Q_{i,j}$. The scalar $\int_B f(x, y) \, d(x, y)$ is called the integral of f over B .

6.3. The Riemann integral of functions of two variables

As for functions of one variable (see [Theorem 6.1.1](#)), continuous functions on $B \subset \mathbb{R}^2$, satisfying the conditions [Items \(I\)](#) and [\(II\)](#) on page 140, are guaranteed to be Riemann integrable. In particular, since the function $f(x) = 1$ for all $x \in B$ is continuous, it is Riemann integrable. This leads to an exact definition of the *area* of a subset B .

Definition 6.3.2 Area

Consider a subset B in \mathbb{R}^2 , which satisfies the conditions [Items \(I\)](#) and [\(II\)](#) on page 140. Then the *area* of B is defined as

$$\text{area}(B) = \int_B 1 \, dX. \quad (6.43)$$

In general it is complicated to calculate the integrals in [Definition 6.3.1](#) and [\(6.43\)](#) via the described procedure, except if B has a simple geometric structure. We will now consider a particular case, where it can be done in a similar way as to what we did in [Theorem 6.3.1](#).

Theorem 6.3.2

- (i) Consider two continuously differentiable functions

$$\alpha_1 : [a, b] \rightarrow \mathbb{R}, \quad \alpha_2 : [a, b] \rightarrow \mathbb{R},$$

and assume that $\alpha_1(x) \leq \alpha_2(x)$ for all $x \in [a, b]$. Consider the subset B defined as

$$B = \{(x, y) \in \mathbb{R}^2 \mid a \leq x \leq b \wedge \alpha_1(x) \leq y \leq \alpha_2(x)\}. \quad (6.44)$$

Then, for any continuous function $f : B \rightarrow \mathbb{R}$,

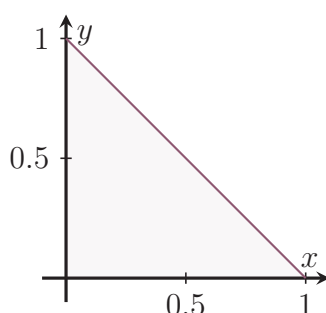
$$\int_B f(x, y) \, d(x, y) = \int_a^b \left(\int_{\alpha_1(x)}^{\alpha_2(x)} f(x, y) \, dy \right) dx. \quad (6.45)$$

- (ii) Similarly, if the subset B can be described as

$$B = \{(x, y) \in \mathbb{R}^2 \mid c \leq y \leq d \wedge \beta_1(y) \leq x \leq \beta_2(y)\} \quad (6.46)$$

for some continuously differentiable functions $\beta_1, \beta_2 : [c, d] \rightarrow \mathbb{R}$ with $\beta_1(y) \leq \beta_2(y)$ for all $y \in [c, d]$, then for any continuous function $f : B \rightarrow \mathbb{R}$,

$$\int_B f(x, y) \, d(x, y) = \int_c^d \left(\int_{\beta_1(y)}^{\beta_2(y)} f(x, y) \, dx \right) dy. \quad (6.47)$$

Figure 6.2: The domain B in (6.48).

In words, [Theorem 6.3.2](#) says that integration over a subset in \mathbb{R}^2 of the special types (6.44) and (6.46) can be reduced to calculation of *two* “standard integrals” in \mathbb{R} . In particular, if both descriptions (6.44) and (6.46) are available, we are free to choose the order of integration in the way that is most convenient; this observation generalizes what we saw in [Theorem 6.3.1](#). Let us illustrate this by an example.

Example 6.3.2

Let

$$B = \{(x, y) \in \mathbb{R}^2 \mid 0 \leq x \leq 1 \wedge 0 \leq y \leq 1 - x\}, \quad (6.48)$$

and consider the function $f(x, y) = 2xy$. Then, according to [Theorem 6.3.2\(i\)](#),

$$\begin{aligned} \int_B f(x, y) \, d(x, y) &= \int_0^1 \left(\int_0^{1-x} 2xy \, dy \right) dx \\ &= \int_0^1 [xy^2]_{y=0}^{y=1-x} dx \\ &= \int_0^1 x(1-x)^2 dx \\ &= \int_0^1 (x^3 - 2x^2 + x) dx = \frac{1}{12}. \end{aligned}$$

The set B in (6.48) is illustrated on [Figure 6.2](#). It also can be described as

$$B = \{(x, y) \in \mathbb{R}^2 \mid 0 \leq y \leq 1 \wedge 0 \leq x \leq 1 - y\}.$$

Thus, we can also apply [Theorem 6.3.2\(ii\)](#), which yield that

$$\int_B f(x, y) \, d(x, y) = \int_0^1 \left(\int_0^{1-y} 2xy \, dx \right) dy.$$

From here we can perform the same operations as above, leading to (we skip some of the details)

$$\begin{aligned}\int_B f(x, y) \, d(x, y) &= \int_0^1 [x^2 y]_{x=0}^{x=1-y} \, dy \\ &= \int_0^1 y(1-y)^2 \, dy = \frac{1}{12}.\end{aligned}$$

In this particular example, the two versions of [Theorem 6.3.2](#) lead to the same calculation, but in many cases one of the results is more convenient to apply than the other.

6.4 Change of variables in \mathbb{R}^2

[Section 6.3](#) provides us with a very general theory for integrations over subsets in \mathbb{R}^2 . However, except in the case of integration over rectangles or cases where results of the type in [Theorem 6.3.2](#) apply, concrete calculations of integrals will be tedious. The purpose of this section is to introduce a technique that in certain cases can transform integrals over complicated subsets into integrals over “simpler subsets”, e.g., rectangles. We will consider an important concrete case in [Section 6.5](#).

The main tool to transfer integration between “complicated” subsets and rectangles is vector functions of the following form.

Definition 6.4.1 Jacobian determinant

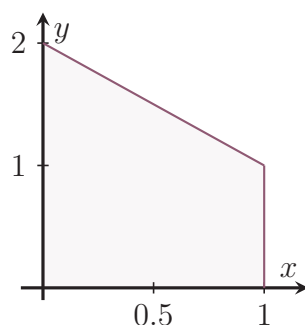
Let A be a subset in \mathbb{R}^2 . Consider a function $\mathbf{r} : A \rightarrow \mathbb{R}^2$, i.e., $\mathbf{r}(u, v) = (r_1(u, v), r_2(u, v))$ with coordinate functions $r_i : A \rightarrow \mathbb{R}$. Suppose \mathbf{r} is a C^1 vector function on A° , i.e., the partial derivatives of the coordinate functions exist and are continuous on the interior of A . Let $\mathbf{J}_r(u, v) \in M_{2 \times 2}(\mathbb{R})$ denote the Jacobian matrix of \mathbf{r} at $(u, v) \in U$. The determinant of $\mathbf{J}_r(u, v)$ is called the *Jacobian determinant* $\det(\mathbf{J}_r(u, v))$ of \mathbf{r} at $(u, v) \in U$.

We can be completely explicitly in [Definition 6.4.1](#). The Jacobian matrix of $\mathbf{r} = (r_1, r_2)$ at (u, v) is:

$$\mathbf{J}_r(u, v) = \begin{bmatrix} \frac{\partial r_1}{\partial u}(u, v) & \frac{\partial r_1}{\partial v}(u, v) \\ \frac{\partial r_2}{\partial u}(u, v) & \frac{\partial r_2}{\partial v}(u, v) \end{bmatrix} \quad (6.49)$$

and the Jacobian determinant is:

$$\det(\mathbf{J}_r(u, v)) = \begin{vmatrix} \frac{\partial r_1}{\partial u}(u, v) & \frac{\partial r_1}{\partial v}(u, v) \\ \frac{\partial r_2}{\partial u}(u, v) & \frac{\partial r_2}{\partial v}(u, v) \end{vmatrix} \quad (6.50)$$

Figure 6.3: The domain B in (6.52).

$$= \frac{\partial r_1}{\partial u}(u, v) \frac{\partial r_2}{\partial v}(u, v) - \frac{\partial r_2}{\partial u}(u, v) \frac{\partial r_1}{\partial v}(u, v). \quad (6.51)$$

Recall that the row vectors in the Jacobi matrix consist precisely of the gradient vectors for the functions r_1 and r_2 .

Example 6.4.1

The function

$$\mathbf{r}(u, v) = \begin{bmatrix} u^2v \\ u + v^2 \end{bmatrix}, \quad (u, v) \in \mathbb{R}^2,$$

is a continuously differentiable map from \mathbb{R}^2 to \mathbb{R}^2 . Its Jacobian matrix is

$$\mathbf{J}_r(u, v) = \begin{bmatrix} \frac{\partial r_1}{\partial u}(u, v) & \frac{\partial r_1}{\partial v}(u, v) \\ \frac{\partial r_2}{\partial u}(u, v) & \frac{\partial r_2}{\partial v}(u, v) \end{bmatrix} = \begin{bmatrix} 2uv & u^2 \\ 1 & 2v \end{bmatrix},$$

and its Jacobian determinant is

$$\det(\mathbf{J}_r(u, v)) = \det \begin{bmatrix} 2uv & u^2 \\ 1 & 2v \end{bmatrix} = 4uv^2 - u^2.$$

Let us now turn to the question of integrating a function $f : B \rightarrow \mathbb{R}$ over a potentially complicated subset $B \subset \mathbb{R}^2$. The key assumption in the subsequent [Theorem 6.4.1](#) is that we can find an injective and surjective map or function $\mathbf{r} : \Gamma \rightarrow B$ from a “simple subset” Γ in \mathbb{R}^2 onto B and then turn the “complicated” integral over B into a “simple integral” over Γ . Before stating [Theorem 6.4.1](#), let us consider an example of such sets B, Γ , and a corresponding map \mathbf{r} . Before reading this example, it is a good idea to remind yourself how we parametrized a line between two points in [Example 3.1.5](#) on page 70.

Example 6.4.2 Parametrization of a subset

Consider the subset

$$B = \{(x, y) \in \mathbb{R}^2 \mid 0 \leq x \leq 1 \wedge y \geq 0 \wedge x + y \leq 2\}, \quad (6.52)$$

which is illustrated on Figure 6.3.

We will now argue that the map

$$\mathbf{r} : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \quad \mathbf{r}(u, v) = \begin{bmatrix} u \\ v(2 - u) \end{bmatrix} \quad (6.53)$$

is mapping the rectangle

$$\Gamma = [0, 1]^2 = \{(u, v) \in \mathbb{R}^2 \mid 0 \leq u \leq 1 \wedge 0 \leq v \leq 1\}$$

bijectionally onto the set B .

In order to do so, consider first a fixed point $(x, y) \in B$ that is located on the “upper” boundary, that is, $x \in [0, 1]$ and $x + y = 2$. Then, for $u = x$

$$\mathbf{r}(x, 1) = (x, 2 - x) = (x, y),$$

i.e., the given point (x, y) is indeed the image of the point $(u, 1) \in \Gamma$. Also, if we consider a point $(x, y) \in B$ that is located on the “lower” boundary, then $x \in [0, 1]$ and $y = 0$; such a point is also in the image of the set Γ under the map \mathbf{r} because

$$\mathbf{r}(u, 0) = (u, 0).$$

Fixing now any $u \in [0, 1]$, and letting the parameter $v \in [0, 1]$ traverse from $v = 0$ to $v = 1$, the point $\mathbf{r}(u, v)$ will initiate in $\mathbf{r}(u, 0) = (u, 0)$ and traverse through a horizontal line to its terminating point $\mathbf{r}(u, 1) = (u, 2 - u)$. That is, for any fixed $u \in [0, 1]$ the function $\mathbf{r}(u, v)$ traverse from the point $(u, 0)$ on the “lower” boundary to the point $(u, 2 - u)$ on the “upper” boundary. When we now let $u \in [0, 1]$ vary, we thus parametrize all horizontal lines connecting the “lower” and “upper” boundary. In other words: each point (x, y) has a representation as $(x, y) = \mathbf{r}(u, v)$ for some $(u, v) \in \Gamma$. This shows that \mathbf{r} indeed maps the set Γ onto B .

When we consider \mathbf{r} as a map on Γ , i.e., $\mathbf{r} : \Gamma \rightarrow B$, then \mathbf{r} is also injective. In order to see this, assume that $\mathbf{r}(u_1, v_1) = \mathbf{r}(u_2, v_2)$ for some $(u_1, v_1), (u_2, v_2) \in \Gamma$. Then, by writing out how \mathbf{r} acts on the vectors,

$$u_1 = u_2 \text{ and } v_1(2 - u_1) = v_2(2 - u_2). \quad (6.54)$$

Using now that $(u_1, v_1), (u_2, v_2) \in \Gamma$, we know that $0 \leq u_1 \leq 1$, and therefore $2 - u_1 = 2 - u_2 \neq 0$. Thus the second equality in (6.54) implies that $v_1 = v_2$,

and therefore $(u_1, v_1) = (u_2, v_2)$. This proves that \mathbf{r} indeed is injective on the set Γ .

Altogether we have now proved that \mathbf{r} maps Γ *bijectively* onto B . We call $(u, v) \in \Gamma$ parameters, and we say that the map \mathbf{r} in (6.53) provides a *parametrization* of the set B .

We are now able to state the important *change of variables theorem*. The result is quite involved, and we will explain how to apply it after having stated it.

Theorem 6.4.1 Change of variables in \mathbb{R}^2

Let Γ be a bounded subset in \mathbb{R}^2 with boundary formed by a finite number of continuously differentiable curves. Suppose that there exists a C^1 vector function $\mathbf{r} : \Gamma \rightarrow \mathbb{R}^2$ that is injective on the interior Γ° and such that

$$\det(\mathbf{J}_{\mathbf{r}}(u, v)) \neq 0 \text{ for all } (u, v) \in \Gamma^\circ. \quad (6.55)$$

Then, for any continuous function $f : \mathbf{r}(\Gamma) \rightarrow \mathbb{R}$,

$$\int_{\mathbf{r}(\Gamma)} f(x, y) \, d\mathbf{x} = \int_{\Gamma} f(r_1(u, v), r_2(u, v)) |\det(\mathbf{J}_{\mathbf{r}}(u, v))| \, d\mathbf{u}, \quad (6.56)$$

where $\mathbf{x} = (x, y)$ and $\mathbf{u} = (u, v)$.

Note that (6.56) involves the *absolute value* of the Jacobian determinant of the map \mathbf{r} . The factor $|\det(\mathbf{J}_{\mathbf{r}}(u, v))|$ describes the local *change in area* of a small region around $(u, v) \in \Gamma^\circ$ under the mapping \mathbf{r} . We will make this statement precise and outline the proof of Theorem 6.4.1 when we consider the change of variables for functions of n variables in Theorem 6.6.3 on page 155. Let us here just consider a very simple example that shows that an “area correcting” factor is needed.

Example 6.4.3

Let $B = [0, 2] \times [-3, 0]$. The area of B is 6. Let us also compute this trivial statement using both (6.43) and (6.56). By Definition 6.3.2 on page 142 we get:

$$\text{area}(B) = \int 1 \, d\mathbf{x} = \int_{-3}^0 \left(\int_0^2 1 \, dx \right) dy = \int_{-3}^0 [x]_0^2 dy = \int_{-3}^0 2 \, dy = [2y]_{-3}^0 = 6$$

We can parametrize B using $\mathbf{r} : \Gamma \rightarrow \mathbb{R}^2$ given by $\mathbf{r}(u, v) = (2u, -3v)$ and $\Gamma = [0, 1]^2$. The Jacobian matrix is $\mathbf{J}_{\mathbf{r}}(u, v) = \text{diag}(2, -3)$, and its

determinant is -6 . Using (6.56) we get

$$\begin{aligned}\text{area}(B) &= \int_{\mathbf{r}(\Gamma)} 1 \, d\mathbf{x} = \int_{\Gamma} |-6| \, d\mathbf{u} = \int_0^1 \left(\int_0^1 6 \, du \right) dv \\ &= \int_0^1 6[u]_0^1 \, dv = \int_0^1 6 \, dv = [6v]_0^1 = 6.\end{aligned}$$

Without the factor $|\det(\mathbf{J}_{\mathbf{r}}(u, v))| = 6$, the latter integral would have given us 1. The Jacobian factor exactly describes the factor by which we “stretch” the unit square Γ to “fit” onto $B = \mathbf{r}(\Gamma)$. Note also that we need the absolute value of the Jacobian determinant as we would have otherwise arrived at a negative area.

For the computations in Example 6.4.3 we certainly did not need the “change of variables” theorem since B was a simple rectangle. Let us return the goal of this section: to obtain a tool to calculate integrals $\int_B f(x, y) \, d\mathbf{x}$, when B is a *complicated* subset. The idea is to *identify* a map $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ and a *simple* domain $\Gamma \subset \mathbb{R}^2$ satisfying the conditions in Theorem 6.4.1 such that $\mathbf{r}(\Gamma) = B$; in that case, we obtain that

$$\begin{aligned}\int_B f(x, y) \, d(x, y) &= \int_{\mathbf{r}(\Gamma)} f(x, y) \, d(x, y) \\ &= \int_{\Gamma} f(r_1(u, v), r_2(u, v)) |\det(\mathbf{J}_{\mathbf{r}}(u, v))| \, d(x, y).\end{aligned}\quad (6.57)$$

The hope is now that the simpler structure of the set Γ compared with B makes it possible to calculate the integral (6.57).

In practice, the difficult task is to identify an appropriate map \mathbf{r} and the set Γ . Let us consider a concrete case, where we can get through with the technical conditions, based on what we already know from Example 6.4.2.

Example 6.4.4

Let us consider the subset B from Example 6.4.2 on page 146. Our aim is to calculate the integral

$$\int_B 2xy \, d(x, y).\quad (6.58)$$

The function \mathbf{r} in (6.53) maps the rectangle $\Gamma = [0, 1]^2$ injectively and surjectively onto B . Hence, $\mathbf{r}(\Gamma) = B$. In order to apply Theorem 6.4.1, we therefore calculate the Jacobian determinant

$$\mathbf{J}_{\mathbf{r}}(u, v) = \begin{vmatrix} 1 & -v \\ 0 & 2 - u \end{vmatrix} = 2 - u$$

and note that $\mathbf{J}_r(u, v) = 2 - u > 0$ on $(u, v) \in \Gamma$. Parametrizations for which the Jacobian determinant is nonzero on Γ are sometimes called *regular*, cf. regular parametrizations [Definition 3.1.4](#) on page 70.

Thus, applying (6.57) with $f(x, y) = 2xy$,

$$\begin{aligned} \int_B 2xy \, d(x, y) &= \int_{\Gamma} f(r_1(u, v), r_2(u, v)) |\det(\mathbf{J}_r(u, v))| \, d(u, v) \\ &= \int_{[0,1]^2} 2uv(2-u)|2-u| \, d(u, v). \end{aligned}$$

Since the Jacobian determinant $2 - u$ is positive on the considered interval $u \in [0, 1]$, we can remove the absolute sign on the Jacobian determinant, and we obtain that

$$\begin{aligned} \int_B 2xy \, d(x, y) &= \int_0^1 \left(\int_0^1 2uv(2-u)(2-u) \, du \right) dv \\ &= \int_0^1 \int_0^1 2v(4u - 4u^2 + u^3) \, du \, dv \\ &= \int_0^1 2v \left[2u^2 - \frac{4}{3}u^3 + \frac{1}{4}u^4 \right]_{u=0}^{u=1} dv \\ &= \frac{11}{12} \int_0^1 2v \, dv = \frac{11}{12} [v^2]_0^1 = \frac{11}{12}. \end{aligned}$$

where we used [Theorem 6.3.1](#) on page 139 to compute the integral over $\Gamma = [0, 1]^2$.

Note that in this particular case, we could have obtained the result in [Example 6.4.4](#) in a much simpler way using the same procedure as in [Example 6.3.2](#). In [Section 6.5](#) we will see an application, where [Theorem 6.4.1](#) is of substantial importance.

6.5 Polar coordinates

In this section we will discuss an important application of [Theorem 6.4.1](#). For subsets B in \mathbb{R}^2 having a convenient representation in polar coordinates, it leads to a useful and explicit formula for integrals of the type $\int_B f(x, y) \, dx \, dy$.

Recall that a point in \mathbb{R}^2 can be represented either by its *rectangular coordinates*

$$(x, y), \quad x, y \in \mathbb{R},$$

or its *polar coordinates*

$$(r, \theta), \quad r \in [0, \infty[, \theta \in]-\pi, \pi].$$

The relation between these two types of coordinates is that

$$x = r \cos(\theta), \quad y = r \sin(\theta).$$

trigonometric functions are 2π -periodic in the angle θ hence we can use any interval of length 2π , e.g., $\theta \in]-\pi, \pi]$ as above, or $\theta \in [0, 2\pi[$.

Example 6.5.1 Parametrization of a subset

Consider the subset

$$B = \{(x, y) \in \mathbb{R}^2 \mid x \geq 0 \wedge y \geq x \wedge x^2 + y^2 \leq 2\}. \quad (6.59)$$

The region B is an eighth of a circular disc with radius $\sqrt{2}$, exactly, the piece above the line $y = x$ and to the right of $x = 0$ (including both lines). This region is easy to describe using polar coordinates as all pairs (r, θ) with $r \in [0, \sqrt{2}]$ and $\theta \in [\pi/4, \pi/2]$. This description leads to the following parametrization of B with $\Gamma = [0, \sqrt{2}] \times [\pi/4, \pi/2]$:

$$\mathbf{p} : \Gamma \rightarrow \mathbb{R}^2, \quad \mathbf{p}(r, \theta) = (r \cos(\theta), r \sin(\theta)).$$

The function \mathbf{p} is a transformation from polar coordinates (r, θ) to Cartesian coordinates (x, y) . The function is surjective. It is also injective, but only on the interior of Γ° . It fails to be injective on Γ since $\mathbf{p}(0, \theta) = (0, 0)$ for all values of θ .

In Section 6.4 we use parameters (u, v) and denoted the vector function \mathbf{r} . We use \mathbf{p} here since it is natural to use (r, θ) as parameters when working with polar coordinates. This is, of course, just a matter of notation. Moreover, we do not need to use $[0, \sqrt{2}] \times [\pi/4, \pi/2]$ as domain for our parameters. We are always able to re-scale the parameter domains of the form $[a, b] \times [c, d]$ to $[0, 1]^2$, e.g., the function

$$\mathbf{r} : [0, 1]^2 \rightarrow \mathbb{R}^2, \quad \mathbf{r}(u, v) = \left(\sqrt{2}u \cos\left(\frac{\pi}{4}(1+v)\right), \sqrt{2}u \sin\left(\frac{\pi}{4}(1+v)\right) \right)$$

is an alternative parametrization of B . In particular, parametrizations are never unique, but, and this is the important fact, the value of $\int_B f(x, y) dx dy$ is *independent* of the chosen parametrization.

Assume now that we have given $\alpha, \beta \in \mathbb{R}$ such that $0 \leq \alpha < \beta \leq 2\pi$, and two functions

$$\varphi_1 : [\alpha, \beta] \rightarrow \mathbb{R}, \quad \varphi_2 : [\alpha, \beta] \rightarrow \mathbb{R} \quad (6.60)$$

such that

$$0 \leq \varphi_1(\theta) \leq \varphi_2(\theta), \text{ for all } \theta \in [\alpha, \beta]. \quad (6.61)$$

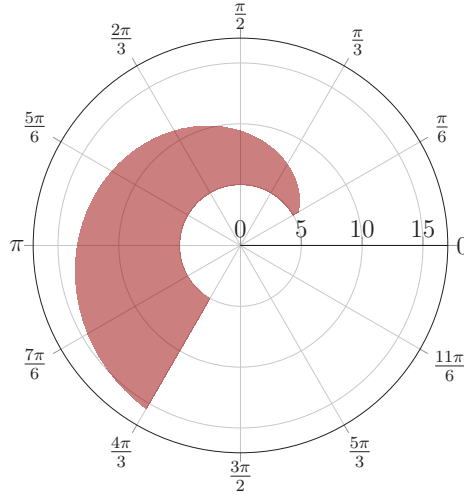


Figure 6.4: The subset $B \subset \mathbb{R}^2$ given by Equation (6.62) with $\alpha = \frac{\pi}{6}$, $\beta = \frac{4\pi}{3}$, $\phi_1(\theta) = 5$ and $\phi_2(\theta) = \sqrt{\frac{180\theta}{\pi}}$.

Associated with such functions, we consider the subset B in \mathbb{R}^2 which is given by

$$B = \left\{ (x, y) = (r \cos(\theta), r \sin(\theta)) \mid \alpha \leq \theta \leq \beta \wedge \varphi_1(\theta) \leq r \leq \varphi_2(\theta) \right\}. \quad (6.62)$$

An example of such a set is shown in Figure 6.4. The subset B considered Example 6.5.1 is also a set of this form.

In general, a subset B of this form will have a complicated representation in terms of rectangular (x, y) -coordinates which makes it cumbersome to calculate an integral of the form $\int_B f(x, y) dx dy$ directly. This is the point, where Theorem 6.4.1 turns out to be helpful. Indeed, consider the map

$$\mathbf{p} : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \quad \mathbf{p}(r, \theta) = \begin{bmatrix} r \cos(\theta) \\ r \sin(\theta) \end{bmatrix}. \quad (6.63)$$

Then \mathbf{p} is mapping the subset

$$\Gamma = \left\{ (r, \theta) \in \mathbb{R}^2 \mid \alpha \leq \theta \leq \beta \wedge \varphi_1(\theta) \leq r \leq \varphi_2(\theta) \right\} \quad (6.64)$$

bijectively onto the subset B ; that is, \mathbf{p} is injective and $\mathbf{p}(\Gamma) = B$. The Jacobi matrix of \mathbf{p} is

$$\mathbf{J}_{\mathbf{p}}(r, \theta) = \begin{bmatrix} \frac{\partial p_1}{\partial r}(r, \theta) & \frac{\partial p_1}{\partial \theta}(r, \theta) \\ \frac{\partial p_2}{\partial r}(r, \theta) & \frac{\partial p_2}{\partial \theta}(r, \theta) \end{bmatrix} = \begin{bmatrix} \cos(\theta) & -r \sin(\theta) \\ \sin(\theta) & r \cos(\theta) \end{bmatrix},$$

so its Jacobian determinant is

$$\det(\mathbf{J}_p(r, \theta)) = \begin{vmatrix} \cos(\theta) & -r \sin(\theta) \\ \sin(\theta) & r \cos(\theta) \end{vmatrix} = r \cos^2(\theta) + r \sin^2(\theta) = r.$$

Note that $\det(\mathbf{J}_p(r, \theta)) > 0$ for $(r, \theta) \in \Gamma^\circ$. Hence, the technical condition (6.55) is satisfied.

Hence, by Theorem 6.4.1 (see the version stated in (6.57)), for any continuous function $f : B \rightarrow \mathbb{R}$,

$$\begin{aligned} \int_B f(x, y) \, d(x, y) &= \int_\Gamma f(p_1(r, \theta), p_2(r, \theta)) |\det(\mathbf{J}_p(r, \theta))| \, d(r, \theta) \\ &= \int_\Gamma f(r \cos(\theta), r \sin(\theta)) r \, d(r, \theta) \\ &= \int_\alpha^\beta \int_{\varphi_1(\theta)}^{\varphi_2(\theta)} f(r \cos(\theta), r \sin(\theta)) r \, dr \, d\theta. \end{aligned}$$

Let us formulate the obtained result as a theorem.

Theorem 6.5.1 Integration in polar coordinates

Consider a subset B in \mathbb{R}^2 of the form (6.62). Then, for any continuous function $f : B \rightarrow \mathbb{R}$,

$$\int_B f(x, y) \, d(x, y) = \int_\alpha^\beta \int_{\varphi_1(\theta)}^{\varphi_2(\theta)} f(r \cos(\theta), r \sin(\theta)) r \, dr \, d\theta. \quad (6.65)$$

Example 6.5.2

Let us continue Example 6.5.1 with B given as Equation (6.59). We want to calculate the integral $\int_B f(x, y) \, dx \, dy$ for the function $f : B \rightarrow \mathbb{R}$ given by $f(x, y) = x^2 + y^2$ using Theorem 6.5.1. Hence, we first need to express f in polar coordinates:

$$f(r \cos(\theta), r \sin(\theta)) = (r \cos(\theta))^2 + (r \sin(\theta))^2 = r^2.$$

Then we need to express B in the form (6.62). We see that $\varphi_1(\theta) = 0$ and $\varphi_2(\theta) = \sqrt{2}$ for $\theta \in [\alpha, \beta] = [\pi/4, \pi/2]$. By Theorem 6.5.1, we finally compute

$$\begin{aligned} \int_B f(x, y) \, d(x, y) &= \int_{\pi/4}^{\pi/2} \int_0^{\sqrt{2}} r^2 r \, dr \, d\theta = \left(\frac{\pi}{2} - \frac{\pi}{4} \right) \int_0^{\sqrt{2}} r^3 \, dr \\ &= \frac{\pi}{4} \left[\frac{1}{4} r^4 \right]_0^{\sqrt{2}} = \frac{\pi}{4}. \end{aligned}$$

6.6 The Riemann integral of functions of several variables

Integration in \mathbb{R}^n is completely analogue to the procedure we described in Section 6.3 for \mathbb{R}^2 so we will not repeat all the details. We will follow the steps in Section 6.3 closely and first consider integration over rectangles.

A *rectangle in \mathbb{R}^n* is a set of the form

$$Q = \{(x_1, x_2, \dots, x_n) \in \mathbb{R}^n \mid a_i \leq x_i \leq b_i, \forall i = 1, 2, \dots, n\}. \quad (6.66)$$

In order to define the integral of a continuous function $f : Q \rightarrow \mathbb{R}$, we will first split the rectangle Q into subrectangles. Precisely as in Section 6.3, this is done by splitting each of the intervals $[a_i, b_i]$ into J intervals of equal length $\Delta x_i = \frac{b_i - a_i}{J}$ and forming the corresponding subrectangles. Recall that in the case of \mathbb{R}^2 the procedure gave us subrectangles $Q_{i,j}, i, j = 1, \dots, J$, that is, we were dealing with J^2 subrectangles. In the case of \mathbb{R}^n , we obtain J^n subrectangles; in order to avoid a cumbersome multi-indexing, we will simply denote them by $Q_i, i = 1, 2, \dots, J^n$. The n -dimensional volume of each of these rectangles $Q_i, i = 1, 2, \dots, J^n$, is $\Delta x_1 \Delta x_2 \cdots \Delta x_n$. For each of these subrectangles $Q_i \subset \mathbb{R}^n$, we *define* its n -dimensional volume as the product of its side lengths

$$\text{vol}_n(Q_i) := \Delta x_1 \Delta x_2 \cdots \Delta x_n,$$

and we pick a point $\xi_i \in Q_i$. Similarly to what we did in (6.31), the corresponding *Riemann sum* S_J is then defined by

$$S_J = \sum_{i=1}^{J^n} f(\xi_i) \text{vol}_n(Q_i) \quad (6.67)$$

as $J \rightarrow \infty$, the Riemann sum S_J tends to a limit, which we denote by

$$\int_Q f(x_1, x_2, \dots, x_n) d(x_1, x_2, \dots, x_n) \quad (6.68)$$

Similar to Theorem 6.3.1, integration over a rectangle in \mathbb{R}^n can be reduced to n consecutive integrations over bounded intervals in \mathbb{R} .

Theorem 6.6.1

Consider a rectangle Q in \mathbb{R}^n of the form (6.66), and let $f : Q \rightarrow \mathbb{R}$ denote

6.6. The Riemann integral of functions of several variables

a continuous function. Then

$$\begin{aligned} & \int_Q f(x_1, x_2, \dots, x_n) \, d(x_1, x_2, \dots, x_n) \\ &= \int_{a_n}^{b_n} \cdots \left(\int_{a_2}^{b_2} \left(\int_{a_1}^{b_1} f(x_1, x_2, \dots, x_n) \, dx_1 \right) dx_2 \right) \cdots dx_n. \end{aligned} \quad (6.69)$$

Calculations using (6.69) are always done starting from *inside*. We first calculate the integral $\int_{a_1}^{b_1} f(x_1, x_2, \dots, x_n) \, dx_1$, then $\int_{a_2}^{b_2} \left(\int_{a_1}^{b_1} f(x_1, x_2, \dots, x_n) \, dx_1 \right) dx_2$, and so on.

The Riemann integral is extended to more general subsets B in \mathbb{R}^n precisely as we did for \mathbb{R}^2 on page 140. To be more precise, we will assume that

(I) B is bounded, i.e.,

$$B \subset [a_1, b_1] \times [a_2, b_2] \times \cdots \times [a_n, b_n] \quad (6.70)$$

for some $a_i, b_i, \in \mathbb{R}, i = 1, \dots, n$.

(II) The boundary ∂B of B is *smooth*, i.e., it is formed as the union of the graphs for a finite number of continuously differentiable functions from \mathbb{R}^{n-1} to \mathbb{R} .

Fix again an integer $J \in \mathbb{N}$, split each of the intervals $[a_1, b_1], [a_2, b_2], \dots, [a_n, b_n]$ into J intervals of equal length $\Delta x_i = \frac{b_i - a_i}{J}$, and form the corresponding subrectangles $Q_i \subset \mathbb{R}^n, i = 1, \dots, \dots, J^n$. Among the rectangles Q_i , we will *exclusively consider* the ones that are completely contained in B . For these sets Q_i , we pick arbitrary points $\xi_i \in Q_i$, and form the corresponding *Riemann sum*

$$S_J = \sum_{\{i|Q_i \subset B\}} f(\xi_i) \operatorname{vol}_n(Q_i). \quad (6.71)$$

where the notation $\sum_{\{i|Q_i \subset B\}}$ means that we only include the i th term in the sum if Q_i is contained in B .

One can again show that if the function f is continuous, the Riemann sums S_J tend to a limit as $J \rightarrow \infty$. The limit is denoted by

$$\int_B f(x_1, x_2, \dots, x_n) \, d(x_1, x_2, \dots, x_n), \quad \int_B f(x_1, x_2, \dots, x_n) \, d\mathbf{x}, \quad (6.72)$$

or by

$$\int_B f(x_1, x_2, \dots, x_n) \, dX, \quad (6.73)$$

6.6. The Riemann integral of functions of several variables

and it is called the *Riemann integral of f over the subset B* . Note that the Riemann integral of f over the subset B is a *real scalar*, and we will sometimes use the compact notation $\int_B f(\mathbf{x}) \, d\mathbf{x}$.

Similar to the definition of *area* in \mathbb{R}^2 , we define the *n -dimensional volume* of $B \subset \mathbb{R}^n$ as

$$\text{vol}_n(B) = \int_B 1 \, d\mathbf{x}.$$

Exercise 6.6.2

A parallelotope P in \mathbb{R}^n “spanned” by n linearly independent vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \subset \mathbb{R}^n$ is defined by:

$$P = \{\mathbf{y} \in \mathbb{R}^n \mid \mathbf{y} = A\mathbf{x}, \quad \text{where } x_i \in [0, 1] \text{ for } i = 1, 2, \dots, n\}$$

where A is the $n \times n$ matrix whose i th column is \mathbf{a}_i . The set of points P can also be written as $P = A([0, 1]^n)$. The n -dimensional volume of P is given by the formula:

$$\text{vol}_n(P) = |\det(A)|. \tag{6.74}$$

We want to prove this fact *without* using tools from integration theory.

- (a) Give a proof of (6.74) in the case $n = 2$.
- (b) Give a proof of (6.74) in the case $n = 3$.
- (c) Give a proof of (6.74) in the general case $n \in \mathbb{N}$ (Note: this is a difficult exercise).

Finally, the definition of a (not necessarily continuous) function $f : B \rightarrow \mathbb{R}$, $B \subset \mathbb{R}^n$, being *Riemann integrable* is similar to how we did in Definition 6.3.1 on page 141, and we leave the details to the reader.

The technique for changing variables is a straightforward generalization of the result we saw for \mathbb{R}^2 in Theorem 6.4.1.

Theorem 6.6.3 Change of variables in \mathbb{R}^n

Let Γ be a bounded subset in \mathbb{R}^n having a smooth boundary. Suppose that there exists a C^1 vector function $\mathbf{r} : \Gamma \rightarrow \mathbb{R}^n$ which is injective on the interior Γ° and such that

$$\det(\mathbf{J}_r(\mathbf{u})) \neq 0 \text{ for all } \mathbf{u} = (u_1, \dots, u_n) \in \Gamma^\circ. \tag{6.75}$$

Then, for any continuous function $f : \mathbf{r}(\Gamma) \rightarrow \mathbb{R}$,

$$\int_{\mathbf{r}(\Gamma)} f(\mathbf{x}) \, d\mathbf{x} = \int_{\Gamma} f(\mathbf{r}(\mathbf{u})) |\det(\mathbf{J}_r(\mathbf{u}))| \, d\mathbf{u}. \tag{6.76}$$

6.6. The Riemann integral of functions of several variables

that is,

$$\int_{\mathbf{r}(\Gamma)} f(x_1, x_2, \dots, x_n) d(x_1, x_2, \dots, x_n) = \int_{\Gamma} f(\mathbf{r}(u_1, \dots, u_n)) |\det(\mathbf{J}_{\mathbf{r}}(u_1, \dots, u_n))| d(u_1, u_2, \dots, u_n).$$

Proof. A full proof of [Theorem 6.6.3](#) is beyond the scope of this book, but it is useful to do a heuristic argument in order to see where the factor $|\det(\mathbf{J}_{\mathbf{r}}(\mathbf{u}))|$ comes from. It is clear from [Example 6.4.3](#) on page 147 that some “Jacobian factor” is needed.

Let $r_1, r_2, \dots, r_k : \Gamma \rightarrow \mathbb{R}$ be the coordinate functions, so that $\mathbf{r}(\mathbf{u}) = (r_1(\mathbf{u}), \dots, r_k(\mathbf{u}))$. Let $\mathbf{u} \in \Gamma$ be fixed. Since \mathbf{r} is continuously differentiable, it follows, e.g., using Taylor’s formula, that

$$\mathbf{r}(\mathbf{u} + \mathbf{h}) = \mathbf{r}(\mathbf{u}) + \mathbf{J}_{\mathbf{r}}(\mathbf{u}) \mathbf{h} + \boldsymbol{\varepsilon}(\mathbf{h}) \|\mathbf{h}\|,$$

for $\|\mathbf{h}\| < \delta$, where δ is chosen so small that $\mathbf{u} + \mathbf{h} \in \Gamma^\circ$. For brevity we write the Jacobian matrix at \mathbf{u} as $\mathbf{J} = \mathbf{J}_{\mathbf{r}}(\mathbf{u})$.

Let $\epsilon_i > 0, i = 1, 2, \dots, n$, be given (small) real scalars. The rectangle, depending on \mathbf{u} and $\boldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)$,

$$\begin{aligned} Q &= [u_1, u_1 + \epsilon_1] \times [u_2, u_2 + \epsilon_2] \times \cdots \times [u_n, u_n + \epsilon_n] \\ &= \{\mathbf{h} \mid u_i \leq h_i \leq u_i + \epsilon_i \text{ for } i = 1, 2, \dots, n\} \end{aligned}$$

is under \mathbf{r} mapped to $\mathbf{r}(Q)$. By Taylor’s formula this rectangle is approximately mapped to

$$\{\mathbf{r}(\mathbf{u}) + \mathbf{J} \mathbf{h} \mid 0 \leq h_1 \leq \epsilon_1, \dots, 0 \leq h_n \leq \epsilon_n\},$$

which can also be written as $\mathbf{r}(\mathbf{u}) + \mathbf{J}([0, \epsilon_1] \times \cdots \times [0, \epsilon_n])$. However, $\mathbf{J}([0, \epsilon_1] \times \cdots \times [0, \epsilon_n])$ is a parallelotope given by $A([0, 1]^n)$, where A is the matrix whose i th column is ϵ_i times the i th column of $\mathbf{J} = \mathbf{J}_{\mathbf{r}}(\mathbf{u})$.

The n -dimensional volume of the parallelotope $\mathbf{J}([0, \epsilon_1] \times \cdots \times [0, \epsilon_n])$ is

$$\epsilon_1 \cdot \epsilon_2 \cdots \epsilon_n \cdot |\det(\mathbf{J}_{\mathbf{r}}(\mathbf{u}))|$$

while the n -dimensional volume of the rectangle Q is

$$\text{vol}_n(Q) = \epsilon_1 \cdot \epsilon_2 \cdots \epsilon_n$$

Intuitively, we thus have that

$$\frac{\text{vol}_n(\mathbf{J}([0, \epsilon_1] \times \cdots \times [0, \epsilon_n]))}{\text{vol}_n(Q)} \rightarrow |\det(\mathbf{J}_{\mathbf{r}}(\mathbf{u}))| \quad \text{as } \boldsymbol{\epsilon} \rightarrow \mathbf{0},$$

6.6. The Riemann integral of functions of several variables

i.e., $|\det(\mathbf{J}_r(\mathbf{u}))|$ is the local expansion factor at the point \mathbf{u} under the transformation $\mathbf{u} \mapsto \mathbf{x} = \mathbf{r}(\mathbf{u})$.

In other words, the “infinitesimal” rectangle $[u_1, u_1 + \Delta u_1] \times \cdots \times [u_n, u_n + \Delta u_n]$ is a rectangle with side lengths $\Delta u_1, \dots, \Delta u_n$, and the image of this “infinitesimal” rectangle under \mathbf{r} is a parallelotope with volume $|\det(\mathbf{J}_r(\mathbf{u}))| \Delta u_1 \cdots \Delta u_n$. Hence, if we consider a Riemann sum (6.71) for $f \circ \mathbf{r}$ over $\mathbf{u} \in \Gamma$, we need to replace $f(\boldsymbol{\xi}_i) \text{vol}_n(Q_i)$ with the corrected volume

$$f(\mathbf{r}(\boldsymbol{\gamma}_i)) \text{vol}_n(Q_i) |\det(\mathbf{J}_r(\mathbf{u}_i))|,$$

where $\boldsymbol{\xi}_i = \mathbf{r}(\boldsymbol{\gamma}_i)$ and $\mathbf{u}_i \in Q_i$. ■

Remark 6.6.1 Jacobian

The “Jacobian factor” $|\det(\mathbf{J}_r(\mathbf{u}))|$ appearing in Equation (6.76) is often simply called *the Jacobian* of the change of variables. When using this terminology, one has to be careful not to confuse the Jacobian with the Jacobian matrix and the Jacobian determinant.

Exercise 6.6.4

Consider again the parallelotope from Exercise 6.6.2. Show that the formula $\text{vol}_n(P) = |\det(A)|$ follows immediately from Theorem 6.6.3.

Let us end this section with a few examples of the important special case of \mathbb{R}^3 .

Example 6.6.1 Coordinate systems in three-dimensional space

There are three commonly used coordinate systems in \mathbb{R}^3 : the well-known Cartesian coordinates (x, y, z) , the cylindrical coordinates (ρ, ϕ, z) , and the spherical coordinates (r, θ, ϕ) . See Figure 6.5 on the next page for an illustration.

In cylindrical coordinates, the position of a point in three-dimensional space is described by the parameters (ρ, ϕ, z) , where ρ is the radial distance from the z -axis, ϕ is the so-called azimuthal angle measured in the xy -plane from the positive x -axis (similar to θ in polar coordinates), and z is the height above the xy -plane. The parametrization of a point in space using cylindrical coordinates can be expressed as the vector function

$$\mathbf{c}(\rho, \phi, z) = \begin{bmatrix} \rho \cos(\phi) \\ \rho \sin(\phi) \\ z \end{bmatrix}, \quad (6.77)$$

where $\rho \geq 0$, $\phi \in [0, 2\pi[$ and $z \in \mathbb{R}$. The Jacobian determinant of the

6.6. The Riemann integral of functions of several variables

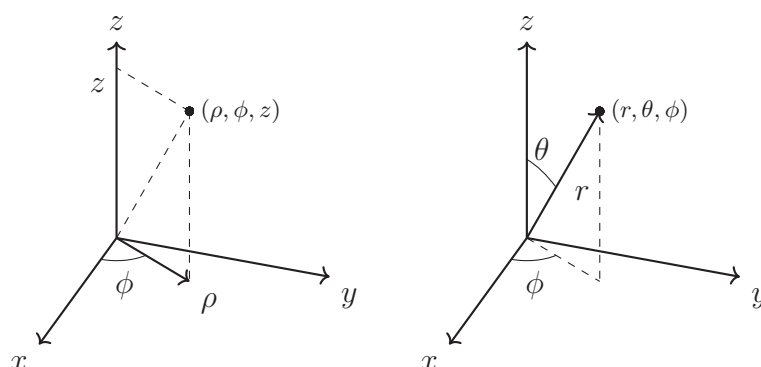


Figure 6.5: Illustration of cylindrical and spherical coordinates of a point in \mathbb{R}^3 .

transformation from cylindrical to Cartesian coordinates is given by the formula $\mathbf{J}_c(\rho, \phi, z) = \rho$.

In spherical coordinates, the position of a point is given by three parameters (r, θ, ϕ) , where r is the radial distance from the origin, θ is the polar angle measured from the positive z -axis, and ϕ is the azimuthal angle in the xy -plane from the positive x -axis. The parametrization of a point in space using spherical coordinates can be expressed as the vector function

$$\mathbf{s}(r, \theta, \phi) = \begin{bmatrix} r \sin(\theta) \cos(\phi) \\ r \sin(\theta) \sin(\phi) \\ r \cos(\theta) \end{bmatrix}, \quad (6.78)$$

where $r \geq 0$, $\theta \in [0, \pi]$ and $\phi \in [0, 2\pi[$. The Jacobian determinant for the transformation from spherical to Cartesian coordinates is $\mathbf{J}_s(r, \theta, \phi) = r^2 \sin(\theta)$.

Example 6.6.2 Volume of a ball

Let us compute the volume of a closed^a ball B in \mathbb{R}^3 of radius $R > 0$. We need to compute

$$\text{vol}_3(B) = \int_B 1 \, d(x, y, z).$$

A point belongs to B if, in spherical coordinates (r, θ, ϕ) , the radial distance r is smaller than or equal to R , or, in other words, if $r \in [0, R]$, $\theta \in [0, \pi]$ and $\phi \in [0, 2\pi[$. Thus, we let $\Gamma = [0, R] \times [0, \pi] \times [0, 2\pi[$ so that $B = \mathbf{s}(\Gamma)$, where \mathbf{s} is given by (6.78). On Γ° , the function \mathbf{s} is injective and its Jacobian

determinant is positive. Hence, by Theorem 6.6.3 on page 155, it follows

$$\begin{aligned} \text{vol}_n(B) &= \int_{s(\Gamma)} 1 \, d(x, y, z) = \int_{\Gamma} \det(\mathbf{J}_s(r, \theta, \phi)) \, d(r, \theta, \phi) \\ &= \int_0^{2\pi} \int_0^{\pi} \int_0^R r^2 \sin(\theta) \, dr \, d\theta \, d\phi = \int_0^{2\pi} \int_0^{\pi} \left[\frac{1}{3} r^3 \right]_0^R \sin(\theta) \, d\theta \, d\phi \\ &= \frac{1}{3} R^3 \int_0^{2\pi} [-\cos(\theta)]_0^{\pi} \, d\phi = \frac{1}{3} R^3 (-1 - (-1)) \int_0^{2\pi} d\phi \\ &= \frac{1}{3} R^3 \cdot 2 \cdot 2\pi = \frac{4\pi}{3} R^3 \end{aligned}$$

^aWhether the ball is open or closed is not essential for the computation and the result is the same.

6.7 The Riemann integral of vector functions

We now turn to the case of integrating vector functions of several variables. The situation is simple, in the sense that, we will just integrate each coordinate of the vector function separately using the techniques we have just developed in Section 6.6 on page 153.

Let us consider a vector function $\mathbf{f} = (f_1, \dots, f_k) : \text{dom}(\mathbf{f}) \rightarrow \mathbb{R}^k$. Given a subset $B \subset \text{dom}(\mathbf{f})$ in \mathbb{R}^n that satisfies our standard assumptions Items (I) and (II) on page 154, we define the Riemann integral of \mathbf{f} over B

$$\int_B \mathbf{f}(\mathbf{x}) \, d\mathbf{x} \tag{6.79}$$

as a vector in \mathbb{R}^k whose i th coordinate is $\int_B f_i(\mathbf{x}) \, d\mathbf{x}$ for $i = 1, \dots, k$.

Example 6.7.1 Center of mass

Let $f : B \rightarrow \mathbb{R}$ be a continuous function on a bounded, connected set $B \subset \mathbb{R}^n$ describing the density (kg m^{-n}). In other words, we consider a “solid” $B \subset \mathbb{R}^n$, whose mass distribution is continuous with density $f(\mathbf{x})$ for $\mathbf{x} \in B$. Typically $n = 3$, but we can easily consider the general case.

The total mass M [kg] of B is given by

$$M = \int_B f(\mathbf{x}) \, d\mathbf{x}. \tag{6.80}$$

Let us argue that the unit of M is indeed kg. The integral is a limit of Riemann sums S_J , and S_J is a sum of terms of the form $f(\xi_i) \text{vol}_n(Q_i)$. Since $f(\xi_i)$ and $\text{vol}_n(Q_i)$ have units kg m^{-n} and m^n , respectively, the Riemann sum and, therefore also, the integral in (6.80) have unit kg. The integral as

6.7. The Riemann integral of vector functions

a limiting process of Riemann sums also explains *why* the mass M is given by (6.80).

By similar considerations, the center of mass $\mathbf{x}^{\text{CM}} \in \mathbb{R}^n$ is the integral of the vector-valued function $\mathbf{x} \mapsto \mathbf{x}f(\mathbf{x})$, $B \rightarrow \mathbb{R}^n$:

$$\mathbf{x}^{\text{CM}} = \frac{1}{M} \int_B \mathbf{x}f(\mathbf{x}) \, d\mathbf{x}, \quad (6.81)$$

i.e.,

$$x_i^{\text{CM}} = \frac{1}{M} \int_B x_i f(\mathbf{x}) \, d\mathbf{x}$$

for each $i = 1, \dots, n$. The unit of each coordinate x_i^{CM} is m.

CHAPTER 7

Vector Fields

Let U be an open subset of \mathbb{R}^n . A *vector field* is a vector function $\mathbf{f} : U \rightarrow \mathbb{R}^k$, where $k = n$, i.e., where the dimension of the domain and co-domain are equal. We have already seen vector fields several times, cf. [Example 1.4.2](#) on page 17. Vector fields are often denoted \mathbf{V} in place of \mathbf{f} , and we will follow this tradition. When dealing with vector fields, we often think of \mathbf{V} as a *field* of vectors on U , i.e., we attach the vector $\mathbf{V}(\mathbf{x})$ to the point $\mathbf{x} \in U$, see [Figure 1.4](#) on page 17.

The vectors could represent the velocity of a fluid, the inner forces or strains in a beam, or the electric field around an antenna. Since a vector field is a special case of a vector function, we can integrate them over subsets of \mathbb{R}^n as outlined in [Section 6.7](#). However, in applications, we are often interested in another type of integration. Let us give two examples of such integrals.

If $\mathbf{V} : U \rightarrow \mathbb{R}^3$ is a force field, we may want to calculate the work done by the force. Recall that work is done when a force acts on something that undergoes a displacement from one position to another. This could be a displacement along a space curve $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^3$ connecting the starting point $\mathbf{r}(a) \in \mathbb{R}^3$ to the end-point $\mathbf{r}(b) \in \mathbb{R}^3$. The curve could be, e.g., the helix considered in [\(1.7\)](#) on [Page 14](#). To calculate the work done, we need to integrate the “length” of the force vectors in the direction of the tangent vector of the curve over the entire curve from the start-point $\mathbf{r}(a) \in \mathbb{R}^3$ to the end-point $\mathbf{r}(b) \in \mathbb{R}^3$.

If $\mathbf{V} : U \rightarrow \mathbb{R}^3$ represents the velocity of a fluid, we may want to calculate the flow rate, i.e., the flowing per unit of time through a surface \mathcal{F} in \mathbb{R}^3 . In this case, the function we are integrating is the length of the velocity vectors in the direction of the *normal* vector to the surface, i.e., in the direction of the orthogonal complement of the tangent plane. A surface could be the upper half of the unit sphere; cf. the parametrization considered in [\(1.8\)](#) on

Page 14 with $u \in [0, 2\pi[$ and $v \in [0, \pi/2]$.

In both these cases we need to integrate scalar quantities $f : U \rightarrow \mathbb{R}$, $U \subseteq \mathbb{R}^3$, over “thin” subsets of U such as curves or surfaces in \mathbb{R}^3 . However, if we use the integration techniques from the previous chapter, we do not get a useful outcome. E.g., let $f(x, y, z) = 1$ be the constant function for $(x, y, z) \in \mathbb{R}^3$. The Riemann integral of f over the upper half of the unit sphere \mathcal{F} in \mathbb{R}^3 is

$$\int_{\mathcal{F}} f(x, y, z) \, dX = \int_{\mathcal{F}} 1 \, dX = \text{vol}_3(\mathcal{F}) = 0$$

since the 3-dimensional volume of the sphere is zero. A sphere of radius $r > 0$ has *surface area* $4\pi r^2$, and the desired outcome of the integration of $f = 1$ over \mathcal{F} , where $r = 1$, is indeed $\frac{1}{2}4\pi = 2\pi$.

However, no matter which continuous function we are trying to integrate, the Riemann integral in \mathbb{R}^n will give 0 if we integrate over “thin” subsets such as the surface considered above. The aim of this chapter is to develop a useful integration over “thin” subsets, and to develop the notion of curve (or line, as we will call it) and surface integrals of vector fields.

7.1 Parametric curves and surfaces

Our starting point when considering curves and surfaces in \mathbb{R}^n will always be a parametrization. So throughout this chapter we let $\mathbf{r} : \Gamma \rightarrow \mathbb{R}^n$ be a continuous vector function, where Γ is a connected subset of \mathbb{R}^m ($m \leq n$) that satisfies [Items \(I\) and \(II\)](#) on page 154. Usually, we will take Γ to be a rectangle of the form

$$\begin{aligned} \Gamma &= [a_1, b_1] \times [a_2, b_2] \times \cdots \times [a_m, b_m] \\ &= \{\mathbf{u} \in \mathbb{R}^m \mid a_j \leq u_j \leq b_j \, \forall j = 1, 2, \dots, m\}, \end{aligned}$$

where $a_j < b_j$ for $j = 1, \dots, m$.

The number of parameters in the parametrization \mathbf{r} is given by m . For curves we have $m = 1$ (and $n \geq 2$), and we then often take Γ to be a bounded and closed interval, i.e., $\Gamma = [a, b]$. For surfaces we have $m = 2$ (and $n \geq 3$), and we then usually take Γ to be a bounded and closed rectangle, i.e., $\Gamma = [a_1, b_1] \times [a_2, b_2]$. A *closed curve* is a curve $\mathbf{r}([a, b])$ with no “endpoints”, i.e., $\mathbf{r}(a) = \mathbf{r}(b)$. It is called a *simple closed curve*, if the curve does not intersect itself, i.e., if \mathbf{r} is injective on the interior of Γ , e.g., on $]a, b[$. A *closed surface* is a surface with no “boundary points”. It is beyond the scope of this text to precisely define what it means for a surface to have no

“boundary points”, but the intuitive meaning should be clear: the boundary $\delta\mathcal{F}$ of a surface \mathcal{F} is where the surface stops. Let us consider some examples.

Example 7.1.1 Closed surfaces

The sphere with center $\mathbf{a} \in \mathbb{R}^3$ and radius $r \geq 0$:

$$\mathcal{F}_1 = \{\mathbf{x} \in \mathbb{R}^3 \mid \|\mathbf{x} - \mathbf{a}\| = r\}$$

is a closed surface with no “boundary points”, i.e., $\delta\mathcal{F}_1 = \emptyset$. The same holds for the torus. Note the potential confusion of terminology here. According to Definition 2.2.4 on page 38 the sphere is equal to its boundary, i.e., $\mathcal{F}_1 = \partial\mathcal{F}_1$. We use $\delta\mathcal{F}_1$ to denote the boundary of the surface, while $\partial\mathcal{F}_1$ denotes the boundary of \mathcal{F}_1 as a subset of \mathbb{R}^n .

The half sphere given by

$$\mathcal{F}_2 = \{\mathbf{x} \in \mathbb{R}^3 \mid \|\mathbf{x}\| = r \wedge x_3 \geq 0\}$$

has “surface boundary” given by

$$\delta\mathcal{F}_2 = \{\mathbf{x} \in \mathbb{R}^3 \mid x_1^2 + x_2^2 = r^2 \wedge x_3 = 0\}.$$

Since \mathcal{F}_2 has a boundary $\delta\mathcal{F}_2 \neq \emptyset$, it is not a closed surface.

Let $h_1, h_2 \in \mathbb{R}$. A hollow cylinder surface without top and bottom given by

$$\mathcal{F}_3 = \{\mathbf{x} \in \mathbb{R}^3 \mid x_1^2 + x_2^2 = r^2 \wedge x_3 \in [h_1, h_2]\}$$

has boundary

$$\delta\mathcal{F}_3 = \{\mathbf{x} \in \mathbb{R}^3 \mid x_1^2 + x_2^2 = r^2 \wedge (x_3 = h_1 \vee x_3 = h_2)\}.$$

Definition 7.1.1 Regular parametrization

Let $m \leq n$, and let $\mathbf{r} : \Gamma \rightarrow \mathbb{R}^n$, $\Gamma \subset \mathbb{R}^m$, be a C^1 vector function on an open set containing Γ , i.e., all partial derivatives of the coordinate functions $\mathbf{r} = (r_1, r_2, \dots, r_m)$ exist and are continuous. The *Jacobian* (also called the *geometric tensor*) of \mathbf{r} is defined to be the function given by:

$$\Gamma \rightarrow \mathbb{R}, \mathbf{u} \mapsto \sqrt{\det(\mathbf{J}^T \mathbf{J})}, \quad (7.1)$$

where $\mathbf{J} := \mathbf{J}_r(\mathbf{u}) \in \mathbb{M}_{n \times m}(\mathbb{R})$ is the Jacobian matrix of \mathbf{r} at $\mathbf{u} \in \Gamma$. If

$$\det\left((\mathbf{J}_r(\mathbf{u}))^T \mathbf{J}_r(\mathbf{u})\right) \neq 0 \quad (7.2)$$

for all $\mathbf{u} \in \Gamma^\circ$, the parametrization \mathbf{r} is said to be *regular*, and $\mathbf{r}(\Gamma)$ is said to be a *smooth m -fold* (a 1-fold is a curve and a 2-fold is a surface) .

Let us study the matrix product $\mathbf{G} := \mathbf{J}^T \mathbf{J}$ in (7.1) in more detail. Recall that the i th row, $i = 1, 2, \dots, n$, of the Jacobian matrix $\mathbf{J}_r(\mathbf{u}) \in \mathbb{R}^{n \times m}$ is the (transpose of the) gradient vector of the j th coordinate function, i.e., $(\nabla r_j(\mathbf{u}))^T \in \mathbb{R}^m$. On the other hand, we can also consider the columns of the Jacobian matrix. The j th column, $j = 1, \dots, m$, of $\mathbf{J}_r(\mathbf{u})$ at $\mathbf{u} = (u_1, u_2, \dots, u_m) \in \Gamma$ is the *tangent vector* $\frac{\partial \mathbf{r}}{\partial u_j}(\mathbf{u})$, also denoted by

$$\mathbf{r}'_{u_j}(\mathbf{u}) := \left(\frac{\partial r_1}{\partial u_j}(\mathbf{u}), \frac{\partial r_2}{\partial u_j}(\mathbf{u}), \dots, \frac{\partial r_n}{\partial u_j}(\mathbf{u}) \right) \in \mathbb{R}^n.$$

Hence, the (i, j) entry of $\mathbf{G} = \mathbf{J}^T \mathbf{J} \in \mathbb{R}^{m \times m}$ is the inner product¹ $\langle \mathbf{r}'_{u_i}, \mathbf{r}'_{u_j} \rangle = (\mathbf{r}'_{u_j})^T \mathbf{r}'_{u_i}$. Such matrices are often called a Gram matrix or the *Gramian* of the tangent vectors. One can show that the condition (7.2) is equivalent with $\mathbf{J}_r(\mathbf{u})$ being of rank m for every $\mathbf{u} \in \Gamma^\circ$, i.e., the tangent vectors \mathbf{r}'_{u_j} , $j = 1, 2, \dots, m$, being linearly independent and, hence, the tangent vectors span the tangent plane of \mathbf{r} at \mathbf{u} .

From Example 2.9.1 on page 63 we have that the \mathbf{G} is a positive semi-definite matrix, and therefore $\det(\mathbf{G}) = \lambda_1 \cdots \lambda_n \geq 0$. Hence, the condition in (7.2) simply says that \mathbf{G} is positive definite, i.e., the eigenvalues of \mathbf{G} cannot become zero.

We introduced the Jacobian in Remark 6.6.1 on page 157 in Chapter 6 for parametrizations where $n = m$. We will show in Lemma 7.1.1 that this form of the Jacobian is identical to the Jacobian introduced in Definition 7.1.1, that is, there is *one* expression of the Jacobian “factor”, namely (7.1), that can be used in all cases. However, it is often cumbersome to compute the Jacobian by the formula $(\det((\mathbf{J}_r(\mathbf{u}))^T \mathbf{J}_r(\mathbf{u})))^{1/2}$. Hence, in Lemma 7.1.1, we will also give explicit formulas in the most common situations. For this formulation we need to introduce the *cross product* of two vectors \mathbf{x}, \mathbf{y} in \mathbb{R}^3 :

$$\begin{aligned} \mathbf{x} \times \mathbf{y} &= \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \times \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \\ &= \begin{bmatrix} x_2 y_3 - x_3 y_2 \\ x_3 y_1 - x_1 y_3 \\ x_1 y_2 - x_2 y_1 \end{bmatrix}. \end{aligned} \tag{7.3}$$

¹In case of complex-valued vector functions \mathbf{r} , one needs to replace the transpose \mathbf{J}^T by the adjoint \mathbf{J}^* in the definition $\mathbf{G} := \mathbf{J}^T \mathbf{J}$.

The cross product is not commutative as it satisfies

$$\mathbf{x} \times \mathbf{y} = -\mathbf{y} \times \mathbf{x}.$$

The cross product is a mapping from $\mathbb{R}^3 \times \mathbb{R}^3$ to \mathbb{R}^3 for which $\mathbf{x} \times \mathbf{y}$ is orthogonal to \mathbf{x} and \mathbf{y} , that is $\mathbf{x} \times \mathbf{y} \in \{\mathbf{x}, \mathbf{y}\}^\perp$. The length of the cross product can be computed by Lagrange's identity:

$$\|\mathbf{x} \times \mathbf{y}\|^2 = \|\mathbf{x}\|^2\|\mathbf{y}\|^2 - \langle \mathbf{x}, \mathbf{y} \rangle^2 \quad (7.4)$$

which, in turn, follows from the Cauchy-Binet formula.

Lemma 7.1.1 Jacobian in special cases

Consider regular parametrizations as in Definition 7.1.1 on page 163.

- (i) (Curves in \mathbb{R}^n). Let $m = 1$ and let $I = \Gamma \subset \mathbb{R}$ be an interval. The tangent vector is $\mathbf{r}'(u) = (r'_1(u), r'_2(u), \dots, r'_m(u))$, and the Jacobian (7.1) is the norm of the tangent vector:

$$\mathbf{u} \mapsto \|\mathbf{r}'(u)\|. \quad (7.5)$$

- (ii) (Surfaces in \mathbb{R}^3). Let $m = 2$ and let $n = 3$. For any $\mathbf{u} = (u_1, u_2) \in \Gamma^\circ$ the tangent plane is spanned by the two tangent vectors $\mathbf{r}'_{u_1}(\mathbf{u})$ and $\mathbf{r}'_{u_2}(\mathbf{u})$, and the Jacobian (7.1) is the norm of the cross product of the two tangent vectors:

$$\mathbf{u} \mapsto \|\mathbf{r}'_{u_1}(\mathbf{u}) \times \mathbf{r}'_{u_2}(\mathbf{u})\|. \quad (7.6)$$

- (iii) ("Thick" subsets of \mathbb{R}^n). Let $m = n$. The Jacobian in (7.1) is identical to the Jacobian introduced in Remark 6.6.1 on page 157, that is, the Jacobian satisfies

$$\sqrt{\det((\mathbf{J}_r(\mathbf{u}))^T \mathbf{J}_r(\mathbf{u}))} = |\det(\mathbf{J}_r(\mathbf{u}))|.$$

Proof. (i): The matrix $\mathbf{G} = \mathbf{J}^T \mathbf{J} \in \mathbb{R}^{1 \times 1}$ is a 1×1 . By the discussion above, the (1, 1) entry of \mathbf{G} is the inner product $\langle \mathbf{r}'(u), \mathbf{r}'(u) \rangle = \|\mathbf{r}'(u)\|^2$. The determinant of a 1×1 matrix is just the entry at (1, 1) itself. Taking the square root of $\|\mathbf{r}'(u)\|^2$ we end up with

$$\sqrt{\det((\mathbf{J}_r(u))^T \mathbf{J}_r(u))} = \|\mathbf{r}'(u)\|.$$

(ii): The matrix $\mathbf{G} = \mathbf{J}^T \mathbf{J} \in \mathbb{R}^{2 \times 2}$ is in the form of a Gram matrix given by:

$$\mathbf{G} = \begin{bmatrix} \langle \mathbf{r}'_{u_1}, \mathbf{r}'_{u_1} \rangle & \langle \mathbf{r}'_{u_1}, \mathbf{r}'_{u_2} \rangle \\ \langle \mathbf{r}'_{u_2}, \mathbf{r}'_{u_1} \rangle & \langle \mathbf{r}'_{u_2}, \mathbf{r}'_{u_2} \rangle \end{bmatrix},$$

where each tangent vector should be evaluated at $\mathbf{u} = (u_1, u_2)$. The determinant of this two-by-two matrix is easily found as:

$$\det(\mathbf{G}) = \langle \mathbf{r}'_{u_1}, \mathbf{r}'_{u_1} \rangle \langle \mathbf{r}'_{u_2}, \mathbf{r}'_{u_2} \rangle - \langle \mathbf{r}'_{u_1}, \mathbf{r}'_{u_2} \rangle^2.$$

By Lagrange's identity (7.4) we conclude that

$$\det(\mathbf{G}) = \|\mathbf{r}'_{u_1}\|^2 \|\mathbf{r}'_{u_2}\|^2 - \langle \mathbf{r}'_{u_1}, \mathbf{r}'_{u_2} \rangle^2 = \|\mathbf{r}'_{u_1} \times \mathbf{r}'_{u_2}\|^2.$$

(iii): Recall that $\det(AB) = \det(A) \det(B)$ and $\det(A) = \det(A^T)$ for square matrices A and B . For brevity, denote $\mathbf{J} := \mathbf{J}_r(\mathbf{u})$. Then

$$\sqrt{\det(\mathbf{J}^T \mathbf{J})} = \sqrt{\det(\mathbf{J}) \det(\mathbf{J})} = |\det(\mathbf{J})|$$

which was what we wanted to show. ■

Note that in this chapter we are interested in the situation of “thin” subsets where $m < n$, hence we will use Items (i) and (ii) from Lemma 7.1.1 in the next sections, while Item (iii) relates to integration as in Chapter 6.

For surfaces as in Lemma 7.1.1, the vector $\mathbf{n}_{\mathcal{F}}(\mathbf{u}) := \mathbf{r}'_{u_1}(\mathbf{u}) \times \mathbf{r}'_{u_2}(\mathbf{u})$ is called the *normal vector* of the surface $\mathcal{F} = \mathbf{r}(\Gamma)$ at the point $\mathbf{r}(\mathbf{u}) \in \mathbb{R}^3$. The tangent plane at $\mathbf{r}(\mathbf{u}) \in \mathbb{R}^3$ can be expressed as

$$\{\mathbf{x} \in \mathbb{R}^3 \mid \langle \mathbf{x} - \mathbf{r}(\mathbf{u}), \mathbf{n}_{\mathcal{F}}(\mathbf{u}) \rangle = 0\},$$

i.e., the set of all vectors “from $\mathbf{r}(\mathbf{u})$ to \mathbf{x} ” orthogonal to the normal vector.

In the remainder of this chapter, we will assume that parametrizations \mathbf{r} are always *regular* and *injective* on Γ° . However, we will allow \mathbf{r} to be continuous and only *piecewise* C^1 . When $m = 1$, this simply means that we can subdivide $[a, b]$ in k pieces $a = u_0 < u_1 \dots u_k = b$ such that \mathbf{r} is C^1 on each interval $[u_{i-1}, u_i]$. Such a curve $\mathbf{r}([a, b])$ can be considered as a joining of k curves, each being C^1 . For piecewise C^1 curves, the tangent vector is not well-defined at the points u_i . For piecewise C^1 surfaces, the normal vector is not well-defined along curves where two surfaces meet, but this will not be an issue when it comes to defining integration since these problematic subsets are “small”.

7.2 Line and surface integrals

It is possible to introduce integration over curves and surface directly from Riemann sums in much the same manner as we did in [Chapter 6](#). However, we will here simply *define* integration of continuous functions over thin sets (m -folds) using the change of variables technique.

Definition 7.2.1 Integral over m -folds

Let $\mathbf{r} : \Gamma \rightarrow \mathbb{R}^n$, $\Gamma \subset \mathbb{R}^m$, be a regular parametrization that is injective on Γ° . Suppose $f : A \rightarrow \mathbb{R}$, $A \subseteq \mathbb{R}^n$ is a scalar, continuous function on the piecewise smooth m -fold $\mathcal{D} := \mathbf{r}(\Gamma)$. Then the Riemann integral of f over \mathcal{D} is defined by

$$\int_{\mathcal{D}} f(\mathbf{x}) \, dS := \int_{\Gamma} f(\mathbf{r}(\mathbf{u})) \sqrt{\det((\mathbf{J}_{\mathbf{r}}(\mathbf{u}))^T \mathbf{J}_{\mathbf{r}}(\mathbf{u}))} \, d\mathbf{u}, \quad (7.7)$$

where the integral on the right-hand-side of (7.7) was introduced in [Section 6.6](#)

Note that the Jacobian (7.1) appearing in (7.7) is again an m -dimensional volume correcting factor. The argument for why this factor is correct will not be given here, but we mention that the argument is similar to the proof sketch of [Theorem 6.6.3](#) on page 155, e.g., in case of the surface integral, $m = 2$, one needs to show that the *area* of the parallelogram spanned by two linearly independent vectors \mathbf{x} and \mathbf{y} in \mathbb{R}^n on the two-dimensional subspace $\text{span}\{\mathbf{x}, \mathbf{y}\}$ of \mathbb{R}^n is given by $\sqrt{\det(J^T J)}$, where $J = [\mathbf{x} \ \mathbf{y}] \in \mathbb{R}^{n \times 2}$.

Remark 7.2.1 Re-parametrization

The integral introduced in [Definition 7.2.1](#) is independent of the chosen parametrization of the m -fold. That is, the value of the integral does not change under re-parametrization of the m -fold. A re-parametrization of $\mathbf{r} : \Gamma_1 \rightarrow \mathbb{R}^n$ is a mapping of the form $\mathbf{r} \circ \mathbf{p}$, where \mathbf{p} is a *bijective* C^1 vector function $\mathbf{p} : \Gamma_2 \rightarrow \Gamma_1$, whose Jacobian determinant is never zero. Here, as usual, Γ_1 and Γ_2 are connected subsets of \mathbb{R}^m satisfying [Items \(I\) and \(II\)](#) on page 154. The domain of $\mathbf{r} \circ \mathbf{p}$ is Γ_2 , hence we have changed parameter domain from Γ_1 to Γ_2 , while not changing the m -fold itself $\mathbf{r}(\Gamma_1) = \mathbf{r} \circ \mathbf{p}(\Gamma_2)$.

Let us explicitly write out the definition in (7.7) for the cases $m = 1$ and $m = 2$. Thus, we assume the setup used in [Definition 7.2.1](#).

The line integral. Let $m = 1$ and $\Gamma = [a, b]$. The “infinitesimal element” dS for curves is written ds . Let us denote the curve by $\mathcal{C} := \mathbf{r}(\Gamma)$. The

integral is called a *line integral* (also called path or curve integral), and the definition in (7.7) becomes

$$\int_{\mathcal{C}} f(\mathbf{x}) \, ds = \int_a^b f(\mathbf{r}(u)) \|\mathbf{r}'(u)\| \, du. \quad (7.8)$$

The length of the curve \mathcal{C} is found by setting $f(\mathbf{x}) = 1$ for all $\mathbf{x} \in \text{dom}(f)$:

$$\text{length}(\mathcal{C}) = \int_{\mathcal{C}} 1 \, ds = \int_a^b \|\mathbf{r}'(u)\| \, du. \quad (7.9)$$

Example 7.2.1 Arc of semi-circle

Let us return to the curve

$$\mathbf{r} : [0, \pi] \rightarrow \mathbb{R}^2, \quad \mathbf{r}(u) = \begin{bmatrix} \cos(u) \\ \sin(u) \end{bmatrix},$$

considered in Example 3.1.4. According to (7.9) the length of the curve $\mathcal{C} = \mathbf{r}([0, \pi])$ is

$$\text{length}(\mathcal{C}) = \int_0^{\pi} \|\mathbf{r}'(u)\| \, du = \int_0^{\pi} \left\| \begin{bmatrix} -\sin(u) \\ \cos(u) \end{bmatrix} \right\| \, du = \int_0^{\pi} 1 \, du = \pi$$

which is as expected half the circumference of a unit circle.

Example 7.2.2 Length of a helix

Let us return to the curve

$$\mathbf{r}(u) = \begin{bmatrix} \cos(u) \\ \sin(u) \\ u \end{bmatrix} \quad \text{for } u \in [0, 2\pi] \quad (7.10)$$

introduced in Equation (1.7) on page 14. The tangent vector is

$$\mathbf{r}'(u) = \begin{bmatrix} -\sin(u) \\ \cos(u) \\ 1 \end{bmatrix}$$

Hence, the length of the helix $\mathcal{C} = \mathbf{r}([0, 2\pi])$ is

$$\begin{aligned} \text{length}(\mathcal{C}) &= \int_0^{2\pi} \|\mathbf{r}'(u)\| \, du = \int_0^{2\pi} \left((-\sin(u))^2 + (\cos(u))^2 + 1 \right)^{1/2} \, du \\ &= \int_0^{2\pi} (1 + 1)^{1/2} \, du = 2\sqrt{2}\pi \end{aligned}$$

which is a factor $\sqrt{2}$ longer than the unit circle in the plane.

The surface integral. Let $m = 2$, $\Gamma \subset \mathbb{R}^2$, and let us denote the surface by $\mathcal{F} := \mathbf{r}(\Gamma)$, and denote the vectors in Γ by $\mathbf{u} = (u_1, u_2)$. The *surface integral* defined in (7.7) reads explicitly, using (7.6),

$$\int_{\mathcal{F}} f(\mathbf{x}) \, dS = \int_{\Gamma} f(\mathbf{r}(\mathbf{u})) \|\mathbf{r}'_{u_1}(\mathbf{u}) \times \mathbf{r}'_{u_2}(\mathbf{u})\| \, d(u_1, u_2). \quad (7.11)$$

The area of the surface \mathcal{F} is found by setting $f(\mathbf{x}) = 1$ for all $\mathbf{x} \in \text{dom}(f)$:

$$\text{area}(\mathcal{F}) = \int_{\mathcal{F}} 1 \, dS = \int_{\Gamma} \|\mathbf{r}'_{u_1}(\mathbf{u}) \times \mathbf{r}'_{u_2}(\mathbf{u})\| \, d(u_1, u_2). \quad (7.12)$$

Example 7.2.3

A “parameter” subset $\Gamma \subset \mathbb{R}^2$ can itself be identified with a surface \mathcal{F} in \mathbb{R}^3 , simply by taking

$$\mathbf{r}(u, v) = (u, v, 0), \quad (u, v) \in \Gamma.$$

Then, the tangent vectors are

$$\mathbf{r}'_u(u, v) = \frac{\partial \mathbf{r}}{\partial u}(u, v) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{r}'_v(u, v) = \frac{\partial \mathbf{r}}{\partial v}(u, v) = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix},$$

so

$$\frac{\partial \mathbf{r}}{\partial u}(u, v) \times \frac{\partial \mathbf{r}}{\partial v}(u, v) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

Thus (7.12) yields

$$\text{area}(\mathcal{F}) = \int_{\mathcal{F}} 1 \, dS = \int_{\Gamma} \|\mathbf{r}'_u(u, v) \times \mathbf{r}'_v(u, v)\| \, d(u, v) = \int_B 1 \, d(u, v)$$

which coincides with the definition of the area of $\Gamma \subset \mathbb{R}^2$ in Definition 6.3.2 on page 142.

Example 7.2.4 Area of cylinder

With an appropriate positioning of the coordinate system, a hollow cylinder \mathcal{Z} with height h and diameter $2r$ can be described as

$$\mathcal{Z} = \{(x, y, z) \in \mathbb{R}^3 \mid x^2 + y^2 = r^2 \wedge z \in [0, h]\}. \quad (7.13)$$

The surface \mathcal{Z} can easily be parameterized using cylindrical coordinates,

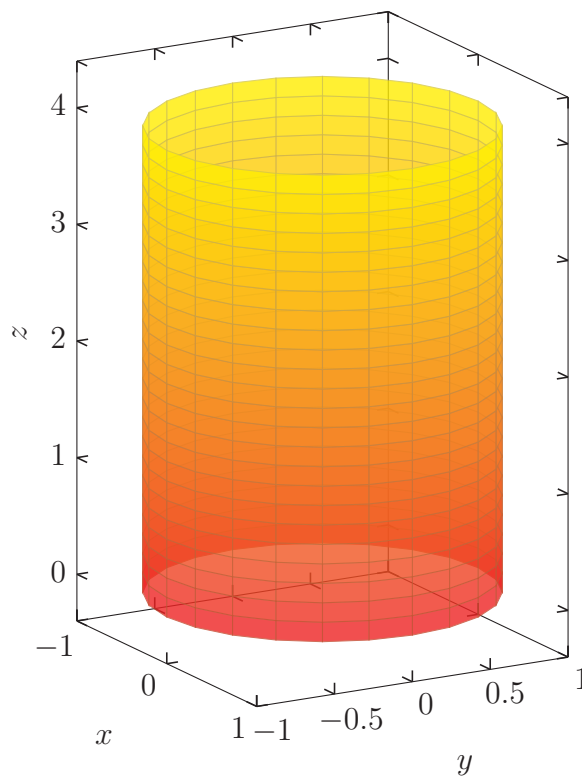


Figure 7.1: The cylinder as in Equation (7.13) with $r = 1$ and $h = 4$.

see Equation (6.77) on page 157, as

$$\mathbf{r}(u, v) = \begin{bmatrix} r \cos(u) \\ r \sin(u) \\ v \end{bmatrix}, \quad (u, v) \in \Gamma = [0, 2\pi] \times [0, h]. \quad (7.14)$$

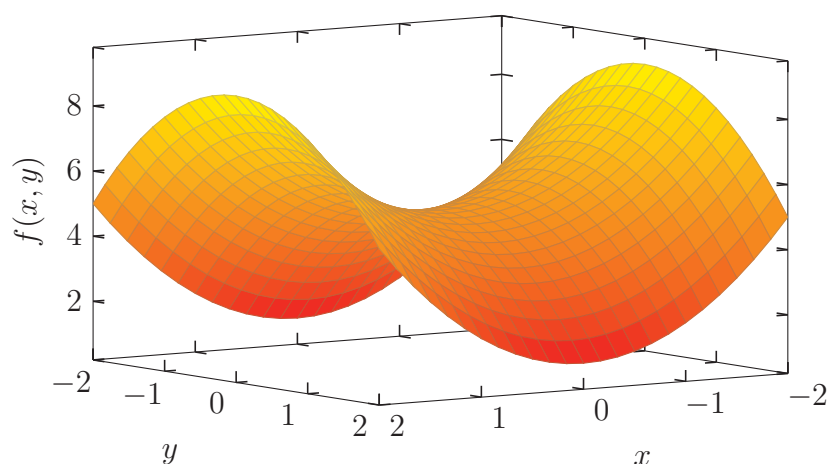
Then

$$\mathbf{r}'_u(u, v) = \frac{\partial \mathbf{r}}{\partial u}(u, v) = \begin{bmatrix} -r \sin(u) \\ r \cos(u) \\ 0 \end{bmatrix} \quad \text{and} \quad \mathbf{r}'_v(u, v) = \frac{\partial \mathbf{r}}{\partial v}(u, v) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

so

$$\frac{\partial \mathbf{r}}{\partial u}(u, v) \times \frac{\partial \mathbf{r}}{\partial v}(u, v) = \begin{bmatrix} -r \sin(u) \\ r \cos(u) \\ 0 \end{bmatrix} \times \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} r \cos(u) \\ r \sin(u) \\ 0 \end{bmatrix}.$$

Note that the vector $\mathbf{n}_{\mathcal{Z}}(u, v) = [r \cos(u), r \sin(u), 0]^T$ is indeed a normal vector to the surface \mathcal{Z} . The Jacobian of the parametrization of \mathcal{Z} is the

Figure 7.2: The graph surface associated with the function h in (7.17).

norm of the normal vector:

$$\|\mathbf{n}_{\mathcal{Z}}(u, v)\| = \|\mathbf{r}'_u(u, v) \times \mathbf{r}'_v(u, v)\| = \sqrt{(r \cos(u))^2 + (r \sin(u))^2} = r,$$

hence, the area-correcting factor does not depend on the parameters u, v , but is constant over the entire surface.

More importantly, the Jacobian $\|\mathbf{n}_{\mathcal{Z}}(u)\|$ is positive, so (7.12) finally yields

$$\begin{aligned} \text{area}(\mathcal{Z}) &= \int_{\mathcal{Z}} 1 \, dS = \int_{\Gamma} \|\mathbf{r}'_u(u, v) \times \mathbf{r}'_v(u, v)\| \, d(u, v) \\ &= \int_{\Gamma} r \, d(u, v) \\ &= \int_0^h \int_0^{2\pi} r \, du \, dv = 2\pi r h. \end{aligned}$$

This coincides with the standard definition of the surface area of a cylinder.

An important class of surfaces occur by considering the graph of a continuous function of two variables.

Definition 7.2.2 Graph surface

Let B be a connected subset of \mathbb{R}^2 which satisfies the conditions Items (I) and (II) on page 140. Given a continuous function

$$h : B \rightarrow \mathbb{R}, \tag{7.15}$$

a surface of the form

$$\mathcal{S} = \{(u, v, h(u, v)) \in \mathbb{R}^3 \mid (u, v) \in B\} \quad (7.16)$$

is called a *graph surface*.

Figure 7.2 illustrates the concept of a graph surface for the function

$$h : [-2, 2] \times [-2, 2] \rightarrow \mathbb{R}, \quad h(x, y) = x^2 - y^2 + 5. \quad (7.17)$$

To compute surface integrals over graph surfaces, we need to specify a Jacobian, cf. Definition 7.1.1 on page 163. So, let us consider a graph surface of $h : B \rightarrow \mathbb{R}$ as in Definition 7.2.2. Let $\Gamma = B$ and define the parametrization

$$\mathbf{r}(u, v) = [u, v, h(u, v)]^T \quad (u, v) \in \Gamma. \quad (7.18)$$

The tangent vectors are

$$\mathbf{r}'_u(u, v) = \begin{bmatrix} 1 \\ 0 \\ \frac{\partial h}{\partial u}(u, v) \end{bmatrix} \quad \text{and} \quad \mathbf{r}'_v(u, v) = \begin{bmatrix} 0 \\ 1 \\ \frac{\partial h}{\partial v}(u, v) \end{bmatrix}.$$

The normal vector $\mathbf{n}_{\mathcal{S}}(u, v)$ is given as the cross product of the two tangent vectors:

$$\mathbf{n}_{\mathcal{S}}(u, v) = \mathbf{r}'_u(u, v) \times \mathbf{r}'_v(u, v) = \begin{bmatrix} -\frac{\partial h}{\partial u}(u, v) \\ -\frac{\partial h}{\partial v}(u, v) \\ 1 \end{bmatrix}$$

Hence, the Jacobian (7.6) is given by

$$(u, v) \mapsto \|\mathbf{n}_{\mathcal{S}}(u, v)\| = \sqrt{\left(\frac{\partial h}{\partial u}(u, v)\right)^2 + \left(\frac{\partial h}{\partial v}(u, v)\right)^2 + 1}.$$

7.3 Vector fields and gradient fields

We already have discussed the notion of vector fields several times before. Recall that this is just vector functions of the form $\mathbf{V} : U \rightarrow \mathbb{R}^n$, where the domain U is an open subset of \mathbb{R}^n . The vector field expressed in terms of its coordinate functions is:

$$\mathbf{V}(x_1, x_2, \dots, x_n) = \begin{bmatrix} V_1(x_1, x_2, \dots, x_n) \\ V_2(x_1, x_2, \dots, x_n) \\ \vdots \\ V_n(x_1, x_2, \dots, x_n) \end{bmatrix}, \quad (x_1, x_2, \dots, x_n) \in U,$$

7.3. Vector fields and gradient fields

where each coordinate function $V_i : U \rightarrow \mathbb{R}$, $i = 1, \dots, n$, is a scalar function of n variables.

Let us give some examples of vector fields.

Example 7.3.1

The function

$$\mathbf{V}(x, y) = \begin{bmatrix} x^2y \\ \cos(x + y) + e^x \end{bmatrix}$$

defines a vector field $\mathbf{V} : U \rightarrow \mathbb{R}^2$ where the domain is $U = \mathbb{R}^2$. Similarly, the function

$$\mathbf{V}(x_1, x_2, x_3) = \begin{bmatrix} x_1x_3 \\ \sin(x_1x_2) \\ (x_1 + x_2 + x_3)^3 \end{bmatrix}$$

defines a vector field $\mathbf{V} : U \rightarrow \mathbb{R}^3$, where $U = \mathbb{R}^3$.

Also, any $n \times n$ matrix with real entries in a natural way corresponds to a vector field $\mathbf{V} : \mathbb{R}^n \rightarrow \mathbb{R}^n$.

Example 7.3.2

Let A denote an arbitrary $n \times n$ matrix with real entries. Then

$$\mathbf{V} : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad \mathbf{V}(\mathbf{x}) := A\mathbf{x} \quad (7.19)$$

defines a so-called *linear vector field*. As a concrete example,

$$\mathbf{V} : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \quad \mathbf{V}(x, y) = \begin{bmatrix} 2x + 4y \\ 3x + 6y \end{bmatrix} \quad (7.20)$$

defines a vector field, and it corresponds to (7.19) with the choice

$$A = \begin{bmatrix} 2 & 4 \\ 3 & 6 \end{bmatrix}.$$

Example 7.3.3

In Example 3.3.4 on page 77 we considered the function

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x, y) = xy^2 + x^2 + y^3,$$

and we computed the gradient vector as

$$\nabla f(x, y) = (y^2 + 2x, 2xy + 3y^2), \quad (x, y) \in \mathbb{R}^2. \quad (7.21)$$

Thus, the gradient ∇f is actually vector field as a function from \mathbb{R}^2 to \mathbb{R}^2 , i.e., $\nabla f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$.

Similarly, in [Example 3.3.5](#), we saw that for the function

$$f : \mathbb{R}^3 \rightarrow \mathbb{R}, \quad f(x_1, x_2, x_3) = x_1x_2 + x_2^2x_3^3,$$

the gradient vector at the point $(x_1, x_2, x_3) \in \mathbb{R}^3$ is

$$\nabla f(x_1, x_2, x_3) = (x_2, x_1 + 2x_2x_3^3, 3x_2^2x_3^2);$$

that is, the gradient is a vector field $\nabla f : \mathbb{R}^3 \rightarrow \mathbb{R}^3$.

The vector fields in the examples above are examples of very nice vector fields in the sense that they are C^∞ vector functions, that is, vector functions whose coordinate functions are infinitely often differentiable in the entire domain of the vector function.

We have already seen that the gradient of a scalar C^1 function $f : U \rightarrow \mathbb{R}$, $U \subseteq \mathbb{R}^n$, gives rise to a vector field $\nabla f : U \rightarrow \mathbb{R}^n$. In [Theorem 7.4.2](#) we will see that vector fields arising from a differentiable function f in this way, have certain particularly pleasant properties. This motivates the following definition.

Definition 7.3.1 Gradient field and anti-derivative

A continuous vector field $\mathbf{V} : U \rightarrow \mathbb{R}^n$, where $U \subseteq \mathbb{R}^n$ is an open set, is called a *gradient field* if there exists a C^1 function $f : U \rightarrow \mathbb{R}$ such that $\mathbf{V} = \nabla f$, i.e.,

$$\mathbf{V}(\mathbf{x}) = (V_1(\mathbf{x}), \dots, V_n(\mathbf{x})) = \left(\frac{\partial f}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}) \right) \quad (7.22)$$

for all $\mathbf{x} \in U$. The function f is called the *anti-derivative* of the vector field \mathbf{V} .

In physics, gradient fields are also known under the name *conservative fields*, and the anti-derivative are known as the (scalar) *potential*. Note that, in the definition of potentials f of vector fields \mathbf{V} , it is not uncommon to use the convention $\mathbf{V} = -\nabla f$ (the anti-derivative is multiplied by -1).

Let $U \subseteq \mathbb{R}^n$ be an open set, and let $\mathbf{V} : U \rightarrow \mathbb{R}^n$ be a C^1 vector function. If \mathbf{V} is a gradient field, i.e., $\mathbf{V} = \nabla f$ for some C^2 function $f : U \rightarrow \mathbb{R}$. Note that we assume \mathbf{V} is C^1 , not only C^0 as in [Definition 7.3.1](#), so we can compute the derivatives of \mathbf{V} . Since $V_i = \frac{\partial f}{\partial x_i}$ the Jacobian matrix of \mathbf{V} is

$$\begin{aligned}
 \mathbf{J}_V(\mathbf{x}) &= \begin{bmatrix} \frac{\partial V_1}{\partial x_1}(\mathbf{x}) & \frac{\partial V_1}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial V_1}{\partial x_n}(\mathbf{x}) \\ \frac{\partial V_2}{\partial x_1}(\mathbf{x}) & \frac{\partial V_2}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial V_2}{\partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial V_n}{\partial x_1}(\mathbf{x}) & \frac{\partial V_n}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial V_n}{\partial x_n}(\mathbf{x}) \end{bmatrix} \\
 &= \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{bmatrix} = \mathbf{H}_f(\mathbf{x}).
 \end{aligned}$$

Hence, for a C^1 gradient field $\mathbf{V} = \nabla f$, its Jacobian matrix is the Hessian matrix of the anti-derivative f , i.e., $\mathbf{J}_V(\mathbf{x}) = \mathbf{H}_f(\mathbf{x})$. Since the Hessian matrix is symmetric for C^2 functions, we conclude that the Jacobian matrix of a gradient field is symmetric. Note that a symmetric Jacobian matrix is a rather restrictive property that even the simplest vector fields need not have.

Exercise 7.3.1

Show that the Jacobian matrix of the linear vector field \mathbf{V} defined in (7.20) is not symmetric. Conclude that the vector field is *not* a gradient field, i.e., conclude that \mathbf{V} has no anti-derivative.

Let us formulate our findings as a lemma:

Lemma 7.3.2

Let $U \subseteq \mathbb{R}^n$ be an open subset, and let $\mathbf{V} : U \rightarrow \mathbb{R}^n$ be a C^1 vector field. If \mathbf{V} is a gradient field, then $\mathbf{J}_V(\mathbf{x})$ is symmetric for all $\mathbf{x} \in U$, i.e.,

$$\frac{\partial V_i}{\partial x_j}(\mathbf{x}) = \frac{\partial V_j}{\partial x_i}(\mathbf{x}) \quad \text{for all } i, j = 1, 2, \dots, n. \quad (7.23)$$

Example 7.3.4

For the vector field

$$\mathbf{V}(x, y) = \begin{bmatrix} x^2 y \\ \cos(x + y) + e^x \end{bmatrix}$$

in Example 7.3.1 we have

$$\frac{\partial V_1}{\partial y}(x, y) = x^2$$

and

$$\frac{\partial V_2}{\partial x}(x, y) = -\sin(x + y) + e^x.$$

Thus, the condition (7.23) is not satisfied, i.e., \mathbf{V} is not a gradient vector field.

We will see in [Theorem 7.4.2](#) on page 181 that the converse of [Lemma 7.3.2](#) is also true if we assume the domain U is simply connected. A subset U is *simple connected* if it is a connected set (recall [Definition 5.2.1](#) on page 121) and if U is “without holes” that goes completely through the set. Intuitively, “without holes” means that any closed curve can be continuously transformed into a single point in U without leaving the subset U . Hence, the sphere in \mathbb{R}^3 is simply connected since any closed curve, i.e., a loop, on the sphere can be shrunk to a point. On the other hand, the “sphere” in \mathbb{R}^2 and the torus in \mathbb{R}^3 are not simply connected since a loop can be “trapped” around a hole.

Definition 7.3.2 Simply connected set

Let $U \subseteq \mathbb{R}^n$ be a connected set. Let \mathbf{x} and \mathbf{y} be any two points in U , and $\mathbf{r}_i : [0, 1] \rightarrow U$, $i = 1, 2$, be two continuous curves from \mathbf{x} to \mathbf{y} in U , that is, $\mathbf{r}_i([0, 1]) \subset U$, $i = 1, 2$, $\mathbf{x} = \mathbf{r}_1(0) = \mathbf{r}_2(0)$ and $\mathbf{y} = \mathbf{r}_1(1) = \mathbf{r}_2(1)$. If for any two such curves, \mathbf{r}_1 can be continuously deformed into \mathbf{r}_2 without leaving U , the set U is *simply connected*. (It is beyond the scope of this text to explain precisely what is meant by “continuously deformed”.)

Let us introduce a large class of subsets, the so-called star-shaped sets, that are straightforward to define and that are simply connected.

Definition 7.3.3 Star-shaped set

A subset $S \subseteq \mathbb{R}^n$ is called a star-shaped set if there exists a point (a “center” of S) $\mathbf{x} \in S$ such that the line segment between \mathbf{x} and any other point in S is contained in S .

Example 7.3.5

Let $n \in \mathbb{N}$, and let us consider subsets of \mathbb{R}^n .

- (a) The sets \mathbb{R}^n itself is star-shaped and therefore also simply connected. More generally, any convex set is star-shaped and therefore also simply connected. Open (and closed) balls $B(\mathbf{x}, r)$ and rectangles in \mathbb{R}^n are examples of convex sets. Convex sets are, in fact, star-shaped where any point in the set can be chosen as the center \mathbf{x} according to [Definition 7.3.3](#).

(b) Let $r_2 > r_1$. The annulus

$$\{\mathbf{x} \in \mathbb{R}^n \mid r_1 \leq \|\mathbf{x}\| \leq r_2\} = \overline{B(\mathbf{0}, r_2)} \setminus B(\mathbf{0}, r_1)$$

is a ring-shaped figure between two concentric balls. It is not star-shaped for any value of $n \in \mathbb{N}$ as it has no “center”. If $n = 1$, the set, being a union of two disjoint intervals $[-r_2, -r_1] \cup [r_1, r_2]$, is not even connected so it cannot be simply connected. If $n = 2$, the annulus is connected, but it is not simply connected as any closed curve looping around the inner circle cannot be contracted to a point. However, if $n \geq 3$, the annulus is simply connected as any closed curve can be continuously transformed into a single point.

(c) Let $\mathbf{x} \in \mathbb{R}^n$ be any point; often we take $\mathbf{x} = \mathbf{0}$. The set $\mathbb{R}^n \setminus \{\mathbf{x}\}$ has the same properties as the annulus when it comes to being star-shaped and simply connected (or not), e.g., for $n = 2$ the set is not simply connected, while for $n \geq 3$ the set is simply connected. The arguments are the same as for the annulus (e.g., by considering $r_1 \rightarrow 0$ and $r_2 \rightarrow \infty$).

We can now state a converse of Lemma 7.3.2. It is beyond the scope of this text to give a proof of the result.

Lemma 7.3.3

Let $U \subseteq \mathbb{R}^n$ be an open and *simply connected* subset, and let $\mathbf{V} : U \rightarrow \mathbb{R}^n$ be a C^1 vector field. If $\mathbf{J}_{\mathbf{V}}(\mathbf{x})$ is symmetric for all $\mathbf{x} \in U$, i.e., if (7.23) holds, then \mathbf{V} is a gradient field.

Example 7.3.6

For the vector field

$$\mathbf{V}(x, y) = \begin{bmatrix} y^2 + 2x \\ 2xy + 3y^2 \end{bmatrix}, \quad (x, y) \in \mathbb{R}^2$$

we have

$$\frac{\partial V_1}{\partial y}(x, y) = 2y$$

and

$$\frac{\partial V_2}{\partial x}(x, y) = 2y.$$

Since \mathbb{R}^2 is a star-shaped set, and therefore simply connected, and since Equation (7.23) on page 175 is satisfied, it follows from Lemma 7.3.3 that

\mathbf{V} is a gradient vector field. This is, of course, what we already knew from Example 7.3.3.

7.4 Line integrals of vector fields and computing anti-derivatives

In (7.8) we defined the line integral of scalar functions of n variables. It is straightforward to extend this integral to a vector-valued integral of vector functions using the same definition as in Equation (6.79) on page 159. However, for vector *fields* $\mathbf{V} : U \rightarrow \mathbb{R}^n$, $U \subseteq \mathbb{R}^n$, it is often more useful to integrate the projection of the vector fields onto the tangent vector of the curve. Let $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ be a regular C^1 parametrization that is injective on $]a, b[$. Recall that the orthogonal projection of the vector $\mathbf{V}(\mathbf{r}(u))$ onto $\mathbf{r}'(u)$ is given by

$$\text{proj}_{\mathbf{r}'(u)}(\mathbf{V}(\mathbf{r}(u))) = \left\langle \mathbf{V}(\mathbf{r}(u)), \frac{\mathbf{r}'(u)}{\|\mathbf{r}'(u)\|} \right\rangle \frac{\mathbf{r}'(u)}{\|\mathbf{r}'(u)\|}$$

where $\frac{\mathbf{r}'(u)}{\|\mathbf{r}'(u)\|}$ is a unit norm tangent vector. The norm of the projection is the absolute value of the inner product of $\mathbf{V}(\mathbf{r}(u))$ and the unit norm tangent vector. If we think of the vector field as a force field, and the curve $\mathcal{C} = \mathbf{r}([a, b])$ the trajectory along which we move, e.g., a particle, the work done by traversing the particle along the curve needs to take sign of the inner product of $\mathbf{V}(\mathbf{r}(u))$ and the unit norm tangent vector into account. Hence, the combined work is calculated as line integral along the curve \mathcal{C} of the *scalar-valued* function

$$\mathbf{u} \mapsto \left\langle \mathbf{V}(\mathbf{r}(u)), \frac{\mathbf{r}'(u)}{\|\mathbf{r}'(u)\|} \right\rangle \quad (7.24)$$

Hence, by Equation (7.8), the integral is:

$$\begin{aligned} \int_{\mathcal{C}} \left\langle \mathbf{V}(\mathbf{r}(u)), \frac{\mathbf{r}'(u)}{\|\mathbf{r}'(u)\|} \right\rangle ds &= \int_a^b \left\langle \mathbf{V}(\mathbf{r}(u)), \frac{\mathbf{r}'(u)}{\|\mathbf{r}'(u)\|} \right\rangle \|\mathbf{r}'(u)\| du \\ &= \int_a^b \langle \mathbf{V}(\mathbf{r}(u)), \mathbf{r}'(u) \rangle du \end{aligned}$$

Let us turn these ideas into a definition.

Definition 7.4.1 Line integral of a vector field

Let $U \subset \mathbb{R}^n$ be an open set, and let $\mathbf{V} : U \rightarrow \mathbb{R}^n$ be a continuous vector field. Let $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ denote a regular, piecewise C^1 parametrization of the curve $\mathcal{C} := \mathbf{r}([a, b])$. The *line integral of the vector field* is defined as

$$\int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s} = \int_a^b \langle \mathbf{V}(\mathbf{r}(u)), \mathbf{r}'(u) \rangle du \quad (7.25)$$

Note that $\int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s}$ is just our notation for the line integral of \mathbf{V} along $\mathcal{C} = \mathbf{r}([a, b])$. In particular, we cannot take the dot product of \mathbf{V} and $d\mathbf{r}$. The notation is used since (7.25) can be written as

$$\int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s} = \int_a^b \mathbf{V}(\mathbf{r}(u)) \cdot \mathbf{r}'(u) du.$$

Example 7.4.1

Consider the vector field

$$\mathbf{V}(x, y) = \begin{bmatrix} x + y \\ xy \end{bmatrix}, \quad (x, y) \in \mathbb{R}^2.$$

(a) The line integral of \mathbf{V} over the curve \mathcal{C}_1 parametrized as

$$\mathbf{r}_1(u) = \begin{bmatrix} u \\ u \end{bmatrix}, \quad u \in [0, 1],$$

is

$$\begin{aligned} \int_{\mathcal{C}_1} \mathbf{V} \cdot d\mathbf{s} &= \int_a^b \mathbf{V}(\mathbf{r}_1(u)) \cdot \mathbf{r}'_1(u) du \\ &= \int_0^1 \begin{bmatrix} 2u \\ u^2 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} du \\ &= \int_0^1 (2u + u^2) du = \frac{4}{3}. \end{aligned}$$

(b) The line integral of \mathbf{V} over the curve \mathcal{C}_2 parametrized as

$$\mathbf{r}_2(u) = \begin{bmatrix} u \\ u^2 \end{bmatrix}, \quad u \in [0, 1],$$

is

$$\int_{\mathcal{C}_2} \mathbf{V} \cdot d\mathbf{s} = \int_a^b \mathbf{V}(\mathbf{r}_2(u)) \cdot \mathbf{r}'_2(u) du$$

7.4. Line integrals of vector fields and computing anti-derivatives

$$\begin{aligned} &= \int_0^1 \begin{bmatrix} u + u^2 \\ u^3 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2u \end{bmatrix} du \\ &= \int_0^1 (u + u^2 + 2u^4) du = \frac{37}{30}. \end{aligned}$$

Note that both the curves \mathcal{C}_1 and \mathcal{C}_2 start at the point $(0, 0)$ and end at the point $(1, 1)$. In particular, we see that, in general, different paths that connect two points also lead to different values of the line integrals.

Given a vector field \mathbf{V} and a parametrization of a curve \mathcal{C} , we can directly calculate the line integral (7.25). However, the calculation can, of course, be cumbersome if the parametrization $\mathbf{r}(t)$ is given by a complicated expression. A significant simplification occurs if the vector field \mathbf{V} is a *gradient field*. Then, it turns out, the line integral of \mathbf{V} over a curve \mathcal{C} only depends on the *starting point* and the *ending point* of the curve, but it is *independent* of the curve itself. This is the content of the following result.

Lemma 7.4.1

Let $U \subset \mathbb{R}^n$ be an open set, and let $\mathbf{V} : U \rightarrow \mathbb{R}^n$ be a continuous vector field. Let $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ denote a regular, piecewise C^1 parametrization of the curve $\mathcal{C} := \mathbf{r}([a, b]) \subset U$. Suppose \mathbf{V} has an anti-derivative $f : U \rightarrow \mathbb{R}$, i.e., suppose \mathbf{V} is a gradient field $\mathbf{V} = \nabla f$. Then

$$f(\mathbf{r}(b)) - f(\mathbf{r}(a)) = \int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s}. \quad (7.26)$$

Proof. Consider a continuously differentiable parametrization $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ of the curve \mathcal{C} . Using $\mathbf{V} = \nabla f$, we can rewrite the integrand of the line integral as

$$\mathbf{V}(\mathbf{r}(u)) \cdot \mathbf{r}'(u) = \nabla f(\mathbf{r}(u)) \cdot \mathbf{r}'(u).$$

Define $h = f \circ \mathbf{r}$, i.e., $h(u) = f(\mathbf{r}(u))$ for $u \in [a, b]$. By Corollary 3.7.2 on page 89, the function h is differentiable, and the derivative of h is given by $h'(u) = \langle \mathbf{r}'(u), \nabla f(\mathbf{r}(u)) \rangle = \nabla f(\mathbf{r}(u)) \cdot \mathbf{r}'(u)$. Therefore, by inserting in the definition of the line integral, we get

$$\begin{aligned} \int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s} &= \int_a^b \nabla f(\mathbf{r}(u)) \cdot \mathbf{r}'(u) du = \int_a^b h'(u) du \\ &= [h(u)]_a^b = h(b) - h(a) = f(\mathbf{r}(b)) - f(\mathbf{r}(a)) \end{aligned}$$

as claimed.

7.4. Line integrals of vector fields and computing anti-derivatives

In case the continuous curve $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ is *piecewise* C^1 , we use the above argument on each C^1 piece $u \in [u_{i-1}, u_i]$. We then piece together each of these finitely many pieces using the insertion rule of integration stated in Equation (6.8) on page 132. This completes the proof. ■

Lemma 7.4.1 explains why gradient fields are also called conservative vector fields. They represent forces of physical systems in which energy is conserved, where the amount of work done as one moves through a path in configuration space is determined solely by the path's starting and ending points. This allows for the definition of potential energy (given by f) that does not depend on the specific route traveled.

We now turn to the question of *computing* the anti-derivatives (i.e., the potential) of a vector field that is *already known* to be a gradient field. The formula for the anti-derivatives is given as a line integral of the vector field:

Theorem 7.4.2 Anti-derivative of a vector field

Let $U \subset \mathbb{R}^n$ be an open and *connected set*, and let $\mathbf{x}_0 \in U$. Let $\mathbf{V} : U \rightarrow \mathbb{R}^n$ be a continuous vector field. Suppose \mathbf{V} has an anti-derivative, i.e., suppose \mathbf{V} is a gradient field. Then the anti-derivative $f : U \rightarrow \mathbb{R}$ that satisfies $f(\mathbf{x}_0) = 0$ is given by

$$f(\mathbf{x}) = \int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s}, \quad \mathbf{x} \in U, \quad (7.27)$$

where $\mathcal{C} \subset U$ is *any* piecewise C^1 curve from \mathbf{x}_0 to \mathbf{x} , that is, there is a piecewise C^1 parametrization $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ such that \mathbf{r} is regular, $\mathbf{x}_0 = \mathbf{r}(a)$, $\mathbf{x} = \mathbf{r}(b)$, and $\mathcal{C} = \mathbf{r}([a, b]) \subset U$.

Moreover, all anti-derivatives are of the form $f(\mathbf{x}) + c$ for some constant $c \in \mathbb{R}$.

The choice $f(\mathbf{x}_0) = 0$ in Theorem 7.4.2 is arbitrary in the sense that if $f : U \rightarrow \mathbb{R}$ is any anti-derivative of $\mathbf{V} : U \rightarrow \mathbb{R}^n$ and if $\mathcal{C} \subset U$ is *any* piecewise C^1 curve from $\mathbf{x}_0 \in U$ to $\mathbf{x} \in U$ that runs entirely inside U , then f satisfies

$$f(\mathbf{x}) - f(\mathbf{x}_0) = \int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s}, \quad \mathbf{x} \in U. \quad (7.28)$$

The results in (7.27) and in (7.28) are just reformulations of the result in Lemma 7.4.1 on the previous page. As a consequence, the integral in (7.27) is (also) independent of the chosen curve as long as the curve is continuous, starts in \mathbf{x}_0 , and ends in \mathbf{x} .

Note that U is assumed to be connected, hence any two points can, by definition, be connected by a continuous curve. However, Theorem 7.4.2

7.4. Line integrals of vector fields and computing anti-derivatives

uses a *piecewise* C^1 continuous curve. It can be shown that, using U is also assumed to be open, that one can always modify the continuous curve so that it becomes piecewise C^1 . It is beyond the scope of this text to include the proof. However, we do prove the final piece of new information in [Theorem 7.4.2](#) the “Moreover”-part: if $f : U \rightarrow \mathbb{R}$ is an anti-derivative of \mathbf{V} , then clearly so is $f + c$, where $c \in \mathbb{R}$ is a constant since

$$\nabla(f + c) = \nabla f + \nabla c = \mathbf{V} + \mathbf{0} = \mathbf{V}.$$

On the other hand, if f_1 and f_2 are two anti-derivatives, then for $g := f_1 - f_2$

$$\nabla g = \nabla f_1 - \nabla f_2 = \mathbf{V} - \mathbf{V} = \mathbf{0}.$$

Now, (7.28) shows that g is constant since

$$g(\mathbf{x}) - g(\mathbf{x}_0) = \int_{\mathcal{C}} \nabla g \cdot d\mathbf{s} = \int_{\mathcal{C}} \mathbf{0} \cdot d\mathbf{s} = 0$$

for all values of $\mathbf{x}, \mathbf{x}_0 \in U$. Thus, $f_1 = f_2 + c$ for some constant $c \in \mathbb{R}$.

In order to apply [Theorem 7.4.2](#), it is an advantage to choose the curve \mathcal{C} as simple as possible. One possibility is to take $\mathbf{x}_0 = \mathbf{0}$ and simply let \mathcal{C} be the straight line from $\mathbf{0}$ to \mathbf{x} . Let us see how this works in practice.

Example 7.4.2

Consider the vector field

$$\mathbf{V}(x, y) = \begin{bmatrix} y^2 + 2x \\ 2xy + 3y^2 \end{bmatrix}, \quad (x, y) \in \mathbb{R}^2.$$

In order to decide whether \mathbf{V} is a gradient vector field, we will now apply [Equation \(7.27\)](#). The simplest curve connecting the starting point $\mathbf{x}_0 = \mathbf{0}$ and the ending point $\mathbf{x} = (x, y)$ is the straight line parametrized by

$$\mathbf{r} : [0, 1] \rightarrow \mathbb{R}^2, \quad \mathbf{r}(u) = u\mathbf{x} = u \begin{bmatrix} x \\ y \end{bmatrix}.$$

With this choice (7.27) yields the function

$$\begin{aligned} f(x, y) &= \int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s} = \int_0^1 \mathbf{V}(\mathbf{r}(u)) \cdot \mathbf{r}'(u) \, du \\ &= \int_0^1 \begin{bmatrix} u^2 y^2 + 2xu \\ 2u^2 xy + 3u^2 y^2 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} \, du \\ &= \int_0^1 (u^2 y^2 x + 2x^2 u + 2u^2 xy^2 + 3u^2 y^3) \, du \end{aligned}$$

7.4. Line integrals of vector fields and computing anti-derivatives

$$= xy^2 + x^2 + y^3.$$

Thus, the function $f(x, y) = xy^2 + x^2 + y^3$ is our *candidate* for an anti-derivative of \mathbf{V} . A fast check, as we already did in [Example 7.3.3](#), confirms that $\mathbf{V} = \nabla f$ indeed is satisfied, i.e., \mathbf{V} is a gradient field, and functions $f(x, y) = xy^2 + x^2 + y^3 + c$, where $c \in \mathbb{R}$, are the anti-derivatives of \mathbf{V} .

For applications of [Theorem 7.4.2](#), it is important to stress that (7.27) defines a function even for vector fields without anti-derivatives. Hence, we should only use the formula (7.27) *after* we have proved existence of an anti-derivative *or* we should check that the computed function f indeed satisfies $\nabla f = \mathbf{V}$. To prove that a vector field has an anti-derivative, that is, to prove the vector field is a gradient field *without* computing the line integral (7.27), we usually apply [Lemma 7.3.3](#) on page 177.

As an alternative, the existence of anti-derivatives is actually characterized in the following result. However, the result is difficult to use in practice as one has to check a condition for infinitely many curves.

Corollary 7.4.3 The circulation theorem

Let $U \subset \mathbb{R}^n$ be an open and *connected set*. Suppose $\mathbf{V} : U \rightarrow \mathbb{R}^n$ is a continuous vector field. Then the following assertions are equivalent:

- (i) The vector field \mathbf{V} is a gradient field, i.e., it has an anti-derivative.
- (ii) For any *closed* and piecewise C^1 curve \mathcal{C} in U it holds

$$\int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s} = 0. \quad (7.29)$$

Proof. Assume (i) holds. Since \mathcal{C} is assumed to be closed, any parametrization $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ of \mathcal{C} satisfies $\mathbf{r}(b) = \mathbf{r}(a)$. Since the anti-derivative f satisfies (7.26), it follows that

$$0 = f(\mathbf{r}(b)) - f(\mathbf{r}(b)) = f(\mathbf{r}(b)) - f(\mathbf{r}(a)) = \int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s},$$

which is the assertion in (ii). We leave the proof of (ii) \Rightarrow (i) to the reader. ■

Example 7.4.3

Consider the curve parametrized as

$$\mathbf{r} : [0, 2\pi] \rightarrow \mathbb{R}^2, \mathbf{r}(u) = \begin{bmatrix} \cos u \\ \sin u \end{bmatrix}.$$

Then $\mathbf{r}(0) = \mathbf{r}(2\pi)$, so the curve is closed. Thus, for any gradient field \mathbf{V} ,

$$\int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s} = 0.$$

For *closed* curves \mathcal{C} the line integral $\int_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s}$ is called the *circulation* of the vector field, and the integral is sometimes denoted by

$$\oint_{\mathcal{C}} \mathbf{V} \cdot d\mathbf{s}$$

to stress the fact that \mathcal{C} is assumed to be closed.

7.5 Surface integrals of vector fields

We introduced the integral of scalar functions over m -folds in \mathbb{R}^n and, in particular, over surfaces in \mathbb{R}^3 in Section 7.2. The scalar function of interest in this section is the projection of a vector field onto the (unit) normal vector of a surface. The integral of this scalar function over the surface is called the *flux* of a vector field. For simplicity, we restrict ourselves to vector fields in \mathbb{R}^3 .

Let $\mathbf{V} : U \rightarrow \mathbb{R}^3$, $U \subseteq \mathbb{R}^3$ be a vector field, and let \mathcal{F} be a piecewise C^1 surface in \mathbb{R}^3 parametrized by $\mathbf{r} : \Gamma \rightarrow \mathbb{R}^3$, where $\Gamma \subset \mathbb{R}^2$. For $(u_1, u_2) \in \Gamma$ the normal vector is $\mathbf{n}_{\mathcal{F}}(u_1, u_2) = \mathbf{r}'_{u_1}(u_1, u_2) \times \mathbf{r}'_{u_2}(u_1, u_2)$. Similar to (7.24), the integrand, i.e., the scalar function, is

$$(u_1, u_2) \mapsto \left\langle \mathbf{V}(\mathbf{r}(u_1, u_2)), \frac{\mathbf{n}_{\mathcal{F}}(u_1, u_2)}{\|\mathbf{n}_{\mathcal{F}}(u_1, u_2)\|} \right\rangle, \quad (7.30)$$

where $\frac{\mathbf{n}_{\mathcal{F}}(u_1, u_2)}{\|\mathbf{n}_{\mathcal{F}}(u_1, u_2)\|}$ is a unit normal vector. Note that, the direction of the normal vector depends on the chosen parametrization, e.g., $-\mathbf{r}'_{u_1}(u_1, u_2) \times \mathbf{r}'_{u_2}(u_1, u_2)$ is also a normal vector to the surface at the point $\mathbf{r}(u_1, u_2) \in \mathcal{F}$ and the integrand will change sign if we use this definition of the normal vector.

By (7.11), the surface integral of the function (7.30) reads:

$$\begin{aligned} \int_{\Gamma} \left\langle \mathbf{V}(\mathbf{r}(u_1, u_2)), \frac{\mathbf{n}_{\mathcal{F}}(u_1, u_2)}{\|\mathbf{n}_{\mathcal{F}}(u_1, u_2)\|} \right\rangle \|\mathbf{n}_{\mathcal{F}}(u_1, u_2)\| d(u_1, u_2) \\ = \int_{\Gamma} \left\langle \mathbf{V}(\mathbf{r}(u_1, u_2)), \mathbf{n}_{\mathcal{F}}(u_1, u_2) \right\rangle d(u_1, u_2). \end{aligned}$$

Let us turn these ideas into a definition.

Definition 7.5.1 Flux: surface integral of a vector field

Let $U \subset \mathbb{R}^3$ be an open set, and let $\mathbf{V} : U \rightarrow \mathbb{R}^3$ be a continuous vector field. Let $\mathbf{r} : \Gamma \rightarrow \mathbb{R}^3$ denote a regular, piecewise C^1 parametrization of the surface $\mathcal{F} := \mathbf{r}(\Gamma)$. The *surface integral of the vector field* is defined as

$$\int_{\mathcal{F}} \mathbf{V} \cdot d\mathbf{S} = \int_{\Gamma} \langle \mathbf{V}(\mathbf{r}(u_1, u_2)), \mathbf{n}_{\mathcal{F}}(u_1, u_2) \rangle d(u_1, u_2). \quad (7.31)$$

Note that $\int_{\mathcal{F}} \mathbf{V} \cdot d\mathbf{S}$ is just our notation for the integral of the vector field \mathbf{V} over the surface \mathcal{F} . The value of this integral is also called the *flux* of the vector field through the surface. Note that the sign of the flux changes if we use $-\mathbf{n}_{\mathcal{F}}(u_1, u_2)$ as normal vector.

Example 7.5.1

Consider the cylinder surface \mathcal{Z} from Example 7.2.4 of radius $r > 0$ and height $h > 0$ parametrized by

$$\mathbf{r}(u, v) = \begin{bmatrix} r \cos(u) \\ r \sin(u) \\ v \end{bmatrix}, \quad (u, v) \in \Gamma = [0, 2\pi] \times [0, h].$$

Let $\mathbf{V} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be given by

$$\mathbf{V}(x, y, z) = \begin{bmatrix} x + 2 \sinh(z) \\ y \\ z(h - z) \end{bmatrix}.$$

We want to compute the flux of \mathbf{V} out through the cylinder, where we by “out” understand the direction away from the z -axis.

We already computed the normal vector in the previous example:

$$\mathbf{n}_{\mathcal{Z}}(u, v) = \mathbf{r}'_u(u, v) \times \mathbf{r}'_v(u, v) = \begin{bmatrix} r \cos(u) \\ r \sin(u) \\ 0 \end{bmatrix},$$

and we note that it points away from the z -axis which is the desired direction.

It is not necessary to compute the Jacobian of the parametrization as the integrand in the flux integral simply is:

$$\langle \mathbf{V}(\mathbf{r}(u, v)), \mathbf{n}_{\mathcal{Z}}(u, v) \rangle = (r \cos(u) + 2 \sinh(v))r \cos(u)$$

$$+ r \sin(u)r \sin(u) + v(h - v) 0 = r^2 + 2r \sinh(v) \cos(u)$$

Hence, the flux is

$$\begin{aligned} \int_{\mathcal{Z}} \mathbf{V} \cdot d\mathbf{S} &= \int_{\Gamma} \langle \mathbf{V}(\mathbf{r}(u, v)), \mathbf{n}_{\mathcal{Z}}(u, v) \rangle d(u, v) \\ &= \int_0^h \int_0^{2\pi} r^2 + 2r \sinh(v) \cos(u) du dv \end{aligned}$$

We split the integral in two integrals and compute first

$$\begin{aligned} \int_0^h \int_0^{2\pi} 2r \sinh(v) \cos(u) du dv &= \int_0^h 2r \sinh(v) [\sin(u)]_0^{2\pi} dv \\ &= \int_0^h 2r \sinh(v) (0 - 0) dv = 0. \end{aligned}$$

We then compute

$$\int_0^h \int_0^{2\pi} r^2 du dv = 2\pi r^2 h$$

Using linearity of the integral, we conclude that the flux is

$$\int_{\mathcal{Z}} \mathbf{V} \cdot d\mathbf{S} = 2\pi r^2 h.$$

Appendices

APPENDIX A

Additional Proofs

A.1 Proof of Taylor's theorem

Before we prove Taylor's theorem we state some results which will be used in the proof. We begin with *Rolle's theorem*:

Theorem A.1.1

Assume that the function $f : [a, b] \rightarrow \mathbb{R}$ is differentiable and that $f(a) = f(b) = 0$. Then there exists a point $\xi \in]a, b[$, for which $f'(\xi) = 0$.

Proof. The function f has a maximal value as well as a minimal value in the interval $[a, b]$. If the maximal value equals the minimal value, then f is constant, and $f'(x) = 0$ for all $x \in]a, b[$. On the other hand, if the minimal value is different from the maximal value, the assumption $f(a) = f(b)$ implies that there exists $\xi \in]a, b[$ where f assumes an extremal value; but this implies that $f'(\xi) = 0$. ■

Lemma A.1.2

Assume that $f : [a, b] \rightarrow \mathbb{R}$ is arbitrarily often differentiable, and let x, x_0 be arbitrary points in $[a, b]$, with $x > x_0$. Assume that for some $N \in \mathbb{N}$,

$$f(x) = 0 \text{ and } f^{(j)}(x_0) = 0 \text{ for all } j \in \{0, 1, \dots, N\}. \quad (\text{A.1})$$

Then there exists $\xi \in]x_0, x[$ such that $f^{(N+1)}(\xi) = 0$.

Proof. The proof consists of a continued application of Rolle's theorem. We first prove the result for $N = 0$, and then move on with $N = 1$, etc. Observe that when we want to prove the theorem for one value of N , the assumption (A.1) is available for all $j \in \{0, 1, \dots, N\}$; this will play a key role in the proof.

A.1. Proof of Taylor's theorem

For $N = 0$, our assumption is that $f(x) = f(x_0) = 0$. We shall prove the existence of a point $\xi_1 \in]x_0, x[$ for which $f'(\xi_1) = 0$; but this is exactly the conclusion in Rolle's theorem.

Now consider $N = 1$; in this case, the assumption (A.1) is available for $j = 0$ and $j = 1$. Via the assumption for $j = 0$ we just proved the existence of $\xi_1 \in]x_0, x[$, for which $f'(\xi_1) = 0$. Now, if we use the assumption for $j = 1$ we also have that $f'(x_0) = 0$. Applying Rolle's theorem again, this time to the function f' and the interval $[x_0, \xi_1]$, we conclude that there exists $\xi_2 \in]x_0, \xi_1[$ such that $f''(\xi_2) = 0$. This is exactly the conclusion we wanted to obtain for $N = 1$.

We now consider the general case, where $N > 1$. Thus, the assumption (A.1) means that

$$f(x) = f(x_0) = f'(x_0) = f''(x_0) = \dots = f^{(N)}(x_0) = 0.$$

Applying Rolle's theorem to the function f and the interval $[x_0, x]$, we can find $\xi_1 \in]x_0, x[$ for which $f'(\xi_1) = 0$. Applying Rolle's theorem to f' and the interval $[x_0, \xi_1]$ leads to a point $\xi_2 \in]x_0, \xi_1[$ for which $f''(\xi_2) = 0$. Applying this procedure $N + 1$ times finally leads to a point $\xi_{N+1} \in]x_0, \xi_N[$ for which $f^{(N+1)}(\xi_{N+1}) = 0$. ■

Now we only need one step before we are ready to prove Taylor's theorem. As before, we let P_N denote the N th Taylor polynomial for f at x_0 .

Lemma A.1.3

Assume that $f : [a, b] \rightarrow \mathbb{R}$ is arbitrarily often differentiable, and let $x_0 \in]a, b[$, $N \in \mathbb{N}$. Then, for each $x \in [a, b]$ we can find ξ between x and x_0 such that

$$f(x) = P_N(x_0) + \frac{f^{(N+1)}(\xi)}{(N+1)!} (x - x_0)^{N+1}.$$

Proof. Fix $x \in [a, b]$. The statement certainly holds if $x = x_0$, so we assume that $x_0 < x$ (the proof is similar if $x_0 > x$). Consider the function ϕ defined via

$$\phi(t) = f(t) - P_N(t) - K(t - x_0)^{N+1}; \tag{A.2}$$

here K is a constant, which we choose such that $\phi(x) = 0$, i.e.,

$$K = \frac{f(x) - P_N(x)}{(x - x_0)^{N+1}}.$$

Differentiating leads to

$$\phi'(t) = f'(t) - P'_N(t) - (N+1)K(t - x_0)^N,$$

and more generally, for each integer $j \in \{1, 2, 3, \dots, N + 1\}$,

$$\begin{aligned}\phi^{(j)}(t) &= f^{(j)}(t) - P_N^{(j)}(t) \\ &\quad - (N + 1)N(N - 1) \cdots (N - j + 2)K(t - x_0)^{N-j+1}.\end{aligned}$$

In particular,

$$\phi^{(N+1)}(t) = f^{(N+1)}(t) - P_N^{(N+1)}(t) - (N + 1)!K. \quad (\text{A.3})$$

By definition,

$$\begin{aligned}P_N(t) &= f(x_0) + \frac{f'(x_0)}{1!}(t - x_0) + \frac{f''(x_0)}{2!}(t - x_0)^2 \\ &\quad + \frac{f^{(3)}(x_0)}{3!}(t - x_0)^3 + \cdots + \frac{f^{(N)}(x_0)}{N!}(t - x_0)^N.\end{aligned}$$

In this expression the first term is a constant, so

$$\begin{aligned}P'_N(t) &= f'(x_0) + \frac{f''(x_0)}{2!}2(t - x_0) \\ &\quad + \frac{f^{(3)}(x_0)}{3!}3(t - x_0)^2 + \cdots + \frac{f^{(N)}(x_0)}{N!}N(t - x_0)^{N-1} \\ &= f'(x_0) + f''(x_0)(t - x_0) \\ &\quad + \frac{f^{(3)}(x_0)}{2!}(t - x_0)^2 + \cdots + \frac{f^{(N)}(x_0)}{(N - 1)!}(t - x_0)^{N-1}.\end{aligned}$$

Differentiating once more,

$$\begin{aligned}P''_N(t) &= f''(x_0) + \frac{2}{2!}f^{(3)}(x_0)(t - x_0) \\ &\quad + \cdots + \frac{(N - 1)}{(N - 1)!}f^{(N)}(x_0)(t - x_0)^{N-2} \\ &= f''(x_0) + f^{(3)}(x_0)(t - x_0) + \cdots + \frac{f^{(N)}(x_0)}{(N - 2)!}(t - x_0)^{N-2}.\end{aligned}$$

More generally, for $j \in \{1, 2, 3, \dots, N\}$,

$$\begin{aligned}P_N^{(j)}(t) &= f^{(j)}(x_0) + f^{(j+1)}(x_0)(t - x_0) \\ &\quad + \frac{f^{(j+2)}(x_0)}{2!}(t - x_0)^2 + \cdots + \frac{f^{(N)}(x_0)}{(N - j)!}(t - x_0)^{N-j}.\end{aligned}$$

Note that

$$P_N^{(N+1)}(t) = 0, \quad \forall t \in \mathbb{R}.$$

A.1. Proof of Taylor's theorem

For $j \in \{0, 1, \dots, N\}$, we have $P_n^{(j)}(x_0) = f^{(j)}(x_0)$, and therefore $\phi^{(j)}(x_0) = 0$. This means that the function ϕ satisfies the conditions in [Lemma A.1.2](#). Thus, there exists $\xi \in]x_0, x[$ for which $\phi^{(N+1)}(\xi) = 0$. But according to [\(A.3\)](#),

$$\phi^{(N+1)}(\xi) = f^{(N+1)}(\xi) - (N+1)!K$$

so

$$f^{(N+1)}(\xi) = (N+1)!K \text{ and } K = \frac{f^{(N+1)}(\xi)}{(N+1)!}.$$

Using the expression [\(A.2\)](#) for ϕ and that $\phi(x) = 0$, we see that

$$f(x) = P_N(x) + \frac{f^{(N+1)}(\xi)}{(N+1)!}(x-x_0)^{N+1}.$$

■

Here is finally the proof for Taylor's theorem, in the case $I = [a, b]$:

Proof of Theorem 4.3.3 on page 103. Let $N \in \mathbb{N}$ and $x \in]a, b[$ be given. According to [Lemma A.1.3](#) we can find $\xi \in]a, b[$ such that

$$f(x) = \sum_{n=0}^N \frac{f^{(n)}(x_0)}{n!}(x-x_0)^n + \frac{f^{(N+1)}(\xi)}{(N+1)!}(x-x_0)^{N+1}.$$

Thus,

$$f(x) - \sum_{n=0}^N \frac{f^{(n)}(x_0)}{n!}(x-x_0)^n = \frac{f^{(N+1)}(\xi)}{(N+1)!}(x-x_0)^{N+1}.$$

By assumption we have that $|f^{(N+1)}(x)| \leq C$ for all $x \in]a, b[$, so this implies that

$$\begin{aligned} \left| f(x) - \sum_{n=0}^N \frac{f^{(n)}(x_0)}{n!}(x-x_0)^n \right| &= \left| \frac{f^{(N+1)}(\xi)}{(N+1)!}(x-x_0)^{N+1} \right| \\ &\leq \frac{C}{(N+1)!} |x-x_0|^{N+1}. \end{aligned}$$

This finally proves [\(4.14\)](#). ■

APPENDIX B

List of symbols

\mathbb{R} : the real numbers

\mathbb{C} : the complex numbers

\mathbb{Z} : the integers, i.e., $0, \pm 1, \pm 2, \dots$

\mathbb{N} : the positive integers, i.e., $1, 2, \dots$

\mathbb{Q} : the fractions

$x \in A$: x is an element in the set A

$A \cup B$: the union of the sets A and B

$A \cap B$: the intersection of the sets A and B

$A \subseteq B$: the set A is a subset of the set B , possible that $A = B$.

$A \subset B$: the set A is a subset of the set B , not possible that $A = B$

\forall : for all

\exists : there exist

$f : A \rightarrow B$: function f with domain A and taking values in B

A : matrix

A^{-1} : inverse of matrix, defined when $\det(A) \neq 0$.

\mathbf{x} : vector

\mathbf{f} : vector function

Some useful formulas:

$$e^{(a+i\omega)t} = e^{at}(\cos \omega t + i \sin \omega t), \quad a, \omega \in \mathbb{R};$$

$$\cos \omega t = \frac{1}{2}(e^{i\omega t} + e^{-i\omega t}), \quad \omega \in \mathbb{R};$$

$$\sin \omega t = \frac{1}{2i}(e^{i\omega t} - e^{-i\omega t}), \quad \omega \in \mathbb{R};$$

$$a + ib = r e^{i\theta}, \quad r = |a + ib|, \quad \theta = \text{Arg}(a + ib),$$

$$a = r \cos \theta, \quad b = r \sin \theta;$$

$$|a + ib| = \sqrt{a^2 + b^2}, \quad a, b \in \mathbb{R} \text{ (Pythagoras);}$$

$$|a + b| \leq |a| + |b|, \quad a, b \in \mathbb{C} \text{ (the triangle inequality);}$$

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

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \text{ when } \det(A) = ad - bc \neq 0$$

$$\det \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{21} \begin{vmatrix} a_{12} & a_{13} \\ a_{32} & a_{33} \end{vmatrix} + a_{31} \begin{vmatrix} a_{12} & a_{13} \\ a_{22} & a_{23} \end{vmatrix}$$

$$\sin(x + y) = \sin x \cos y + \cos x \sin y$$

$$\cos(x + y) = \cos x \cos y - \sin x \sin y$$

Relationship between growth of classical functions:

$$\frac{x^a}{b^x} \rightarrow 0 \text{ for } x \rightarrow \infty \text{ for all } a > 0, b > 1$$

$$\frac{\ln x}{x^a} \rightarrow 0 \text{ for } x \rightarrow \infty \text{ for all } a > 0.$$

Partial integration:

$$\int f(x)g(x) dx = F(x)g(x) - \int F(x)g'(x) dx, \text{ where } F(x) = \int f(x) dx.$$

Integration by substitution:

$$\int f(g(x))g'(x) dx = F(g(x)), \text{ where } F(x) = \int f(x) dx.$$

Some useful integrals:

$$\int x^n dx = \frac{1}{n+1} x^{n+1}, n \neq -1;$$

$$\int \frac{1}{x} dx = \ln(x);$$

$$\int \cos(x) dx = \sin(x);$$

$$\int \sin(x) dx = -\cos(x);$$

$$\int e^x dx = e^x.$$

Bibliography

- [1] Garcia, S. and Horn, R. *A Second Course in Linear Algebra*. Cambridge Mathematical Textbooks. Cambridge University Press, 2017. ISBN: 9781107103818.

Index

- C^0 vector function, 73
- C^1 vector function, 92

- affine map, 13
- anti-derivative, 132
 - vector field, 174
- approximating polynomial, *see*
Taylor polynomial
- area, 142
 - surface, 169
- attains its maximum, 116
- attains its minimum, 116

- ball
 - in V , 36
- basis, 42
 - orthonormal, 43
- bijjective, 4
- boundary
 - of a set, 38
- bounded set, 40

- chain rule
 - generalized, 92
 - simple, 89
- change of variable, 155
- circle, 19
- circulation, 184
- closed set, 38

- closure of set, 39
- codomain, 2
- complement, 37
- connected set, 121
- conservative field, *see* gradient
field
- continuity
 - vector function, 71
- continuous function
 - of one variable, 66
- continuously differentiable
function, 67
- coordinate function, 13
- cross product, 164
- curve
 - parametrization, 69
 - boundary, 39
 - connected set, 121
 - in \mathbb{R}^k , 69
 - regular, 70
 - tangent vector, 70
- curve integral, *see* line integral

- definite integral, 132
- diagonalizable, 53
- differentiability
 - vector function of several
variables, 90
- differentiable function

- of one variable, 67
 - of several variables, 85
- differential, 90, 92
 - scalar function of several variables, 85
- directional derivative, 79
- Discrete Fourier Transform, 52
- domain, 2
- dot product, 24
- double integration, 139

- ellipsoid, 19
- ending point, 180
- extremum, 116

- flux, 185
- Fourier matrix, 52
- function, 2
 - continuous, 71
 - continuously differentiable, 67
 - directional derivative, 79
 - gradient, 77
 - of one variable, 9, 65
 - of several variables, 10
 - partial derivative, 76
 - rectified linear unit, 68
 - ReLU, 68
 - scalar functions of several variables, 14
 - vector field, 14
 - vector functions of several variables, 14

- generalized chain rule, 92
- geometric tensor, 163
- gradient field, 174
- gradient vector, 77
- Gram-Schmidt process, 45
- Gramian, 164
- graph
 - for function of one variable, 16
 - of function of several variables, 16
- graph surface, 172

- Hessian matrix, 81
- hyperbola, 19

- image, 3
- indefinite integral, 132
- injective, 4
- inner product, 26
 - Frobenius, 29
 - on \mathbb{F}^n , 28
 - on \mathbb{C}^n , 26
 - on L^2 , 30
- inner product space, 26, 27
- integration
 - by substitution, 135
 - linearity of, 134
 - partial, 135
 - rules, one variable, 131, 134
- interior of set, 40
- invertible matrix, 5

- Jacobian, 163
- Jacobian determinant, 144
- Jacobian matrix, 90

- length
 - curve, 168
- level sets, 17
- line integral, 168
 - vector field, 179
- linear vector field, 173
- linear map
 - differential, 91
- local minimum, 119
- local maximum, 119, 124
- local minimum, 124

- matrix, 4
 - diagonalizable, 53

- Fourier, 52
- Hermitian, 6, 54
- Hessian matrix, 81
- idempotent, 6
- inverse, 5
- Jacobian, 90
- normal, 6, 54
- orthogonal projection, 6
- positive definite, 6
- positive semi-definite, 6
- real diagonalizable, 53
- real orthogonal, 6
- similar to, 53
- symmetric, 6
- unitarily diagonalizable, 53
- unitary, 6
- maximum, 116
- minimum, 116
- norm, 24, 27
- normal vector, 166
- open set, 37
- orthogonal, 30
- orthogonal projection, 40
- orthogonal projection matrix, 6
- orthonormal, 31
- orthonormal basis, 43
- parametrization
 - curve, 14, 39, 69
 - regular, 164
 - surface, 14
- parametrization of a set B , 147
- parametrized
 - curve, 123
- partial integration, 135
- partial derivative, 76
- polar coordinates, 150
- positive definite, 6
- positive semi-definite, 6
- potential, 174
- quadratic form, 11
 - definition, 11
- quadratic forms
 - symmetric matrix, 60
- quadratic form
 - differential, 87
- range, 3
- real orthogonally diagonalizable, 53
- rectangle
 - in \mathbb{R}^n , 153
- rectangle in \mathbb{R}^2 , 137
- rectangular coordinates, 149
- ReLU, 68
- remainder term, 101
- Riemann integrable, 129, 141
- Riemann sum, 129, 138, 141, 154
- Rolle's theorem, 188
- saddle point, 119, 125
- scalar, 1
- set
 - boundary, 38
 - bounded, 40
 - closed, 38
 - closure, 39
 - complement, 37
 - connected, 121
 - interior, 40
 - open, 37
- similar, 53
- simple connected, 176
- simply connected, 176
- smooth set, 154
- spectral decomposition
 - normal matrix, 62
 - real, symmetric matrix, 59
- Spectral theorem
 - normal matrix, 61
 - real, symmetric matrix, 57

- spectrum, 23
- sphere
in \mathbb{R}^3 , 19
- standard inner product
on \mathbb{R}^n , 24
- standard inner product, *see* inner product
- starting point, 180
- stationary point, 122
- strict local maximum, 124
- strict local minimum, 124
- surface integral, 169
vector field, 185
- surjective, 3
- tangent
curve, 70
scalar function of one variable, 95
scalar function of several variables, 106
vector function of one variable, 70
- tangent plane, *see* tangent
- tangent vector, 70
surface, 164
- Taylor polynomial
one variable, first degree, 95
- Taylor polynomial
notation, 95
one variable, K th degree, 98
several variables, first degree, 106
- several variables, second degree, 108
- vector function of several variables, 113
- Taylor's formula
remainder term, 101
scalar function of one variable, 102
scalar function of several variables, 111
- Taylor's theorem
scalar function of one variable, 103
- the Jacobian, 157
- trace, 7
- unit norm vector, 30
- unit vector, 30
- unitarily diagonalizable, 53
- vector, 161
notation, 13
- vector field, 14, 161
- vector function
 C^0 , 73
 C^1 , 92
continuous, 71
continuously differentiable, 92
differential, 90
- vector function of several variables, 13
- volume
in \mathbb{R}^n , 155