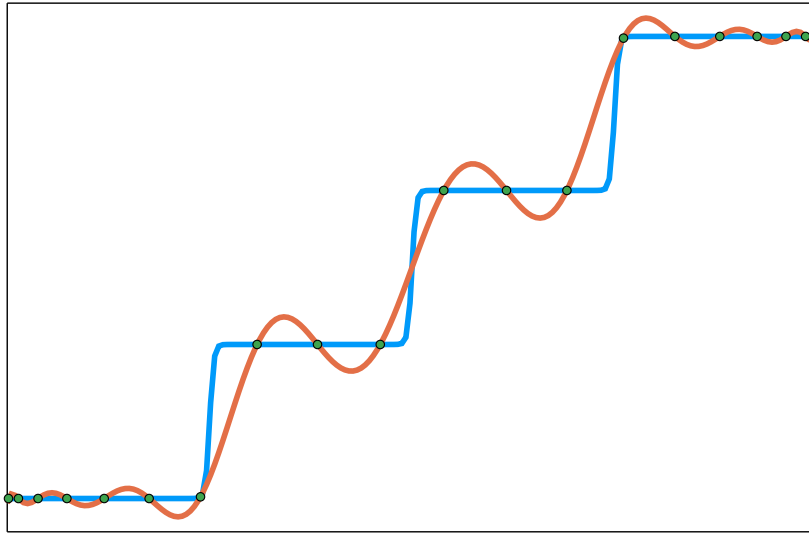


## MATH-UA 9252: Numerical Analysis



Urbain VAES  
urbain.vaes@nyu.edu

NYU PARIS, Fall term 2024

### Weekly schedule:

- Lectures on Monday and Wednesday from 15:00 to 16:15 in room 406;
- Recitation on Monday and Wednesday from 16:20 to 17:05 in room 406;
- Office hour on Monday from 17:15 to 18:15 (Paris time).

# License

The copyright of these notes rests with the author and their contents are made available under a Creative Commons [“Attribution-ShareAlike 4.0 International”](#) license. You are free to copy, distribute, transform and build upon the course material under the following terms:

- **Attribution.** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
- **ShareAlike.** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.



# Course syllabus

**Course content.** This course is aimed at giving a first introduction to classical topics in numerical analysis, including floating point arithmetics and round-off errors, the numerical solution of linear and nonlinear equations, iterative methods for eigenvalue problems, interpolation and approximation of functions, and numerical quadrature. If time permits, we will also cover numerical methods for solving ordinary differential equations.

**Prerequisites.** The course assumes a basic knowledge of linear algebra and calculus. Prior programming experience in Julia, Python or a similar language is desirable but not required.

**Study goals.** After the course, the students will be familiar with the key concepts of *stability*, *convergence* and *computational complexity* in the context of numerical algorithms. They will have gained a broad understanding of the classical numerical methods available for performing fundamental computational tasks, and be able to produce efficient computer implementations of these methods.

**Education method.** The weekly schedule comprises two lectures ( $2 \times 1\text{h}15$  per week) and an exercise session ( $1\text{h}30$  per week). The course material includes rigorous proofs as well as illustrative numerical examples in the Julia programming language, and the weekly exercises blend theoretical questions and practical computer implementation tasks.

**Assessment.** Computational homework will be handed out on a weekly or biweekly basis, each of them focusing on one of the main topics covered in the course. The homework assignments will count towards 70% of the final grade, and the final exam will count towards 30%.

**Literature and study material.** A comprehensive reference for this course is the following textbook: A. QUARTERONI, R. SACCO, and F. SALERI. *Numerical mathematics*, volume **37** of *Texts in Applied Mathematics*. Springer-Verlag, Berlin, 2007. Other pointers to the literature will be given within each chapter.

# Acknowledgments

I am grateful to Vincent Legat, Tony Lelièvre, Gabriel Stoltz and Paul Van Dooren for allowing me to draw inspiration from their lectures notes in numerical analysis.

- Chapter 5 of these notes follows closely the structure of [15, Chapter 3].
- The presentation of the material in Chapter 2 is based on [7].
- Chapter 7 is based on [2, Chapter 2] and [15, Chapter 5].
- Chapter 8 is based on [2, Chapter 3].



I would also like to thank Jean-François Barthélémy and Khôi Nguyễn, who made contributions to several chapters, as well as the students who found several errors and typos in these notes.

# Notations

Unless otherwise specified, we use the following notation throughout these notes.

- Lower case bold is used to denote vectors, e.g.  $\mathbf{x} \in \mathbf{C}^n$ , and upper case sans serif is used to denote matrices, e.g.  $\mathbf{A} \in \mathbf{C}^{m \times n}$ . The entries of a vector  $\mathbf{x} \in \mathbf{C}^n$  are denoted by  $(x_i)$ , and those of a matrix  $\mathbf{A} \in \mathbf{C}^{m \times n}$  are denoted by  $(a_{ij})$  or  $(a_{i,j})$ .
- The notations  $\langle \bullet, \bullet \rangle$  and  $\|\bullet\|$  without a subscript always refer to the Euclidean inner product (A.1) and induced norm.
- The sequence  $x_1, x_2, \dots$  is denoted by  $(x_n)_{n \in \mathbf{N}}$ , or sometimes just  $(x_n)$ .
- The notation  $B_R(\mathbf{y})$  refers to the open ball of radius  $R$  centered at  $\mathbf{y}$ .

# Contents

<b>Notations</b>	<b>iv</b>
<b>1 Floating point arithmetic</b>	<b>4</b>
1.1 Binary representation of real numbers . . . . .	5
1.2 Set of values representable in floating point formats . . . . .	7
1.3 Arithmetic operations between floating point formats . . . . .	10
1.4 Encoding of floating point numbers  . . . . .	13
1.5 Integer formats  . . . . .	15
1.6 Exercises . . . . .	16
1.7 Discussion and bibliography . . . . .	24
<b>2 Interpolation and approximation</b>	<b>25</b>
2.1 Interpolation . . . . .	26
2.2 Approximation . . . . .	40
2.3 Exercises . . . . .	48
2.4 Discussion and bibliography . . . . .	54
<b>3 Numerical integration</b>	<b>55</b>
3.1 The closed Newton–Cotes method . . . . .	56
3.2 Composite methods with equidistant nodes . . . . .	57
3.3 Richardson extrapolation and Romberg’s method . . . . .	63
3.4 Methods with non-equidistant nodes . . . . .	67
3.5 Introduction to probabilistic integration methods . . . . .	71
3.6 Exercises . . . . .	73
3.7 Discussion and bibliography . . . . .	79
<b>4 Solution of linear systems of equation</b>	<b>80</b>
4.1 Conditioning . . . . .	81
4.2 Direct solution method . . . . .	85
4.3 Iterative methods for linear systems . . . . .	97
4.4 Exercises . . . . .	117
4.5 Discussion and bibliography . . . . .	126

<b>5</b>	<b>Solution of nonlinear systems</b>	<b>127</b>
5.1	The bisection method . . . . .	128
5.2	Fixed point methods . . . . .	129
5.3	Convergence of fixed point methods . . . . .	130
5.4	Examples of fixed point methods . . . . .	134
5.5	A numerical experiment . . . . .	142
5.6	Exercises . . . . .	144
5.7	Discussion and bibliography . . . . .	147
<b>6</b>	<b>Numerical computation of eigenvalues</b>	<b>148</b>
6.1	Numerical methods for eigenvalue problems: general remarks . . . . .	149
6.2	Simple vector iterations . . . . .	149
6.3	Methods based on a subspace iteration . . . . .	153
6.4	Projection methods . . . . .	158
6.5	Exercises . . . . .	163
6.6	Discussion and bibliography . . . . .	168
<b>7</b>	<b>Numerical ordinary differential equations</b>	<b>169</b>
7.1	Analysis of the continuous problem . . . . .	170
7.2	One-step methods . . . . .	174
7.3	Multistep methods . . . . .	183
7.4	Absolute stability . . . . .	187
7.5	Exercises . . . . .	193
<b>8</b>	<b>Optimization</b>	<b>194</b>
8.1	Definition and characterization of convexity . . . . .	195
8.2	Unconstrained optimization . . . . .	197
8.3	Constrained optimization . . . . .	199
<b>A</b>	<b>Background material</b>	<b>202</b>
A.1	Inner products and norms . . . . .	202
A.2	Completeness . . . . .	205
A.3	Contraction mappings and the Banach fixed point theorem . . . . .	206
A.4	Vector norms . . . . .	207
A.5	Matrix norms . . . . .	207
A.6	Diagonalization and spectral theorem . . . . .	209
A.7	Similarity transformation and Jordan normal form . . . . .	212
A.8	Oldenburger's theorem and Gelfand's formula . . . . .	213
<b>B</b>	<b>Brief introduction to Julia</b>	<b>215</b>
<b>C</b>	<b>Chebyshev polynomials</b>	<b>223</b>

# Introduction

## Goals of computer simulation

In a wide variety of scientific disciplines, ranging from physics to biology and economics, the phenomena under consideration are well-described by mathematical equations. More often than not, it is too difficult to solve these equations analytically, and so one has to recur to *computer simulation* in order to obtain approximate solutions. Computer simulation enables to gain understanding of the phenomena examined, to explain observations and to make predictions. It plays a crucial role in a number of practical applications including weather forecasting, drug discovery through molecular modeling, flight simulation, and structural engineering, to mention just a few.

Numerical simulation may also be employed in order to calibrate mathematical models of physical phenomena, particularly when observation through experiment is impractical or too costly. For example, it is frequently the case that the parameters in mathematical models for turbulence are estimated not from real data, but from synthetic data generated by computer simulation of the fundamental equations of fluid mechanics. Relying on “computer experiments” is attractive in this context because these enable to perform accurate measurements without disturbing the system being observed. Numerical simulation is also very useful to understand and build simplified models for physical phenomena at very small scales, if direct observation is beyond the capabilities of experimental physics.

## The definition of numerical analysis

Numerical analysis sits at the interface between mathematics and computer science. Nick Trefethen, author of several influential works in mathematics, defines numerical analysis as *the study of algorithms for the problems of continuous mathematics* [14]. Devising and studying algorithms to solve mathematical problems is the central concern of numerical analysis and our main focus in this course. The word *continuous* in the definition is used to indicate that the problems in the realm of numerical analysis involve real or complex variables. Discrete problems, which involve variables that take finitely or countably many values, are usually studied in other fields of mathematics or computer science.



## Sources of error in computational science

It is important for practitioners of computer simulation to be aware of the different sources of error likely to affect numerical results obtained in applications, which may be classified as follows:

- **Modeling error.** There may be a discrepancy between the mathematical model and the underlying physical phenomenon.
- **Data error.** The data of the problem, such as the initial conditions or the parameters entering the equations, are usually known only approximately.
- **Discretization error.** The *discretization* of mathematical equations, i.e. turning them into finite-dimensional problems amenable to computer simulation, adds another source of error.
- **Discrete solver error.** The method employed to solve the discretized equations, especially if it is of iterative nature, may also introduce an error.
- **Round-off errors.** Finally, the limited accuracy of computer arithmetics causes additional errors.

Of these, only the last three are in the domain of numerical analysis, and in this course we focus mainly on the *solver* and *round-off* errors. The order of magnitude of the overall error is dictated by the largest among the above sources of error.

## Aims of this course

The aim of this course is to present the standard numerical methods for performing the tasks most commonly encountered in applications: the solution of linear and nonlinear systems of equations, the solution of eigenvalue problems, interpolation and approximation of functions, and numerical integration. For a given task, there are usually several numerical methods to choose from, and these often include parameters which must be fixed appropriately in order to guarantee a good efficiency. In order to guide these choices, we study carefully the *convergence* and *stability* of the various methods we present. Six topics will be covered in these lecture notes.

- **Floating point arithmetic.** In [Chapter 1](#), we discuss how real numbers are represented, manipulated and stored on a computer. There is an uncountable infinity of real numbers, but only a finite subset of these can be represented exactly on a machine. This subset is specified in the *IEEE 754* standard, which is widely accepted today and employed in most programming languages, including [Julia](#).
- **Interpolation and extrapolation of functions.** In [Chapter 2](#), we focus on the topics of interpolation and approximation. *Interpolation* is concerned with the construction of a function within a given set, for example that of polynomials, that takes given values when evaluated at a discrete set of points. The aim of *approximation*, on the other hand,

is usually to determine, within a class of simple functions, which one is closest to a given function. Depending on the metric employed to measure closeness, this may or may not be a well-defined problem.

- **Numerical integration.** In [Chapter 3](#), we study numerical methods for computing definite integrals. This chapter is strongly related to the previous one, as numerical approximations of the integral of a function are often obtained by first approximating the function, say by a polynomial, and then integrating this approximation exactly.
- **Solution of linear systems.** In [Chapter 4](#), we study the standard numerical methods for solving linear systems. Linear systems are ubiquitous in science, often arising from the discretization of linear elliptic partial differential equations, which themselves govern a large number of physical phenomena including heat propagation, electromagnetism, gravitation and the deformation of solids.
- **Solution of nonlinear equations.** In [Chapter 5](#), we present widely used methods for solving nonlinear equations. Like linear equations, nonlinear equations are omnipresent in science, a prime example being the Navier–Stokes equation describing the motion of fluid flows. Nonlinear equations are usually much more difficult to solve and require dedicated techniques.
- **Solution of eigenvalue problems.** In [Chapter 6](#), we present and study the standard iterative methods for calculating the eigenfunctions and eigenvalues of a matrix. Eigenvalue problems have a large number of applications, for instance in quantum physics and vibration analysis. They are also at the root of the PageRank algorithm for ranking web pages, which played a key role in the early success of Google search.

## Why Julia?

Throughout the course, the `Julia` programming language is employed to exemplify some of the methods and key concepts. In the author’s opinion, the `Julia` language has several advantages compared to other popular languages in the context of scientific computing, such as `Matlab` or `Python`.

- Its main advantage over `Matlab` is that it is free and open source, with the byproduct that it benefits from contributions from a large number of contributors around the world. Additionally, `Julia` is a fully-fledged programming language that can be used for applications unrelated to mathematics.
- Its main advantages over `Python` are significantly better performance and a more concise syntax for mathematical operations, especially those involving vectors and matrices. It should be recognized, however, that although adoption of `Julia` is rapidly increasing, `Python` still enjoys a more mature ecosystem and is much more widely used.

# Chapter 1

## Floating point arithmetic

1.1	Binary representation of real numbers . . . . .	5
1.1.1	Conversion between binary and decimal formats . . . . .	6
1.2	Set of values representable in floating point formats . . . . .	7
1.2.1	Denormalized floating point numbers . . . . .	8
1.2.2	Relative error and machine epsilon . . . . .	8
1.3	Arithmetic operations between floating point formats . . . . .	10
1.4	Encoding of floating point numbers <sup>ⓐ</sup> . . . . .	13
1.5	Integer formats <sup>ⓐ</sup> . . . . .	15
1.6	Exercises . . . . .	16
1.7	Discussion and bibliography . . . . .	24

### Introduction

When we study numerical algorithms in the next chapters, we assume implicitly that the operations involved are performed exactly. On a computer, however, only a subset of the real numbers can be stored and, consequently, many arithmetic operations are performed only approximately. This is the source of the so-called *round-off errors*. The rest of this chapter is organized as follows.

- In [Section 1.1](#), we discuss the binary representation of real numbers.
- In [Section 1.2](#), we describe the set of floating point numbers that can be represented in the usual floating point formats;
- In [Section 1.3](#) we explain how arithmetic operations between floating point numbers behave. We insist in particular on the fact that, in a calculation involving several successive arithmetic operations, the result of each intermediate operation is stored as a floating point number, with a possible error.

- In [Section 1.4](#), we briefly present how floating point numbers are encoded according to the IEEE 754 standard, which is widely adopted today. We discuss also the encoding of special values such as `Inf`, `-Inf` and `NaN`.
- Finally, in [Section 1.5](#), we present the standard integer formats and their encoding.

In order to completely describe computer arithmetic, one would in principle need to also discuss the conversion mechanisms between different number formats, as well as a number of edge cases. A comprehensive discussion of the subject is beyond the scope of this course; our aim in this chapter is only to introduce the key concepts.

## 1.1 Binary representation of real numbers

Given any integer number  $\beta > 0$ , called the *base*, a real number  $x$  can always be expressed as a finite or infinite series of the form

$$x = \pm \sum_{k=-n}^{\infty} a_k \beta^{-k}, \quad a_k \in \{0, \dots, \beta - 1\}. \quad (1.1)$$

The number  $x$  may then be denoted as  $\pm(a_{-n}a_{-n+1} \dots a_{-1}a_0.a_1a_2 \dots)_\beta$ , where the subscript  $\beta$  indicates the base. This numeral system is called the *positional notation* and is universally used today, both by humans (usually with  $\beta = 10$ ) and machines (usually with  $\beta = 2$ ). If the base  $\beta$  is omitted, it is always assumed in this course that  $\beta = 10$  unless otherwise specified – this is the *decimal* representation. The *digits*  $a_{-n}, a_{-n+1}, \dots$  are also called *bits* if  $\beta = 2$ . In computer science, several bases other than 10 are regularly employed, for example the following:

- Base 2 (binary) is the usual choice for storing numbers on a machine. The binary format is convenient because the digits have only two possible values, 0 or 1, and so they can be stored using simple electrical circuits with two states. We employ the binary notation extensively in the rest of this chapter. Notice that, just like multiplying and dividing by 10 is easy in base 10, multiplying and dividing by 2 is very simple in base 2: these operations amount to shifting all the bits by one position to the left or right, respectively.
- Base 16 (hexadecimal) is sometimes convenient to represent numbers in a compact manner. In order to represent the values 0-15 with a single digit, 16 different symbols are required, which are conventionally denoted by  $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F\}$ . With this notation, we have  $(FF)_{16} = (255)_{10}$ , for example.

The hexadecimal notation is often used in programming languages for describing colors specified by a triplet  $(r, g, b)$  of values between 0 and 255, corresponding to the primary colors *red*, *green* and *blue*. The number of possible values for each component is  $256 = 16^2$ , and so only 2 digits are required to represent these in the hexadecimal notation. Hexadecimal numbers are also employed in IPv6 addresses, which are used to identify computers connected to a network.

### 1.1.1 Conversion between binary and decimal formats

Obtaining the decimal representation of a binary number can be achieved from (1.1), using the decimal representations of the powers of 2. Since all the positive and negative powers of 2 have a finite decimal representation, any real number with a finite representation in base 2 has a finite representation also in base 10. For example,  $(0.01)_2 = (0.25)_{10}$  and  $(0.111)_2 = (0.875)_{10}$ .

*Example 1.1* (Converting a binary number to decimal notation). Let us calculate the decimal representation of  $x = (0.\overline{10})_2$ , where the horizontal line indicates repetition:  $x = (0.101010\dots)_2$ . By definition, it holds that

$$x = \sum_{k=0}^{\infty} a_k 2^{-k},$$

where  $a_k = 0$  if  $k$  is even and 1 otherwise. Thus, the series may be rewritten as

$$x = \sum_{k=0}^{\infty} 2^{-(2k+1)} = \frac{1}{2} \sum_{k=0}^{\infty} (2^{-2})^k.$$

We recognize on the right-hand side a geometric series with common ratio  $r = 2^{-2} = \frac{1}{4}$ , and so we obtain

$$x = \frac{1}{2} \left( \frac{1}{1-r} \right) = \frac{2}{3} = (0.\overline{6})_{10}.$$

Obtaining the binary representation of a decimal number is more difficult, because negative powers of 10 have *infinite* binary representations, as [Exercise 1.4](#) demonstrates. There is, however, a simple procedure to perform the conversion, which we present for the specific case of a real number  $x$  with decimal representation of the form  $x = (0.a_1 \dots a_n)_{10}$ . In this setting, the bits  $(b_1, b_2, \dots)$  in the binary representation of  $x = (0.b_1 b_2 b_2 \dots)_2$  may be obtained as follows:

---

**Algorithm 1** Conversion of a number to binary format

---

```

1:  $i \leftarrow 1$ 
2: while  $x \neq 0$  do
3:    $x \leftarrow 2x$ 
4:   if  $x \geq 1$  then
5:      $b_i \leftarrow 1$ 
6:   else
7:      $b_i \leftarrow 0$ 
8:   end if
9:    $x \leftarrow x - b_i$ 
10:   $i \leftarrow i + 1$ 
11: end while

```

---

*Example 1.2* (Converting a decimal number to binary notation). Let us calculate the binary representation of  $x = \frac{1}{3} = (0.\overline{3})_{10}$ . We apply [Algorithm 1](#) and collate the values of  $i$  and  $x$  obtained at the beginning of each iteration, i.e. just before [Line 3](#), in the table below.

$i$	$x$	Result
1	$\frac{1}{3}$	0.0000...
2	$\frac{2}{3}$	0.0100...
3	$\frac{1}{3}$	0.0000...

Since  $x$  in the last row is again  $\frac{1}{3}$ , successive bits alternate between 0 and 1, and the binary representation of  $x$  is given by  $(0.\overline{01})_2$ . This is not surprising since  $2x = (0.66)_{10} = (0.\overline{10})_2$ , as we saw in [Example 1.1](#).

## 1.2 Set of values representable in floating point formats

We mentioned in the introduction that, because of memory limitations, only a subset of the real numbers can be stored exactly in a computer. Nowadays, the vast majority of programming languages comply with the IEEE 754 standard, which requires that the set of representable numbers be of the form

$$\mathbf{F}(p, E_{\min}, E_{\max}) = \left\{ (-1)^s 2^E (b_0.b_1b_2 \dots b_{p-1})_2 : \right. \\ \left. s \in \{0, 1\}, b_i \in \{0, 1\} \text{ and } E_{\min} \leq E \leq E_{\max} \right\}. \quad (1.2)$$

In addition to these, floating number formats provide the special entities `Inf`, `-Inf` and `NaN`, the latter being an abbreviation for *Not a Number*. Three parameters appear in the set definition (1.2). The parameter  $p \in \mathbf{N}_{>0}$  is the number of significant bits (also called the precision), and  $(E_{\min}, E_{\max}) \in \mathbf{Z}^2$  are respectively the minimum and maximum exponents. From the precision, the *machine epsilon* is defined as  $\varepsilon_M = 2^{-(p-1)}$ ; its significance is discussed in [Section 1.2.2](#).

For a number  $x \in \mathbf{F}(p, E_{\min}, E_{\max})$ ,  $s$  is called the *sign*,  $E$  is the *exponent* and  $b_0.b_1b_2 \dots b_{p-1}$  is the *significand*. The latter can be divided into a *leading bit*  $b_0$  and the *fraction*  $b_1b_2 \dots b_{p-1}$ , to the right of the binary point. The most widely used floating point formats are the *single* and *double precision* formats, which are called respectively **Float32** and **Float64** in Julia. Their parameters, together with those of the lesser-known half-precision format, are summarized in [Table 1.1](#). In the rest of this section we use the shorthand notation  $\mathbf{F}_{16}$ ,  $\mathbf{F}_{32}$  and  $\mathbf{F}_{64}$ . Note that  $\mathbf{F}_{16} \subset \mathbf{F}_{32} \subset \mathbf{F}_{64}$ .

	Half precision	Single precision	Double precision
$p$	11	24	53
$E_{\min}$	-14	-126	-1022
$E_{\max}$	15	127	1023

Table 1.1: Floating point formats. The first column corresponds to the *half-precision* format. This format, which is available through the **Float16** type in Julia, is more recent than the single and double precision formats. It was introduced in the 2008 revision to the IEEE 754 standard of 1985, a revision known as IEEE 754-2008.

*Remark 1.1.* Some definitions, notably that in [[10](#), Section 2.5.2], include a general base  $\beta$  instead of the base 2 as an additional parameter in the definition of the number format (1.2).

Since the binary format ( $\beta = 2$ ) is always employed in practice, we focus on this case for simplicity in most of this chapter.

*Remark 1.2.* Given a real number  $x \in \mathbf{F}(p, E_{\min}, E_{\max})$ , the exponent  $E$  and significand are generally not uniquely defined. For example, the number  $2.0 \in \mathbf{F}_{64}$  may be expressed as  $(-1)^0 2^1 (1.00 \dots 00)_2$  or, equivalently, as  $(-1)^0 2^2 (0.100 \dots 00)_2$ .

In **Julia**, non-integer number literals are interpreted as **Float64** by default, which can be verified by using the `typeof` function. For example, the instruction “`a = 0.1`” is equivalent to “`a = Float64(0.1)`”. In order to define a number of type **Float32**, the suffix `f0` must be appended to the decimal expansion. For instance, the instruction “`a = 4.0f0`” defines a floating point number `a` of type **Float32**; it is equivalent to writing “`a = Float32(4.0)`”.

### 1.2.1 Denormalized floating point numbers

We can decompose the set  $\mathbf{F}(p, E_{\min}, E_{\max})$  in two disjoint parts:

$$\begin{aligned} \mathbf{F}(p, E_{\min}, E_{\max}) = & \left\{ (-1)^s 2^E (1.b_1 b_2 \dots b_{p-1})_2 : \right. \\ & \left. s \in \{0, 1\}, b_i \in \{0, 1\} \text{ and } E_{\min} \leq E \leq E_{\max} \right\} \\ & \cup \left\{ (-1)^s 2^{E_{\min}} (0.b_1 b_2 \dots b_{p-1})_2 : s \in \{0, 1\}, b_i \in \{0, 1\} \right\}. \end{aligned}$$

The numbers in the second set are called *subnormal* or *denormalized*.

### 1.2.2 Relative error and machine epsilon

Let  $x$  be a nonzero real number and  $\hat{x}$  be an approximation. We define the absolute and relative errors of the approximation as follows.

**Definition 1.1** (Absolute and relative error). The absolute error is given by  $|x - \hat{x}|$ , whereas the relative error is

$$\frac{|x - \hat{x}|}{|x|}$$

The following result establishes a link between the machine  $\varepsilon_M$  and the relative error between a real number and the closest member of a floating point format.

**Proposition 1.1.** Let  $x_{\min}$  and  $x_{\max}$  denote the smallest and largest non-denormalized positive numbers in a format  $F = \mathbf{F}(p, E_{\min}, E_{\max})$ . If  $x \in [-x_{\max}, -x_{\min}] \cup [x_{\min}, x_{\max}]$ , then

$$\min_{\hat{x} \in F} \frac{|x - \hat{x}|}{|x|} \leq \frac{1}{2} 2^{-(p-1)} = \frac{1}{2} \varepsilon_M. \quad (1.3)$$

*Proof.* For simplicity, we assume that  $x > 0$ . Let  $n = \lfloor \log_2(x) \rfloor$  and  $y := 2^{-n}x$ . Since  $y \in [1, 2)$ , it admits a binary representation of the form  $(1.b_1 b_2 \dots)_2 \neq (1.\bar{1})_2$ . Thus  $x = 2^n (1.b_1 b_2 \dots)_2$ ,

## Absolute spacing between Float64 numbers

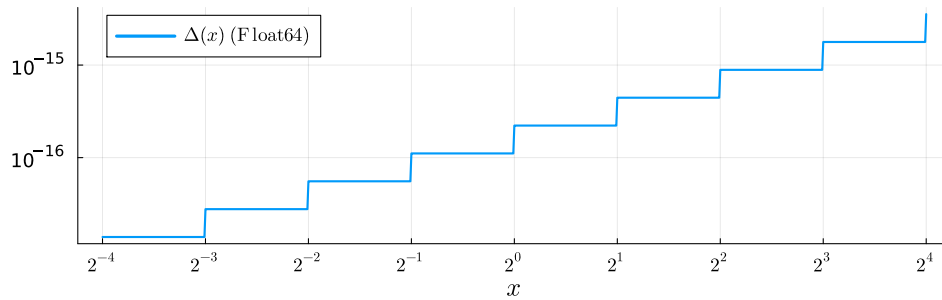


Figure 1.1: Absolute spacing between double-precision floating point numbers, for  $x \in \mathbf{F}_{64}$ . In this figure,  $\Delta(x)$  denotes the distance between  $x$  and its successor in  $\mathbf{F}_{64}$ .

and from the assumption that  $x_{\min} \leq x \leq x_{\max}$  we deduce that  $E_{\min} \leq n \leq E_{\max}$ . We now define the number  $x_- \in F$  by truncating the binary expansion of  $x$  as follows:

$$x_- = 2^n(1.b_1 \dots b_{p-1})_2.$$

The distance between  $x_-$  and its successor in  $F$ , which we denote by  $x_+$ , is given by  $2^{n-p+1}$ . Consequently, it holds that

$$(x_+ - x) + (x - x_-) = x_+ - x_- = 2^{n-p+1}.$$

Since both summands on the left-hand side are positive, this implies that either  $x_+ - x$  or  $x - x_-$  is bounded from above by  $\frac{1}{2}2^{n-p+1} \leq \frac{1}{2}2^{-p+1}x$ , which concludes the proof.  $\square$

The machine epsilon, which was defined as  $\varepsilon_M = 2^{-(p-1)}$ , coincides with the maximum relative spacing between a non-denormalized floating point number  $x$  and its successor in the floating point format, defined as the smallest number in the format that is strictly larger than  $x$ .

Figure 1.1 depicts the density of double-precision floating point numbers, i.e. the number of  $\mathbf{F}_{64}$  members per unit on the real line. The figure shows that the density decreases as the absolute value of  $x$  increases. We also notice that the density is piecewise constant with discontinuities at powers of 2. Figure 1.2 illustrates the relative spacing between successive floating point numbers. Although the absolute spacing increases with the absolute value of  $x$ , the relative spacing oscillates between  $\frac{1}{2}\varepsilon_M$  and  $\varepsilon_M$ .

The picture of the relative spacing between successive floating point numbers looks quite different for denormalized numbers. This is illustrated in Figure 1.3, which shows that the relative spacing increases beyond the machine epsilon in the denormalized range. Fortunately, in the usual  $\mathbf{F}_{32}$  and  $\mathbf{F}_{64}$  formats, the transition between denormalized and non-denormalized numbers occurs at such a small value that it rarely needs worrying about.



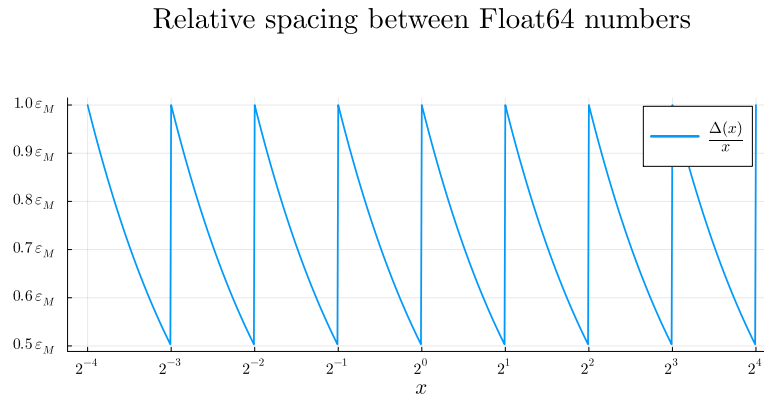


Figure 1.2: Relative spacing between successive double-precision floating point numbers in the “normal range”. The relative spacing oscillates between  $\frac{1}{2}\varepsilon_M$  and  $\varepsilon_M$ .

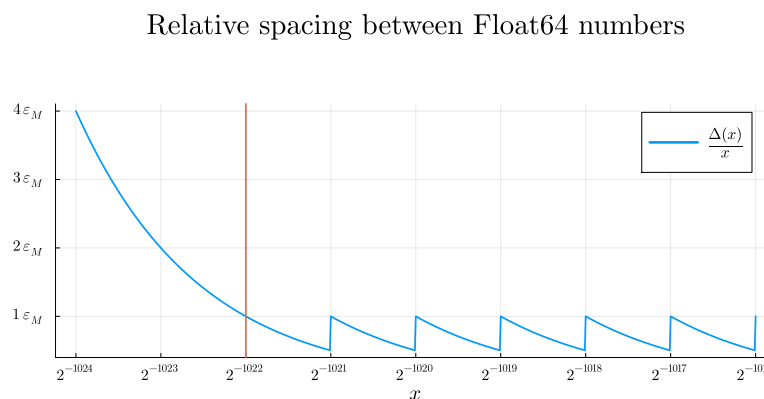


Figure 1.3: Relative spacing between successive double-precision floating point numbers, over a range which includes denormalized number. The vertical red line indicates the transition from denormalized to non-denormalized numbers.

*Example 1.3.* In Julia, the machine epsilon can be obtained using the `eps` function. For example, the instruction `eps(Float16)` returns  $\varepsilon_M$  for the half-precision format.

### 1.3 Arithmetic operations between floating point formats

Now that we have presented the set of values representable on a computer, we attempt in this section to understand precisely how arithmetic operations between floating point formats are performed. The key mechanism governing arithmetic operations on a computer is that of *rounding*, the action of approximating a real number regarded as infinitely precise by a number in a floating point format  $\mathbf{F}(p, E_{\min}, E_{\max})$ . The IEEE 754 standard stipulates that the default mechanism for rounding a real number  $x$ , called *round to nearest*, should behave as follows:

- **Standard case:** The number  $x$  is rounded to the *nearest representable number*, if this number is unique.
- **Edge case:** When there are two equally near representable numbers in the floating point format, the one with the least significant bit equal to zero is delivered.

- **Infinities:** If the real number  $x$  is larger than the largest representable number in the format, that is larger than or equal to  $x_{\max} = 2^{E_{\max}}(2 - \varepsilon)$ , then there are two cases,
  - If  $x < 2^{E_{\max}}(2 - 2^{-p})$ , then  $x_{\max}$  is delivered;
  - Otherwise, the special value **Inf** is delivered.

In other words,  $x_{\max}$  is delivered if it would be delivered by following the rules of the first two bullet points in a different floating point format with the same precision but a larger exponent  $E_{\max}$ . A similar rule applies for large negative numbers.

When a binary arithmetic operation ( $+$ ,  $-$ ,  $\times$ ,  $/$ ) is performed on floating point numbers in format  $\mathbf{F}$ , the result delivered by the computer is obtained by rounding the exact result of the operation according to the rules given above. In other words, the arithmetic operation is performed as if the computer first calculated an intermediate exact result, and then rounded this intermediate result in order to provide a final result in  $\mathbf{F}$ .

Mathematically, arithmetic operations between floating point numbers in a given format  $\mathbf{F}$  may be formalized by introducing the rounding operator  $\text{fl} : \mathbf{R} \rightarrow \mathbf{F}$  and by defining, for any binary operation  $\circ \in \{+, -, \times, /\}$ , the corresponding machine operation

$$\hat{\circ} : \mathbf{F} \times \mathbf{F} \rightarrow \mathbf{F}; (x, y) \mapsto \text{fl}(x \circ y).$$

We defined this operator for arguments in the same floating point format  $\mathbf{F}$ . If the arguments of a binary arithmetic operation are of different types, the format of the end result, known as the *destination format*, depends on that of the arguments: as a rule of thumb, it is given by the most precise among the formats of the arguments. In addition, recall that a floating point literal whose format is not explicitly specified is rounded to double-precision format and so, for example, the addition  $0.1 + 0.1$  produces the result  $\text{fl}_{64}(\text{fl}_{64}(0.1) + \text{fl}_{64}(0.1))$ , where  $\text{fl}_{64}$  is the rounding operator to the double-precision format.

*Example 1.4.* Using the `typeof` function, we check that the floating point literal `1.0` is indeed interpreted as a double-precision number:

```
julia> a = 1.0; typeof(a)
```

```
Float64
```

When two numbers in different floating point formats are passed to a binary operation, the result is in the more precise format.

```
julia> typeof(Float16(1) + Float32(1))
```

```
Float32
```

```
julia> typeof(Float32(1) + Float64(1))
```

```
Float64
```

If a mathematical expression contains several binary arithmetic operations to be performed in succession, the result of each intermediate calculation is stored in a floating point format

dictated by the formats of its argument, and this floating point number is employed in the next binary operation. A consequence of this mechanism is that the machine operands  $\hat{+}$  and  $\hat{*}$  are generally *not associative*. For example, in general

$$(x \hat{+} y) \hat{+} z \neq x \hat{+} (y \hat{+} z)$$

*Example 1.5.* Let  $x = 1$  and  $y = 3 \times 2^{-13}$ . Both of these numbers belong to  $\mathbf{F}_{16}$  and, denoting by  $\hat{+}$  machine addition in  $\mathbf{F}_{16}$ , we have

$$(x \hat{+} y) \hat{+} y = 1 \tag{1.4}$$

but

$$x \hat{+} (y \hat{+} y) = 1 + 2^{-10}. \tag{1.5}$$

To explain this somewhat surprising result, we begin by writing the normalized representations of  $x$  and  $y$  in the  $\mathbf{F}_{16}$  format:

$$\begin{aligned} x &= (-1)^0 \times 2^0 \times (1.0000000000)_2 \\ y &= (-1)^0 \times 2^{-12} \times (1.1000000000)_2. \end{aligned}$$

The exact result of the addition  $x + y$  is given by  $r = 1 + 3 \times 2^{-13}$ , which in binary notation is

$$r = (1.\underbrace{0000000000}_{11 \text{ zeros}}11)_2.$$

Since the length of the significand in the half-precision ( $\mathbf{F}_{16}$ ) format is only  $p = 11$ , this number is not part of  $\mathbf{F}_{16}$ . The result of the machine addition  $\hat{+}$  is therefore obtained by rounding  $r$  to the nearest member of  $\mathbf{F}_{16}$ , which is 1. This reasoning can then be repeated in order to conclude that, indeed,

$$(x \hat{+} y) \hat{+} y = x \hat{+} y = 1.$$

In order to explain the result of (1.5), note that the exact result of the addition  $y + y$  is  $r = 3 \times 2^{-12}$ , which belongs to the floating point format, so it also holds that  $y \hat{+} y = 3 \times 2^{-12}$ . Therefore,

$$x \hat{+} (y \hat{+} y) = 1 \hat{+} 3 \times 2^{-12} = \text{fl}_{16}(1 + 3 \times 2^{-12}).$$

The argument of the  $\mathbf{F}_{16}$  rounding operator does not belong to  $\mathbf{F}_{16}$ , since its binary representation is given by

$$(1.\underbrace{0000000000}_{10 \text{ zeros}}11)_2.$$

This time the nearest member of  $\mathbf{F}_{16}$  is given by  $1 + 2^{-10}$ .

When a numerical computation unexpectedly returns **Inf** or **-Inf**, we say that an *overflow error* occurred. Similarly, *underflow* is said to occur when a number is smaller than the smallest

representable number in a floating point format.

## 1.4 Encoding of floating point numbers

Once a number format is specified through parameters  $(p, E_{\min}, E_{\max})$ , the choice of encoding, i.e. the machine representation of numbers in this format, has no bearing on the magnitude and propagation of round-off errors. Studying encoding is, therefore, not essential for our purposes in this course, but we opted to cover the topic anyway in the hope that it will help the students build intuition on floating point numbers. We focus mainly on the single precision format, but the following discussion applies *mutatis mutandis* to the double and half-precision formats. The material in this section is for information purposes only.

We already mentioned in [Remark 1.2](#) that a number in a floating point format may have several representations. On a computer, however, a floating point number is always stored in the same manner (except for the number 0, see [Remark 1.4](#)). The values of the exponent and significand which are selected by the computer, in the case where there are several possible choices, are determined from the following rules:

- Either  $E > E_{\min}$  and  $b_0 = 1$ ;
- Or  $E = E_{\min}$ , in which case the leading bit may be 0.

The following result proves that these rules enable to define the exponent and significand of a number in a set of floating point numbers uniquely.

**Proposition 1.2.** *Assume that*

$$(-1)^s (2^E b_0 . b_1 \dots b_{p-1})_2 = (-1)^{\tilde{s}} (2^{\tilde{E}} \tilde{b}_0 . \tilde{b}_1 \dots \tilde{b}_{p-1})_2, \quad (1.6)$$

where the parameter sets  $(s, E, b_0, \dots, b_{p-1})$  and  $(\tilde{s}, \tilde{E}, \tilde{b}_0, \dots, \tilde{b}_{p-1})$  both satisfy the above rule. Then  $E = \tilde{E}$  and  $b_i = \tilde{b}_i$  for  $i \in \{0, \dots, p-1\}$ .

*Proof.* We show that  $E = \tilde{E}$ , after which the equality of significands follows trivially. Let us assume for contradiction that  $E > \tilde{E}$  and denote the left and right-hand sides of (1.6) by  $x$  and  $\tilde{x}$ , respectively. Then  $E > E_{\min}$ , implying that  $b_0 = 1$  and so  $2^E \leq |x| < 2^{E+1}$ . On the other hand, it holds that  $|\tilde{x}| < 2^{\tilde{E}+1}$  regardless of whether  $\tilde{E} = E_{\min}$  or not. Since  $E \geq \tilde{E} + 1$  by assumption, we deduce that  $|\tilde{x}| < 2^E \leq |x|$ , which contradicts the equality  $x = \tilde{x}$ .  $\square$

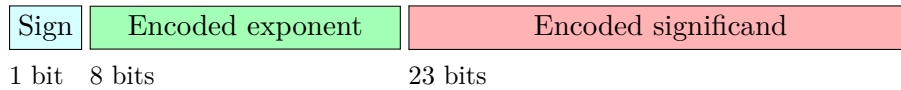
Now that we have explained how a unique set of parameters (sign, exponent, significand) can be assigned to any floating point number, we describe how these parameters are stored on the computer in practice. As their names suggest, the **Float16**, **Float32** and **Float64** formats use 16, 32 and 64 bits of memory, respectively. A naive approach for encoding these number formats would be to store the full binary representations of the sign, exponent and significand.

For the **Float32** format, this approach would require 1 bit for the sign, 8 bits to cover the 254 possible values of the exponent, and 24 bits for the significand, i.e. for storing  $b_0, \dots, b_{p-1}$ . This leads to a total number of 33 bits, which is one more than is available, and this is without

the special values **NaN**, **Inf** and **-Inf**. So how are numbers in the  $\mathbf{F}_{32}$  format actually stored? To answer this question, we begin with two observations:

- If  $E > E_{\min}$ , then necessarily  $b_0 = 1$  in the unique representation of the significand. Consequently, the leading bit need not be explicitly specified in the case; it is said to be *implicit*. We will see, as a consequence, that  $p - 1$  instead of  $p$  bits are in fact sufficient for the significand.
- In the  $\mathbf{F}_{32}$  format, 8 bits at minimum need to be reserved for the exponent, which enables the representation of  $2^8 = 256$  different values, but there are only 254 possible values for the exponent. This suggests that  $256 - 254 = 2$  combinations of the 8 bits can be exploited in order to represent the special values **Inf**, **-Inf** and **NaN**.

Simplifying a little bit, we may view single precision floating point number as an array of 32 bits as illustrated below:



According to the IEEE 754 standard, the first bit is the sign  $s$ , the next 8 bits  $e_0e_1 \dots e_6e_7$  encode the exponent, and the last 23 bits  $b_1b_2 \dots b_{p-2}b_{p-1}$  encode the significand. Let us emphasize that when we say “encode the exponent” here, we just mean that the bits contain information from which the exponent can be uniquely determined, but we have not yet described how this is achieved. Let us introduce the integer number  $e = (e_0e_1 \dots e_6e_7)_2$ ; that is to say,  $0 \leq e \leq 2^8 - 1$  is the integer number whose binary representation is given by  $e_0e_1 \dots e_6e_7$ . One may determine the exponent and significand of a floating point number from the following rules.

- **Denormalized numbers:** If  $e = 0$ , then the implicit leading bit  $b_0$  is zero, the fraction is  $b_1b_2 \dots b_{p-2}b_{p-1}$ , and the exponent is  $E = E_{\min}$ . In other words, using the notation of [Section 1.2](#), we have  $x = (-1)^s 2^{E_{\min}} (0.b_1b_2 \dots b_{p-2}b_{p-1})_2$ . In particular, if  $b_1b_2 \dots b_{p-2}b_{p-1} = 00 \dots 00$ , then it holds that  $x = 0$ .
- **Non-denormalized numbers:** If  $0 < e < 255$ , then the implicit leading bit  $b_0$  of the significand is 1 and the fraction is given by  $b_1b_2 \dots b_{p-2}b_{p-1}$ . The exponent is given by

$$E = e - \text{bias} = E_{\min} + e - 1.$$

where the exponent bias for the single and double precision formats are given in [Table 1.2](#). In this case  $x = (-1)^s 2^{e - \text{bias}} 1.b_1b_2 \dots b_{p-2}b_{p-1}$ . Notice that  $E = E_{\min}$  if  $e = 1$ , as in the case of denormalized numbers.

- **Infinities:** If  $e = 255$  and  $b_1b_2 \dots b_{p-2}b_{p-1} = 00 \dots 00$ , then  $x = \text{Inf}$  if  $s = 0$  and **-Inf** otherwise.
- **Not a Number:** If  $e = 255$  and  $b_1b_2 \dots b_{p-2}b_{p-1} \neq 00 \dots 00$ , then  $x = \text{NaN}$ . Notice that the special value **NaN** can be encoded in many different manners. These extra degrees of

freedom were reserved for passing information on the reason for the occurrence of NaN, which is usually an indication that something has gone wrong in the calculation.

	Half precision	Single precision	Double precision
Exponent bias ( $-E_{\min} + 1$ )	15	127	1023
Exponent encoding (bits)	5	8	11
Significand encoding (bits)	10	23	52

Table 1.2: Encoding parameters for floating point formats

*Remark 1.3* (Encoding efficiency). With 32 bits, at most  $2^{32}$  different numbers could in principle be represented. In practice, as we saw in [Exercise 1.10](#), the **Float32** format enables to represent

$$(E_{\max} - E_{\min})2^p + 2^{p+1} - 1 = 253 \times 2^{23} + 2^{25} - 1 = 2^{32} - 2^{24} - 1 \approx 99.6\% \times 2^{32},$$

different real numbers, indicating a very good encoding efficiency.

*Remark 1.4* (Nonuniqueness of the floating point representation of 0.0). The sign  $s$  is clearly unique for any number in a floating point format, except for 0.0, which could in principle be represented as

$$(-1)^0 2^{E_{\min}} (0.00 \dots 00)_2 \quad \text{or} \quad (-1)^1 2^{E_{\min}} (0.00 \dots 00)_2.$$

In practice, both representations of 0.0 are available on most machines, and these behave slightly differently. For example  $1/(0.0) = \text{Inf}$  but  $1/(-0.0) = -\text{Inf}$ .

## 1.5 Integer formats

The machine representation of integer formats is much simpler than that of floating point numbers. In this short section, we give a few orders of magnitude for common integer formats and briefly discuss overflow issues. Programming languages typically provide integer formats based on 16, 32 and 64 bits. In Julia, these correspond to the types **Int16**, **Int32** and **Int64**, the latter being the default for integer literals.

The most common encoding for integer numbers, which is used in Julia, is known as *two's complement*: a number encoded with  $p$  bits given by  $b_{p-1}b_{p-2} \dots b_0$  corresponds to

$$x = -b_{p-1}2^{p-1} + \sum_{i=0}^{p-2} b_i 2^i.$$

This encoding enables to represent uniquely all the integers from  $N_{\min} = -2^{p-1}$  to  $N_{\max} = 2^{p-1} - 1$ . In contrast with floating point formats, integer formats do not provide special values like **Inf** and **NaN**. The number delivered by the machine when a calculation exceeds the maxi-

num representable value in the format, called the *overflow behavior*, generally depends on the programming language.

Since the overflow behavior of integer formats is not universal across programming languages, a detailed discussion is of little interest. We only mention that Julia uses a *wraparound* behavior, where  $N_{\max} + 1$  silently returns  $N_{\min}$  and, similarly,  $-N_{\min} - 1$  gives  $N_{\max}$ ; the numbers loop back. This can lead to unexpected results, such as  $2^{64}$  evaluating to 0.

## 1.6 Exercises

⚙️ **Exercise 1.1.** Show that if a number  $x \in \mathbf{R}$  admits a finite representation (1.1) in base  $\beta$ , then it also admits an infinite representation in the same base. **Hint:** You may have learned before that  $(0.\overline{9})_{10} = 1$ .

⚙️ **Exercise 1.2.** How many digits does it take to represent all the integers from 0 to  $10^{10} - 1$  in decimal and binary formats? What about the hexadecimal format?

⚙️ **Exercise 1.3.** Find the decimal representation of  $(0.000\overline{1100})_2$ .

⚙️ **Exercise 1.4.** Find the binary representation of  $(0.1)_{10}$ .

□ **Exercise 1.5.** Implement [Algorithm 1](#) on a computer and verify that it works. Your function should take two arguments: an array of integers `[a_1, ..., a_n]` containing the digits after the decimal point and the  $m$  number of bits to return. The bits should be returned as an array `[b_1, ..., b_m]`.

□ **Exercise 1.6.** As mentioned above, [Algorithm 1](#) works only for decimal numbers of the specific form  $x = (0.a_1 \dots a_n)_{10}$ . Find and implement a similar algorithm for integer numbers. More precisely, write a function that takes an integer  $n$  as argument and returns an array containing the bits of the binary expansion  $(b_m \dots b_0)_2$  of  $n$ , from the least significant  $b_0$  to the most significant  $b_m$ . That is to say, your code should return `[b_0, b_1, ...]`.

```
function to_binary(n)
    # Your code comes here
end

# Check that it works
number = 123456789
bits = to_binary(number)
pows2 = 2 .^ range(0, length(bits) - 1)
@assert sum(bits*pows2) == number
```

⚙️ **Exercise 1.7.** Show that the successor of 1 in  $\mathbf{F}_{64}$  is  $1 + \varepsilon_{64}$ , where  $\varepsilon_{64}$  is the machine epsilon for the double-precision format.

⚙️ **Exercise 1.8.** Write down the values of the smallest and largest, in absolute value, positive real numbers representable in the  $\mathbf{F}_{32}$  and  $\mathbf{F}_{64}$  formats.

❄ **Exercise 1.9** (Relative error and machine epsilon). *Prove that the inequality (1.3) is sharp. That is to say, find  $x \in \mathbf{R}$  such that the inequality is an equality.*

❄ **Exercise 1.10** (Cardinality of the set of floating point numbers). *Show that, if  $E_{\max} \geq E_{\min}$ , then  $\mathbf{F}(p, E_{\min}, E_{\max})$  contains exactly*

$$(E_{\max} - E_{\min})2^p + 2^{p+1} - 1$$

distinct real numbers. (In particular, the special values `Inf`, `-Inf` and `NaN` are not counted.) *Hint: Count first the numbers with  $E > E_{\min}$  and then those with  $E = E_{\min}$ .*

❄ **Exercise 1.11.** *Calculate the machine epsilon  $\varepsilon_{16}$  for the  $\mathbf{F}_{16}$  format. Write the results of the arithmetic operations  $1 \hat{+} \varepsilon_{16}$  and  $1 \hat{-} \varepsilon_{16}$  in the form*

$$2^E(1.b_1 \dots b_{p-1})_2.$$

❄ **Exercise 1.12.** *Let  $\varepsilon_{16}$  be the machine epsilon for the  $\mathbf{F}_{16}$  format, and define  $y = \frac{4}{3}\varepsilon_{16}$ . What is the relative error between  $\Delta = (1 + y) - 1$ , and the machine approximation  $\hat{\Delta} = (1 \hat{+} y) \hat{-} 1$ ?*

□ **Exercise 1.13** (Numerical differentiation). *Let  $f(x) = \exp(x)$ . By definition, the derivative of  $f$  at 0 is given by*

$$f'(0) = \lim_{\delta \rightarrow 0} \left( \frac{f(\delta) - f(0)}{\delta} \right).$$

*The expression within brackets on the right-hand side may be used with a small but nonzero  $\delta$  as an approximation for  $f'(0)$ . Implement this approach using double-precision numbers and the same values for  $\delta$  as in the table below. Explain the results you obtain.*

$\delta$	$\frac{\varepsilon_{64}}{4}$	$\frac{\varepsilon_{64}}{2}$	$\varepsilon_{64}$
Approximation of $f'(0)$	0	2	1

□ **Exercise 1.14** (Avoiding overflow). *Write a code to calculate the weighted average*

$$S := \frac{\sum_{j=0}^J w_j j}{\sum_{j=0}^J w_j}, \quad w_j = \exp(j), \quad J = 1000.$$

*You may need to first rewrite  $S$  differently.*

□ **Exercise 1.15.** *Plot the function  $x \mapsto \log(e^{e^x} - 1)$  over the interval  $[0, 10]$ .*

□ **Exercise 1.16** (Calculating the sample variance). *Assume that  $(x_n)_{1 \leq n \leq N}$ , with  $N = 10^6$ , are independent random variables distributed according to the uniform distribution  $\mathcal{U}(L, L + 1)$ . That is, each  $x_n$  takes a random value uniformly distributed between  $L$  and  $L + 1$  where  $L = 10^9$ . In Julia, these samples can be generated with the following lines of code:*

```
N, L = 10^6, 10^9
x = L .+ rand(N)
```



It is well known that the variance of  $x_n \sim \mathcal{U}(L, L + 1)$  is given by  $\sigma^2 = \frac{1}{12}$ . Numerically, the variance can be estimated from the sample variance:

$$s^2 = \frac{1}{N-1} \left( \left( \sum_{n=1}^N x_n^2 \right) - N\bar{x}^2 \right), \quad \bar{x} = \frac{1}{N} \sum_{n=1}^N x_n. \quad (1.7)$$

Write a computer code to calculate  $s^2$  with the best possible accuracy. Can you find a formula that enables better accuracy than (1.7)?

**Remark 1.5.** In order to estimate the true value of  $s^2$  for your samples, you can use the **BigFloat** format, to which the array  $x$  can be converted by using the instruction

```
x = BigFloat.(x)
```

**Exercise 1.17.** Euler proved that

$$\frac{\pi^2}{6} = \lim_{N \rightarrow \infty} \sum_{n=1}^N \frac{1}{n^2}.$$

Using the default **Float64** format, estimate the error obtained when the series on the right-hand side is truncated after  $10^{10}$  terms. Can you rearrange the sum for best accuracy?

**Exercise 1.18.** Let  $x$  and  $y$  be positive real numbers in the interval  $[2^{-10}, 2^{10}]$  (so that we do not need to worry about denormalized numbers, assuming we are working in single or double precision), and let us define the machine addition operator  $\hat{+}$  for arguments in real numbers as

$$\hat{+} : \mathbf{R} \times \mathbf{R} \rightarrow \mathbf{R}; (x, y) \mapsto \text{fl}(\text{fl}(x) + \text{fl}(y)).$$

Prove the following bound on the relative error between the sum  $x + y$  and its machine approximation  $x \hat{+} y$ :

$$\frac{|(x + y) - (x \hat{+} y)|}{|x + y|} \leq \frac{\varepsilon_M}{2} \left( 2 + \frac{\varepsilon_M}{2} \right).$$

**Hint:** decompose the numerator as

$$(x + y) - (x \hat{+} y) = (x - \text{fl}(x)) + (y - \text{fl}(y)) + (\text{fl}(x) + \text{fl}(y) - (x \hat{+} y)),$$

and then use [Proposition 1.1](#).

**Exercise 1.19.** Is `Float32(0.1) * Float32(10) == 1` equal to **true** or **false** given the default rounding rule defined by the IEEE standard? Explain.

**Solution.** By default, real numbers are rounded to the nearest floating point number. This can be checked in Julia with the command `rounding(Float32)`, which prints the default rounding mode. The

exact binary representation of the real number  $x = 0.1$  is

$$\begin{aligned} x &= (0.000\overline{1100})_2 \\ &= 2^{-4} \times \underbrace{(1.1001100110011001100\overline{1100})_2}_{24 \text{ bits}} \end{aligned}$$

The first task is to determine the member of  $\mathbf{F}_{32}$  that is nearest  $x$ . We have

$$\begin{aligned} x^- &:= \max\{y : y \in \mathbf{F}_{32} \text{ and } y \leq x\} = 2^{-4} \times (1.1001100110011001100)_2 \\ x^+ &:= \min\{y : y \in \mathbf{F}_{32} \text{ and } y \geq x\} = 2^{-4} \times (1.1001100110011001101)_2. \end{aligned}$$

Since the number  $(0.\overline{1100})_2$  is closer to 1 than to 0, the number  $x$  is closer to  $x^+$  than to  $x^-$ . Therefore, the number obtained when writing `Float32(0.1)` is  $x^+$ . To conclude the exercise, we need to calculate  $\text{fl}(10 \times x^+)$ , and to this end we first write the exact binary representation of the real number  $10 \times x^+ = (1010)_2 \times x^+$ . We have

$$\begin{aligned} (1010)_2 \times x^+ &= (1000)_2 \times x^+ + (10)_2 \times x^+ = 2^{-4} \times (1100.11001100110011001101)_2 \\ &\quad + 2^{-4} \times (11.0011001100110011001101)_2 \\ &= 2^{-4} \times \underbrace{(10000.00000000000000000001)_2}_{24 \text{ bits}}. \end{aligned}$$

This can be checked in Julia by writing `bitstring(Float32(0.1) * Float64(10.0))`. Clearly, when rounding to the nearest  $\mathbf{F}_{32}$  number, the number  $2^{-4}(10000)_2 = 1$  is obtained; the boolean expression in the question is thus `true`.  $\triangle$

*Remark 1.6.* It should not be inferred from [Exercise 1.19](#) that `Float32(1/i) * i` is always exact in floating point arithmetic. For example `Float32(1/41) * 41` does not evaluate to 1, and neither do `Float16(1/11) * 11` and `Float64(1/49) * 49`.

⚙️ **Exercise 1.20.** Given that the default rounding rule specified by the IEEE 754 standard is “round to nearest, tie to even”, does `Float16(0.1) + Float16(0.2) == Float16(0.3)` evaluate to `true` or `false`. Explain.

⚙️ **Exercise 1.21.** Explain why `Float32(sqrt(2))^2 - 2` is not zero in Julia.

*Solution.* The exact binary representation of  $x := \sqrt{2}$ , found using the `bitstring` in the `Float64` format, is given by

$$x = \underbrace{(1.01101010000010011110011001100\dots)_2}_{24 \text{ bits}}$$

The first task is to determine the member of  $\mathbf{F}_{32}$  that is nearest  $x$ . We have

$$\begin{aligned} x^- &:= \max\{x : x \in \mathbf{F}_{32} \text{ and } x \leq \sqrt{2}\} = \underbrace{(1.01101010000010011110011)_2}_{24 \text{ bits}} \\ x^+ &:= \min\{x : x \in \mathbf{F}_{32} \text{ and } x \geq \sqrt{2}\} = \underbrace{(1.01101010000010011110100)_2}_{24 \text{ bits}}, \end{aligned}$$

and we calculate

$$\begin{aligned}x - x^- &= 2^{-24}(0.01100\dots)_2, \\x^+ - x &= 2^{-21}(1 - (0.11001100\dots)_2) \geq 2^{-21}(1 - (0.11001101)_2) = 2^{-21}(0.00110011)_2.\end{aligned}$$

We deduce that  $x - x^- < x^+ - x$ , and so  $\text{fl}(x) = x^-$ . To conclude the exercise, we need to show that  $\text{fl}((x^-)^2)$  is not equal to 2. The exact binary expansion of  $(x^-)^2$  is

$$(x^-)^2 = \underbrace{(1.11111111111111111111111111111111)}_{24 \text{ bits}}011011001111111010101001_2.$$

The member of  $\mathbf{F}_{32}$  nearest this number is

$$(1.111111111111111111111111111111)_2 = 2 - 2^{-23},$$

which is precisely the result returned by Julia. △

□ **Exercise 1.22** (Numerical differentiation). Let  $f(x) = \exp(x)$  and let  $d(\delta)$  be the approximation of  $f'(x)$  obtained from the following piece of code:

```
f, x = exp, 1
d(δ) = (f(x+δ) - f(x))/δ
```

Plot for fixed  $x = 1$  the error  $\text{abs}(d(\delta) - \exp(x))$  as a function of  $\delta$  in logarithmic scale, and explain the result.

*Solution.* We assume that  $\delta \in \mathbf{F}_{64}$  for simplicity. This is not a restrictive assumption as  $\delta$  can only take floating point values in computer programs. The proof is rather technical, and so it is given for information purposes only. We rewrite  $d(\delta)$  mathematically as

$$d(\delta) = \left( \widehat{f}(x + \delta) - \widehat{f}(x) \right) \widehat{\delta},$$

where  $\widehat{f}(x) = \text{fl}(f(x))$ . We wish to bound  $|f'(x) - d(\delta)|$ . By the triangle inequality,

$$|f'(x) - d(\delta)| \leq \left| f'(x) - \frac{f(x + \delta) - f(x)}{\delta} \right| + \left| \frac{f(x + \delta) - f(x)}{\delta} - \left( \widehat{f}(x + \delta) - \widehat{f}(x) \right) \widehat{\delta} \right|. \quad (1.8)$$

**Bounding the first term in (1.8).** By Taylor's theorem, there exists  $\xi \in [x, x + \delta]$  such that

$$f(x + \delta) = f(x) + \delta f'(x) + \frac{\delta^2}{2} f''(\xi).$$

Therefore, the first term in (1.8) satisfies

$$\left| f'(x) - \frac{f(x + \delta) - f(x)}{\delta} \right| = \frac{\delta}{2} |f''(\xi)| = \frac{\delta}{2} |f''(x)| + \mathcal{O}(\delta^2).$$

**Bounding the second term in (1.8) – the roundoff error.** This is more tedious but not difficult; the main ingredient is the triangle inequality. Specifically, we will use the bound

$$\begin{aligned}
 & \left| \frac{f(x+\delta) - f(x)}{\delta} - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \widehat{\nearrow} \delta \right| \\
 & \leq \delta^{-1} |f(x+\delta) - f(x \widehat{+} \delta)| + \delta^{-1} |f(x \widehat{+} \delta) - \widehat{f}(x \widehat{+} \delta)| + \delta^{-1} |\widehat{f}(x) - f(x)| \\
 & \quad + \delta^{-1} \left| \left( \widehat{f}(x \widehat{+} \delta) - \widehat{f}(x) \right) - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \right| \\
 & \quad + \left| \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) / \delta - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \widehat{\nearrow} \delta \right|. \tag{1.9}
 \end{aligned}$$

Note that without absolute values, both sides are indeed equal. The first three terms on the right-hand side are together a bound from above for

$$\left| \frac{f(x+\delta) - f(x)}{\delta} - \frac{\widehat{f}(x \widehat{+} \delta) - \widehat{f}(x)}{\delta} \right|, \tag{1.10}$$

while the two other terms account for the roundoff errors associated with the machine subtraction and division operators, respectively. We will show that the dominant terms in (1.9) are the first three; the latter two are negligible in comparison. Employing [Proposition 1.1](#) with  $x = f(a)$ , with  $x = a \pm b$  and with  $x = a/b$ , we deduce that the following inequalities are satisfied:

$$\left| \widehat{f}(a) - f(a) \right| \leq \varepsilon |f(a)|, \quad |(a \pm b) - (a \widehat{\pm} b)| \leq \varepsilon |a \pm b|, \quad |a/b - a \widehat{\nearrow} b| \leq \varepsilon |a/b|. \tag{1.11}$$

The first two inequalities are valid for all  $(a, b) \in \mathbf{F}_{64} \times \mathbf{F}_{64}$ , while the third inequality is valid for all  $(a, b) \in \mathbf{F}_{64} \times \mathbf{F}_{64} \setminus \{0\}$ . The first inequality in (1.11) can be employed in order to bound the second and third term on the right-hand side of (1.9):

$$\begin{aligned}
 \delta^{-1} |f(x \widehat{+} \delta) - \widehat{f}(x \widehat{+} \delta)| + \delta^{-1} |\widehat{f}(x) - f(x)| & \leq \delta^{-1} \varepsilon (|f(x+\delta)| + |f(x)|) \\
 & = 2\delta^{-1} \varepsilon |f(x)| + \mathcal{O}(\varepsilon). \tag{1.12}
 \end{aligned}$$

Using Taylor's theorem and then the second inequality in (1.11) (with  $+$ ), we then bound the first term on right-hand side of (1.9) as

$$\begin{aligned}
 \delta^{-1} |f(x+\delta) - f(x \widehat{+} \delta)| & = \delta^{-1} |f'(\xi)((x+\delta) - (x \widehat{+} \delta))| \\
 & \leq \delta^{-1} \varepsilon |f'(\xi)| |x+\delta| = \delta^{-1} \varepsilon |f'(x)| |x| + \mathcal{O}(\varepsilon). \tag{1.13}
 \end{aligned}$$

Combining (1.12) and (1.13), we obtain that the expression in (1.10) scales as  $\delta^{-1} \varepsilon$ :

$$\left| \frac{f(x+\delta) - f(x)}{\delta} - \frac{\widehat{f}(x \widehat{+} \delta) - \widehat{f}(x)}{\delta} \right| = \mathcal{O}(\delta^{-1} \varepsilon). \tag{1.14}$$

Next, using the second inequality in (1.11) (with  $-$ ) and then (1.14), we bound the fourth term in (1.9) as follows:

$$\begin{aligned}
 \delta^{-1} \left| \left( \widehat{f}(x \widehat{+} \delta) - \widehat{f}(x) \right) - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \right| \\
 \leq \delta^{-1} \varepsilon |\widehat{f}(x \widehat{+} \delta) - \widehat{f}(x)| = \delta^{-1} \varepsilon \underbrace{|f(x+\delta) - f(x)|}_{=\mathcal{O}(\delta)} + \mathcal{O}(\delta^{-1} \varepsilon^2) = \mathcal{O}(\varepsilon + \delta^{-1} \varepsilon^2). \tag{1.15}
 \end{aligned}$$

Notice that this term is smaller than the first three in (1.9) when  $\delta \ll 1$ . We deduce by a triangle inequality from (1.14) and (1.15) that

$$\delta^{-1} \left| (f(x + \delta) - f(x)) - (\widehat{f}(x + \delta) \widehat{-} \widehat{f}(x)) \right| = \mathcal{O}(\delta^{-1}\varepsilon).$$

Finally, using the third inequality in (1.11) together with this equation, we bound the fifth term on the right-hand side of (1.9):

$$\begin{aligned} & \left| (\widehat{f}(x + \delta) \widehat{-} \widehat{f}(x)) / \delta - (\widehat{f}(x + \delta) \widehat{-} \widehat{f}(x)) \widehat{/} \delta \right| \\ & \leq \varepsilon |\widehat{f}(x + \delta) \widehat{-} \widehat{f}(x)| / \delta = \frac{\varepsilon}{\delta} \underbrace{|f(x + \delta) - f(x)|}_{\mathcal{O}(\delta)} + \mathcal{O}(\delta^{-1}\varepsilon^2) = \mathcal{O}(\varepsilon + \delta^{-1}\varepsilon^2). \end{aligned}$$

This term is also negligible in front of the other dominant contributions to the roundoff error given in (1.12) and (1.13). Going back to (1.9), we conclude that

$$\left| \frac{f(x + \delta) - f(x)}{\delta} - (\widehat{f}(x + \delta) \widehat{-} \widehat{f}(x)) \widehat{/} \delta \right| \leq \delta^{-1}\varepsilon (2|f(x)| + |f'(x)x|) + \mathcal{O}(\varepsilon + \delta^{-1}\varepsilon^2).$$

**Concluding the proof.** Going back to (1.8), we conclude that

$$|f'(x) - d(\delta)| \leq \frac{\delta}{2}|f''(x)| + \frac{\varepsilon}{\delta} (2|f(x)| + |f'(x)x|),$$

up to higher order terms. For fixed  $x$ , the right-hand side is minimized when

$$\delta = \sqrt{\varepsilon} \sqrt{\frac{2|f(x)| + |f'(x)x|}{|f''(x)|}},$$

which is a well-known formula for the optimal step size in numerical differentiation. The error as a function of  $\delta$  for  $x = 1$  is depicted in Figure 1.4. △

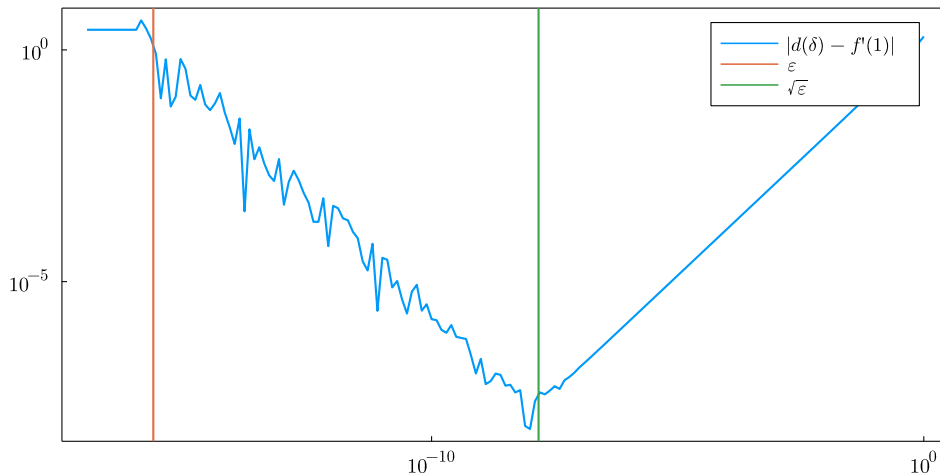


Figure 1.4: Solution to Exercise 1.22.

⚙️ **Exercise 1.23.** Explain why `exp(log(Float16(7))) == 7` is **false**.

*Solution.* We begin by finding the binary representation of  $\log(7)$ . In the  $\mathbf{F}_{64}$  format, the sign and exponent are encoded over 1 and 11 bits respectively, and so the fraction is given by the last 52 bits, which can be obtained from the command

```
julia> bitstring(log(7))[13:end]
"1111001000100111001010101110001100100101101001010111"
```

Therefore, since `exponent(log(7))` returns 0, we have

$$\log(7) = (\underbrace{1.1111001000}_{11 \text{ bits}} 1001100101010111000110010010110100101011\dots)_2 \quad (1.16)$$

The number returned by the command `log(Float16(7))` is given by  $\text{fl}_{16}(\log(7))$ , where  $\text{fl}_{16}$  denotes the half-precision rounding operator. Rounding the right-hand side of (1.16) to 11 bits, we obtain

$$\text{fl}_{16}(\log(7)) = (1.1111001001)_2 = 1.9462890625.$$

The number returned by the code `exp(log(Float16(7)))` is

$$\text{fl}_{16}\left(\exp\left(\text{fl}_{16}(\log(7))\right)\right) = \text{fl}_{16}(\exp(1.9462890625)).$$

The rounding operator appears twice on the left-hand side, because the computer rounds *after every operation*. To explain the result of the rounding operation, we begin by calculating the binary expansion of  $\exp(1.9462890625)$ .

```
julia> bitstring(exp(1.9462890625))[13:end]
"11000000001010110111011100001110001000011000100011110"
```

```
julia> exponent(exp(1.9462890625))
2
```

Therefore,

$$\exp(1.9462890625) = 2^2(\underbrace{1.1100000000}_{11 \text{ bits}} 1010110111011100001110001000011000100011\dots)_2,$$

and, rounding to 11 bits, we finally obtain

$$\begin{aligned} \text{fl}_{16}(\exp(1.9462890625)) &= 2^2(1.110000000\mathbf{1})_2 = 4\left(1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{2^{10}}\right) \\ &= 7 + \frac{1}{2^8} = 7.00390625, \end{aligned}$$

which is different from 7. △

⚙️ **Exercise 1.24.** Determine the encoding of the following **Float32** numbers:

- $x_1 = 2.0^{E_{\min}}$
- $x_2 = -2.0^{E_{\min}-p-1} = -2.0^{-149}$
- $x_3 = 2.0^{E_{\max}}(2 - 2^{-p+1})$

Check your results using the Julia function `bitstring`.

⚙️ **Exercise 1.25** (Summary). *True or false?*

1. Let  $(\bullet)_2$  denote binary representation. It holds that  $(0.1111)_2 + (0.0001)_2 = 1$ .

2. It holds that  $(1000)_2 \times (0.001)_2 = 1$ .

3. It holds that

$$(0.\overline{1})_3 = \frac{1}{2}.$$

4. In base 16, all the natural numbers from 1 to 200 can be represented using 2 digits.

5. In Julia, `Float64(.1) == Float32(.1)` evaluates to `true`.

6. The spacing (in absolute value) between successive double-precision (`Float64`) floating point numbers is constant.

7. It holds that  $(0.\overline{10101})_2 = (1.2345)_{10}$ .

8. Machine addition  $\hat{+}$  is an associative operation. More precisely, given any three double-precision floating point numbers  $x$ ,  $y$  and  $z$ , the following equality holds:

$$(x \hat{+} y) \hat{+} z = x \hat{+} (y \hat{+} z).$$

9. The machine epsilon is the smallest strictly positive number that can be represented in a floating point format.

10. Let  $\varepsilon$  denote the machine epsilon for the double-precision format. Let also  $\hat{+}$  and  $\hat{/}$  denote respectively the machine addition and the machine division operators for the double-precision format. It holds that  $1 \hat{+} (\varepsilon \hat{/} 64) = 1$  and that  $\varepsilon \hat{/} 64 \neq 0$ .


11. Assume that  $x \in \mathbf{R}$  belongs to the double-precision floating point format (that is, assume that  $x \in \mathbf{F}_{64}$ ). Then  $-x \in \mathbf{F}_{64}$ .

## 1.7 Discussion and bibliography

This chapter is mostly based on the original 1985 IEEE 754 standard [4] and the reference book [10]. A significant revision to the 1985 IEEE standard was published in 2008 [5], adding for example specifications for the half precision and quad precision formats, and a minor revision was published in 2019 [6]. The original IEEE standard and its revisions constitute the authoritative guide on floating point formats. It was intended to be widely disseminated and is written very clearly and concisely, but is not available for free online. Another excellent source for learning about floating point numbers and round-off errors is D. Goldberg’s paper “*What every computer scientist should know about floating-point arithmetic*” [3], freely available online.

# Chapter 2

## Interpolation and approximation

2.1	Interpolation . . . . .	26
2.1.1	Vandermonde matrix . . . . .	27
2.1.2	Lagrange interpolation formula . . . . .	27
2.1.3	Gregory–Newton interpolation . . . . .	28
2.1.4	Interpolation error . . . . .	33
2.1.5	Interpolation at Chebyshev nodes . . . . .	35
2.1.6	Hermite interpolation . . . . .	38
2.1.7	Piecewise interpolation . . . . .	39
2.2	Approximation . . . . .	40
2.2.1	Least squares approximation of data points . . . . .	40
2.2.2	Mean square approximation of functions . . . . .	42
2.2.3	Orthogonal polynomials . . . . .	43
2.2.4	Orthogonal polynomials and numerical integration: an introduction 	46
2.3	Exercises . . . . .	48
2.4	Discussion and bibliography . . . . .	54

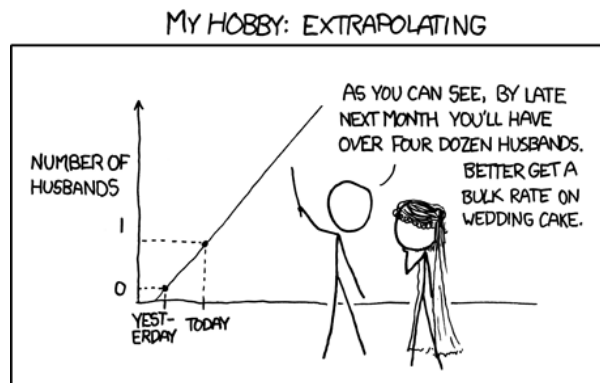


Figure 2.1: Source: <https://xkcd.com/605/>



## Introduction

In this chapter, we study numerical methods for interpolating and approximating functions. The Cambridge dictionary defines interpolation as *the addition of something different in the middle of a text, piece of music, etc. or the thing that is added*. The concept of interpolation in mathematics is consistent with this definition; interpolation consists in finding, given a set of points  $(x_i, y_i)$ , a function  $f$  in a finite-dimensional space that goes through these points. Throughout this course, you use the `plot` function in Julia, which performs piecewise linear interpolation for drawing functions, but there are a number of other standard interpolation methods. Our first goal in this chapter is to present an overview of these methods and the associated error estimates.

In the second part of this chapter, we focus on *function approximation*, which is closely related to the subject of mathematical interpolation. Indeed, a simple manner for approximating a general function by another one in a finite-dimensional space is to select a set of real numbers on the  $x$  axis, called *nodes*, and find the associated interpolant. As we shall demonstrate, not all sets of interpolation nodes are equal, and special care is required in order to avoid undesired oscillations. The field of function approximation is vast, so our aim in this chapter is to present only an introduction to the subject. In order to quantify the quality of an approximation, a metric on the space of functions, or a subset thereof, must be specified in order to measure errors. Without a metric, saying that two functions are close is almost meaningless!

## 2.1 Interpolation

Assume that we are given  $n + 1$  nodes  $x_0, \dots, x_n$  on the  $x$  axis, together with values  $u_0, \dots, u_n$ , which may be the values taken by an unknown function  $u(x)$  when evaluated at these points. Suppose that we are looking for an interpolation  $\hat{u}(x)$  in a subspace  $\text{Span}\{\varphi_0, \dots, \varphi_n\}$  of the vector space of continuous functions, i.e. an interpolating function of the form

$$\hat{u}(x) = \alpha_0\varphi_0(x) + \dots + \alpha_n\varphi_n(x),$$

where  $\alpha_0, \dots, \alpha_n$  are real coefficients. In order for  $\hat{u}(x)$  to be an interpolating function, we must require that

$$\forall i \in \{0, \dots, n\}, \quad \hat{u}(x_i) = u_i.$$

This leads to a linear system of  $n + 1$  equations and  $n + 1$  unknowns, the latter being the coefficients  $\alpha_0, \dots, \alpha_n$ . This system of equations in matrix form reads

$$\begin{pmatrix} \varphi_0(x_0) & \varphi_1(x_0) & \dots & \varphi_n(x_0) \\ \varphi_0(x_1) & \varphi_1(x_1) & \dots & \varphi_n(x_1) \\ \vdots & \vdots & & \vdots \\ \varphi_0(x_n) & \varphi_1(x_n) & \dots & \varphi_n(x_n) \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} = \begin{pmatrix} u_0 \\ u_1 \\ \vdots \\ u_n \end{pmatrix}. \quad (2.1)$$

### 2.1.1 Vandermonde matrix

Since polynomials are very convenient for evaluation, integration, and differentiation, they are a natural choice for interpolation purposes. The simplest basis of the subspace of polynomials of degree less than or equal to  $n$  is given by the monomials:

$$\varphi_0(x) = 1, \quad \varphi_1(x) = x, \quad \dots, \quad \varphi_n(x) = x^n.$$

In this case, the linear system (2.1) for determining the coefficients of the interpolant reads

$$\begin{pmatrix} 1 & x_0 & \dots & x_0^n \\ 1 & x_1 & \dots & x_1^n \\ \vdots & \vdots & & \vdots \\ 1 & x_n & \dots & x_n^n \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} = \begin{pmatrix} u_0 \\ u_1 \\ \vdots \\ u_n \end{pmatrix}. \quad (2.2)$$

The matrix on the left-hand side is called a *Vandermonde* matrix. If the abscissae  $x_0, \dots, x_n$  are distinct, then this is a full rank matrix, and so (2.2) admits a unique solution, implying as a corollary that the interpolating polynomial exists and is unique. It is possible to show that the condition number of the Vandermonde increases dramatically with  $n$ . Consequently, solving (2.2) is not a viable method in practice for calculating the interpolating polynomial.

### 2.1.2 Lagrange interpolation formula

One may wonder whether polynomial basis functions  $\varphi_0, \dots, \varphi_n$  can be defined in such a manner that the matrix in (2.1) is the identity matrix. The answer to this question is positive; it suffices to take as a basis the *Lagrange polynomials*, which are given by

$$\varphi_i(x) = \frac{(x - x_0)(x - x_1) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_n)}{(x_i - x_0)(x_i - x_1) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_n)} = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}.$$

It is simple to check that

$$\varphi_i(x_j) = \delta_{i,j} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases}$$

Finding the interpolant in this basis is immediate:

$$\widehat{u}(x) = u_0\varphi_0(x) + \dots + u_n\varphi_n(x).$$

While simple, this approach to polynomial interpolation has a few disadvantages:

- First, evaluating  $\widehat{u}(x)$  is computationally costly when  $n$  is large.
- Second, all the basis functions change when adding new interpolation nodes.
- Finally, Lagrange interpolation is numerically unstable because of cancellations between large terms. Indeed, it is often the case that Lagrange polynomials take very large values

over the interpolation intervals; this occurs, for example, when many equidistant interpolation nodes are employed, as illustrated in Figure 2.2.

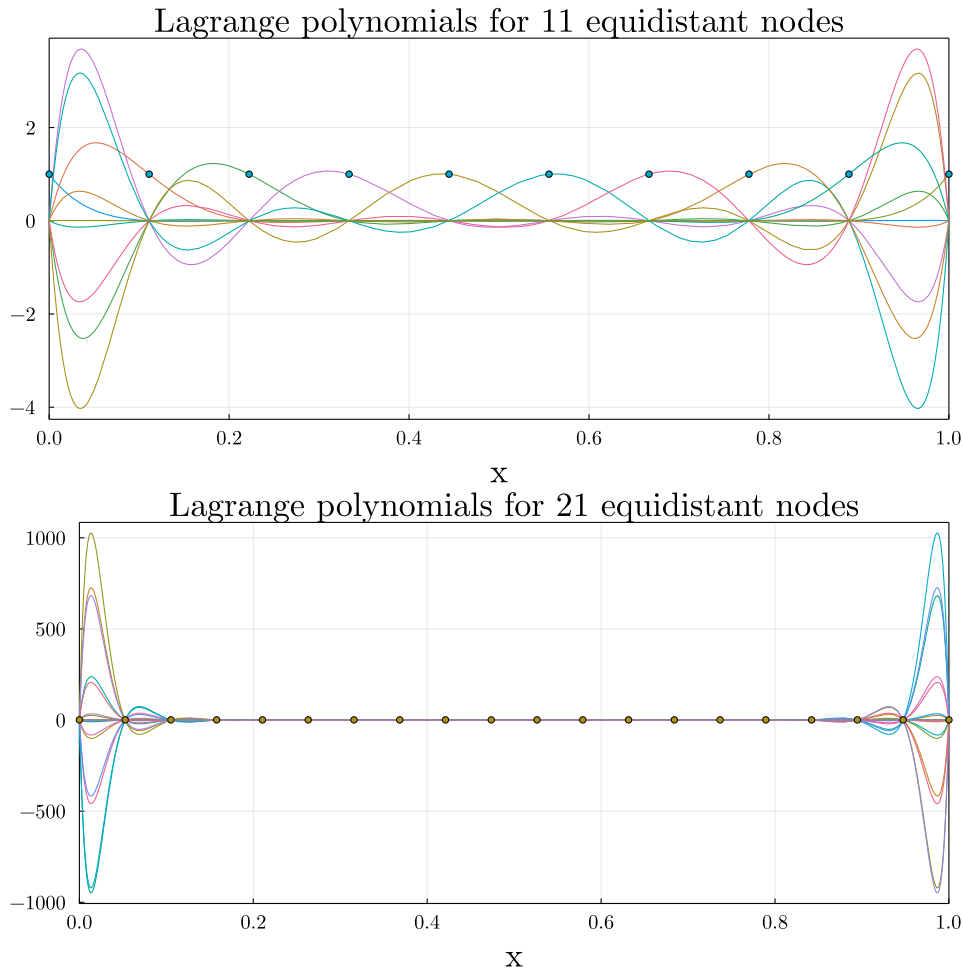


Figure 2.2: Lagrange polynomials associated with equidistant nodes over the  $(0, 1)$  interval.

### 2.1.3 Gregory–Newton interpolation

By Taylor's formula, any polynomial  $p$  of degree  $n$  may be expressed as

$$p(x) = p(0) + p'(0)x + \frac{p''(0)}{2}x^2 + \dots + \frac{p^{(n)}(0)}{n!}x^n. \quad (2.3)$$

The constant coefficient can be obtained by evaluating the polynomial at 0, the linear coefficient can be identified by evaluating the first derivative at 0, and so on. Assume now that we are given the values taken by  $p$  when evaluated at the integer numbers  $\{0, \dots, n\}$ . We ask the following question: can we find a formula similar in spirit to (2.3), but including only evaluations of  $p$  and not of its derivatives? To answer this question, we introduce the difference operator  $\Delta$  which acts on functions as follows:

$$\Delta f(x) = f(x+1) - f(x).$$

The operator  $\Delta$  is a linear operator on the space of continuous functions. It maps constant functions to 0, and the linear function  $x$  to the constant function 1, suggesting a resemblance

with the differentiation operator. In order to further understand this connection, let us define the *falling power* of a real number  $x$  as

$$x^{\underline{k}} = x(x-1)(x-2)\dots(x-k+1). \quad (2.4)$$

We then have that

$$\begin{aligned} \Delta x^{\underline{k}} &= (x+1)x(x-1)\dots(x-k+2) - x(x-1)(x-2)\dots(x-k+1) \\ &= ((x+1) - (x-k+1))(x(x-1)\dots(x-k+2)) = kx^{\underline{k-1}} \end{aligned} \quad (2.5)$$

In other words, the action of the difference operator on falling powers mirrors that of the differentiation operator on monomials. The falling powers form a basis of the space of polynomials, and so any polynomial in  $\mathbf{P}(n)$ , i.e. of degree less than or equal to  $n$ , can be expressed as

$$p(x) = \alpha_0 + \alpha_1 x^{\underline{1}} + \alpha_2 x^{\underline{2}} + \dots + \alpha_n x^{\underline{n}}. \quad (2.6)$$

It is immediate to show that  $\alpha_i = \Delta^i p(0)/i!$ , where  $\Delta^i p$  denotes the function obtained after  $i$  applications of the operator  $\Delta$ . Therefore, any polynomial of degree less than or equal to  $n$  may be expressed as

$$p(x) = p(0) + \Delta p(0)x^{\underline{1}} + \frac{\Delta^2 p(0)}{2}x^{\underline{2}} + \dots + \frac{\Delta^n p(0)}{n!}x^{\underline{n}}. \quad (2.7)$$

An expansion of the form (2.7) is called a *Newton series*, which is the discrete analog of the continuous Taylor series. From the definition of  $\Delta$ , it is clear that the coefficients in (2.7) depend only on  $p(0), \dots, p(n)$ . We conclude that, given points  $n+1$  points  $(i, u_i)$  for  $i \in \{0, \dots, n\}$ , the unique interpolating polynomial is given by (2.7), after replacing  $p(i)$  by  $u_i$ .

*Example 2.1.* Let us use (2.6) in order to calculate the value of

$$S(n) := \sum_{i=0}^n i^2.$$

Since  $\Delta S(n) = (n+1)^2$ , which is a second degree polynomial in  $n$ , we deduce that  $S(n)$  is a polynomial of degree 3. Let us now determine its coefficients.

$n$	0	1	2	3
$\Delta^0 S(n)$	<b>0</b>	1	5	14
$\Delta^1 S(n)$	<b>1</b>	4	9	
$\Delta^2 S(n)$	<b>3</b>	5		
$\Delta^3 S(n)$	<b>2</b>			

We conclude that

$$S(n) = \mathbf{1}n + \frac{\mathbf{3}}{2!}n(n-1) + \frac{\mathbf{2}}{3!}n(n-1)(n-2) = \frac{n(2n+1)(n+1)}{6}$$

Notice that when falling powers are employed as polynomial basis, the matrix in (2.1) is lower triangular, and so the algorithm described in [Example 2.1](#) could be replaced by the forward

substitution method. Whereas the coefficients of the Lagrange interpolant can be obtained immediately from the values of  $u$  at the nodes, calculating the coefficients of the expansion in (2.6) requires  $\mathcal{O}(n^2)$  operations. However, Gregory–Newton interpolation has several advantages over Lagrange interpolation:

- If a point  $(n + 1, p_{n+1})$  is added to the set of interpolation points, only one additional term, corresponding to the falling power  $x^{\underline{n+1}}$ , needs to be calculated in (2.7). All the other coefficients are unchanged. Therefore, the Gregory–Newton approach is well-suited for incremental interpolation.
- The Gregory–Newton interpolation method is more numerically stable than Lagrange interpolation, because the basis functions do not take very large values.
- A polynomial in the form of a Newton series can be evaluated efficiently using Horner’s method, which is based on rewriting the polynomial as

$$p(x) = \alpha_0 + x \left( \alpha_1 + (x - 1) \left( \alpha_2 + (x - 2) \left( \alpha_3 + (x - 3) \dots \right) \right) \right).$$

Evaluating this expression starting from the innermost bracket leads to an algorithm with a cost scaling linearly with the degree of the polynomial.

### Non-equidistant nodes

So far, we have described the Gregory–Newton method in the simple setting where interpolation nodes are just a sequence of successive natural numbers. The method can be generalized to the setting of nodes  $x_0 \neq \dots \neq x_n$  which are not necessarily equidistant. In this case, we take as basis the following functions instead of the falling powers:

$$\varphi_i(x) = (x - x_0)(x - x_1) \dots (x - x_{i-1}), \tag{2.8}$$

with the convention that the empty product is 1. By (2.1), the coefficients of the interpolating polynomial in this basis solve the following linear system:

$$\begin{pmatrix} 1 & & & & & 0 \\ 1 & x_1 - x_0 & & & & \\ 1 & x_2 - x_0 & (x_2 - x_0)(x_2 - x_1) & & & \vdots \\ \vdots & \vdots & & \ddots & & \\ 1 & x_n - x_0 & \dots & \dots & \prod_{j=0}^{n-1} (x_n - x_j) & \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{pmatrix} = \begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}. \tag{2.9}$$

This system could be solved using, for example, forward substitution. Clearly  $\alpha_0 = u_0$  from the first equation, and then from the second equation we obtain

$$\alpha_1 = \frac{u_1 - u_0}{x_1 - x_0} =: [u_0, u_1],$$

which may be viewed as an approximation of the slope of  $u$  at  $x_0$ . The right-hand side of this equation is an example of a *divided difference*. In general, divided differences are defined

recursively as follows:

$$[u_0, u_1, \dots, u_d] := \frac{[u_1, \dots, u_d] - [u_0, \dots, u_{d-1}]}{x_d - x_0}, \quad [u_i] = u_i. \quad (2.10)$$

Let us start by observing that our divided differences coincide with the coefficients in (2.6).

**Lemma 2.1.** *Assume that  $(0, u_0), \dots, (n, u_n) \in \mathbf{R}^2$ . Then, provided that  $(i, k) \in \mathbf{N}^2$  satisfies  $i + k \leq n$ , we have*

$$\frac{1}{k!} \Delta^k u(i) = [u_i, \dots, u_{i+k}].$$

*Proof.* The result is clearly true for  $k = 0$ , and we shall assume by induction that it holds for  $k - 1$ . Using our induction hypothesis, we get

$$\begin{aligned} \frac{1}{k!} \Delta^k u(i) &= \frac{1}{k} \left( \frac{\Delta^{k-1} u(i+1)}{(k-1)!} - \frac{\Delta^{k-1} u(i)}{(k-1)!} \right) \\ &= \frac{1}{k} \left( [u_{i+1}, \dots, u_{i+k}] - [u_i, \dots, u_{i+k-1}] \right). \end{aligned}$$

By the definition of divided differences in (2.10), it follows that we have

$$\frac{1}{k!} \Delta^k u(i) = \frac{1}{k} \underbrace{(x_{i+k} - x_i)}_{i+k-i} [u_i, \dots, u_{i+k}] = [u_i, \dots, u_{i+k}],$$

which allows us to conclude the proof.  $\square$

In light of Lemma 2.1, the reader will not be surprised that the coefficients in the  $\varphi_0, \dots, \varphi_n$  basis in (2.8) are given by the divided differences.

**Proposition 2.2.** *Assume that  $(x_0, u_0), \dots, (x_n, u_n)$  are  $n+1$  points in the plane with distinct abscissae. Then the interpolating polynomial of degree  $n$  may be expressed as*

$$p(x) = \sum_{i=0}^n [u_0, \dots, u_i] \varphi_i(x),$$

where  $\varphi_i(x)$ , for  $i = 0, \dots, n$ , are the basis functions defined in (2.8).

*Proof.* The statement is true for  $n = 0$ . Reasoning by induction, we assume that it holds true for polynomials of degree up to  $n - 1$ . Let  $p_1(x)$  and  $p_2(x)$  be the interpolating polynomials at the points  $x_0, x_1, \dots, x_{n-2}, x_{n-1}$  and  $x_0, x_1, \dots, x_{n-2}, x_n$ , respectively. Then

$$p(x) = p_1(x) + \frac{x - x_{n-1}}{x_n - x_{n-1}} (p_2(x) - p_1(x)) \quad (2.11)$$

is a polynomial of degree  $n$  that interpolates all the data points. By the induction hypothesis,

it holds that

$$p_1(x) = u_0 + [u_0, u_1](x - x_0) + \dots + [u_0, u_1, \dots, u_{n-2}, \mathbf{u}_{n-1}] \prod_{i=0}^{n-2} (x - x_i),$$

$$p_2(x) = u_0 + [u_0, u_1](x - x_0) + \dots + [u_0, u_1, \dots, u_{n-2}, \mathbf{u}_n] \prod_{i=0}^{n-2} (x - x_i).$$

Here we used a bold font in order to emphasize the difference between the two expressions. Substituting these expressions in (2.11), we obtain

$$p(x) = u_0 + [u_0, u_1](x - x_0) + \dots + [u_0, \dots, u_{n-2}] \prod_{i=0}^{n-2} (x - x_i) + \frac{[u_0, u_1, \dots, u_{n-2}, u_n] - [u_0, u_1, \dots, u_{n-2}, u_{n-1}]}{x_n - x_{n-1}} \prod_{i=0}^{n-1} (x - x_i).$$

In Exercise 2.4, we show that divided differences are invariant under permutations of the data points, and so we have that

$$\frac{[u_0, u_1, \dots, u_{n-2}, u_n] - [u_0, u_1, \dots, u_{n-2}, u_{n-1}]}{x_n - x_{n-1}} = [u_0, \dots, u_n],$$

which enables to conclude. □

*Example 2.2.* Assume that we are looking for the third degree polynomial going through the following points:

$$(-1, 10), \quad (0, 4), \quad (2, -2), \quad (4, -40).$$

We have to calculate the divided difference  $\alpha_i = [u_0, \dots, u_i]$  for  $i \in \{0, 1, 2, 3\}$ . To this end, it is convenient to use a table:

$i$	0	1	2	3
$[u_i]$	<b>10</b>	4	-2	-40
$x_{i+1} - x_i$	1	2	2	
$[u_i, u_{i+1}]$	<b>-6</b>	-3	-19	
$x_{i+2} - x_i$	<b>3</b>	4		
$[u_i, u_{i+1}, u_{i+2}]$	<b>1</b>	-4		
$x_{i+3} - x_i$	5			
$[u_i, u_{i+1}, u_{i+2}, u_{i+3}]$	<b>-1</b>			

We deduce that the expression of the interpolating polynomial is

$$p(x) = \mathbf{10} + (-\mathbf{6})(x + 1) + \mathbf{1}(x + 1)x + (-\mathbf{1})(x + 1)x(x - 2) = -x^3 + 2x^2 + -3x + 4.$$

### 2.1.4 Interpolation error

Assume that  $u(x)$  is a continuous function and denote by  $\widehat{u}(x)$  its interpolation through the points  $(x_i, u_i)$ , where  $u_i = u(x_i)$  for  $i = 0, \dots, n$ . In this section, we study the behavior of the error in the limit as  $n \rightarrow \infty$ .

**Theorem 2.3** (Interpolation error). *Assume that  $u: [a, b] \rightarrow \mathbf{R}$  is a function in  $C^{n+1}([a, b])$  and let  $x_0, \dots, x_n$  denote  $n + 1$  distinct interpolation nodes. Then for all  $x \in [a, b]$ , there exists  $\xi = \xi(x)$  in the interval  $[a, b]$  such that*

$$e_n(x) := u(x) - \widehat{u}(x) = \frac{u^{(n+1)}(\xi)}{(n+1)!} (x - x_0) \dots (x - x_n).$$

*Proof.* The statement is obvious if  $x \in \{x_0, \dots, x_n\}$ , so we assume from now on that  $x$  does not coincide with an interpolation node. Let us use the compact notation  $\omega_n = \prod_{i=0}^n (x - x_i)$  and introduce the function

$$g(t) = e_n(t)\omega_n(x) - e_n(x)\omega_n(t). \quad (2.12)$$

The function  $g$  is smooth and takes the value 0 when evaluated at  $x_0, \dots, x_n, x$ . Since  $g$  has  $n + 2$  roots in the interval  $[a, b]$ , Rolle's theorem implies that  $g'$  has at least  $n + 1$  distinct roots in  $(a, b)$ . Another application of Rolle's theorem then yields that  $g''$  has at least  $n$  distinct roots in  $(a, b)$ . Iterating this reasoning, we deduce that  $g^{(n+1)}$  has one root  $t_*$  in  $(a, b)$ . From (2.12), we calculate that

$$g^{(n+1)}(t) = e_n^{(n+1)}(t)\omega_n(x) - e_n(x)\omega_n^{(n+1)}(t) = u^{(n+1)}(t)\omega_n(x) - e_n(x)(n+1)!. \quad (2.13)$$

Here we employed the fact that  $\widehat{u}^{(n+1)}(t) = 0$ , because  $\widehat{u}$  is a polynomial of degree at most  $n$ . Evaluating (2.13) at  $t_*$  and rearranging, we obtain that

$$e_n(x) = \frac{u^{(n+1)}(t_*)}{(n+1)!} \omega_n(x),$$

which completes the proof. □

As a corollary to Theorem 2.3, we deduce the following error bound.

**Corollary 2.4** (Upper bound on the interpolation error). *Assume that  $u$  is smooth in the interval  $[a, b]$  and let*

$$C_{n+1} = \sup_{x \in [a, b]} |u^{(n+1)}(x)|.$$

*Then*

$$E_n := \sup_{x \in [a, b]} |e_n(x)| \leq \frac{C_{n+1}}{4(n+1)} h^{n+1} \quad (2.14)$$

*where  $h$  is the maximum spacing between two successive interpolation nodes.*



*Proof.* By Theorem 2.3, it holds that

$$\forall x \in [a, b], \quad |e_n(x)| \leq \frac{C_{n+1}}{(n+1)!} |(x-x_0)\dots(x-x_n)|. \quad (2.15)$$

The product on the right-hand side is bounded from above by

$$\frac{h^2}{4} \times 2h \times 3h \times 4h \times \dots \times nh = \frac{n!h^{n+1}}{4}. \quad (2.16)$$

The first factor comes from the fact that, if  $x \in [x_i, x_{i+1}]$ , then

$$\left| (x-x_i)(x-x_{i+1}) \right| \leq \frac{(x_{i+1}-x_i)^2}{4},$$

because the left-hand side is maximized when  $x$  is the midpoint of the interval  $[x_i, x_{i+1}]$ . Substituting (2.16) into (2.15), we deduce the statement.  $\square$

We now ask the following natural question: does  $E_n$  given in (2.14) tend to zero as the maximum spacing between successive nodes tends to 0? By Corollary 2.4, the answer to this question is positive when  $C_n$  does not grow too quickly as  $n \rightarrow \infty$ . For example the interpolation error for the function  $u(x) = \sin(x)$  decreases very quickly as  $n \rightarrow \infty$  when equidistant interpolation nodes are employed, as illustrated in Figure 2.3.

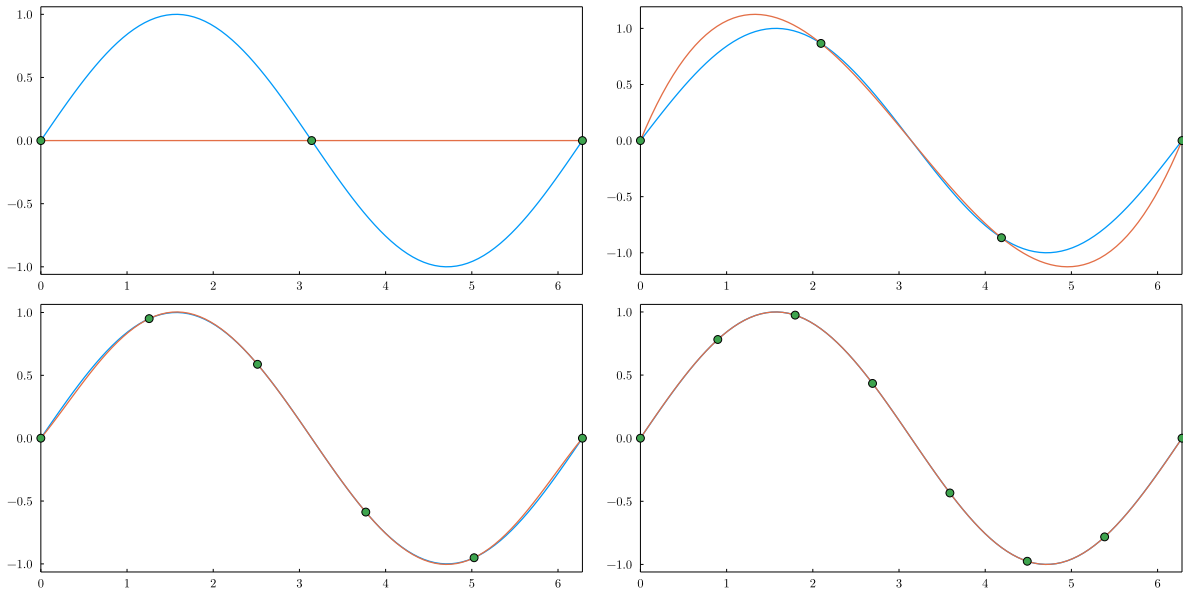


Figure 2.3: Interpolation (in orange) of the function  $u(x) = \sin(x)$  (in blue) using 3, 4, 6, and 8 equidistant nodes.

In some cases, however, the constant  $C_n$  grows quickly with  $n$ , to the extent that  $E_n$  may increase with  $n$ ; in this case, the maximum interpolation error grows when nodes are added! The classic example illustrating this potential issue is that of the Runge function:

$$u(x) = \frac{1}{1+25x^2}. \quad (2.17)$$

It is possible to show that, for this function, the upper bound in (2.14) tends to  $\infty$  in the

limit as the number  $n$  of interpolation nodes increases. We emphasize that this does not prove that  $E_n \rightarrow \infty$  in the limit as  $n \rightarrow \infty$ , because (2.14) provides only an *upper bound* on the error. In fact, the interpolation error for the Runge function can either grow or decrease, depending on the choice of interpolation nodes. With equidistant nodes, it turns out that  $E_n \rightarrow \infty$ , as illustrated in Figure 2.4.

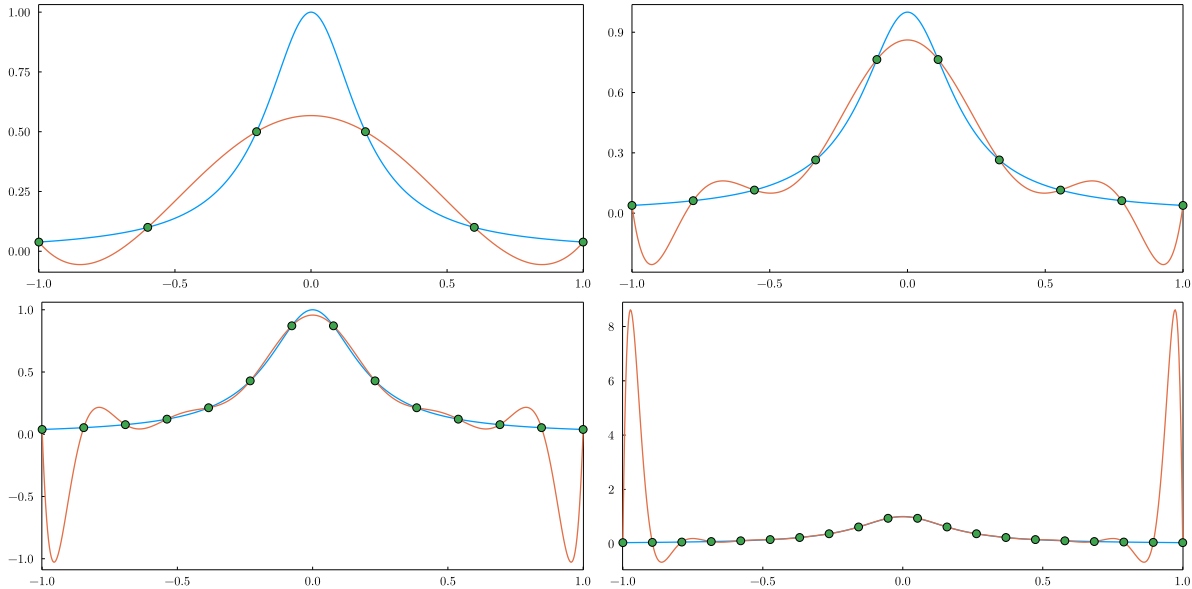


Figure 2.4: Interpolation (in orange) of the Runge function (2.17) (in blue) using 6, 10, 14, and 20 equidistant nodes.

### 2.1.5 Interpolation at Chebyshev nodes

Sometimes, interpolation is employed as a technique for approximating functions. The spectral collocation method, for example, is a technique for solving partial differential equations where a discrete solution is first calculated, and then a continuous solution is constructed using polynomial or Fourier interpolation. When the interpolation nodes are not given a priori as data, it is natural to wonder whether these can be picked in such a manner that the interpolation error, measured in a function norm, is minimized. For example, given a continuous function  $u(x)$  and a number of nodes  $n + 1$ , is it possible to choose nodes  $x_0, \dots, x_n$  such that

$$E := \sup_{x \in [a, b]} |u(x) - \hat{u}(x)|$$

is minimized? Here  $\hat{u}$  is the polynomial interpolating  $u$  at the nodes. Achieving this goal in general is a difficult task, because  $\xi = \xi(x)$  is unknown in the expression of the interpolation error from Theorem 2.3:

$$e_n(x) = \frac{u^{(n+1)}(\xi)}{(n+1)!} (x - x_0) \dots (x - x_n).$$

In view of this difficulty, we will focus on the simpler problem of finding interpolation nodes such that the product  $(x - x_0) \dots (x - x_n)$  is minimized in the supremum norm. This problem

amounts to finding the optimal interpolation nodes, in the sense that  $E$  is minimized, in the particular case where  $u$  is a polynomial of degree  $n+1$ , because in this case  $u^{(n+1)}(\xi)$  is a constant factor. It turns out that this problem admits an explicit solution, which we will deduce from the following theorem.

**Theorem 2.5** (Minimum  $\infty$  norm). *Assume that  $p$  is a monic polynomial of degree  $n \geq 1$ :*

$$p(x) = \alpha_0 + \alpha_1 x + \cdots + \alpha_{n-1} x^{n-1} + x^n.$$

*Then it holds that*

$$\sup_{x \in [-1, 1]} |p(x)| \geq \frac{1}{2^{n-1}} =: E. \quad (2.18)$$

*In addition, the lower bound is achieved for  $p_*(x) = 2^{-(n-1)}T_n(x)$ , where  $T_n$  is the Chebyshev polynomial of degree  $n$ :*

$$T_n(x) = \cos(n \arccos x) \quad (-1 \leq x \leq 1). \quad (2.19)$$

*Proof.* By Exercise C.5, the polynomial  $x \mapsto 2^{-(n-1)}T_n(x)$  is indeed monic, and it is clear that it achieves the lower bound (2.18) since  $|T_n(x)| \leq 1$  for all  $x \in [-1, 1]$ .

In order to prove (2.18), we assume by contradiction that there is a different monic polynomial  $q$  of degree  $n$  such that

$$\sup_{x \in [-1, 1]} |q(x)| < E. \quad (2.20)$$

Let us introduce  $x_i = \cos(i\pi/n)$ , for  $i = 0, \dots, n$ , and observe that

$$p(x_i) = 2^{-(n-1)}T_n(x_i) = (-1)^i E.$$

The function  $h(x) := p(x) - q(x)$  is a polynomial of degree at most  $n - 1$  which, by the assumption (2.20), is strictly positive at  $x_0, x_2, x_4, \dots$  and strictly negative at  $x_1, x_3, x_5, \dots$ . Therefore, the polynomial  $h(x)$  changes sign at least  $n$  times and so, by the intermediate value theorem, it has at least  $n$  roots. But this is impossible, because  $h(x) \neq 0$  and  $h(x)$  is of degree at most  $n - 1$ .  $\square$

*Remark 2.1* (Derivation of Chebyshev polynomials). The polynomial  $p_*$  achieving the lower bound in (2.18) oscillates between the values  $-E$  and  $E$ , which are respectively its minimum and maximum values over the interval  $[-1, 1]$ . It attains the values  $E$  or  $-E$  at  $n + 1$  distinct points  $x_0 < \dots < x_n$ , with  $x_0 = -1$  and  $x_n = 1$ . It turns out that these properties, which can be shown to hold a priori using Chebyshev's *equioscillation theorem*, are sufficient to derive an explicit expression for the polynomial  $p_*$ , as we formally demonstrate hereafter.

The points  $x_1, \dots, x_{n-1}$  are local extrema of  $p_*$ , and so  $p'_*(x) = 0$  at these nodes. We therefore deduce that  $p_*$  satisfies the differential equation

$$n^2 (E^2 - p_*(x)^2) = p'_*(x)^2 (1 - x^2). \quad (2.21)$$

Indeed, both sides are polynomials of degree  $2n$  with single roots at  $-1$  and  $1$ , with double roots at  $x_1, \dots, x_{n-1}$ , and with the same coefficient of the leading power. In order to solve (2.21), we rearrange the equation and take the square root:

$$\frac{\frac{p'_*(x)}{E}}{\sqrt{1 - \frac{p_*(x)^2}{E^2}}} = \pm \frac{n}{\sqrt{1 - x^2}} \quad \Leftrightarrow \quad \frac{d}{dx} \left( \arccos \left( \frac{p_*(x)}{E} \right) \right) = \pm n \frac{d}{dx} \arccos(x).$$

Integrating both sides and taking the cosine, we obtain

$$p_*(x) = E \cos(C + n \arccos(x)).$$

Requiring that  $|p_*(-1)| = E$ , we deduce  $C = 0$ .

From Theorem 2.5, we deduce the following corollary.

**Corollary 2.6** (Chebyshev nodes). *Assume that  $x_0 < x_1 < \dots < x_n$  are in the interval  $[a, b]$ . The supremum norm of the product  $\omega(x) := (x - x_0) \cdots (x - x_n)$  over  $[a, b]$  is minimized when*

$$x_i = a + (b - a) \frac{1 + \cos \left( \frac{(2i+1)\pi}{2(n+1)} \right)}{2} \quad (2.22)$$

*Proof.* We consider the affine change of variable

$$\begin{aligned} \zeta &: [-1, 1] \rightarrow [a, b]; \\ y &\mapsto \frac{a + b + y(b - a)}{2}. \end{aligned}$$

The function

$$\begin{aligned} p(y) &:= \frac{2^{n+1}}{(b - a)^{n+1}} \omega(\zeta(y)) = \frac{2^{n+1}}{(b - a)^{n+1}} (\zeta(y) - x_0) \cdots (\zeta(y) - x_n) \\ &= (y - y_0) \cdots (y - y_n), \quad y_i = \zeta^{-1}(x_i), \end{aligned}$$

is a monic polynomial of degree  $n + 1$  such that

$$\sup_{y \in [-1, 1]} |p(y)| = \frac{2^{n+1}}{(b - a)^{n+1}} \sup_{x \in [a, b]} |(x - x_0) \cdots (x - x_n)|. \quad (2.23)$$

By Theorem 2.5, the left-hand side is minimized when  $p$  is equal to  $2^{-n}T_{n+1}$ , i.e. when the roots of  $p$  coincide with the roots of  $T_{n+1}$ . This occurs when

$$y_i = \zeta^{-1}(x_i) = \cos \left( \frac{(2i + 1)\pi}{2(n + 1)} \right).$$

Applying the inverse change of variable  $x_i = \zeta(y_i)$ , we deduce the result.  $\square$

Corollary 2.6 is useful for interpolation. The nodes

$$x_i = a + (b - a) \frac{1 + \cos\left(\frac{(2i+1)\pi}{2(n+1)}\right)}{2}, \quad i = 0, \dots, n, \quad (2.24)$$

are known as Chebyshev nodes and, more often than not, employing these nodes for interpolation produces much better results than using equidistant nodes, both in the case where  $u$  is a polynomial of degree  $n + 1$ , as we just proved, but also for general  $u$ . As an example we plot in Figure 2.5 the interpolation of the Runge function using Chebyshev nodes. In this case, the interpolating polynomial converges uniformly to the Runge function as we increase the number of interpolation nodes!

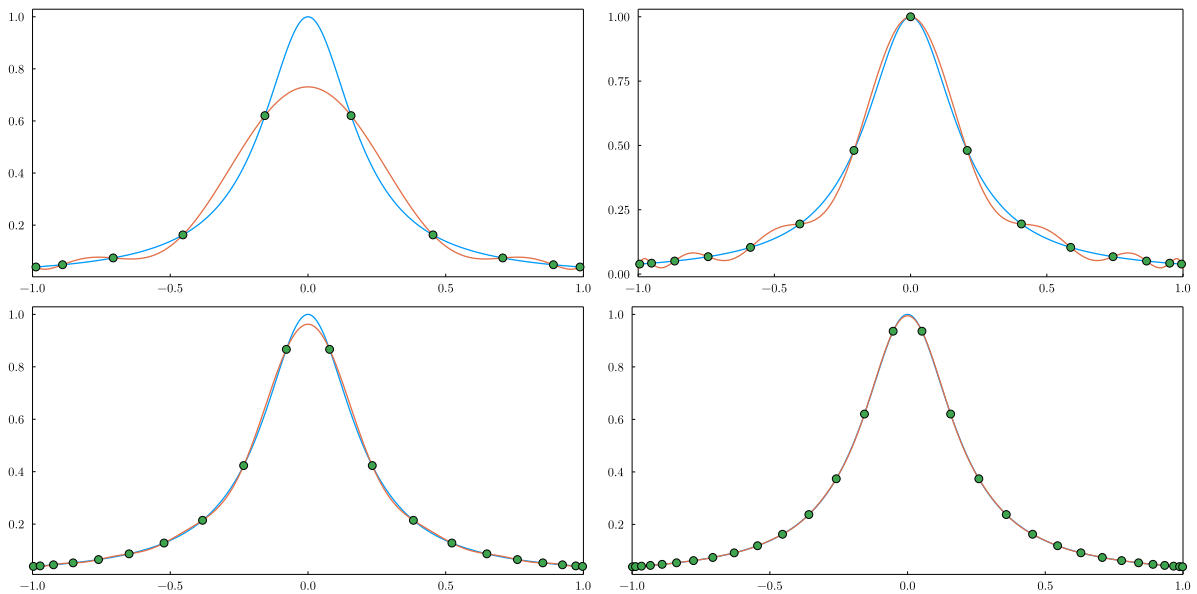


Figure 2.5: Interpolation (in orange) of the Runge function (2.17) (in blue) using 10, 15, 20, and 30 Chebyshev nodes.

### 2.1.6 Hermite interpolation

Hermite interpolation, sometimes also called Hermite–Birkoff interpolation, generalizes Lagrange interpolation to the case where, in addition to the function values  $u_0, \dots, u_n$ , the values of some of the derivatives are given at the interpolation nodes. For simplicity, we assume in this section that only the first derivative is specified. In this case, the aim of Hermite interpolation is to find, given data  $(x_i, u_i, u'_i)$  for  $i \in \{0, \dots, n\}$ , a polynomial  $\hat{u}$  of degree at most  $2n + 1$  such that

$$\forall i \in \{0, \dots, n\}, \quad \hat{u}(x_i) = u_i, \quad \hat{u}'(x_i) = u'_i.$$

In order to construct the interpolating polynomial, it is useful to define the functions

$$\psi_i(x) = \prod_{j=0, j \neq i}^n \left( \frac{x - x_j}{x_i - x_j} \right)^2, \quad i = 0, \dots, n.$$

The function  $\psi_i$  is the square of the usual Lagrange polynomials associated with  $x_i$ , and it satisfies

$$\psi_i(x_i) = 1, \quad \psi_i'(x_i) = \sum_{j=0, j \neq i}^n \frac{2}{x_i - x_j}, \quad \forall j \neq i \quad \psi_i(x_j) = \psi_i'(x_j) = 0.$$

We consider the following ansatz for  $\hat{u}$ :

$$\hat{u}(x) = \sum_{i=0}^n \psi_i(x) q_i(x),$$

where  $q_i$  are polynomials to be determined of degree at most one, so that  $\hat{u}$  is of degree at most  $2n + 1$ . We then require

$$\hat{u}(x_i) = q_i(x_i), \quad \hat{u}'(x_i) = \psi_i'(x_i) q_i(x_i) + q_i'(x_i).$$

From the first equation, we deduce that  $q_i(x_i) = u_i$ , and from the second equation we then have  $q_i'(x_i) = \hat{u}'(x_i) - \psi_i'(x_i) u_i$ . We conclude that the interpolating polynomial is given by

$$\hat{u}(x) = \sum_{i=0}^n \psi_i(x) \left( u_i + (u_i' - \psi_i'(x_i) u_i)(x - x_i) \right).$$

The following theorem gives an expression of the error.

**Theorem 2.7** (Hermite interpolation error). *Assume that  $u: [a, b] \rightarrow \mathbf{R}$  is a function in  $C^{2n+2}([a, b])$  and let  $\hat{u}$  denote the Hermite interpolation of  $u$  at the nodes  $x_0, \dots, x_n$ . Then for all  $x \in [a, b]$ , there exists  $\xi = \xi(x)$  in the interval  $[a, b]$  such that*

$$u(x) - \hat{u}(x) = \frac{u^{(2n+2)}(\xi)}{(2n+2)!} (x - x_0)^2 \dots (x - x_n)^2.$$

*Proof.* See Exercise 2.8. □

### 2.1.7 Piecewise interpolation

The interpolation methods we discussed so far are in some sense global; they aim to construct one polynomial that goes through all the data points. This approach is attractive because the interpolant is infinitely smooth but, as we observed, it is not always fruitful, in particular when equidistant interpolation nodes are employed. An alternative approach is to divide the domain in a number of small intervals and perform polynomial interpolation within each interval. Although the resulting interpolating function is usually not smooth over the full domain, this “local” approach to interpolation is more robust.

Several methods belong in the category of piecewise interpolation. We mention, for instance, piecewise Lagrange interpolation and cubic splines interpolation. In this section, we briefly describe the former method, which is widely used in the context of the *finite element method*. Information on the latter method is available in [10, Section 8.7.1].

For simplicity, we illustrate the method in dimension 1, but piecewise Lagrange interpolation can be extended to several dimensions. Assume that we wish to approximate a function  $u: [a, b] \rightarrow \mathbf{R}$ . We consider a subdivision  $a = x_0 < x_1 < \dots < x_n = b$  of the interval  $[a, b]$  and let  $h$  denote the maximum spacing:

$$h = \max_{i \in \{0, \dots, n-1\}} |x_{i+1} - x_i|.$$

Within each subinterval  $I_i = [x_i, x_{i+1}]$ , we consider a further subdivision

$$x_i = x_i^{(0)} < x_i^{(1)} < \dots < x_i^{(m)} = x_{i+1},$$

where the nodes  $x_i^{(0)}, \dots, x_i^{(m)}$  are equally spaced with distance  $h/m$ . The idea of piecewise Lagrange interpolation is to calculate, for each interval  $I_i$  in the partition, the interpolating polynomial  $p_i$  at the nodes  $x_i^{(0)}, \dots, x_i^{(m)}$ . The interpolant is then defined as

$$\widehat{u}(x) = p_\iota(x), \quad (2.25)$$

where  $\iota = \iota(x)$  is the index of the interval to which  $x$  belongs. Since  $x_i^{(m)} = x_{i+1} = x_{i+1}^{(0)}$ , the interpolant defined by (2.25) is continuous. If the function  $u$  belongs to  $C^{m+1}([a, b])$ , then by Corollary 2.4 the interpolation error within each subinterval may be bounded from above as follows:

$$\sup_{x \in I_i} |u(x) - \widehat{u}(x)| \leq \frac{C_{m+1}(h/m)^{m+1}}{4(m+1)}, \quad C_{m+1} := \sup_{x \in [a, b]} |u^{(m+1)}(x)|, \quad (2.26)$$

and so we deduce

$$\sup_{x \in [a, b]} |u(x) - \widehat{u}(x)| \leq Ch^{m+1},$$

for an appropriate constant  $C$  independent of  $h$ . This equation shows that the error is guaranteed to decrease to 0 in the limit as  $h \rightarrow 0$ . In practice, the number  $m$  of interpolation nodes within each interval can be small.

## 2.2 Approximation

In this section, we focus on the subject of *approximation*, both of discrete data points and of continuous functions. We begin, in Section 2.2.1 with a discussion of least squares approximation for data points, and in Section 2.2.2 we focus on function approximation in the mean square sense.

### 2.2.1 Least squares approximation of data points

Consider  $n + 1$  distinct  $x$  values  $x_0 < \dots < x_n$ , and suppose that we know the values  $u_0, \dots, u_n$  taken by an unknown function  $u$  when evaluated at these points. We wish to approximate the function  $u$  by a function of the form

$$\widehat{u}(x) = \sum_{i=0}^m \alpha_i \varphi_i(x) \in \text{Span}\{\varphi_0, \dots, \varphi_m\}, \quad (2.27)$$

for some  $m < n$ . In many cases of practical interest, the basis functions  $\varphi_0, \dots, \varphi_m$  are polynomials. In contrast with interpolation, here we seek a function  $\hat{u}$  in a finite-dimensional function space of dimension  $m$  strictly lower than the number of data points. In order for  $\hat{u}$  to be a good approximation, we wish to find coefficients  $\alpha_0, \dots, \alpha_m$  such that the following linear system is approximately satisfied.

$$\mathbf{A}\boldsymbol{\alpha} := \begin{pmatrix} \varphi_0(x_0) & \varphi_1(x_0) & \dots & \varphi_m(x_0) \\ \varphi_0(x_1) & \varphi_1(x_1) & \dots & \varphi_m(x_1) \\ \varphi_0(x_2) & \varphi_1(x_2) & \dots & \varphi_m(x_2) \\ \vdots & \vdots & & \vdots \\ \varphi_0(x_{n-2}) & \varphi_1(x_{n-2}) & \dots & \varphi_m(x_{n-2}) \\ \varphi_0(x_{n-1}) & \varphi_1(x_{n-1}) & \dots & \varphi_m(x_{n-1}) \\ \varphi_0(x_n) & \varphi_1(x_n) & \dots & \varphi_m(x_n) \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_m \end{pmatrix} \approx \begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ \vdots \\ u_{n-2} \\ u_{n-1} \\ u_n \end{pmatrix} =: \mathbf{b}.$$

In general, since the matrix on the left-hand side has more lines than columns, there does not exist an exact solution to this equation. In order to find an approximate solution, a natural approach is to find coefficients  $\alpha_0, \dots, \alpha_m$  such that the residual vector  $\mathbf{r} = \mathbf{A}\boldsymbol{\alpha} - \mathbf{b}$  is small in some vector norm. A particularly popular approach, known as least squares approximation, is to minimize the Euclidean norm of  $\mathbf{r}$  or, equivalently, the square of the Euclidean norm:

$$\|\mathbf{r}\|^2 = \sum_{i=0}^n |u_i - \hat{u}(x_i)|^2 = \sum_{i=0}^n \left( u_i - \sum_{j=0}^m \alpha_j \varphi_j(x_i) \right)^2.$$

Let us denote the right-hand side of this equation by  $J(\boldsymbol{\alpha})$ , which we view as a function of the vector of coefficients  $\boldsymbol{\alpha}$ . A necessary condition for  $\boldsymbol{\alpha}_*$  to be a minimizer is that  $\nabla J(\boldsymbol{\alpha}_*) = 0$ . The gradient of  $J$ , written as a column vector, is given by

$$\begin{aligned} \nabla J(\boldsymbol{\alpha}) &= \nabla \left( (\mathbf{A}\boldsymbol{\alpha} - \mathbf{b})^T (\mathbf{A}\boldsymbol{\alpha} - \mathbf{b}) \right) \\ &= \nabla \left( \boldsymbol{\alpha}^T (\mathbf{A}^T \mathbf{A}) \boldsymbol{\alpha} - \mathbf{b}^T \mathbf{A} \boldsymbol{\alpha} - \boldsymbol{\alpha}^T \mathbf{A}^T \mathbf{b} + \mathbf{b}^T \mathbf{b} \right) \\ &= 2(\mathbf{A}^T \mathbf{A}) \boldsymbol{\alpha} - 2\mathbf{A}^T \mathbf{b}. \end{aligned}$$

We deduce that  $\boldsymbol{\alpha}_*$  solves the linear system

$$\mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}_* = \mathbf{A}^T \mathbf{b}, \quad (2.28)$$

where the matrix on the left-hand side is given by:

$$\mathbf{A}^T \mathbf{A} := \begin{pmatrix} \sum_{i=0}^n \varphi_0(x_i) \varphi_0(x_i) & \sum_{i=0}^n \varphi_0(x_i) \varphi_1(x_i) & \dots & \sum_{i=0}^n \varphi_0(x_i) \varphi_m(x_i) \\ \sum_{i=0}^n \varphi_1(x_i) \varphi_0(x_i) & \sum_{i=0}^n \varphi_1(x_i) \varphi_1(x_i) & \dots & \sum_{i=0}^n \varphi_1(x_i) \varphi_m(x_i) \\ \vdots & \vdots & & \vdots \\ \sum_{i=0}^n \varphi_m(x_i) \varphi_0(x_i) & \sum_{i=0}^n \varphi_m(x_i) \varphi_1(x_i) & \dots & \sum_{i=0}^n \varphi_m(x_i) \varphi_m(x_i) \end{pmatrix}.$$



Equation (2.28) is a system of  $m$  equations with  $m$  unknowns, which admits a unique solution provided that  $\mathbf{A}^T\mathbf{A}$  is full rank or, equivalently, the columns of  $\mathbf{A}$  are linearly independent. The linear equations (2.28) are known as the *normal equations*. As a side note, we mention that the solution  $\boldsymbol{\alpha}_* = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}$  coincides with the maximum likelihood estimator for  $\alpha$  under the assumption that the data is generated according to  $u_i = u(x_i) + \varepsilon_i$ , for some function  $u \in \text{Span}\{\varphi_0, \dots, \varphi_m\}$  and random noise  $\varepsilon_i \sim \mathcal{N}(0, 1)$ .

*Remark 2.2.* From equation (2.28) we deduce that

$$\mathbf{A}\boldsymbol{\alpha}_* = \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{b}.$$

The matrix  $\Pi_{\mathbf{A}} := \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T$  on the right-hand side is the orthogonal projection operator onto  $\text{col}(\mathbf{A}) \subset \mathbf{R}^n$ , the subspace spanned by the columns of  $\mathbf{A}$ . Indeed, it holds that  $\Pi_{\mathbf{A}}^2 = \Pi_{\mathbf{A}}$ , which is the defining property of a projection operator.

To conclude this section, we note that the matrix  $\mathbf{A}^+ = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T$  is a left inverse of the matrix  $\mathbf{A}$ , because  $\mathbf{A}^+\mathbf{A} = \mathbf{I}$ . It is also called the Moore–Penrose inverse or pseudoinverse of the matrix  $\mathbf{A}$ , which generalizes the usual inverse matrix. In Julia, the backslash operator silently uses the Moore–Penrose inverse when employed with a rectangular matrix. Therefore, solving the normal equations (2.28) can be achieved by just writing  $\boldsymbol{\alpha} = \mathbf{A} \setminus \mathbf{b}$ .

## 2.2.2 Mean square approximation of functions

The approach described in Section 2.2.1 can be generalized to the setting where the actual function  $u$ , rather than just discrete evaluations of it, is available. In this section, we seek an approximation of the form (2.27) such that the error  $\widehat{u}(x) - u(x)$ , measured in some function norm, is minimized. Of course, the solution to this minimization problem depends in general on the norm employed, and may in some cases not even be unique. Instead of specifying a particular norm, as done in Section 2.2.1, in this section we retain some generality and assume only that the norm is induced by an inner product on the space of real-valued continuous functions:

$$\langle \bullet, \bullet \rangle : C([a, b]) \times C([a, b]) \rightarrow \mathbf{R}. \quad (2.29)$$

In other words, we seek to minimize

$$J(\boldsymbol{\alpha}) := \|\widehat{u} - u\|^2 = \langle \widehat{u} - u, \widehat{u} - u \rangle.$$

This is again a function of the  $m + 1$  variables  $\alpha_0, \dots, \alpha_m$ . Before calculating its gradient, we rewrite the function  $J(\boldsymbol{\alpha})$  in a simpler form:

$$\begin{aligned} J(\boldsymbol{\alpha}) &= \left\langle u - \sum_{j=0}^m \alpha_j \varphi_j, u - \sum_{k=0}^m \alpha_k \varphi_k \right\rangle \\ &= \sum_{j=0}^m \sum_{k=0}^m \alpha_j \alpha_k \langle \varphi_j, \varphi_k \rangle - 2 \sum_{j=0}^m \alpha_j \langle u, \varphi_j \rangle + \langle u, u \rangle = \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\alpha} - 2\mathbf{b}^T \boldsymbol{\alpha} + \langle u, u \rangle, \end{aligned}$$

where we introduced

$$\mathbf{G} := \begin{pmatrix} \langle \varphi_0, \varphi_0 \rangle & \langle \varphi_0, \varphi_1 \rangle & \cdots & \langle \varphi_0, \varphi_m \rangle \\ \langle \varphi_1, \varphi_0 \rangle & \langle \varphi_1, \varphi_1 \rangle & \cdots & \langle \varphi_1, \varphi_m \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \varphi_m, \varphi_0 \rangle & \langle \varphi_m, \varphi_1 \rangle & \cdots & \langle \varphi_m, \varphi_m \rangle \end{pmatrix}, \quad \mathbf{b} := \begin{pmatrix} \langle u, \varphi_0 \rangle \\ \langle u, \varphi_1 \rangle \\ \vdots \\ \langle u, \varphi_m \rangle \end{pmatrix}. \quad (2.30)$$

Employing the same approach as in the previous section, we then obtain  $\nabla J(\boldsymbol{\alpha}) = \mathbf{G}\boldsymbol{\alpha} - \mathbf{b}$ , and so the minimizer of  $J(\boldsymbol{\alpha})$  is the solution to the equation

$$\mathbf{G}\boldsymbol{\alpha} = \mathbf{b}. \quad (2.31)$$

The matrix  $\mathbf{G}$ , known as the *Gram matrix*, is positive semi-definite and nonsingular provided that the basis functions are linearly independent, see [Exercise 2.9](#). Therefore, the solution  $\boldsymbol{\alpha}_*$  exists and is unique. In addition, since the Hessian of  $J$  is equal to  $\mathbf{G}$ , the vector  $\boldsymbol{\alpha}_*$  is indeed a minimizer. Note that if  $\langle \bullet, \bullet \rangle$  is defined as a finite sum of the form

$$\langle f, g \rangle = \sum_{i=0}^n f(x_i)g(x_i), \quad (2.32)$$

then (2.31) coincides with the normal equations (2.28) from the previous section. We remark that (2.32) is in fact not an inner product on the space of continuous functions, but it is an inner product on the space of polynomials of degree less than or equal to  $n$ .

In practice, the matrix and right-hand side of the linear system (2.31) can usually not be calculated exactly, because the inner product  $\langle \bullet, \bullet \rangle$  is defined through an integral; see (2.33) in the next section.

*Remark 2.3.* Rewriting the normal equations (2.31) in terms of  $\hat{u}$  we obtain

$$\langle \hat{u} - u, \varphi_0 \rangle = 0, \quad \dots, \quad \langle \hat{u} - u, \varphi_m \rangle = 0.$$

Therefore, the optimal approximation  $\hat{u} \in \text{Span}\{\varphi_0, \dots, \varphi_m\}$  satisfies

$$\forall v \in \text{Span}\{\varphi_0, \dots, \varphi_m\}, \quad \langle \hat{u} - u, v \rangle = 0.$$

This shows that the optimal approximation  $\hat{u}$ , in the sense of the norm  $\|\bullet\|$ , is the orthogonal projection of  $u$  onto  $\text{Span}\{\varphi_0, \dots, \varphi_m\}$ .

### 2.2.3 Orthogonal polynomials

The Gram matrix  $\mathbf{G}$  in (2.31) is equal to the identity matrix when the basis functions are orthonormal for the inner product considered. In this case, the solution to the linear system is

$$\alpha_i = \langle u, \varphi_i \rangle, \quad i = 0, \dots, m,$$

and so the best approximation  $\hat{u}$  (for the norm induced by the inner product considered!) is simply given by

$$\hat{u} = \sum_{i=0}^m \langle u, \varphi_i \rangle \varphi_i.$$

The coefficients  $\langle u, \varphi_i \rangle$  of the basis functions in this expansion are called *Fourier coefficients*. Given a finite dimensional subspace  $\mathcal{S}$  of the space of continuous functions, an orthonormal basis can be constructed via the Gram–Schmidt process. In this section, we focus on the particular case where  $\mathcal{S} = \mathbf{P}(n)$  – the subspace of polynomials of degree less than or equal to  $n$ . Another widely used approach, which we do not explore in this course, is to use trigonometric basis functions. We also assume that the inner product (2.29) is of the form

$$\langle f, g \rangle = \int_a^b f(x)g(x)w(x) dx, \quad (2.33)$$

where  $w(x)$  is a given nonnegative weight function such that

$$\int_a^b w(x) dx > 0.$$

Let  $\varphi_0(x), \varphi_1(x), \varphi_2(x) \dots$  denote the orthonormal polynomials obtained by applying the Gram–Schmidt procedure to the monomials  $1, x, x^2, \dots$ . These depend in general on the weight  $w(x)$  and on the interval  $[a, b]$ . A few of the popular classes of orthogonal polynomials are presented in the table below:

Name	$w(x)$	$[a, b]$
Legendre	1	$[-1, 1]$
Chebyshev	$\frac{1}{\sqrt{1-x^2}}$	$(-1, 1)$
Hermite	$\exp\left(-\frac{x^2}{2}\right)$	$[-\infty, \infty]$
Laguerre	$e^{-x}$	$[0, \infty]$

Orthogonal polynomials have a rich structure, and in the rest of this section we prove some of their key properties, one of which will be very useful in the context of numerical integration in Chapter 3. We begin by showing that orthogonal polynomials have distinct real roots.

**Proposition 2.8.** *Assume for simplicity that  $w(x) > 0$  for all  $x \in [a, b]$ , and let  $\varphi_0, \varphi_1, \dots$  denote orthonormal polynomials of increasing degree for the inner product (2.33). Then for all  $n \in \mathbf{N}$ , the polynomial  $\varphi_n$  has  $n$  distinct roots in the open interval  $(a, b)$ .*

*Proof.* Reasoning by contradiction, we assume that  $\varphi_n$  changes sign at only  $k < n$  points of the open interval  $(a, b)$ , which we denote by  $x_1, \dots, x_k$ . Then

$$\varphi_n(x) \times (x - x_0)(x - x_1) \dots (x - x_k)$$

is either everywhere nonnegative or everywhere nonpositive over  $[a, b]$ . But then

$$\int_a^b \varphi_n(x) \times (x - x_1) \dots (x - x_k) w(x) dx$$

is nonzero, which is a contradiction because the product  $(x - x_1) \dots (x - x_k)$  is a polynomial of degree  $k$ , which is orthogonal to  $\varphi_n$  by assumption. Indeed, being orthogonal to  $\varphi_0, \dots, \varphi_{n-1}$ , the polynomial  $\varphi_n$  is also orthogonal to any linear combination of these polynomials.  $\square$

Next, we show that orthogonal polynomials satisfy a three-term recurrence relation.

**Proposition 2.9.** *Assume that  $\varphi_0, \varphi_1, \dots$  are orthonormal polynomials for some inner product of the form (2.33) such that  $\varphi_i$  is of degree  $i$ . Then*

$$\forall n \in \{1, 2, \dots\}, \quad \alpha_{n+1}\varphi_{n+1}(x) = (x - \beta_n)\varphi_n(x) - \alpha_n\varphi_{n-1}(x), \quad (2.34)$$

where

$$\alpha_n = \langle x\varphi_n, \varphi_{n-1} \rangle, \quad \beta_n = \langle x\varphi_n, \varphi_n \rangle.$$

In addition,  $\alpha_1\varphi_1(x) = (x - \beta_0)\varphi_0(x)$ .

*Proof.* Since  $x\varphi_n(x)$  is a polynomial of degree  $n + 1$ , it may be decomposed as

$$x\varphi_n(x) = \sum_{i=0}^{n+1} \gamma_{n,i}\varphi_i(x). \quad (2.35)$$

Taking the inner product of both sides of this equation with  $\varphi_i$  and employing the orthonormality assumption, we obtain an expression for the coefficients:

$$\gamma_{n,i} = \langle x\varphi_n, \varphi_i \rangle.$$

From the expression (2.33) of the inner product, it is clear that  $\langle x\varphi_n, \varphi_i \rangle = \langle \varphi_n, x\varphi_i \rangle$ . Since  $x\varphi_i$  is a polynomial of degree  $i + 1$  and  $\varphi_n$  is orthogonal to all polynomials of degree strictly less than  $n$ , we deduce that  $\gamma_{n,i} = 0$  if  $i < n - 1$ . Consequently, we can rewrite the right-hand side of (2.35) as a sum involving only three terms

$$x\varphi_n(x) = \langle x\varphi_n, \varphi_{n-1} \rangle \varphi_{n-1}(x) + \langle x\varphi_n, \varphi_n \rangle \varphi_n(x) + \langle x\varphi_n, \varphi_{n+1} \rangle \varphi_{n+1}(x). \quad (2.36)$$

Since  $\langle x\varphi_n, \varphi_{n+1} \rangle = \langle x\varphi_{n+1}, \varphi_n \rangle$ , we obtain the statement after rearranging.  $\square$

*Remark 2.4.* Notice that the polynomials in Proposition 2.9 are orthonormal by assumption, and so the coefficient  $\alpha_{n+1}$  is just a normalization constant. We deduce that

$$\varphi_{n+1}(x) = \frac{(x - \beta_n)\varphi_n(x) - \alpha_n\varphi_{n-1}(x)}{\|(x - \beta_n)\varphi_n(x) - \alpha_n\varphi_{n-1}(x)\|},$$

which enables to calculate the orthogonal polynomials recursively.

### 2.2.4 Orthogonal polynomials and numerical integration: an introduction

Equation (2.36) may be rewritten in matrix form as follows:

$$\begin{pmatrix} x\varphi_0(x) \\ x\varphi_1(x) \\ x\varphi_2(x) \\ \vdots \\ x\varphi_{m-1}(x) \\ x\varphi_m(x) \end{pmatrix} = \begin{pmatrix} \beta_0 & \alpha_1 & & & & \\ \alpha_1 & \beta_1 & \alpha_2 & & & \\ & \alpha_2 & \beta_2 & \alpha_3 & & \\ & & \ddots & \ddots & \ddots & \\ & & & \alpha_{m-1} & \beta_{m-1} & \alpha_m \\ & & & & \alpha_m & \beta_m \end{pmatrix} \begin{pmatrix} \varphi_0(x) \\ \varphi_1(x) \\ \varphi_2(x) \\ \vdots \\ \varphi_{m-1}(x) \\ \varphi_m(x) \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ \alpha_{m+1}\varphi_{m+1}(x) \end{pmatrix}.$$

Let  $\mathbf{T}$  denote the matrix on the left-hand side of this equation, and let  $r_0, \dots, r_m$  denote the roots of  $\varphi_{m+1}$ . By [Proposition 2.8](#), these are distinct and all belong to the interval  $(a, b)$ . The second term on the right-hand side cancels out when  $x$  is a root of  $\varphi_{m+1}$ , and so

$$\forall r \in \{r_0, \dots, r_m\}, \quad \begin{pmatrix} r\varphi_0(r) \\ r\varphi_1(r) \\ \vdots \\ r\varphi_{m-1}(r) \\ r\varphi_m(r) \end{pmatrix} = \begin{pmatrix} \beta_0 & \alpha_1 & & & & \\ \alpha_1 & \beta_1 & \alpha_2 & & & \\ & \ddots & \ddots & \ddots & & \\ & & \alpha_{m-1} & \beta_{m-1} & \alpha_m & \\ & & & \alpha_m & \beta_m \end{pmatrix} \begin{pmatrix} \varphi_0(r) \\ \varphi_1(r) \\ \vdots \\ \varphi_{m-1}(r) \\ \varphi_m(r) \end{pmatrix}.$$

In other words, for any root  $r$  of  $\varphi_{m+1}$ , the vector  $(\varphi_0(r) \ \dots \ \varphi_m(r))^T$  is an eigenvector of the matrix  $\mathbf{T}$ , with associated eigenvalue equal to  $r$ . Since  $\mathbf{T}$  is a symmetric matrix, the eigenvectors associated with distinct eigenvalues are orthogonal for the Euclidean inner product of  $\mathbf{R}^{m+1}$ , so given that the eigenvalues of  $\mathbf{T}$  are distinct, we deduce that

$$\forall i \neq j, \quad \sum_{i=0}^m \varphi_i(r_i)\varphi_i(r_j) = 0. \quad (2.37)$$

Let us construct the matrix

$$\mathbf{P} = \begin{pmatrix} \varphi_0(r_0) & \varphi_1(r_0) & \dots & \varphi_m(r_0) \\ \varphi_0(r_1) & \varphi_1(r_1) & \dots & \varphi_m(r_1) \\ \varphi_0(r_2) & \varphi_1(r_2) & \dots & \varphi_m(r_2) \\ \vdots & \vdots & \dots & \vdots \\ \varphi_0(r_m) & \varphi_1(r_m) & \dots & \varphi_m(r_m) \end{pmatrix}.$$

Equation (2.37) indicates that the rows of  $\mathbf{P}$  are orthogonal, and so the matrix  $\mathbf{D} = \mathbf{P}\mathbf{P}^T$  is diagonal with elements given by

$$d_{ii} = \sum_{j=0}^m |\varphi_j(r_i)|^2, \quad i = 0, \dots, m.$$

(Here we start counting the rows from 0 for convenience.) Since  $\mathbf{P}\mathbf{P}^T\mathbf{D}^{-1} = \mathbf{I}$ , we deduce that the inverse of  $\mathbf{P}$  is given by  $\mathbf{P}^{-1} = \mathbf{P}^T\mathbf{D}^{-1}$ . Consequently,

$$\mathbf{P}^T\mathbf{D}^{-1}\mathbf{P} = \mathbf{P}^{-1}\mathbf{P} = \mathbf{I},$$

which means that the *columns* of  $\mathbf{P}$  are orthonormal for the inner product  $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x}^T\mathbf{D}^{-1}\mathbf{y}$ . In other words, the polynomials  $\varphi_1, \dots, \varphi_m$  are orthonormal for the inner product

$$\begin{aligned} \langle \bullet, \bullet \rangle_{m+1} : \mathbf{P}(m) \times \mathbf{P}(m) &\rightarrow \mathbf{R}; \\ (p, q) &\mapsto \sum_{i=0}^m \frac{p(r_i)q(r_i)}{d_{ii}}. \end{aligned}$$

We have thus shown that, if  $\varphi_0, \varphi_1, \varphi_2, \dots$  is a family of orthonormal polynomials for an inner product  $\langle \bullet, \bullet \rangle$ , then these are also orthonormal for the inner product  $\langle \bullet, \bullet \rangle_{m+1}$ . We reformulate our findings in the following result where, since  $m$  was arbitrary in the previous reasoning, we add a superscript to indicate when the quantities involved depend on  $m$ .

**Theorem 2.10.** *Orthonormal polynomials  $\varphi_0, \dots, \varphi_m$  for the inner product*

$$\langle f, g \rangle = \int_a^b f(x)g(x)w(x) \, dx$$

*are also orthonormal for the inner product*

$$\langle f, g \rangle_{m+1} = \sum_{i=0}^m f(r_i^{(m+1)})g(r_i^{(m+1)})w_i^{(m+1)},$$

*where  $r_0^{(m+1)}, \dots, r_m^{(m+1)}$  are the roots of  $\varphi_{m+1}$  and the weights  $w_i^{(m+1)}$  are given by*

$$w_i^{(m+1)} = \frac{1}{\sum_{j=0}^m |\varphi_j(r_i^{(m+1)})|^2}, \quad i = 0, \dots, m.$$

As an immediate corollary, we deduce that

$$\forall (p, q) \in \mathbf{P}(m) \times \mathbf{P}(m), \quad \langle p, q \rangle = \langle p, q \rangle_{m+1}, \quad (2.38)$$

Indeed, denoting by  $p = \alpha_0\varphi_0 + \dots + \alpha_m\varphi_m$  and  $q = \beta_0\varphi_0 + \dots + \beta_m\varphi_m$  the expansions of the polynomials  $p$  and  $q$  in the orthonormal basis, we have

$$\begin{aligned} \langle p, q \rangle &= \langle \alpha_0\varphi_0 + \dots + \alpha_m\varphi_m, \beta_0\varphi_0 + \dots + \beta_m\varphi_m \rangle \\ &= \sum_{i=0}^m \sum_{j=0}^m \alpha_i\beta_j \langle \varphi_i, \varphi_j \rangle = \alpha_0\beta_0 + \dots + \alpha_m\beta_m = \sum_{i=0}^m \sum_{j=0}^m \alpha_i\beta_j \langle \varphi_i, \varphi_j \rangle_{m+1} \\ &= \langle \alpha_0\varphi_0 + \dots + \alpha_m\varphi_m, \beta_0\varphi_0 + \dots + \beta_m\varphi_m \rangle_{m+1} = \langle p, q \rangle_{m+1}. \end{aligned}$$

To conclude this section, we prove the following statement, which is another consequence of [Theorem 2.10](#) and has applications to numerical integration.

**Theorem 2.11.** *It holds that*

$$\forall p \in \mathbf{P}(2m+1), \quad \int_a^b p(x) w(x) dx = \sum_{i=0}^m p(r_i^{(m+1)}) w_i^{(m+1)}. \quad (2.39)$$

*Proof.* Taking  $q = 1$  in (2.38) and employing the definitions of  $\langle \bullet, \bullet \rangle$  and  $\langle \bullet, \bullet \rangle_{m+1}$ , we have that (2.38) is satisfied for any  $p \in \mathbf{P}(m)$ . Next, any polynomial  $p \in \mathbf{P}(2m+1)$  may be decomposed as  $p(x) = \varphi_{m+1}(x)q(x) + \rho(x)$ , for some polynomial  $q$  of degree  $m$  (the quotient of the polynomial division of  $p$  by  $\varphi_{m+1}$ ) and some polynomial  $\rho$  of degree lower than or equal to  $m$  (the remainder of the polynomial division). Therefore, since (2.39) was already shown to hold for polynomials of degree up to  $m$ , we obtain

$$\begin{aligned} \int_a^b p(x)w(x) dx &= \int_a^b \varphi_{m+1}(x)q(x)w(x) dx + \int_a^b \rho(x)w(x) dx \\ &= 0 + \int_a^b \rho(x)w(x) dx = 0 + \sum_{i=0}^m \rho(r_i) w_i \\ &= \sum_{i=0}^m \varphi_{m+1}(r_i) q(r_i) w_i + \sum_{i=0}^m \rho(r_i) w_i = \sum_{i=0}^m p(r_i) w_i, \end{aligned}$$

where we dropped the  $(m+1)$  superscript for conciseness and we used, in the penultimate inequality, the fact that  $r_0, \dots, r_m$  are the roots of the polynomial  $\varphi_{m+1}$ .  $\square$

Since the left-hand side of (2.39) is an integral and the right-hand side is a sum, we have just constructed an integration formula, which enjoys a very nice property: it is exact for polynomials of degree up to  $2m+1$ ! A formula of this type is called a *quadrature formula*, with  $m+1$  nodes  $r_0^{(m+1)}, \dots, r_m^{(m+1)}$  and associated weights  $w_0^{(m+1)}, \dots, w_m^{(m+1)}$ . Note that the nodes and weights of the quadrature depend on the weight  $w(x)$  and on the degree  $m$ . We will revisit this subject in Chapter 3.

## 2.3 Exercises

⚙️ **Exercise 2.1.** Find the polynomial  $p(x) = ax + b$  (a straight line) that goes through the points  $(x_0, u_0)$  and  $(x_1, u_1)$ .

⚙️ **Exercise 2.2.** Find the polynomial  $p(x) = ax^2 + bx + c$  (a parabola) that goes through the points  $(0, 1)$ ,  $(1, 3)$  and  $(2, 7)$ .

⚙️ **Exercise 2.3.** Prove the following recurrence relation for Chebyshev polynomials:

$$T_{i+1}(x) = 2xT_i(x) - T_{i-1}(x), \quad i = 1, 2, \dots$$

⚙️ **Exercise 2.4.** Show by recursion that

$$[u_0, u_1, \dots, u_n] = \sum_{j=0}^n \frac{u_j}{\prod_{k \in \{0, \dots, n\} \setminus \{j\}} (x_j - x_k)}. \quad (2.40)$$

Deduce from this identity that

$$[u_0, u_1, \dots, u_n] = [u_{\sigma_1}, u_{\sigma_2}, \dots, u_{\sigma_n}],$$

for any permutation  $\sigma$  of  $(0, 1, 2, \dots, n)$ .

*Solution.* The first statement (2.40) is clear when  $n = 0$ . Reasoning by induction, we assume that the statement is true up to  $n - 1$  and prove that it then also holds for  $n$ . Using the definition (2.10) and the induction hypothesis, we obtain that

$$\begin{aligned} [u_0, u_1, \dots, u_n] &= \frac{[u_1, \dots, u_n] - [u_0, \dots, u_{n-1}]}{x_n - x_0} \\ &= \frac{1}{x_n - x_0} \left( \sum_{j=1}^n \frac{u_j}{\prod_{k \in \{1, \dots, n\} \setminus \{j\}} (x_j - x_k)} - \sum_{j=0}^{n-1} \frac{u_j}{\prod_{k \in \{0, \dots, n-1\} \setminus \{j\}} (x_j - x_k)} \right) \end{aligned}$$

Rewriting the fractions with a common denominator leads to

$$[u_0, u_1, \dots, u_n] = \frac{1}{x_n - x_0} \sum_{j=0}^n \frac{u_j ((x_j - x_0) - (x_j - x_n))}{\prod_{k \in \{0, \dots, n\} \setminus \{j\}} (x_j - x_k)} = \sum_{j=0}^n \frac{u_j}{\prod_{k \in \{0, \dots, n\} \setminus \{j\}} (x_j - x_k)},$$

which concludes the proof of the first statement. The second statement then follows immediately, because the right-hand side of (2.40) is invariant under permutations.  $\triangle$

**Exercise 2.5.** Using the Gregory–Newton formula, find an expression for

$$\sum_{i=1}^n i^4.$$

**Exercise 2.6.** Let  $(f_0, f_1, f_2, \dots) = (1, 1, 2, \dots)$  denote the Fibonacci sequence. Prove that there does not exist a polynomial  $p$  such that

$$\forall n \in \mathbf{N}, \quad f_n = p(n). \quad (2.41)$$

*Solution.* Assume by contradiction that  $p: \mathbf{R} \rightarrow \mathbf{R}$  is a polynomial such that (2.41) is satisfied, and let  $n$  be the degree of this polynomial. Then it holds that  $\Delta^{n+1}p = 0$  (2.5), where both sides are viewed as functions from  $\mathbf{R}$  to  $\mathbf{R}$ . On the other hand, since  $p(n) = f_n$  for all  $n \in \mathbf{N}$ , we can calculate explicitly the values of taken by the function  $\Delta^m p$  when evaluated at all the natural numbers, for all  $m \in \mathbf{N}$ . We collate a few values in the following table.

$n$	0	1	2	3	4	5	6
$\Delta^0 p(n)$	1	1	2	3	5	8	13
$\Delta^1 p(n)$	0	1	1	2	3	5	8
$\Delta^2 p(n)$	1	0	1	1	2	3	5
$\Delta^3 p(n)$	-1	1	0	1	1	2	3

It appears from these calculations that the Fibonacci sequence is shifted one position to the right with



each additional application of  $\Delta$ . In other words, our calculations suggest that

$$\forall(m, n) \in \mathbf{N} \times \mathbf{N}, \quad \Delta^m p(m+n) = f_n, \quad (2.42)$$

which is a contradiction. To conclude, let us prove (2.42) rigorously. This equation is obvious for  $m = 0$  by assumption. Now, reasoning by contradiction, we assume that (2.42) is true up to  $m$ . Then by definition of the difference operator  $\Delta$ , we have

$$\begin{aligned} \Delta^{m+1} p(m+n+1) &= \Delta^m p(m+n+2) - \Delta^m p(m+n+1) \\ &= f_{n+2} - f_{n+1} = f_n. \end{aligned}$$

Here, we used the induction hypothesis (2.42) in the second equality, and the definition of the Fibonacci series in the third one.  $\triangle$

**⚙️ Exercise 2.7.** Using the Gregory–Newton formula, show that

$$\forall n \in \mathbf{N}, \quad 2^n = 1 + n + \frac{n^2}{2!} + \frac{n^3}{3!} + \frac{n^4}{4!} + \dots \quad (2.43)$$

*Solution.* Equation (2.43) is a particular case of the following more general statement: for any function  $f \in \mathbf{R} \rightarrow \mathbf{R}$ , it holds that

$$\forall n \in \mathbf{N}, \quad f(n) = f(0) + \Delta f(0)n + \Delta^2 f(0) \frac{n^2}{2!} + \Delta^3 f(0) \frac{n^3}{3!} + \Delta^4 f(0) \frac{n^4}{4!} + \dots \quad (2.44)$$

In order to show this equation, it is sufficient to prove that for any  $n_* \in \mathbf{N}$ , the two sides of (2.44) coincide for every  $n \in \{0, \dots, n_*\}$ . Since  $n^p = 0$  for all  $n \in \{0, \dots, p-1\}$  by definition (2.4) of the falling powers, the right-hand side of (2.44) coincides for all  $n \in \{0, \dots, n_*\}$  with

$$g(n) = f(0) + \Delta f(0)n + \Delta^2 f(0) \frac{n^2}{2!} + \dots + \Delta^{n_*} f(0) \frac{n^{n_*}}{n_*!}.$$

We recognize on the right-hand side Newton's expression of the interpolating polynomial through the points  $(0, f(0)), \dots, (n_*, f(n_*))$ , and so  $g(n) = f(n)$  for all  $n \in \{0, \dots, n_*\}$ , which concludes the proof.  $\triangle$

*Remark 2.5.* Remarkably, equation (2.43) holds in fact for any  $n \in \mathbf{R}_{>0}$ . However, showing this more general statement is beyond the scope of this course.

**⚙️ Exercise 2.8.** Prove [Theorem 2.7](#).

**⚙️ Exercise 2.9.** Show that the matrix  $\mathbf{G}$  in (2.30) is positive definite if the basis functions  $\varphi_0, \dots, \varphi_m$  are linearly independent.

**□ Exercise 2.10.** Write a Julia code for interpolating the following function using a polynomial of degree 20 over the interval  $[-1, 1]$ .

$$f(x) = \tanh\left(\frac{x+1/2}{\varepsilon}\right) + \tanh\left(\frac{x}{\varepsilon}\right) + \tanh\left(\frac{x-1/2}{\varepsilon}\right), \quad \varepsilon = .01.$$

Use equidistant and then Chebyshev nodes, and compare the two approaches in terms of accuracy. Plot the function  $f$  together with the approximating polynomials.

□ **Exercise 2.11.** Write from scratch a function to obtain the polynomial interpolating the data points

$$(x_0, u_0), \dots, (x_n, u_n).$$

Your function should return the values taken by the interpolating polynomial when evaluated at the points  $X_0, \dots, X_m$ . You may use the following code to test your function

```
import Plots
function interp(X, x, u)
    # Your code comes here
end
n, m = 10, 100
f(t) = cos(2π * t)
x = LinRange(0, 1, n)
X = LinRange(0, 1, m)
u = f.(x)
U = interp(X, x, u)
Plots.plot(X, f.(X), label="Original function")
Plots.plot!(X, U, label="Interpolation")
Plots.scatter!(x, u, label="Data")
```

□ **Exercise 2.12.** We wish to use interpolation to approximate the following parametric function, called an epitrochoid:

$$x(\theta) = (R + r) \cos \theta + d \cos \left( \frac{R + r}{r} \theta \right) \quad (2.45)$$

$$y(\theta) = (R + r) \sin \theta - d \sin \left( \frac{R + r}{r} \theta \right), \quad (2.46)$$

with  $R = 5$ ,  $r = 2$  and  $d = 3$ , and for  $\theta \in [0, 4\pi]$ . Write a Julia program to interpolate  $x(\theta)$  and  $y(\theta)$  using 40 equidistant points. Use the **BigFloat** format in order to reduce the impact of round-off errors. After constructing the polynomial interpolations  $\hat{x}(\theta)$  and  $\hat{y}(\theta)$ , plot the parametric curve  $\theta \mapsto (\hat{x}(\theta), \hat{y}(\theta))$ . Your plot should look similar to [Figure 2.6](#).

□ **Exercise 2.13** (Solving the Laplace equation using a spectral method). The classical Laplace equation with homogeneous Dirichlet boundary conditions in dimension 1 reads

$$\text{Find } u \in C^2([0, 1]) \text{ such that } \begin{cases} -u''(x) = f(x) & \forall x \in (0, 1), \\ u(0) = u(1) = 0. \end{cases} \quad (2.47)$$

Our goal in this exercise is to approximate the exact solution  $u(x)$  using interpolation. Specifically, we propose to proceed in two steps:

- Interpolate the right-hand side using a polynomial with equidistant nodes. That is, find a

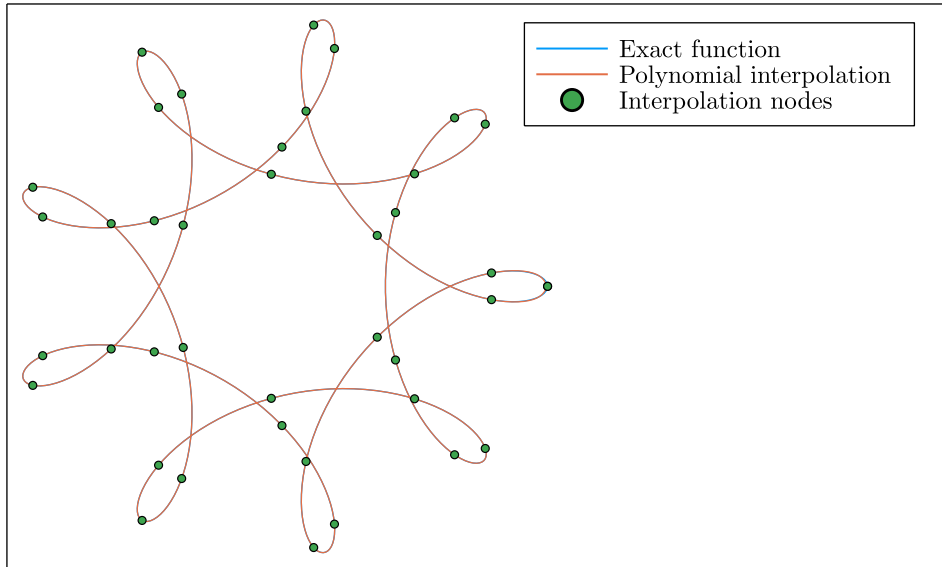


Figure 2.6: Solution for Exercise 2.12.

polynomial  $\hat{f} \in \mathbf{P}(n)$  such that

$$\forall i \in \{0, \dots, n\}, \quad \hat{f}(x_i) = f(x_i), \quad x_i = \frac{i}{n}.$$

- Solve (2.47) with  $\hat{f}$  instead of  $f$ . Since  $\hat{f}$  is a polynomial, this can be achieved analytically.

Implement this program in the case where

$$f(x) = \exp(\sin(2\pi x)) \cos(2\pi x)^2 - \exp(\sin(2\pi x)) \sin(2\pi x),$$

and compare for various values of  $n$  the approximate solution you obtain with the exact solution to (2.47), which is given by  $u(x) = (2\pi)^{-2} \exp((\sin(2\pi x)) - 1)$  in this case.

□ **Exercise 2.14** (Modeling the vapor pressure of mercury). The dataset loaded through the following Julia commands contains data on the vapor pressure of mercury as a function of the temperature.

```
import RDatasets
data = RDatasets.dataset("datasets", "pressure")
```

Find a low-dimensional mathematical model of the form

$$p(T) = \exp(\alpha_0 + \alpha_1 T + \alpha_2 T^2 + \alpha_3 T^3) \quad (2.48)$$

for the pressure as a function of the temperature. Plot the approximation together with the data. An example solution is given in Figure 2.7.

□ **Exercise 2.15.** Let  $u: [0, 2\pi] \rightarrow \mathbf{R}$  and

$$x_k = \left( \frac{2k\pi}{2n+1} \right), \quad k = 0, \dots, 2n. \quad (2.49)$$

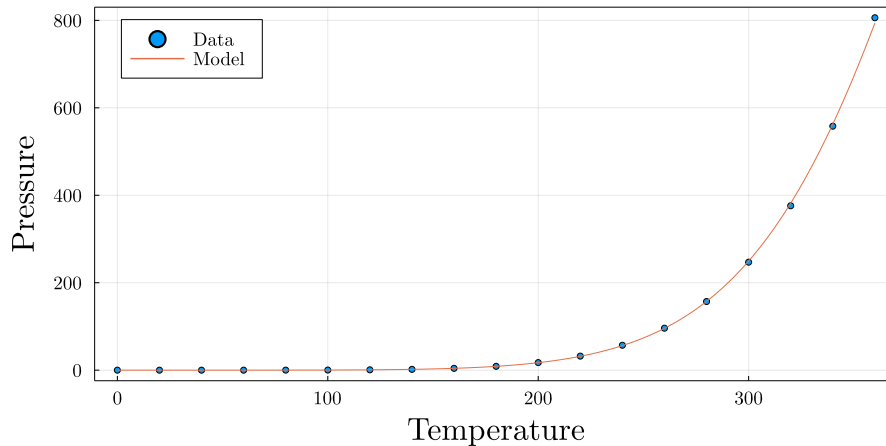


Figure 2.7: Solution for Exercise 2.14.

We wish to interpolate  $u$  at these nodes using complex exponentials:

$$\hat{u} = \sum_{k=-n}^n a_k e^{ikx}.$$

Write a function

```
function fourier_interpolate(u, x, X)
    # Your code comes here ...
end
```

which takes three arguments:

- $u$  is the function to interpolate;
- $x$  are the interpolation nodes, given by (2.49) in the test code below; you can assume that this array contains an odd number of elements.
- $X$  is a one-dimensional array of values on the  $x$  axis.

The function should return a one-dimensional array containing the values that  $\hat{u}$  takes when evaluated at the points contained in  $X$ . You can use the following code to test your function:

```
import Plots
n, m = 5, 1000
x = 2π/(2n+1) * (0:2n)
X = 2π/m * (0:m)

u(x) = sign(x - π)
# u(x) = exp(sin(x) + cos(5x))
# u(x) = x^2 * (x - 2π)^2 / π^4

@time U = fourier_interpolate(u, x, X)
Plots.plot(X, u.(X), label="u(x)", legend=:bottomright)
```

```

Plots.plot!(X, U, label="û(x)")
Plots.scatter!(x, u.(x), label="Interpolation points")
Plots.xlims!(0, 2π)

```

An example of the output plot is illustrated in Figure 2.8.

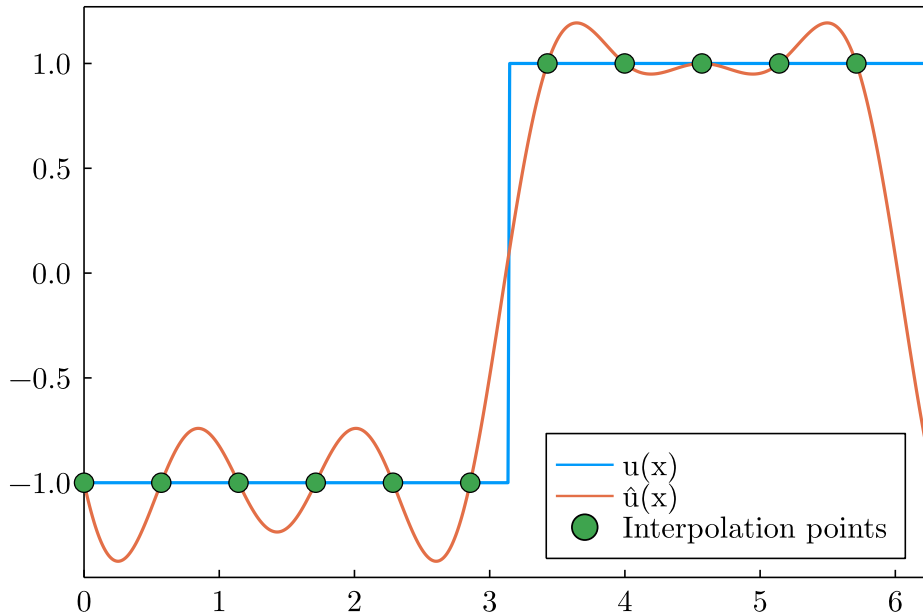


Figure 2.8: Example solution for Exercise 2.15.

## 2.4 Discussion and bibliography

A comprehensive study of approximation theory would require to cover the  $L^\infty$  setting as well as other functional settings. A pillar of  $L^\infty$  approximation theorem is Chebyshev's equioscillation theorem, which we alluded to in Remark 2.1. An excellent introductory reference on approximation theory is [8] (in French). See also [10, Chapter 10] and the references therein.

# Chapter 3

## Numerical integration

3.1	The closed Newton–Cotes method . . . . .	56
3.2	Composite methods with equidistant nodes . . . . .	57
3.3	Richardson extrapolation and Romberg’s method . . . . .	63
3.4	Methods with non-equidistant nodes . . . . .	67
3.5	Introduction to probabilistic integration methods . . . . .	71
3.6	Exercises . . . . .	73
3.7	Discussion and bibliography . . . . .	79

### Introduction

Integrals are ubiquitous in science and mathematics. In this chapter, we are concerned with the problem of calculating numerically integrals of the form

$$I = \int_{\Omega} u(\mathbf{x}) \, d\mathbf{x}, \tag{3.1}$$

Perhaps somewhat surprisingly, the numerical calculation of such integrals when  $n \gg 1$  is still a very active area of research today. In this chapter, however, we will focus for simplicity on the one-dimensional setting where  $\Omega = [a, b] \subset \mathbf{R}$ . We assume throughout this chapter that the function  $u$  is Riemann-integrable. Then, by definition,

$$I = \lim_{h \rightarrow 0} \sum_{i=0}^{n-1} u(t_i)(z_{i+1} - z_i),$$

where  $a = z_0 < \dots < z_n = b$  is a partition of the interval  $[a, b]$  such that the maximum spacing between successive  $x$  values is equal to  $h$ , and with  $t_i \in [x_i, x_{i+1}]$  for all  $i \in \{0, \dots, n - 1\}$ .

All the numerical integration formulas that we present in this chapter are based on a deter-

ministic approximation of the form

$$\widehat{I} = \sum_{i=0}^n w_i u(x_i), \quad (3.2)$$

where  $x_0 < \dots < x_n$  are the *integration points* and  $w_0, \dots, w_n$  are the *integration weights*. In many cases, integration formulas contain a small parameter that can be refined to improve the accuracy of the approximation. In methods based on equidistant interpolation nodes, for example, this parameter encodes the distance between nodes and is typically denoted by  $h$ . We shall often use the notation  $\widehat{I}_h$  to emphasize the dependence of the approximation on  $h$ . The difference  $E_h = I - \widehat{I}_h$  is called the *integration error* or *discretization error*. The *degree of precision*, defined hereafter, is an important measure of the quality of quadrature rule.

**Definition 3.1.** The *degree of precision* of an integration method is the smallest integer number  $d$  such that the integration error is zero for all  $u \in \mathbf{P}(d)$ , i.e. for all the polynomials of degree less than or equal to  $d$ .

We observe that, without loss of generality, we can assume that the integration interval is equal to  $[-1, 1]$ . Indeed, using the change of variable

$$\begin{aligned} \zeta: [-1, 1] &\rightarrow [a, b]; \\ y &\mapsto \frac{b+a}{2} + \frac{b-a}{2}y, \end{aligned} \quad (3.3)$$

we have

$$\int_a^b u(x) dx = \int_{-1}^1 u(\zeta(y)) \zeta'(y) dy = \frac{b-a}{2} \int_{-1}^1 u \circ \zeta(y) dy, \quad (3.4)$$

and the right-hand side is the integral of  $u \circ \zeta$  over the interval  $[-1, 1]$ .

### 3.1 The closed Newton–Cotes method

Given a set of equidistant points  $-1 = x_0 < \dots < x_n = 1$ , a natural method for approximating the integral (3.1) of a function  $u: [-1, 1] \rightarrow \mathbf{R}$  is to first construct the interpolating polynomial  $\widehat{u}$  at the nodes, and then calculate the exact integral of this polynomial. By construction, this method is exact for polynomials of degree up to  $n$ , and so the degree of precision is equal to *at least*  $n$ . Let  $\varphi_0, \dots, \varphi_n$  denote the Lagrange polynomials associated with the integration nodes. Then we have

$$I \approx \int_{-1}^1 \widehat{u}(x) dx = \int_{-1}^1 \sum_{i=0}^n u(x_i) \varphi_i(x) dx = \sum_{i=0}^n u(x_i) \underbrace{\int_{-1}^1 \varphi_i(x) dx}_{w_i}.$$

The weights are independent of the function  $u$ , and so they can be calculated a priori. The class of integration methods obtained using this approach are known as *Newton–Cotes methods*. We present a few particular cases:

- $n = 1, d = 1$  (trapezoidal rule):

$$\int_{-1}^1 u(x) dx \approx u(-1) + u(1). \quad (3.5)$$

- $n = 2, d = 3$  (Simpson's rule):

$$\int_{-1}^1 u(x) dx \approx \frac{1}{3}u(-1) + \frac{4}{3}u(0) + \frac{1}{3}u(1). \quad (3.6)$$

- $n = 3, d = 3$  (Simpson's  $\frac{3}{8}$  rule):

$$\int_{-1}^1 u(x) dx \approx \frac{1}{4}u(-1) + \frac{3}{4}u(-1/3) + \frac{3}{4}u(1/3) + \frac{1}{4}u(1).$$

- $n = 4, d = 5$  (Boole's rule):

$$\int_{-1}^1 u(x) dx \approx \frac{7}{45}u(-1) + \frac{32}{45}u\left(-\frac{1}{2}\right) + \frac{12}{45}u(0) + \frac{32}{45}u\left(\frac{1}{2}\right) + \frac{7}{45}u(1).$$

*Remark 3.1.* Note that, although it is based on a quadratic polynomial interpolation, Simpson's rule (3.6) has a degree of precision equal to 3. This is because any integration rule with nodes and weights symmetric around  $x = 0$  is exact for odd functions, in particular  $x^3$ . Likewise, the degree of precision of Boole's rule is equal to 5.

In principle, this approach could be employed in order to construct integration rules of arbitrary high degree of precision. In practice, however, the weights become more and more imbalanced as the number of interpolation points increases, with some of them becoming negative. As a result, roundoff errors become increasingly detrimental to accuracy. In addition, in cases where the interpolating polynomial does not converge to  $u$ , for example if  $u$  is Runge's function, the approximate integral may not converge to the correct value in the limit as  $n \rightarrow \infty$ , even in exact arithmetic!

The integration rules presented in this section, which are based on equidistant nodes that include the endpoints of the integration interval, are called *closed Newton-Cotes* methods. A similar approach can be employed in order to construct to integration rules based on equidistant nodes that do not include the endpoints; these are called *open Newton-Cotes* methods.

## 3.2 Composite methods with equidistant nodes

A natural alternative to the approach presented in Section 3.1 is to construct an integration rule using piecewise polynomial interpolation, which we studied in Section 2.1.7. After partitioning the integration interval in a number of subintervals, the integral can be approximated by using one of the rules presented in Section 3.1 within each subinterval.



**Composite trapezoidal rule.** Let us illustrate the composite approach with an example. To this end, we introduce a partition  $a = x_0 < \dots < x_n = b$  of the interval  $[a, b]$  and assume that the nodes are equidistant with  $x_{i+1} - x_i = h$ . Using (3.4) with  $a = x_i$  and  $b = x_{i+1}$ , we first generalize (3.5) to an interval  $[x_i, x_{i+1}]$  as follows:

$$\int_{x_i}^{x_{i+1}} u(x) dx = \frac{h}{2} \int_{-1}^1 u \circ \zeta(y) dy \approx \frac{h}{2} (u \circ \zeta(-1) + u \circ \zeta(1)) = \frac{h}{2} (u(x_i) + u(x_{i+1})),$$

where  $\approx$  in this equation indicates approximation using the trapezoidal rule. Applying this approximation to each subinterval of the partition, we obtain the composite trapezoidal rule:

$$\begin{aligned} \int_a^b u(x) dx &= \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} u(x) dx \approx \frac{h}{2} \sum_{i=0}^{n-1} (u(x_i) + u(x_{i+1})) \\ &= \frac{h}{2} (u(x_0) + 2u(x_1) + 2u(x_2) + \dots + 2u(x_{n-2}) + 2u(x_{n-1}) + u(x_n)). \end{aligned} \quad (3.7)$$

Like the trapezoidal rule (3.5), the composite trapezoidal rule (3.7) has a degree of precision equal to 1. However, the integration error of the method depends on the parameter  $h$ , which represents the width of each subinterval: for very small  $h$ , equation (3.7) is expected to provide a good approximation of the integral. An error estimate can be obtained directly from the formula in Theorem 2.3 for the interpolation error, provided that we assume that  $u \in C^2[a, b]$ .

**Theorem 3.1** (Integration error for the composite trapezoidal rule). *Let  $\widehat{I}_h$  denote the approximate integral calculated using (3.7). Then*

$$|I - \widehat{I}_h| \leq \frac{b-a}{12} C_2 h^2, \quad C_2 := \sup_{\xi \in [a, b]} |u''(\xi)|. \quad (3.8)$$

*Proof.* Denoting by  $\widehat{u}_h$  the piecewise linear interpolation of  $u$ , we have

$$\int_{x_i}^{x_{i+1}} u(x) - \widehat{u}_h(x) dx = \frac{1}{2} \int_{x_i}^{x_{i+1}} u''(\xi(x))(x - x_i)(x - x_{i+1}) dx.$$

Since  $(x - x_i)(x - x_{i+1})$  is nonpositive over the interval  $[x_i, x_{i+1}]$ , we deduce that

$$\left| \int_{x_i}^{x_{i+1}} u(x) - \widehat{u}_h(x) dx \right| \leq \frac{1}{2} \left( \sup_{\xi \in [a, b]} |u''(\xi)| \right) \int_{x_i}^{x_{i+1}} (x - x_i)(x_{i+1} - x) dx = C_2 \frac{h^3}{12}.$$

Summing the contributions of all the intervals, we obtain

$$|I - \widehat{I}_h| \leq \sum_{i=0}^{n-1} \left| \int_{x_i}^{x_{i+1}} u(x) - \widehat{u}_h(x) dx \right| \leq n \times C_2 \frac{h^3}{12} = \frac{b-a}{12} C_2 h^2,$$

which concludes the proof.  $\square$

The integration error therefore scales as  $\mathcal{O}(h^2)$ . (Strictly speaking, we have shown only that the integration error admits an upper bound that scales at  $\mathcal{O}(h^2)$ , but it turns out that the dependence on  $h$  of this bound is optimal).

**Example 3.1** (Integration error for the composite trapezoidal rule). Let us approximate the integral of  $x^2$  over  $[0, 1]$  via the trapezoidal rule. In other words, let us calculate the right-hand side of

$$\int_0^1 x^2 dx \approx \frac{h}{2} (x_0^2 + 2x_1^2 + \dots + 2x_{n-1}^2 + x_n^2),$$

which we shall denote  $\widehat{I}_h$  hereafter. Using  $x_i = ih$ , and  $x_n = 1$ , we obtain

$$\widehat{I}_h = \frac{h}{2} \left( 2h^2 \sum_{i=1}^{n-1} i^2 + 1 \right) = \frac{h}{2} \left( \frac{h^2}{3} (n-1)n(2n-1) + 1 \right) = \frac{h}{2} \left( \frac{h^2}{3} (2n^3 - 3n^2 + n) + 1 \right).$$

Now, observing that  $n = 1/h$  yields

$$\widehat{I}_h = \frac{h}{2} \left( \frac{h^2}{3} (2h^{-3} - 3h^{-2} + h^{-1}) + 1 \right) = \frac{h}{2} \left( \frac{2}{3}h^{-1} - 1 + \frac{h}{3} + 1 \right) = \frac{1}{3} + \frac{h^2}{6},$$

which is exactly the error shown by [Theorem 3.1](#).

**Composite Simpson rule.** The composite Simpson rule is derived in [Exercise 3.2](#). Given an odd number  $n + 1$  of equidistant points  $a = x_0 < x_1 < \dots < x_n = b$ , this rule is given by

$$\widehat{I}_h = \frac{h}{3} \left( u(x_0) + 4u(x_1) + 2u(x_2) + 4u(x_3) + 2u(x_4) + \dots + 2u(x_{n-2}) + 4u(x_{n-1}) + u(x_n) \right). \quad (3.9)$$

This approximation is obtained by integrating the piecewise quadratic interpolant over a partition of the integration interval into  $n/2$  subintervals of equal width. Obtaining an optimal error estimate, in terms of the dependence on  $h$ , for this integration formula is slightly more involved.

**Theorem 3.2** (Integration error for the composite Simpson rule). *Let  $\widehat{I}_h$  denote the approximate integral calculated using (3.9). Then*

$$|I - \widehat{I}_h| \leq (b - a) \frac{C_4 h^4}{180}, \quad C_4 := \sup_{\xi \in [a, b]} |u^{(4)}(\xi)|. \quad (3.10)$$

*Proof.* For a given subinterval  $[x_{2i}, x_{2i+2}]$ , let us denote by  $\widehat{u}_2$  the quadratic interpolating polynomial at  $x_{2i}, x_{2i+1}, x_{2i+2}$ , and by  $\widehat{u}_3(\bullet; \alpha)$  the cubic interpolating polynomial relative to the nodes  $x_{2i}, x_{2i+1}, x_{2i+2}, \alpha$ , for some  $\alpha \in [x_{2i}, x_{2i+1}]$  that does not coincide with the integration nodes. We have

$$\int_{x_{2i}}^{x_{2i+2}} u(x) - \widehat{u}_2(x) dx = \int_{x_{2i}}^{x_{2i+2}} u(x) - \widehat{u}_3(x; \alpha) dx + \int_{x_{2i}}^{x_{2i+2}} \widehat{u}_3(x; \alpha) - \widehat{u}_2(x) dx. \quad (3.11)$$

The second term on the right-hand side is zero, because the integrand is a cubic polynomial with zeros at  $x_{2i}, x_{2i+1}$  and  $x_{2i+2}$ , and because

$$\int_{x_{2i}}^{x_{2i+2}} (x - x_{2i})(x - x_{2i+1})(x - x_{2i+2}) dx = 0.$$

By [Theorem 2.3](#), the first term in [\(3.11\)](#) is bounded from above as follows:

$$\begin{aligned} \left| \int_{x_{2i}}^{x_{2i+2}} u(x) - \widehat{u}_3(x; \alpha) \, dx \right| &\leq \int_{x_{2i}}^{x_{2i+2}} \left| \frac{u^{(4)}(\xi(x))}{24} (x - x_{2i})(x - x_{2i+1})(x - x_{2i+2})(x - \alpha) \right| \, dx \\ &\leq \frac{C_4}{24} \int_{x_{2i}}^{x_{2i+2}} |(x - x_{2i})(x - x_{2i+1})(x - x_{2i+2})(x - \alpha)| \, dx. \end{aligned}$$

This inequality is valid for any  $\alpha \in A := [x_{2i}, x_{2i+2}] \setminus \{x_{2i}, x_{2i+1}, x_{2i+2}\}$ . Denoting by  $h(\alpha)$  the integral on the right-hand side, we observe that

$$\lim_{\alpha \rightarrow x_{2i+1}} h(\alpha) = \int_{x_{2i}}^{x_{2i+2}} (x - x_{2i})(x - x_{2i+1})^2(x_{2i+2} - x) \, dx = \frac{4}{15}h^5.$$

Therefore, we conclude that

$$\left| \int_{x_{2i}}^{x_{2i+2}} u(x) - \widehat{u}_2(x) \, dx \right| \leq \inf_{\alpha \in A} \left| \int_{x_{2i}}^{x_{2i+2}} u(x) - \widehat{u}_3(x; \alpha) \, dx \right| \leq \frac{C_4}{90}h^5.$$

Summing the contributions of all the subintervals, we finally obtain

$$|I - \widehat{I}_h| \leq \frac{n}{2} \times \frac{C_4 h^5}{90} = (b - a) \frac{C_4 h^4}{180}, \quad (3.12)$$

which concludes the proof.  $\square$

*Remark 3.2.* The cancellation of the second term in [\(3.11\)](#) also follows from the fact that the degree of precision of the Simpson rule [\(3.6\)](#) is equal to 3, and so

$$\int_{x_{2i}}^{x_{2i+2}} \widehat{u}_3(x) - \widehat{u}_2(x) \, dx = \frac{1}{3}(\widehat{u}_3 - \widehat{u}_2)(x_{2i}) + \frac{4}{3}(\widehat{u}_3 - \widehat{u}_2)(x_{2i+1}) + \frac{1}{3}(\widehat{u}_3 - \widehat{u}_2)(x_{2i+2}) = 0,$$

where we used the short-hand notation  $\widehat{u}_3(x) = \widehat{u}_3(x; \alpha)$ .

**General composite quadrature rules.** In view of [\(3.4\)](#), any quadrature rule with  $N + 1$  integration points of the form

$$\int_{-1}^1 u(x) \, dx \approx \sum_{j=0}^N w_j u(x_j) \quad (3.13)$$

admits a composite version for an arbitrary integration interval  $[a, b]$  obtained by applying the rule locally in  $M$  equally-sized subintervals:

$$\int_a^b u(x) \, dx \approx \frac{h}{2} \sum_{i=0}^{M-1} \sum_{j=0}^N w_j u \left( a + ih + \frac{h}{2} + \frac{x_j h}{2} \right), \quad h := \frac{b - a}{M}. \quad (3.14)$$

Clearly, the single-interval rule [\(3.13\)](#) and the composite rule [\(3.14\)](#) have the same degree of precision. To conclude this section, we prove a relation between this degree of precision of a rule and the convergence rate with respect to  $h$  of the composite rule [\(3.14\)](#).

**Theorem 3.3.** *Suppose that the degree of precision of the single-interval rule (3.13) is equal to  $d$ . Then for any  $u \in C^{(d+1)}[a, b]$ , there is  $c > 0$  such that*

$$\forall M \geq 1, \quad \left| I[u] - \widehat{I}_h[u] \right| \leq ch^{d+1}, \quad c := \frac{C_{d+1}(b-a)}{(d+1)!} \left( 1 + \frac{1}{2} \sum_{i=0}^J |w_j| \right). \quad (3.15)$$

Here  $I[u]$  and  $\widehat{I}_h[u]$  denote respectively the left-hand side and the right-hand side of (3.14), and  $C_{d+1}$  is defined in (3.17).

*Proof.* Fix  $i \in \{0, \dots, M-1\}$ , and let  $p_i$  denote the Taylor expansion of degree  $d$  of the function  $u$  around  $z_i := a + ih$ :

$$p_i(x) = u(z_i) + u'(z_i)(x - z_i) + \dots + \frac{u^{(d)}(z_i)}{d!}(x - z_i)^d. \quad (3.16)$$

By the mean-value form of the remainder, it is simple to prove that

$$\forall x \in [z_i, z_{i+1}], \quad |u(x) - p_i(x)| = \frac{C_{d+1}}{(d+1)!} h^{d+1}, \quad C_{d+1} = \sup_{x \in [a, b]} |u^{(d+1)}(x)|. \quad (3.17)$$

Let us define the local contributions to the total integral as follows:

$$I^{(i)}[u] := \int_{z_i}^{z_{i+1}} u(x) dx, \quad \widehat{I}_h^{(i)}[u] = \frac{h}{2} \sum_{i=0}^N w_i u \left( z_i + \frac{h}{2} + \frac{x_i h}{2} \right),$$

so that  $I[u] = \sum_{i=1}^{M-1} I^{(i)}[u]$  and  $\widehat{I}_h[u] = \sum_{i=1}^{M-1} \widehat{I}_h^{(i)}[u]$ . Since  $p$  is a polynomial of degree at most  $d$ , it follows that

$$\begin{aligned} \left| I^{(i)}[u] - \widehat{I}_h^{(i)}[u] \right| &= \left| I^{(i)}[u] - I^{(i)}[p] + \widehat{I}_h^{(i)}[p] - \widehat{I}_h^{(i)}[u] \right| \\ &\leq \left| I^{(i)}[u] - I^{(i)}[p] \right| + \left| \widehat{I}_h^{(i)}[u] - \widehat{I}_h^{(i)}[p] \right| = \left| I^{(i)}[u - p] \right| + \left| \widehat{I}_h^{(i)}[u - p] \right|, \end{aligned}$$

where we used the triangle inequality and linearity in the second line. Thus, using (3.17) we deduce that

$$\begin{aligned} \left| I^{(i)}[u] - \widehat{I}_h^{(i)}[u] \right| &\leq \frac{C_{d+1}}{(d+1)!} h^{d+1} \int_{z_i}^{z_{i+1}} 1 dx + \frac{h}{2} \frac{C_{d+1}}{(d+1)!} h^{d+1} \sum_{i=0}^J |w_j| \\ &= \frac{C_{d+1}}{(d+1)!} h^{d+2} \left( 1 + \frac{1}{2} \sum_{i=0}^J |w_j| \right). \end{aligned}$$

Summing the contributions of the local errors, we obtain

$$\left| I[u] - \widehat{I}_h[u] \right| \leq \sum_{i=0}^{M-1} \left| I^{(i)}[u] - \widehat{I}_h^{(i)}[u] \right| \leq M \frac{C_{d+1}}{(d+1)!} h^{d+2} \left( 1 + \frac{1}{2} \sum_{i=0}^J |w_j| \right),$$

which leads to the result since  $Mh = b - a$ . □

*Remark 3.3.* A few comments are in order.

- The constant  $c$  in (3.15) is not sharp in general, but the power of  $h$  is optimal.
- This result may be viewed as a generalization of Theorems 3.1 and 3.2, albeit with a worse constant prefactor.
- Instead of the polynomial  $p_i$  in (3.16), we could have used any polynomial approximation of  $u$  such that (3.17) is satisfied, with possibly a different prefactor but the same power of  $h$  on the right-hand side. A natural choice, for example, would have been to define  $p_i$  by interpolation through (possibly a subset or superset of) the local integration points.
- Theorem 3.3 further motivates why the degree of precision of an integration rule is an interesting metric.

**Estimating the error a posteriori.** In practice, it is useful to be able to estimate the integration error so that, if the error is deemed too large, a better approximation of the integral can be calculated by using a smaller value for the step size  $h$ . Calculating the exact error  $I - \widehat{I}_h$  is impossible in general, because this would require to know the exact value of the integral, but it is possible to calculate a rough approximation of the error based on two numerical approximations of the integral, as we illustrate formally hereafter for the composite Simpson rule.

Suppose that  $\widehat{I}_{2h}$  and  $\widehat{I}_h$  are two approximations of the integral, calculated using the composite Simpson rule with step size  $2h$  and  $h$ , respectively. If we assume that the error proportionally to  $\mathcal{O}(h^4)$  as (3.12) suggests, then it holds approximately that

$$I - \widehat{I}_h \approx \frac{1}{24}(I - \widehat{I}_{2h}). \quad (3.18)$$

This implies that

$$I - \widehat{I}_{2h} = (I - \widehat{I}_h) + (\widehat{I}_h - \widehat{I}_{2h}) \approx \frac{1}{16}(I - \widehat{I}_{2h}) + (\widehat{I}_h - \widehat{I}_{2h}).$$

Rearranging this equation gives an approximation of the error for  $\widehat{I}_{2h}$ :

$$I - \widehat{I}_{2h} \approx \frac{16}{15}(\widehat{I}_h - \widehat{I}_{2h}).$$

Using (3.18), we can then derive an error estimate for  $\widehat{I}_h$ :

$$|I - \widehat{I}_h| \approx \frac{1}{15}|\widehat{I}_h - \widehat{I}_{2h}|. \quad (3.19)$$

The right-hand side can be calculated numerically, because it does not depend on the exact value of the integral. In practice, the two sides of (3.19) are often very close for small  $h$ . In the code example below, we approximate the integral

$$I = \int_0^{\frac{\pi}{2}} \cos(x) dx = 1 \quad (3.20)$$

for different step sizes and compare the exact error with the approximate error obtained using (3.19). The results obtained are summarized in Table 3.1, which shows a good match between the two quantities.

Table 3.1: Comparison between the exact integration error and the approximate integration error calculated using (3.19).

$h$	Exact error $ I - \widehat{I}_h $	Approximate error $\frac{1}{15} \widehat{I}_h - \widehat{I}_{2h} $
$2^{-4}$	$5.166847063531321 \times 10^{-7}$	$5.185892840930961 \times 10^{-7}$
$2^{-5}$	$3.226500089326123 \times 10^{-8}$	$3.229464703065806 \times 10^{-8}$
$2^{-6}$	$2.0161285974040766 \times 10^{-9}$	$2.016591486390477 \times 10^{-9}$
$2^{-7}$	$1.2600120946615334 \times 10^{-10}$	$1.260084925291949 \times 10^{-10}$

```
# Composite Simpson's rule
function composite_simpson(u, a, b, n)
    # Integration nodes
    x = LinRange(a, b, n + 1)
    # Evaluation of u at the nodes
    ux = u.(x)
    # Step size
    h = x[2] - x[1]
    # Approximation of the integral
    return (h/3) * sum([ux[1]; ux[end]; 4ux[2:2:end-1]; 2ux[3:2:end-2]])
end
# Function to integrate
u(x) = cos(x)
# Integration bounds
a, b = 0, pi/2
# Exact integral
I = 1.0
# Number of subintervals
ns = [8; 16; 32; 64; 128]
# Approximate integrals
Î = composite_simpson.(u, a, b, ns)
# Calculate exact and approximate errors
for i in 2:length(ns)
    println("Exact error: $(I - Î[i]), ",
           "Approx error: $((Î[i] - Î[i-1])/15)")
end
```

### 3.3 Richardson extrapolation and Romberg's method

In the previous section, we showed how the integration error could be approximated based on two approximations of the integral with different step sizes. The aim of this section is to show that, by cleverly combining two approximations  $\widehat{I}_h$  and  $\widehat{I}_{2h}$  of an integral, an approximation even better than  $\widehat{I}_h$  can be constructed.

This approach is based on *Richardson's extrapolation*, which is a general method for accelerating the convergence of sequences, with applications beyond numerical integration. The idea is the following: assume that  $J(h)$  is an approximation with step size  $h$  of some unknown quantity  $J_* = \lim_{h \rightarrow 0} J(h)$ , and that we have access to evaluations of  $J$  at  $h, h/2, h/4, h/8 \dots$ . If  $J$  extends to a smooth function over  $[0, H]$ , then by Taylor expansion it holds that

$$J(\eta) = J(0) + J'(0)\eta + J''(0)\frac{\eta^2}{2} + J^{(3)}(0)\frac{\eta^3}{3!} + \dots + J^{(k)}(0)\frac{\eta^k}{k!} + \mathcal{O}(\eta^{k+1}).$$

**Elimination of the linear error term.** Let us assume that  $J'(0) \neq 0$ , so that the leading order term after the constant  $J(0)$  scales as  $\eta$ . Then we have

$$\begin{aligned} J(h) &= J(0) + J'(0)h + \mathcal{O}(h^2) \\ J(h/2) &= J(0) + J'(0)\frac{h}{2} + \mathcal{O}(h^2). \end{aligned}$$

We now ask the following question: can we combine linearly  $J(h)$  and  $J(h/2)$  in order to construct an approximation  $J_1(h/2)$  of  $J(0)$  with an error scaling as  $\mathcal{O}(h^2)$ ? Employing the ansatz  $J_1(h/2) = \alpha J(h) + \beta J(h/2)$ , we calculate

$$J_1(h/2) = (\alpha + \beta)J(0) + J'(0)h \left( \alpha + \frac{1 - \alpha}{2} \right) + \mathcal{O}(h^2). \quad (3.21)$$

Since we want this expression to approximate  $J(0)$  for small  $h$ , we need to impose that  $\alpha + \beta = 1$ . Then, in order for the term multiplying  $h$  to cancel out, we require that

$$\alpha + \frac{1 - \alpha}{2} = 0 \quad \Leftrightarrow \quad \alpha = -1.$$

This yields the formula

$$J_1(h/2) = 2J(h/2) - J(h). \quad (3.22)$$

Notice that, in the case where  $J$  is a linear function,  $J_1(h/2)$  is exactly equal to  $J(0)$ . This reveals a geometric interpretation of (3.22): the approximation  $J_1(h/2)$  is simply the  $y$  intercept of the straight line passing through the points  $(h/2, J(h/2))$  and  $(h, J(h))$ .

**Elimination of the quadratic error term.** If we had tracked the coefficient of  $h^2$  in the previous paragraph, we would have obtained instead of (3.21) the following equation:

$$J_1(h/2) = J(0) - J^{(2)}(0)\frac{h^2}{4} + \mathcal{O}(h^3).$$

Provided that we have access also to  $J(h/4)$ , we can also calculate

$$J_1(h/4) = 2J(h/4) - J(h/2) = J(0) - J^{(3)}(0)\frac{h^2}{16} + \mathcal{O}(h^3).$$

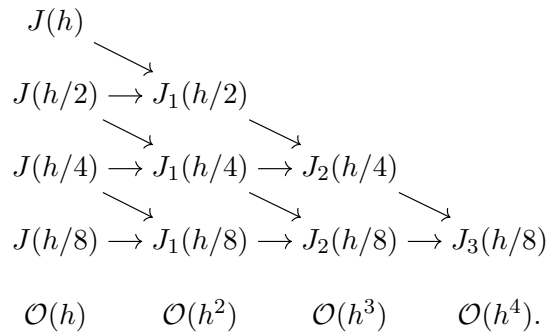
At this point, it is natural to wonder whether we can combine  $J_1(h/2)$  and  $J_1(h/4)$  in order to produce an even better approximation of  $J(0)$ . Applying the same reasoning as in the previous

section leads us to introduce

$$J_2(h/4) = \frac{4J_1(h/4) - J_1(h/2)}{4 - 1} = J(0) + \mathcal{O}(h^3).$$

This is an exact approximation of  $J(0)$  if  $J$  is a quadratic polynomial, indicating that  $J_2(h/4)$  is simply the  $y$  intercept of the quadratic polynomial interpolating the function  $J$  through the three points  $(h/4, J(h/4))$ ,  $(h/2, J(h/2))$  and  $(h, J(h))$ .

**Elimination of higher order terms.** The procedure above can be repeated in order to eliminate terms of higher and higher orders. The following schematic illustrates, for example, the calculation of an approximation  $J_3(h/8) = J(0) + \mathcal{O}(h^4)$ .



Here, the last row indicates the scaling of the error with respect to the parameter  $h$  in the limit as  $h \rightarrow 0$ . The linear combination in order to calculate  $J_i(h/2^i)$  is always of the form

$$J_i(h/2^i) = \frac{2^i J_{i-1}(h/2^i) - J_{i-1}(h/2^{i-1})}{2^i - 1}, \quad J_0 = J.$$

In practice we calculate the values taken by  $J, J_1, J_2, \dots$  at specific values of  $h$ , but these are in fact functions of  $h$ . In Figure 3.1, we plot these functions when  $J(h) = 1 + \sin(h)$ . It appears clearly from the figure that, for sufficiently small  $h$ ,  $J_3(h)$  provides the most precise approximation of  $J(0) = 1$ . Constructing the functions in Julia can be achieved in just a few lines of code.

```

J(h) = 1 + sin(h)
J_1(h) = 2J(h) - J(2h)
J_2(h) = (4J_1(h) - J_1(2h))/3
J_3(h) = (8J_2(h) - J_2(2h))/7

```

**Generalization.** Sometimes, it is known a priori that the Taylor development of the function  $J$  around zero contains only even powers of  $h$ . In this case, the Richardson extrapolation procedure can be slightly modified in order to produce approximations with errors scaling as  $\mathcal{O}(h^4)$ ,



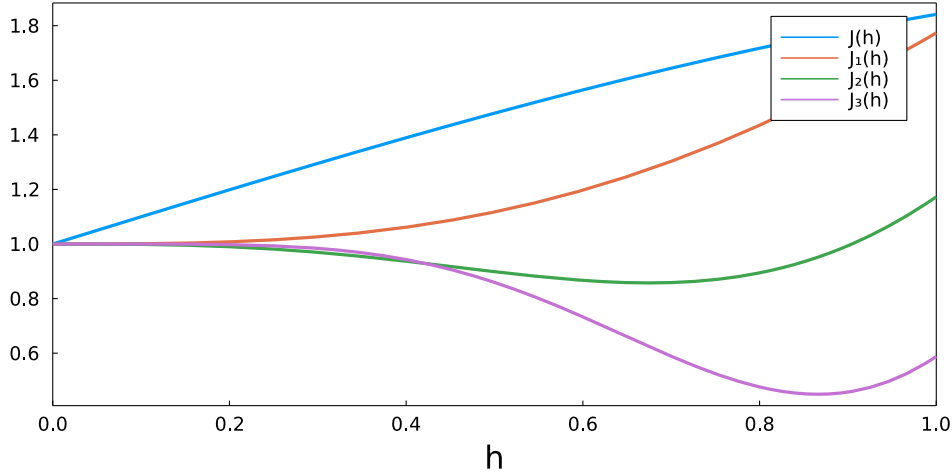
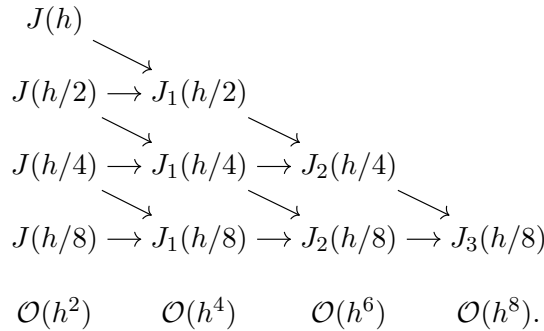


Figure 3.1: Illustration of the functions  $J_1$ ,  $J_2$  and  $J_3$  constructed by Richardson extrapolation.

then  $\mathcal{O}(h^6)$ , then  $\mathcal{O}(h^8)$ , etc. This procedure is illustrated below:



This time, the linear combinations required for populating this table are given by

$$J_i(h/2^i) = \frac{2^{2i} J_{i-1}(h/2^i) - J_{i-1}(h/2^{i-1})}{2^{2i} - 1}. \tag{3.23}$$

**Application to integration: Romberg’s method** Romberg’s integration method consists of applying Richardson’s extrapolation to the function

$$J(h) = \widehat{I}_h = u(x_0) + 2u(x_1) + 2u(x_2) + \cdots + 2u(x_{n-1}) + 2u(x_n), \quad h \in \left\{ \frac{b-a}{n} : n \in \mathbf{N} \right\}.$$

where  $a = x_0 < x_1 < \cdots < x_n = b$  are equidistant nodes. The right-hand side of this equation is simply the composite trapezoidal rule with step size  $h$ . It is possible to show that  $J(h)$  may be expanded as follows:

$$\forall k \in \mathbf{N}, \quad J(h) = I + \alpha_1 h^2 + \alpha_2 h^4 + \cdots + \alpha_k h^{2k} + \mathcal{O}(h^{2k+2}). \tag{3.24}$$

This is the content of the following result.

**Lemma 3.4.** *Let  $J(h)$  denote the approximation of the integral  $I[u]$  by the composite trapezium rule, and assume that  $u \in C^\infty[a, b]$ . Then  $J(h)$  may be expanded as in (3.24).*

*Proof.* Using the same notation as in [Theorem 3.3](#), we introduce

$$I^{(i)}[u] = \int_{x_i}^{x_{i+1}} u(x) \, dx, \quad \widehat{I}_h^{(i)}[u] = \frac{h}{2} \left( u(x_i) + u(x_{i+1}) \right), \quad i = 0, \dots, n-1.$$

Fix  $i \in \{0, \dots, n-1\}$  and let  $p_i$  denote the Taylor expansion of  $u$  around  $x_{i+\frac{1}{2}} := x_i + \frac{h}{2}$ , of degree  $2k+1$ :

$$p_i(x) = u(x_{i+\frac{1}{2}}) + u'(x_{i+\frac{1}{2}}) \left( x - x_{i+\frac{1}{2}} \right) + \dots + \frac{u^{(2k+1)}(x_{i+\frac{1}{2}})}{(2k+1)!} \left( x - x_{i+\frac{1}{2}} \right)^{2k+1}.$$

From the mean-value form of the remainder, it is clear that

$$I^{(i)}[u] = I^{(i)}[p_i] + \mathcal{O}(h^{2k+3}), \quad \widehat{I}_h^{(i)}[u] = \widehat{I}_h^{(i)}[p_i] + \mathcal{O}(h^{2k+3}).$$

Substituting  $p_i$  and noting that odd powers cancel when integrating, we have that

$$\begin{aligned} I^{(i)}[p_i] &= u(x_{i+\frac{1}{2}})h + \omega_2 h^3 + \omega_4 h^5 + \dots + \omega_{2k} h^{2k+1}, \\ \widehat{I}_h^{(i)}[p_i] &= u(x_{i+\frac{1}{2}})h + \eta_2 h^3 + \eta_4 h^5 + \dots + \eta_{2k} h^{2k+1}, \end{aligned}$$

for appropriate coefficients. Thus we obtain that

$$\begin{aligned} I[u] - \widehat{I}_h[u] &= \sum_{i=0}^n I^{(i)}[p_i] - \widehat{I}_h^{(i)}[p_i] \\ &= (\omega_2 - \eta_2)h^2 + (\omega_4 - \eta_4)h^4 + \dots + (\omega_{2k} - \eta_{2k})h^{2k} + \mathcal{O}(h^{2k+2}), \end{aligned}$$

which concludes the proof.  $\square$

Richardson's extrapolation [\(3.23\)](#) can therefore be employed in order to compute approximations of the integral with increasing accuracy. The convergence of Romberg's method for calculating the integral [\(3.20\)](#) is illustrated in [Figure 3.2](#).

### 3.4 Methods with non-equidistant nodes

The Newton–Cotes method relies on equidistant integration nodes, and the only degrees of freedom are the integration weights. If the nodes are not fixed, then additional degrees of freedom are available, and these can be leveraged in order to construct a better integration formula. The total number of degrees of freedom for a general integration rule of the form [\(3.2\)](#) is  $2n+2$  which, in principle, should enable to construct an integration rule with a degree of precision equal to  $2n+1$ .

A necessary condition for an integration rule of the form [\(3.2\)](#) to have a degree of precision equal to  $2n+1$  is that it integrates exactly all the monomials of degree 0 to  $2n+1$ . This condition is also sufficient because, assuming that it is satisfied, we have by linearity of the

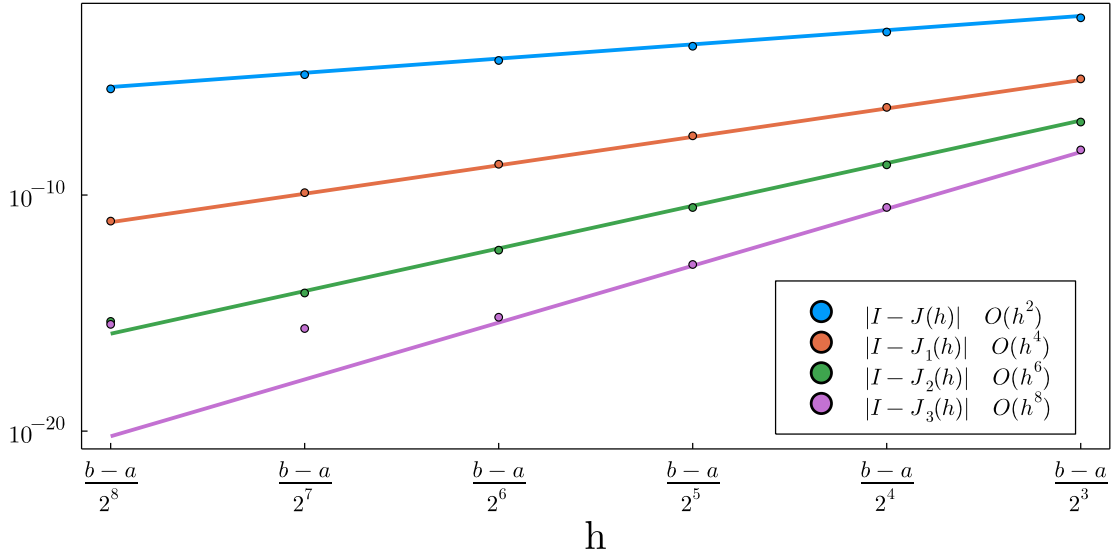


Figure 3.2: Convergence of Romberg’s method. The straight lines correspond to the monomial functions  $f(h) = C_i h^i$ , with  $i = 2, 4, 6, 8$  and for appropriate constants  $C_i$ . We observe a good agreement between the observed and theoretical convergence rates.

functionals  $I$  and  $\widehat{I}$  that

$$\begin{aligned} \widehat{I}(\alpha_0 + \alpha_1 x + \dots + \alpha_{2n+1} x^{2n+1}) &= \alpha_0 \widehat{I}(1) + \alpha_1 \widehat{I}(x) + \dots + \alpha_{2n+1} \widehat{I}(x^{2n+1}) \\ &= \alpha_0 I(1) + \alpha_1 I(x) + \dots + \alpha_{2n+1} I(x^{2n+1}) \\ &= I(\alpha_0 + \alpha_1 x + \dots + \alpha_{2n+1} x^{2n+1}), \end{aligned}$$

Here  $I(u)$  and  $\widehat{I}(u)$  denote respectively the exact integral of  $u$  and its approximate integral using (3.2). In order to find the nodes and weights of the integration rule, we can therefore solve the following nonlinear system of  $2n + 2$  equations with  $2n + 2$  unknowns:

$$\sum_{i=0}^n w_i x_i^d = \int_{-1}^1 x^d dx, \quad d = 0, \dots, 2n + 1. \tag{3.25}$$

The quadrature rule thus obtained is called the *Gauss–Legendre quadrature*.

*Example 3.2.* Let us derive the Gauss–Legendre quadrature with  $n + 1 = 2$  nodes. The system of equations that we need to solve in this case is the following:

$$w_0 + w_1 = 2, \quad w_0 x_0 + w_1 x_1 = 0, \quad w_0 x_0^2 + w_1 x_1^2 = \frac{2}{3}, \quad w_0 x_0^3 + w_1 x_1^3 = 0.$$

The solution to these equations is given by

$$-x_0 = x_1 = \frac{\sqrt{3}}{3}, \quad w_0 = w_1 = 1.$$

**Connection with orthogonal polynomials.** Let  $(L_n)_{n \in \mathbf{N}}$  denote the Legendre polynomials, i.e. the orthogonal polynomials with for the inner product

$$\langle f, g \rangle = \int_{-1}^1 f(x)g(x) dx.$$

The nodes and weights of the Gauss–Legendre quadrature rules can be obtained constructively from Legendre polynomials, as shown in Section 2.2.4. We shall now demonstrate this connection in much more direct manner. Specifically, we prove that the integration nodes are given by the roots of a Legendre polynomial.

**Theorem 3.5.** *For every  $n \in \mathbf{N}$ , there exists a unique solution to the system of equations (3.25). The nodes  $(x_i)_{i \in \{0, \dots, n\}}$  are the roots of  $L_{n+1}$  and the weights are given by*

$$w_i = \int_{-1}^1 \ell_i(x) dx, \quad i \in \{1, \dots, n\}, \quad (3.26)$$

where  $\ell_i$  is the Lagrange polynomial

$$\frac{(x - x_0) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n)}{(x_i - x_0) \cdots (x_i - x_{i-1})(x_i - x_{i+1}) \cdots (x_i - x_n)}.$$

In addition, the weights are all positive.

*Proof.* We show first that the solution to (3.25) exists, then that it is unique, and finally that the weights are positive.

**Existence.** We begin by showing that, if  $(x_i)_{i \in \{0, \dots, n\}}$  are the roots of  $L_{n+1}$  and the weights are defined from (3.26), then the equations (3.25) are satisfied. To this end, it is sufficient to show that for all  $p \in \mathbf{P}(2n + 1)$ ,

$$\int_{-1}^1 p(x) dx = \sum_{i=0}^n w_i p(x_i). \quad (3.27)$$

Take  $p \in \mathbf{P}(2n + 1)$  and let  $q \in \mathbf{P}(n)$  and  $r \in \mathbf{P}(n)$  be the polynomials such that

$$p(x) = q(x)L_{n+1}(x) + r(x).$$

The quotient  $q$  and the remainder  $r$  can be obtained by Euclidean division of  $p$  by  $L_{n+1}$ . Since  $L_{n+1}$  is orthogonal to any polynomial in  $\mathbf{P}(n)$ , in particular  $q$ , and the nodes  $(x_i)_{i \in \{0, \dots, n\}}$  are the roots of  $L_{n+1}$ , it holds that

$$\int_{-1}^1 p(x) dx - \sum_{i=0}^n w_i p(x_i) = \int_{-1}^1 r(x) dx - \sum_{i=0}^n w_i r(x_i). \quad (3.28)$$

Given that  $r \in \mathbf{P}(n)$ , the remainder  $r$  must coincides with its polynomial interpolation at the points  $x_0, \dots, x_m$ . Therefore,

$$r = r(x_0)\ell_0 + \cdots + r(x_n)\ell_n,$$

and so

$$\int_{-1}^1 r(x) dx = \sum_{i=0}^n r(x_i) \int_{-1}^1 \ell_i(x) dx = \sum_{i=0}^n r(x_i) w_i,$$

where we used (3.26) in the last equality. Consequently, the right-hand side of (3.28) is zero, and since  $p$  was arbitrary this implies that (3.27) is satisfied.

**Uniqueness.** Next, we show that the nodes are necessarily the roots of  $L_{n+1}$ . To this end, assume that  $x_0, \dots, x_n$  and weights  $w_0, \dots, w_n$  are such that the equations (3.25) are satisfied, and let

$$q(x) = (x - x_0) \dots (x - x_n).$$

Our goal is to show that  $q(x)$  coincides with  $L_{n+1}$  up to a constant factor. In order to show this, it is sufficient to prove that  $q(x)$  is orthogonal to  $x^d$  for all  $d = 0, \dots, n$ , because the only polynomial in  $\mathbf{P}(n+1)$  that satisfies these orthogonality relations is the Legendre polynomial  $L_{n+1}$ , or a multiple thereof. For the values of  $d$  considered, the polynomial  $q(x)x^d$  belongs to  $\mathbf{P}(2n+1)$ . Given that the integration rule with nodes  $x_0, \dots, x_n$  and weights  $w_0, \dots, w_n$  is exact for all the elements of  $\mathbf{P}(2n+1)$  by assumption, we deduce that

$$\forall d \in \{0, \dots, n\}, \quad \int_{-1}^1 q(x)x^d dx = \sum_{i=0}^n w_i q(x_i)x_i^d = 0.$$

Finally, we show that the weights are necessarily given by (3.26). Since the integration rule must be exact for any Lagrange polynomial  $\ell_j$ , we have

$$\int_{-1}^1 \ell_j(x) dx = \sum_{i=0}^n w_i \ell_j(x_i) = w_j,$$

which concludes the proof of uniqueness.

**Positivity of the weights.** Since the integration rule is exact for all the polynomials in  $\mathbf{P}(2n+1)$  and  $\ell_j(x)^2 \in \mathbf{P}(2n+1)$ , we deduce that

$$\int_{-1}^1 |\ell_j(x)|^2 dx = \sum_{i=0}^n w_i |\ell_j(x_i)|^2 = w_j.$$

The left-hand side is positive, and so  $w_j$  must also be positive. □

Since the integration weights are all positive, the Gauss–Legendre quadrature rules are less susceptible to roundoff errors than the Newton–Cotes methods. In addition, we have the following result.

**Theorem 3.6.** *Assume that  $u \in C([-1, 1])$ , and let  $\widehat{I}_n(u)$  denote the approximation of  $I(u)$  using the Gauss–Legendre quadrature. Then  $\widehat{I}_n(u) \rightarrow I(u)$  in the limit as  $n \rightarrow \infty$ .*

*Proof.* The positivity of the weights is crucial for the proof. By the Weierstrass approximation

theorem, for all  $\varepsilon > 0$  there exists polynomial  $p$  such that

$$E := \max_{x \in [-1, 1]} |p(x) - u(x)| \leq \varepsilon.$$

Since the degree of precision of the Gauss–Legendre quadrature with  $n + 1$  integration nodes is  $2n + 1$ , there exists  $N$  sufficiently large such that  $\widehat{I}_n(p) = I(p)$  for all  $n \geq N$ . Thus

$$\begin{aligned} \forall n \geq N, \quad \left| \widehat{I}_n(u) - I(u) \right| &\leq \left| \widehat{I}_n(u) - \widehat{I}_n(p) \right| + \left| \widehat{I}_n(p) - I(p) \right| + \left| I(p) - I(u) \right| \\ &= \left| \widehat{I}_n(u) - \widehat{I}_n(p) \right| + \left| I(p) - I(u) \right| \\ &\leq \sum_{i=0}^n |w_i| E + \int_{-1}^1 E \, dx = 2\varepsilon + 2\varepsilon = 4\varepsilon. \end{aligned}$$

Indeed, all the weights are positive and the Gauss–Legendre quadrature rule is exact for the constant function  $u(x) = 1$ , which implies that the weights add up to 2. Since  $\varepsilon > 0$  was arbitrary, this concludes the proof.  $\square$

**Computation of the integration nodes.** We proved that the integration nodes are given by the roots of the Legendre polynomials. In practice, calculating these roots can be achieved by calculating the eigenvalues of a tridiagonal matrix, see [Section 2.2.4](#). This is known as the Golub–Welsch algorithm.

**Generalization to higher dimensions.** Gauss–Legendre integration is ubiquitous in numerical methods for partial differential equations, in particular the *finite element method*. Its generalization to higher dimensions is immediate: for a function  $u: [-1, 1] \times [-1, 1] \rightarrow \mathbf{R}$ , we have

$$\int_0^1 \int_0^1 u(x, y) \, dy dx \approx \sum_{i=0}^n \sum_{j=0}^n w_i w_j u(x_i, y_j).$$

The degree of precision of this integration rule is the same as that of the corresponding one-dimensional rule, and this approach can be generalized to any dimension  $d$ . The associated computational cost, however, scales as  $n^d$ , and so it is not a good idea to use a deterministic method of this type in the high dimensional setting. The explosion of the computational cost as the dimension increases is known as *the curse of dimensionality*.

### 3.5 Introduction to probabilistic integration methods

So far in this chapter, we covered only *deterministic* integration formulas. Much of the research around the calculation of high-dimensional integrals today is concerned with *probabilistic* integration methods using probabilistic approaches. These methods are based on the connection between integrals and expectations. To illustrate the simplest probabilistic approach, called the *Monte Carlo method*, we consider the problem of approximating the integral

$$I = \int_0^1 u(x) \, dx.$$

This integral may be expressed as the expectation  $\mathbf{E}[u(X)]$ , where  $\mathbf{E}$  is the expectation operator and  $X \sim \mathcal{U}(0, 1)$  is a uniformly distributed random variable over the interval  $[0, 1]$ . Therefore, in practice, the integral  $I$  may be approximated by generating a large number of *independent* samples  $X_1, X_2, \dots$  drawn from the distribution  $\mathcal{U}(0, 1)$  and averaging  $f(X_i)$  over all of these samples:

$$\hat{I}_N = \frac{1}{N} \sum_{n=1}^N u(X_n).$$

The quantity  $\hat{I}_N$ , where the subscript  $N$  denotes the number of samples employed, is itself a random variable; it is called an *estimator* of the exact integral  $I$ . In Julia, the calculation of the integral using the Monte Carlo method can be achieved with the following code.

```
N = 1000
u(x) = x^2
X = rand(N)
Î = (1/N) * sum(u.(X))
```

Since the expectation operator is linear, we calculate that

$$\mathbf{E}[\hat{I}_N] = \frac{1}{N} \sum_{n=1}^N \mathbf{E}[u(X_n)] = \frac{1}{N} \sum_{n=1}^N \int_0^1 u(x) dx = \frac{1}{N} \sum_{n=1}^N I = I.$$

The estimator  $\hat{I}_N$  is therefore said to be *unbiased*, because the bias  $\mathbf{E}[\hat{I}_N] - I$  is zero. Now assume that the function  $u$  is square integrable over the interval  $[0, 1]$  and let

$$\sigma^2 := \mathbf{V}[u(X)] := \mathbf{E}\left[\left(u(X) - \mathbf{E}[u(X)]\right)^2\right] = \int_0^1 |u(x) - I|^2 dx,$$

where  $\mathbf{V}[X]$  denotes the variance of the random variable  $X$ . We can calculate explicitly the variance of the estimator  $\hat{I}_N$ :

$$\mathbf{V}[\hat{I}_N] = \mathbf{E}\left[\left(\frac{1}{N} \sum_{n=1}^N u(X_n) - I\right)^2\right] = \mathbf{E}\left[\frac{1}{N^2} \sum_{n=1}^N \sum_{m=1}^N (u(X_n) - I)(u(X_m) - I)\right].$$

Since the samples  $X_1, X_2, \dots$  are independent, it holds that

$$\mathbf{E}\left[(u(X_n) - I)(u(X_m) - I)\right] = \delta_{nm}\sigma^2,$$

where  $\delta_{mn}$  is the Kronecker delta. Consequently, we deduce that

$$\mathbf{V}[\hat{I}_N] = \frac{1}{N^2} \sum_{n=1}^N \sum_{m=1}^N \mathbf{E}\left[(u(X_n) - I)(u(X_m) - I)\right] = \frac{1}{N^2} \sum_{n=1}^N \sum_{m=1}^N \delta_{mn}\sigma^2 = \frac{\sigma^2}{N}.$$

Therefore, the variance of  $\hat{I}_N$  decreases as  $1/N$  when the number of samples  $N$  increases, indicating that the estimator becomes increasingly accurate. To state this more precisely, we will use Chebyshev's inequality.

**Theorem 3.7** (Chebyshev's inequality). *Let  $Z$  be a random variable with mean  $m$  and variance  $s^2$ . Then for any real  $k > 0$ ,*

$$\mathbf{P}[|Z - m| \geq ks] \leq \frac{1}{k^2}.$$

Let  $\varepsilon > 0$ . Employing Chebyshev's inequality with  $k = \varepsilon/\sqrt{\mathbf{V}[\widehat{I}_N]}$ , we obtain

$$\mathbf{P}[|\widehat{I}_N - I| \geq \varepsilon] \leq \frac{\mathbf{V}[\widehat{I}_N]}{\varepsilon^2} = \frac{\sigma^2}{N\varepsilon^2}. \quad (3.29)$$

Consequently, for any  $\varepsilon > 0$ , the probability that the integration error  $|\widehat{I}_N - I|$  is greater than or equal to  $\varepsilon$  tends to zero in the limit as  $N \rightarrow \infty$ . In probabilistic jargon, we say that  $\widehat{I}_N$  converges *in probability* to  $I$ . Equation (3.29) can also be employed in order to construct a *confidence interval* for the exact integral, as we demonstrate in the following paragraph.

**Construction of a  $(1 - \alpha)$  confidence interval.** By definition, a  $(1 - \alpha)$  confidence interval for  $I$  is an interval  $[z_1, z_2]$ , the endpoints of which being random variables, such that the probability that  $I \in [z_1, z_2]$  is greater than or equal to  $1 - \alpha$ . To construct such an interval, we begin by finding  $\varepsilon$  such that the right-hand side of (3.29) is equal to  $\alpha$ , which gives  $\varepsilon = \sqrt{\sigma^2/N\alpha}$ . For this value of  $\varepsilon$ , we have

$$\mathbf{P}[|\widehat{I}_N - I| \geq \varepsilon] \leq \alpha.$$

Since  $|\widehat{I}_N - I| \geq \varepsilon$  if and only if  $I \notin (\widehat{I}_N - \varepsilon, \widehat{I}_N + \varepsilon)$ , we conclude that

$$\mathbf{P}[I \in (\widehat{I}_N - \varepsilon, \widehat{I}_N + \varepsilon)] \geq 1 - \alpha.$$

We have thus shown that

$$\left( \widehat{I}_N - \sqrt{\frac{\sigma^2}{N\alpha}}, \widehat{I}_N + \sqrt{\frac{\sigma^2}{N\alpha}} \right)$$

is a  $(1 - \alpha)$  confidence interval for  $I$ .

## 3.6 Exercises

⚙ **Exercise 3.1.** *Derive the Simpson's integration rule (3.6).*

⚙ **Exercise 3.2.** *Derive the composite Simpson integration rule (3.9).*

⚙ **Exercise 3.3.** *Consider the integration rule*

$$\int_0^1 u(x) dx \approx w_1 u(0) + w_2 u(1) + w_3 u'(0).$$

*Find  $w_1$ ,  $w_2$  and  $w_3$  so that this integration rule has the highest possible degree of precision.*



⚙️ **Exercise 3.4.** Consider the integration rule

$$\int_{-1}^1 u(x) \, dx \approx w_1 u(x_1) + w_2 u'(x_1).$$

Find  $w_1$ ,  $w_2$  and  $x_1$  so that this integration rule has the highest possible degree of precision.

⚙️ **Exercise 3.5.** What is the degree of precision of the following quadrature rule?

$$\int_{-1}^1 u(x) \, dx \approx \frac{2}{3} \left( 2u\left(-\frac{1}{2}\right) - u(0) + 2u\left(\frac{1}{2}\right) \right).$$

⚙️ **Exercise 3.6.** The Gauss–Hermite quadrature rule with  $n + 1$  nodes is an approximation of the form

$$\int_{-\infty}^{\infty} u(x) e^{-\frac{x^2}{2}} \, dx \approx \sum_{i=0}^n w_i u(x_i),$$

such that the rule is exact for all polynomials of degree less than or equal to  $2n + 1$ . Find the Gauss–Hermite rule with two nodes.

⚙️ **Exercise 3.7.** Use Romberg’s method to construct an integration rule with an error term scaling as  $\mathcal{O}(h^4)$ . Is there a link between the method you obtained and another integration rule seen in class?

⚙️ **Exercise 3.8** (Improving the error bound for the composite trapezoidal rule). The notation used in this exercise is the same as in Section 3.2. In particular,  $\widehat{I}_h$  denotes the approximate integral obtained by using the composite trapezoidal rule (3.7), and  $\widehat{u}_h$  is the corresponding piecewise linear interpolant.

A version of the mean value theorem states that, if  $g: [a, b] \rightarrow \mathbf{R}$  is a non-negative integrable function and  $f: [a, b] \rightarrow \mathbf{R}$  is continuous, then there exists  $\xi \in (a, b)$  such that

$$\int_a^b f(x)g(x) \, dx = f(\xi) \int_a^b g(x) \, dx. \quad (3.30)$$

- Using (3.30), show that, for all  $i \in \{0, \dots, n - 1\}$ , there exists  $\xi_i \in (x_i, x_{i+1})$  such that

$$\int_{x_i}^{x_{i+1}} u(x) - \widehat{u}_h(x) \, dx = -u''(\xi_i) \frac{h^3}{12}.$$

- Prove, by using the intermediate value theorem, that if  $f: [a, b] \rightarrow \mathbf{R}$  is a continuous function, then for any set  $\xi_0, \dots, \xi_{n-1}$  of points within the interval  $(a, b)$ , there exists  $c \in (a, b)$  such that

$$\frac{1}{n} \sum_{i=0}^{n-1} f(\xi_i) = f(c).$$

- Combining the previous items, conclude that there exists  $\xi \in (a, b)$  such that

$$I - \widehat{I}_h = -u''(\xi)(b - a) \frac{h^2}{12},$$

which is a more precise expression of the error than that obtained in (3.8).

**Remark 3.4.** One may convince oneself of (3.30) by rewriting this equation as

$$\frac{\int_a^b f(x)g(x) \, dx}{\int_a^b g(x) \, dx} = f(\xi).$$

The left-hand side is the average of  $f(x)$  with respect to the probability measure with density given by

$$x \mapsto \frac{g(x)}{\int_a^b g(x) \, dx}.$$

⚙️ **Exercise 3.9** (From the final exam of Spring 2022). Construct an integration rule of the form

$$\int_{-1}^1 u(x) \, dx \approx w_1 u\left(-\frac{1}{2}\right) + w_2 u(0) + w_3 u\left(\frac{1}{2}\right)$$

with a degree of precision equal to at least 2. What is the degree of precision of the rule constructed?

*Solution.* The Lagrange polynomials associated with  $-1/2$ ,  $0$  and  $1/2$  are respectively

$$\begin{aligned} p_1(x) &= 2x \left(x - \frac{1}{2}\right), \\ p_2(x) &= -4 \left(x + \frac{1}{2}\right) \left(x - \frac{1}{2}\right), \\ p_3(x) &= 2 \left(x + \frac{1}{2}\right) x. \end{aligned}$$

We deduce that

$$w_1 = \int_{-1}^1 p_1(x) \, dx = \frac{4}{3}, \quad w_2 = \int_{-1}^1 p_2(x) \, dx = -\frac{2}{3}, \quad w_3 = \int_{-1}^1 p_3(x) \, dx = \frac{4}{3}.$$

By construction, the degree of precision is at least 2. However, the integration rule is exact also when  $u(x) = x^3$ . Since the rule is not exact for  $u(x) = x^4$ , we conclude that the degree of precision is 3.  $\triangle$

⚙️ **Exercise 3.10** (From the final exam of Spring 2022). The Gauss–Laguerre quadrature rule with  $n$  nodes is an approximation of the form

$$\int_0^\infty u(x) e^{-x} \, dx \approx \sum_{i=1}^n w_i u(x_i),$$

such that the rule is exact when  $u$  is a polynomial of degree less than or equal to  $2n - 1$ .

- Find the Gauss–Laguerre rule with one node ( $n = 1$ ).
- Find the Gauss–Laguerre quadrature rule with two nodes ( $n = 2$ ). You may find it useful to first calculate the Laguerre polynomial of degree 2.

*Solution.* Below are the derivations of the Gauss–Laguerre rules with 1 and 2 nodes.

**Gauss–Laguerre rule with 1 node.** We are looking for  $w_1$  and  $x_1$  such that

$$\forall(a, b) \in \mathbf{R}^2, \quad \int_0^\infty (a + bx) e^{-x} dx = w_1(a + bx_1).$$

The left-hand side is equal to

$$a \int_0^\infty e^{-x} dx + b \int_0^\infty x e^{-x} dx = a + b \int_0^\infty x e^{-x} dx.$$

Using integration by parts, we find the value of the remaining integral on the right-hand side:

$$\begin{aligned} \int_0^\infty x e^{-x} dx &= \int_0^\infty -(x e^{-x})' + e^{-x} dx \\ &= -(x e^{-x}) \Big|_{x=\infty} + (x e^{-x}) \Big|_{x=0} + \int_0^\infty e^{-x} dx \\ &= 0 + 0 + 1. \end{aligned}$$

(To be fully rigorous, we would need to write the first term on the second line as a limit.) Therefore, we obtain

$$a + b = w_1(a + bx_1),$$

which implies that  $w_1 = x_1 = 1$ .

**Gauss–Laguerre rule with 2 nodes.** The integration nodes are given by the roots of the Laguerre polynomials, which are the orthogonal polynomials for the inner product

$$\langle f, g \rangle := \int_0^\infty f(x)g(x) e^{-x} dx.$$

The first polynomial is  $\ell_0(x) = 1$ . It is simple to check that the only linear monic polynomial orthogonal to  $\ell_0$  is given by  $\ell_1(x) = x - 1$ . Next, by integration by parts we calculate that

$$\int_0^\infty x^2 e^{-x} dx = \int_0^\infty -(x^2 e^{-x})' + 2x e^{-x} dx = 2.$$

and, similarly,

$$\int_0^\infty x^3 e^{-x} dx = \int_0^\infty -(x^3 e^{-x})' + 3x^2 e^{-x} dx = 6.$$

Consider the ansatz  $\ell_2(x) = x^2 + a\ell_1(x) + b$ . In order for  $\ell_2$  to be orthogonal to  $\ell_0$  and  $\ell_1$ , it is necessary that

$$\begin{aligned} 0 &= \int_0^\infty \ell_2(x) \ell_0(x) e^{-x} dx = 2 + b, \\ 0 &= \int_0^\infty \ell_2(x) \ell_1(x) e^{-x} dx = 4 + a \int_0^\infty \ell_1(x) \ell_1(x) dx = 4 + a. \end{aligned}$$

Therefore, we conclude that  $a = -4$  and  $b = -2$ , which gives

$$\ell_2(x) = x^2 - 4x + 2.$$

The roots are given by  $2 \pm \sqrt{2}$ , so we have  $x_1 = 2 - \sqrt{2}$  and  $x_2 = 2 + \sqrt{2}$ . It remains to find the weights. To this end, we need only two additional equations; it is sufficient to require that, for any  $(a, b) \in \mathbf{R}^2$ ,

$$\begin{aligned} a + b &= \int_0^\infty (a + bx) e^{-x} dx = w_1(a + bx_1) + w_2(a + bx_2) \\ &= a(w_1 + w_2) + 2b(w_1 + w_2) + \sqrt{2}b(w_2 - w_1). \end{aligned}$$

Letting  $a = 1$  and  $b = 0$ , we obtain  $w_1 + w_2 = 1$ . Then, letting  $a = 0$  and  $b = 1$ , we deduce

$$1 = 2 + \sqrt{2}(w_2 - w_1) \quad \Leftrightarrow w_2 - w_1 = -\frac{\sqrt{2}}{2}.$$

Therefore

$$w_1 = \frac{2 + \sqrt{2}}{4}, \quad w_2 = \frac{2 - \sqrt{2}}{4},$$

which concludes the exercise.  $\triangle$

**□ Exercise 3.11** (Calculating the volume of hyperballs). Let  $B_d$  denote the  $d$ -dimensional unit ball for the Euclidean norm:

$$B_d = \left\{ \mathbf{x} \in \mathbf{R}^d : \|\mathbf{x}\| \leq 1 \right\}.$$

The volume of  $B_d$  is defined as the integral of the characteristic function over  $B_d$ :

$$\text{vol}(B_d) = \underbrace{\int_{\mathbf{R}} \cdots \int_{\mathbf{R}}}_{d \text{ times}} \chi(\mathbf{x}) dx_1 \dots dx_d, \quad \chi(\mathbf{x}) := \begin{cases} 1 & \text{if } \mathbf{x} \in B_d \\ 0 & \text{otherwise.} \end{cases}$$

Complete the following tasks:

- Write a function

```
function hyperball_volume(dim, n)
    # Your code comes here
    return vol,  $\sigma$ 
end
```

that calculates the volume of the unit ball  $B_d$  with  $d = \text{dim}$  using a Monte Carlo approach with  $n$  samples drawn from an appropriate distribution. Your function should return an estimation of the volume together with the standard deviation of the estimator (which you should estimate from the samples).

- Using the function `hyperball_volume`, plot the volumes for  $d$  going from 1 to 15, together with a 99% confidence interval. See [Figure 3.3](#) for an example solution with  $n = 10^7$ .

You are allowed to use your knowledge of the fact that  $\text{vol}(B_2) = \pi$  and  $\text{vol}(B_3) = 4\pi/3$ , but do not use the general formula for the volume of  $B_d$ .

**□ Exercise 3.12.** Complete the following tasks:

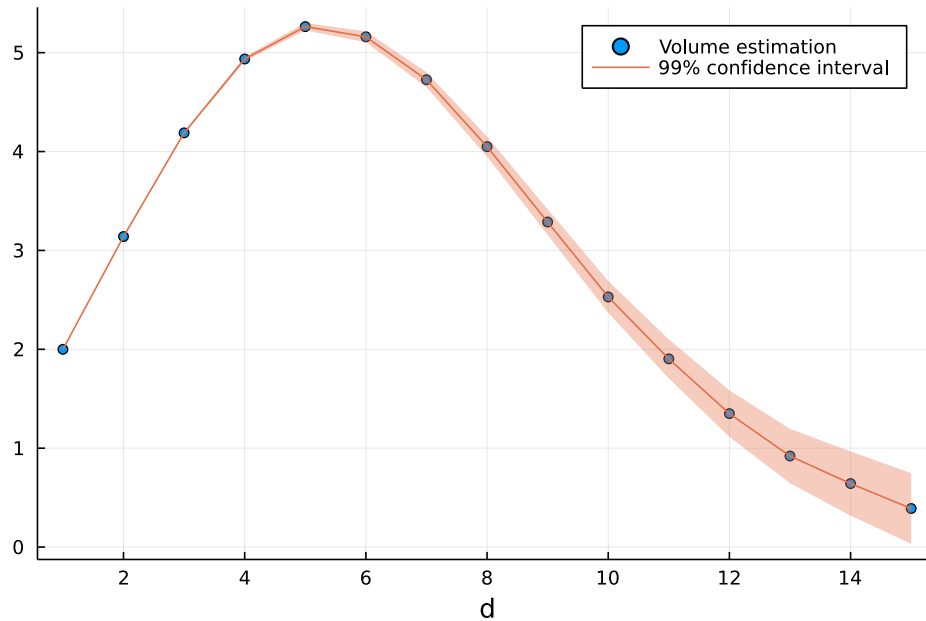


Figure 3.3: Example solution for Exercise 3.11.

- Write a function `legendre(n)` which returns the Legendre polynomial of degree  $n$ . To this end, you may use the `Polynomials` library and Rodrigues' formula:

$$L_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n.$$

- Write a function `get_nodes_and_weights(n)` which returns the nodes and weights of the Gauss–Legendre quadrature with  $n$  nodes. In order to construct Lagrange polynomials, you may find it useful to use the `fromroots` functions.
- Write a function `composite_gauss_legendre(u, a, b, n, N)`, which returns an approximation of the integral

$$\int_a^b u(x) dx$$

obtained by partitioning the integration interval  $[a, b]$  into  $N$  cells, and applying the Gauss–Legendre quadrature within each cell.

- Take  $u(x) = \cos(x)$ ,  $a = -1$  and  $b = 1$ . Illustrate on the same plot the error for the values  $n \in \{1, 2, 3\}$  and  $N$  varying from 1 to 40. For each value of  $n$ , estimate the order of convergence with respect to  $N$ , i.e. find  $\alpha(n)$  such that

$$|\widehat{I}_{n,N} - I| \propto CN^{-\alpha},$$


where  $I$  denotes the exact value of the Integral and  $\widehat{I}_{n,N}$  denotes its approximation.

### 3.7 Discussion and bibliography

The presentation of part of the material follows that in [7], and some exercises come from [10, Chapter 9]. The main advantage of probabilistic integration approaches is that they generalize naturally to high-dimensional and infinite-dimensional settings.

# Chapter 4

## Solution of linear systems of equation

4.1	Conditioning . . . . .	81
4.2	Direct solution method . . . . .	85
4.2.1	LU decomposition . . . . .	86
4.2.2	Backward and forward substitution . . . . .	90
4.2.3	Gaussian elimination with pivoting  . . . . .	91
4.2.4	Direct method for Hermitian positive definite matrices . . . . .	94
4.2.5	Direct methods for banded matrices . . . . .	95
4.3	Iterative methods for linear systems . . . . .	97
4.3.1	Basic iterative methods . . . . .	98
4.3.2	The conjugate gradient method . . . . .	106
4.4	Exercises . . . . .	117
4.5	Discussion and bibliography . . . . .	126

### Introduction

This chapter is devoted to the numerical solution of linear problems of the following form:

$$\text{Find } \mathbf{x} \in \mathbf{C}^n \text{ such that } \quad \mathbf{Ax} = \mathbf{b}, \quad \mathbf{A} \in \mathbf{C}^{n \times n}, \quad \mathbf{b} \in \mathbf{C}^n. \quad (4.1)$$

Systems of this type appear in a variety of applications. They naturally arise in the context of linear partial differential equations, which we use as main motivating example. Partial differential equations govern a wide range of physical phenomena including heat propagation, gravity, and electromagnetism, to mention just a few. Linear systems in this context often have a particular structure: the matrix  $\mathbf{A}$  is generally very sparse, which means that most of the entries are equal to 0, and it is often Hermitian and positive definite, provided that these properties are satisfied by the underlying operator.

There are two main approaches for solving linear systems:

- Direct methods enable to calculate the exact solution to systems of linear equations, up to round-off errors, in a finite number of steps. Although this is an attractive property, direct methods are usually too computationally costly for large systems: The cost of inverting a general  $n \times n$  matrix, measured in number of floating operations, scales as  $n^3$ !
- Iterative methods, on the other hand, enable to progressively calculate increasingly accurate approximations of the solution. Iterations may be stopped once the *the residual* is sufficiently small. These methods are often preferable when the dimension  $n$  of the linear system is very large.

This chapter is organized as follows.

- In [Section 4.1](#), we introduce the concept of *conditioning*. The condition number of a matrix provides information on the sensitivity of the solution to perturbations of the right-hand side  $\mathbf{b}$  or matrix  $\mathbf{A}$ . It is useful, for example, in order to determine the potential impact of round-off errors.
- In [Section 4.2](#), we present the direct method for solving systems of linear equations. We study in particular the LU decomposition for an invertible matrix, as well as its variant for symmetric positive definite matrices, which is called the Cholesky decomposition.
- In [Section 4.3](#), we present iterative methods for solving linear systems. We focus, in particular, on basic iterative methods based on a splitting of the matrix  $\mathbf{A}$  and on the conjugate gradient method.

## 4.1 Conditioning

Before studying the properties of numerical methods for solving linear system, it is crucial to understand the impact of *round-off errors* on the solution, which impose a bound on the accuracy we can hope to achieve. We begin with a motivating example.

*Example 4.1.* Suppose that we wish to solve the following equation

$$\mathbf{Ax} := \begin{pmatrix} 1 & 1 \\ 1 & 1 - 10^{-12} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 10^{-12} \end{pmatrix} =: \mathbf{b}.$$

The exact solution is given by  $(1 \ -1)^T$ . In Julia, this equation can be solved as follows:

```
A = [1 1; 1 (1-1e-12)]
b = [0; 1e-12]
x = A\b
```

The solution returned by the program is the following:

```
1.0000221222095027
-1.0000221222095027
```



The relative error on the solution is of the order of  $10^{-5}$ , which is about 12 orders of magnitude larger than the machine epsilon for the **Float64** type. In order to understand this, note that the final error on  $\mathbf{x}$  arises from three sources:

- First, the machine representation  $\mathbf{b}$  of the right-hand side is only *an approximation* of the true right-hand side  $\mathbf{b}$ .
- Similarly, the machine representation  $\mathbf{A}$  of the matrix is only *an approximation* of the true matrix  $\mathbf{A}$ .
- Finally, the operation  $\mathbf{A} \setminus \mathbf{b}$  itself leads to additional round-off errors, as the computer implementation of the backslash operator is based on elementary arithmetic operations. Understanding the magnitude of the error introduced at this step is delicate, and we shall not address this question.

It follows from the first two items that, when writing  $\mathbf{A} \setminus \mathbf{b}$ , we are in fact asking the computer for a solution  $\mathbf{x} + \Delta \mathbf{x}$  to a perturbed equation

$$(\mathbf{A} + \Delta \mathbf{A})(\mathbf{x} + \Delta \mathbf{x}) = \mathbf{b} + \Delta \mathbf{b},$$

where  $\Delta \mathbf{A}$  and  $\Delta \mathbf{b}$  represent the rounding errors on the matrix and right-hand side, respectively. Understanding the impact of the perturbations  $\Delta \mathbf{b}$  and  $\Delta \mathbf{A}$  is the goal of this section. We shall prove, in particular, that the magnitude of the error  $\Delta \mathbf{x}$  is related to the perturbations  $\Delta \mathbf{b}$  and  $\Delta \mathbf{A}$  through a quantity known as the *condition number* of the matrix  $\mathbf{A}$ .

In general, the condition number for a given problem measures the sensitivity of the solution to the input data. In order to define this concept precisely, we consider a general problem of the form  $F(x, d) = 0$ , with unknown  $x$  and data  $d$ . The linear system (4.1) can be recast in this form, with the input data equal to  $\mathbf{b}$  or  $\mathbf{A}$  or both. We denote the solution corresponding to perturbed input data  $d + \Delta d$  by  $x + \Delta x$ . The absolute and relative condition numbers are defined as follows.

**Definition 4.1** (Condition number for the problem  $F(x, d) = 0$ ). The absolute and relative condition numbers with respect to perturbations of  $d$  are defined as

$$K_{\text{abs}}(d) = \lim_{\varepsilon \rightarrow 0} \left( \sup_{\|\Delta d\| \leq \varepsilon} \frac{\|\Delta x\|}{\|\Delta d\|} \right), \quad K(d) = \lim_{\varepsilon \rightarrow 0} \left( \sup_{\|\Delta d\| \leq \varepsilon} \frac{\|\Delta x\|/\|x\|}{\|\Delta d\|/\|d\|} \right).$$

The short notation  $K$  is reserved for the relative condition number, which is often more useful in applications.

In the rest of this section, we obtain an upper bound on the relative condition number for the linear system (4.1) with respect to perturbations first of  $\mathbf{b}$ , and then of  $\mathbf{A}$ . We use the notation  $\|\bullet\|$  to denote both a vector norm on  $\mathbf{C}^n$  and the induced operator norm on matrices.

**Proposition 4.1** (Perturbation of the right-hand side). *Let  $\mathbf{x} + \Delta\mathbf{x}$  denote the solution to the perturbed equation  $\mathbf{A}(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{b} + \Delta\mathbf{b}$ . Then it holds that*

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\|\|\mathbf{A}^{-1}\| \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|}, \quad (4.2)$$

*Proof.* It holds by definition of  $\Delta\mathbf{x}$  that  $\mathbf{A}\Delta\mathbf{x} = \Delta\mathbf{b}$ . Therefore, we have

$$\|\Delta\mathbf{x}\| = \|\mathbf{A}^{-1}\Delta\mathbf{b}\| \leq \|\mathbf{A}^{-1}\| \|\Delta\mathbf{b}\| = \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{b}\|} \|\mathbf{A}^{-1}\| \|\Delta\mathbf{b}\| \leq \frac{\|\mathbf{A}\|\|\mathbf{x}\|}{\|\mathbf{b}\|} \|\mathbf{A}^{-1}\| \|\Delta\mathbf{b}\|. \quad (4.3)$$

Here we employed (A.11), proved in Appendix A, in the first and last inequalities. Rearranging the inequality (4.3), we obtain (4.2).  $\square$

Proposition 4.1 implies that the relative condition number of (4.1) with respect to perturbations of the right-hand side is bounded from above by  $\|\mathbf{A}\|\|\mathbf{A}^{-1}\|$ . Exercise 4.1 shows that there are values of  $\mathbf{x}$  and  $\Delta\mathbf{b}$  for which the inequality (4.2) is an equality, indicating that the inequality is sharp.

Studying the impact of perturbations of the matrix  $\mathbf{A}$  is slightly more difficult, because this time the variation  $\Delta\mathbf{x}$  of the solution does not depend linearly on the perturbation of the data. Before stating and proving the main result, we show an ancillary lemma.

**Lemma 4.2.** *Let  $\mathbf{B} \in \mathbf{C}^{n \times n}$  be such that  $\|\mathbf{B}\| < 1$  in some submultiplicative matrix norm  $\|\bullet\|$ . Then  $\mathbf{I} - \mathbf{B}$  is invertible and*

$$\|(\mathbf{I} - \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}, \quad (4.4)$$

*where  $\mathbf{I} \in \mathbf{C}^{n \times n}$  is the identity matrix.*

*Proof.* It holds for any matrix  $\mathbf{B} \in \mathbf{C}^{n \times n}$  that

$$\mathbf{I} - \mathbf{B}^{n+1} = (\mathbf{I} - \mathbf{B})(\mathbf{I} + \mathbf{B} + \cdots + \mathbf{B}^n). \quad (4.5)$$

Since  $\|\mathbf{B}\| < 1$  in a submultiplicative matrix norm, both sides of the equation are convergent in the limit as  $n \rightarrow \infty$ . The left-hand side converges to the identity matrix  $\mathbf{I}$ , and the right-hand side converges as  $n \rightarrow \infty$  because  $\{\mathbf{S}_0, \mathbf{S}_1, \dots\}$  with

$$\mathbf{S}_n := \mathbf{I} + \mathbf{B} + \cdots + \mathbf{B}^n$$

is a Cauchy sequence in the vector space of matrices endowed with the norm for which  $\|\mathbf{B}\| < 1$ . Indeed, by the triangle inequality and the submultiplicative property of the norm, it holds that

$$\begin{aligned} \|\mathbf{S}_{n+m} - \mathbf{S}_n\| &= \|\mathbf{B}^{n+1} + \cdots + \mathbf{B}^{n+m}\| \\ &\leq \|\mathbf{B}^{n+1}\| + \cdots + \|\mathbf{B}^{n+m}\| \leq \|\mathbf{B}\|^{n+1} + \cdots + \|\mathbf{B}\|^{n+m} \\ &\leq \frac{\|\mathbf{B}\|^{n+1}}{1 - \|\mathbf{B}\|} \xrightarrow{n \rightarrow \infty} 0, \end{aligned}$$

where we employed the formula for a geometric series in the last inequality. Equating the limits of both sides of (4.5), we obtain

$$\mathbf{I} = (\mathbf{I} - \mathbf{B}) \sum_{i=0}^{\infty} \mathbf{B}^i.$$

This implies that  $(\mathbf{I} - \mathbf{B})$  is invertible with inverse given by a so-called *Neumann series*

$$(\mathbf{I} - \mathbf{B})^{-1} = \sum_{i=0}^{\infty} \mathbf{B}^i.$$

Applying the triangle inequality repeatedly, and then using the submultiplicative property of the norm, we obtain

$$\forall n \in \mathbf{N}, \quad \left\| \sum_{i=0}^n \mathbf{B}^i \right\| \leq \sum_{i=0}^n \|\mathbf{B}^i\| \leq \sum_{i=0}^n \|\mathbf{B}\|^i = \frac{1}{1 - \|\mathbf{B}\|}.$$

where we used the summation formula for geometric series in the last equality. Letting  $n \rightarrow \infty$  in this equation and using the continuity of the norm enables to conclude the proof.  $\square$

**Proposition 4.3** (Perturbation of the matrix). *Let  $\mathbf{x} + \Delta\mathbf{x}$  denote the solution to the perturbed equation  $(\mathbf{A} + \Delta\mathbf{A})(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{b}$ . If  $\mathbf{A}$  is invertible and  $\|\Delta\mathbf{A}\| < \|\mathbf{A}^{-1}\|^{-1}$ , then*

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} \left( \frac{1}{1 - \|\mathbf{A}^{-1}\Delta\mathbf{A}\|} \right). \quad (4.6)$$

*Proof.* Left-multiplying both sides of the perturbed equation with  $\mathbf{A}^{-1}$ , we obtain

$$(\mathbf{I} + \mathbf{A}^{-1}\Delta\mathbf{A})(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{x} \quad \Leftrightarrow \quad (\mathbf{I} + \mathbf{A}^{-1}\Delta\mathbf{A})\Delta\mathbf{x} = -\mathbf{A}^{-1}\Delta\mathbf{A}\mathbf{x}. \quad (4.7)$$

Since  $\|\mathbf{A}^{-1}\Delta\mathbf{A}\| \leq \|\mathbf{A}^{-1}\| \|\Delta\mathbf{A}\| < 1$  by assumption, we deduce from [Lemma 4.2](#) that the matrix on the left-hand side is invertible with a norm bounded as in (4.4). Consequently, using in addition the assumed submultiplicative property of the norm, we obtain that

$$\|\Delta\mathbf{x}\| = \|(\mathbf{I} + \mathbf{A}^{-1}\Delta\mathbf{A})^{-1}\mathbf{A}^{-1}\Delta\mathbf{A}\mathbf{x}\| \leq \frac{\|\mathbf{A}^{-1}\Delta\mathbf{A}\|}{1 - \|\mathbf{A}^{-1}\Delta\mathbf{A}\|} \|\mathbf{x}\|.$$

which enables to conclude the proof.  $\square$

Using [Proposition 4.3](#), we deduce that the relative condition number of (4.1) with respect to perturbations of the matrix  $\mathbf{A}$  is also bounded from above by  $\|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ , because the term between brackets on the right-hand side of (4.6) converges to 1 as  $\|\Delta\mathbf{A}\| \rightarrow 0$ .

[Propositions 4.1](#) and [4.3](#) show that the condition number of the linear system (4.1), with respect to perturbations of either  $\mathbf{b}$  or  $\mathbf{A}$ , depends only on  $\mathbf{A}$ . This motivates the following definition of the condition number of a matrix.

**Definition 4.2** (Condition number of a matrix). The condition number of a matrix  $\mathbf{A}$  asso-

ciated with a vector norm  $\|\bullet\|$  is defined as

$$\kappa(\mathbf{A}) = \|\mathbf{A}\|\|\mathbf{A}^{-1}\|.$$

The condition number for the  $p$ -norm, defined in Definition A.4, is denoted by  $\kappa_p(\mathbf{A})$ .

Note that the condition number  $\kappa(\mathbf{A})$  associated with an induced norm, i.e. a matrix norm induced by a vector norm, is at least one. Indeed, since the identity matrix has induced norm 1, it holds that

$$1 = \|\mathbf{I}\| = \|\mathbf{A}\mathbf{A}^{-1}\| \leq \|\mathbf{A}\|\|\mathbf{A}^{-1}\|.$$

Since the 2-norm of an invertible matrix  $\mathbf{A} \in \mathbf{C}^{n \times n}$  coincides with the spectral radius  $\rho(\mathbf{A}^T\mathbf{A})$ , the condition number  $\kappa_2$  corresponding to the 2-norm is equal to

$$\kappa_2(\mathbf{A}) = \sqrt{\frac{\lambda_{\max}(\mathbf{A}^T\mathbf{A})}{\lambda_{\min}(\mathbf{A}^T\mathbf{A})}},$$

where  $\lambda_{\max}(\mathbf{A}^T\mathbf{A})$  and  $\lambda_{\min}(\mathbf{A}^T\mathbf{A})$  are the maximal and minimal (both real and positive) eigenvalues of the matrix  $\mathbf{A}^T\mathbf{A}$ .

*Example 4.2* (Perturbation of the matrix). Consider the following linear system with perturbed matrix

$$(\mathbf{A} + \Delta\mathbf{A}) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0.01 \end{pmatrix}, \quad \mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 0.01 \end{pmatrix}, \quad \Delta\mathbf{A} = \begin{pmatrix} 0 & 0 \\ 0 & \varepsilon \end{pmatrix},$$

where  $0 < \varepsilon \ll 0.01$ . Here the eigenvalues of  $\mathbf{A}$  are given by  $\lambda_1 = 1$  and  $\lambda_2 = 0.01$ . The solution when  $\varepsilon = 0$  is given by  $(0, 1)^T$ , and the solution to the perturbed equation is

$$\begin{pmatrix} x_1 + \Delta x_1 \\ x_2 + \Delta x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ \frac{1}{1+100\varepsilon} \end{pmatrix}.$$

Consequently, we deduce that, in the 2-norm,

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} = \left| \frac{100\varepsilon}{1+100\varepsilon} \right| \approx 100\varepsilon = 100 \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|}.$$

In this case, the relative impact of perturbations of the matrix is close to  $\kappa_2(\mathbf{A}) = 100$ .

## 4.2 Direct solution method

In this section, we present the *direct method* for solving linear systems of the form (4.1) with a general invertible matrix  $\mathbf{A} \in \mathbf{C}^{n \times n}$ . The direct method can be decomposed into three steps:

- First calculate the so-called LU decomposition of  $\mathbf{A}$ , i.e. find an upper triangular matrix  $\mathbf{U}$  and a *unit* lower triangular matrix  $\mathbf{L}$  such that  $\mathbf{A} = \mathbf{L}\mathbf{U}$ . A unit lower triangular matrix is a lower triangular matrix with only ones on the diagonal.

- Then solve  $\mathbf{L}\mathbf{y} = \mathbf{b}$  using a method called *forward substitution*.
- Finally, solve  $\mathbf{U}\mathbf{x} = \mathbf{y}$  using a method called *backward substitution*.

By construction, the solution  $\mathbf{x}$  thus obtained is a solution to (4.1). Indeed, we have that

$$\mathbf{A}\mathbf{x} = \mathbf{L}\mathbf{U}\mathbf{x} = \mathbf{L}\mathbf{y} = \mathbf{b}.$$

### 4.2.1 LU decomposition

In this section, we first discuss the existence and uniqueness of the LU factorization of a matrix. We then describe a numerical algorithm for calculating the factors  $\mathbf{L}$  and  $\mathbf{U}$ , based on *Gaussian elimination*.

#### Existence and uniqueness of the decomposition

We present a necessary and sufficient condition for the existence of a unique LU decomposition of a matrix. To this end, we define the principal submatrix of order  $i$  of a matrix  $\mathbf{A} \in \mathbf{C}^{n \times n}$  as the matrix  $\mathbf{A}_i = \mathbf{A}[1 : i, 1 : i]$ , in Julia notation.

**Proposition 4.4.** *The LU factorization of a matrix  $\mathbf{A} \in \mathbf{C}^{n \times n}$ , with  $\mathbf{L}$  unit lower triangular and  $\mathbf{U}$  upper triangular, exists and is unique if and only if the principal submatrices of  $\mathbf{A}$  of all orders are nonsingular.*

*Proof.* We prove only the “if” direction; see [10, Theorem 3.4] for the “only if” implication.

The statement is clear if  $n = 1$ . Reasoning by induction, we assume that the result is proved up to  $n - 1$ . Since the matrix  $\mathbf{A}_{n-1}$  and all its principal submatrices are nonsingular by assumption, it holds that  $\mathbf{A}_{n-1} = \mathbf{L}_{n-1}\mathbf{U}_{n-1}$  for a unit lower triangular matrix  $\mathbf{L}_{n-1}$  and an upper triangular matrix  $\mathbf{U}_{n-1}$ . These two matrices are nonsingular, for if either of them were singular then the product  $\mathbf{A}_{n-1} = \mathbf{L}_{n-1}\mathbf{U}_{n-1}$  would be singular as well. Let us decompose  $\mathbf{A}$  as follows:

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{n-1} & \mathbf{c} \\ \mathbf{d}^T & a_{nn} \end{pmatrix}.$$

Let  $\boldsymbol{\ell}$  and  $\mathbf{u}$  denote the solutions to  $\mathbf{L}_{n-1}\mathbf{u} = \mathbf{c}$  and  $\mathbf{U}_{n-1}^T\boldsymbol{\ell} = \mathbf{d}$ . These solutions exist and are unique, because the matrices  $\mathbf{L}_{n-1}$  and  $\mathbf{U}_{n-1}$  are nonsingular. Letting  $u_{nn} = a_{nn} - (\boldsymbol{\ell}^T\mathbf{u})^{-1}$ , we check that  $\mathbf{A}$  factorizes as

$$\begin{pmatrix} \mathbf{A}_{n-1} & \mathbf{c} \\ \mathbf{d}^T & a_{nn} \end{pmatrix} = \begin{pmatrix} \mathbf{L}_{n-1} & \mathbf{0}_{n-1} \\ \boldsymbol{\ell}^T & 1 \end{pmatrix} \begin{pmatrix} \mathbf{U}_{n-1} & \mathbf{u} \\ \mathbf{0}_{n-1}^T & u_{nn} \end{pmatrix}.$$

This completes the proof of the existence of the decomposition. The uniqueness of the factors follows from the uniqueness of  $\boldsymbol{\ell}$ ,  $\mathbf{u}$  and  $u_{nn}$ . □

Proposition 4.4 raises the following question: are there classes of matrices whose principal matrices are all nonsingular? The answer is positive, and we mention, as an important example, the class of positive definite matrices. Proving this is the aim of Exercise 4.4.

### Gaussian elimination algorithm for computing L and U

So far we have presented a condition under which the LU decomposition of a matrix exists and is unique, but not a practical method for calculating the matrices L and U. We describe in this section an algorithm, known as *Gaussian elimination*, for calculating the LU decomposition of a matrix. We begin by introducing the concept of *Gaussian transformation*.

**Definition 4.3.** A Gaussian transformation is a matrix of the form  $M_k = I - \mathbf{c}^{(k)} \mathbf{e}_k^T$ , where  $\mathbf{e}_k$  is the column vector with entry at index  $k$  equal to 1 and all the other entries equal to zero, and  $\mathbf{c}^{(k)}$  is a column vector of the following form:

$$\mathbf{c}^{(k)} = \begin{pmatrix} 0 & 0 & \dots & 0 & c_{k+1}^{(k)} & c_{k+2}^{(k)} & \dots & c_n^{(k)} \end{pmatrix}^T.$$

The action of a Gaussian transformation  $M_k$  left-multiplying a matrix  $A \in \mathbf{C}^{n \times n}$  is to replace the rows from index  $k+1$  to index  $n$  by a linear combination involving themselves and the  $k$ -th row. To see this, let us denote by  $(\mathbf{r}^{(i)})_{1 \leq i \leq n}$  the rows of a matrix  $T \in \mathbf{C}^{n \times n}$ . Then, we have

$$M_k T = (I - \mathbf{c}^{(k)} \mathbf{e}_k^T) T = \begin{pmatrix} 1 & & & & & & & & \\ & 1 & & & & & & & \\ & & \ddots & & & & & & \\ & & & 1 & & & & & \\ & & & -c_{k+1}^{(k)} & 1 & & & & \\ & & & \vdots & & \ddots & & & \\ & & & -c_n^{(k)} & & & \ddots & & \\ & & & & & & & 1 & \end{pmatrix} \begin{pmatrix} \mathbf{r}^{(1)} \\ \mathbf{r}^{(2)} \\ \vdots \\ \mathbf{r}^{(k)} \\ \mathbf{r}^{(k+1)} \\ \vdots \\ \mathbf{r}^{(n)} \end{pmatrix} = \begin{pmatrix} \mathbf{r}^{(1)} \\ \mathbf{r}^{(2)} \\ \vdots \\ \mathbf{r}^{(k)} \\ \mathbf{r}^{(k+1)} - c_{k+1}^{(k)} \mathbf{r}^{(k)} \\ \vdots \\ \mathbf{r}^{(n)} - c_n^{(k)} \mathbf{r}^{(k)} \end{pmatrix}$$

We show in [Exercise 4.2](#) that the inverse of a Gaussian transformation matrix is given by

$$(I - \mathbf{c}^{(k)} \mathbf{e}_k^T)^{-1} = I + \mathbf{c}^{(k)} \mathbf{e}_k^T. \tag{4.8}$$

The idea of the Gaussian elimination algorithm is to successively left-multiply A with Gaussian transformation matrices  $M_1$ , then  $M_2$ , etc. appropriately chosen in such a way that the matrix  $A^{(k)}$ , obtained after  $k$  iterations, is upper triangular up to column  $k$ . That is to say, the Gaussian transformations are constructed so that all the entries in columns 1 to  $k$  under the diagonal of the matrix  $A^{(k)}$  are equal to zero. The resulting matrix  $A^{(n-1)}$  after  $n-1$  iterations is then upper triangular and satisfies

$$A^{(n-1)} = M_{n-1} \dots M_1 A.$$

Rearranging this equation, we deduce that

$$A = (M_1^{-1} \dots M_{n-1}^{-1}) A^{(n-1)}.$$

The first factor is lower triangular by (4.8) and [Exercise 4.3](#). The product in the definition of the matrix L admits a simple explicit expression.







## Computer implementation

The Gaussian elimination procedure is summarized as follows.

```

A(0) ← A, L ← I
for i ∈ {1, ..., n - 1} do
    Construct Mi as in Lemma 4.6.
    A(i) ← MiA(i-1), L ← LMi-1
end for
U ← A(n-1).

```

In practice, it is not necessary to explicitly create the Gaussian transformation matrices, or to perform full matrix multiplications. A more realistic version of the algorithm in Julia is given below. The code exploits the relation (4.9) between L and the parameters of the Gaussian transformations.

```

1  # A is an invertible matrix of size n x n
2  L = [i == j ? 1.0 : 0.0 for i in 1:n, j in 1:n]
3  U = copy(A)
4  for i in 1:n-1
5      for r in i+1:n
6          U[i, i] == 0 && error("Pivotal entry is zero!")
7          ratio = U[r, i] / U[i, i]
8          L[r, i] = ratio
9          U[r, i:end] -= U[i, i:end] * ratio
10     end
11 end
12 # L is unit lower triangular and U is upper triangular

```

## Computational cost

The computational cost of the algorithm, measured as the number of floating point operations (flops) required, is dominated by the Gaussian transformations, in line 9 in the above code. All the other operations amount to a computational cost scaling as  $\mathcal{O}(n^2)$ , which is negligible compared to the cost of the LU factorization when  $n$  is large. This factorization requires

$$\underbrace{2}_{\text{and } * \text{ for } i \text{ in } 1:n-1} \times \underbrace{\sum_{i=1}^{n-1}}_{\text{for } i \text{ in } 1:n-1} \underbrace{(n-i)}_{\text{for } r \text{ in } i+1:n} \underbrace{(n-i+1)}_{\text{indices } [i:\text{end}]} \text{ flops} = \frac{2}{3}n^3 + \mathcal{O}(n^2) \text{ flops.}$$

### 4.2.2 Backward and forward substitution

Once the LU factorization has been completed, the solution to the linear system can be obtained by first using forward, and then backward substitution, which are just bespoke methods for solving linear systems with lower and upper triangular matrices, respectively. Let us consider the case of a lower triangular system:

$$Ly = b$$

Notice that the unknown  $y_1$  may be obtained from the first equation of the system. Then, since  $y_1$  is known, the value of  $y_2$  can be obtained from the second equation, etc. A simple implementation of this algorithm is as follows:

```
# L is unit lower triangular
y = copy(b)
for i in 2:n
    for j in 1:i-1
        y[i] -= L[i, j] * y[j]
    end
end
```

### 4.2.3 Gaussian elimination with pivoting

The Gaussian elimination algorithm that we presented in [Section 4.2.1](#) relies on the existence of an LU factorization. In practice, this assumption may not be satisfied, and in this case a modified algorithm, called Gaussian elimination *with pivoting*, is required.

In fact, pivoting is useful even if the usual LU decomposition of  $A$  exists, as it enables to reduce the condition number of the matrices  $L$  and  $U$ . There are two types of pivoting: partial pivoting, where only the rows are rearranged through a permutation at each iteration, and complete pivoting, where both the rows and the columns are rearranged at each iteration.

Showing rigorously why pivoting is useful is beyond the scope of this course. In this section, we only present the partial pivoting method. Its influence on the condition number of the factors  $L$  and  $U$  is studied empirically in [Exercise 4.6](#). It is useful at this point to introduce the concept of a row permutation matrix.

#### Row permutation matrix

**Definition 4.4.** Let  $\sigma : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  be a permutation, i.e. a bijection on the set  $\{1, \dots, n\}$ . The row permutation matrix associated with  $\sigma$  is the matrix with entries

$$p_{ij} = \begin{cases} 1 & \text{if } i = \sigma(j), \\ 0 & \text{otherwise.} \end{cases}$$

When a row permutation  $P$  left-multiplies a matrix  $B \in \mathbf{C}^{n \times n}$ , row  $i$  of matrix  $B$  is moved to row index  $\sigma(i)$  in the resulting matrix, for all  $i \in \{1, \dots, n\}$ . A permutation matrix has a single entry equal to 1 per row and per column, and its inverse coincides with its transpose:  $P^{-1} = P^T$ .

#### Partial pivoting

Gaussian elimination with partial pivoting applies for any invertible matrix  $A$ , and it outputs 3 matrices: a row permutation  $P$ , a unit triangular matrix  $L$ , and an upper triangular matrix  $U$ . These are related by the relation

$$PA = LU.$$

This is sometimes called a PLU decomposition of the matrix  $\mathbf{A}$ . It is not unique in general but, unlike the usual LU decomposition, it always exists provided that  $\mathbf{A}$  is invertible. We take this for granted in this course.

The idea of partial pivoting is to rearrange the rows at each iteration of the Gaussian elimination procedure in such a manner that the pivotal entry is as large as possible in absolute value. One step of the procedure reads

$$\mathbf{A}^{(k+1)} = \mathbf{M}_{k+1}\mathbf{P}_{k+1}\mathbf{A}^{(k)}. \quad (4.10)$$

Here  $\mathbf{P}_{k+1}$  is a simple row permutation matrix which, when acting on  $\mathbf{A}^{(k)}$ , interchanges row  $k+1$  and row  $\ell$ , for some index  $\ell \geq k+1$ . The row index  $\ell$  is selected in such a way that the absolute value of the pivotal entry, in position  $(k+1, k+1)$  of the product  $\mathbf{P}_{k+1}\mathbf{A}^{(k)}$ , is maximum. The matrix  $\mathbf{M}_{k+1}$  is then the unique Gaussian transformation matrix ensuring that  $\mathbf{A}^{(k+1)}$  is upper triangular up to column  $k+1$ , obtained as in Lemma 4.6. The resulting matrix  $\mathbf{A}^{(n-1)}$  after  $n-1$  steps of the form (4.10) is upper triangular and satisfies

$$\mathbf{A}^{(n-1)} = \mathbf{M}_{n-1}\mathbf{P}_{n-1} \cdots \mathbf{M}_1\mathbf{P}_1\mathbf{A} \quad \Leftrightarrow \quad \mathbf{A} = (\mathbf{M}_{n-1}\mathbf{P}_{n-1} \cdots \mathbf{M}_1\mathbf{P}_1)^{-1}\mathbf{A}^{(n-1)}.$$

The first factor in the decomposition of  $\mathbf{A}$  is not necessarily lower triangular. However, using the notation  $\mathbf{M} = \mathbf{M}_{n-1}\mathbf{P}_{n-1} \cdots \mathbf{M}_1\mathbf{P}_1$  and  $\mathbf{P} = \mathbf{P}_{n-1} \cdots \mathbf{P}_1$ , we have

$$\mathbf{P}\mathbf{A} = \mathbf{P}\mathbf{M}^{-1}\mathbf{U} = (\mathbf{P}\mathbf{M}^{-1})\mathbf{U} =: \mathbf{L}\mathbf{U}. \quad (4.11)$$

Lemma 4.8 below shows that, as the notation  $\mathbf{L}$  suggests, the matrix  $\mathbf{L} = (\mathbf{P}\mathbf{M}^{-1})$  on the right-hand side is indeed lower triangular. Before stating and proving the lemma, we note that  $\mathbf{P}$  is a row permutation matrix, and so the solution to the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be obtained by solving  $\mathbf{L}\mathbf{U}\mathbf{x} = \mathbf{P}^T\mathbf{b}$  by forward and backward substitution. Since  $\mathbf{P}$  is a very sparse matrix, the right-hand side  $\mathbf{P}^T\mathbf{b}$  can be calculated very efficiently.

**Lemma 4.8.** *The matrix  $\mathbf{L} = \mathbf{P}\mathbf{M}^{-1}$  is unit lower triangular with all entries bounded in absolute value from above by 1. It admits the expression*

$$\mathbf{L} = \mathbf{I} + (\mathbf{P}_{n-1} \cdots \mathbf{P}_2\mathbf{c}^{(1)})\mathbf{e}_1^T + (\mathbf{P}_{n-1} \cdots \mathbf{P}_3\mathbf{c}^{(2)})\mathbf{e}_2^T + \cdots + (\mathbf{P}_{n-1}\mathbf{c}^{(n-2)})\mathbf{e}_{n-2}^T + \mathbf{c}^{(n-1)}\mathbf{e}_{n-1}^T.$$

*Proof.* Let  $\mathbf{M}^{(k)} = \mathbf{M}_k\mathbf{P}_k \cdots \mathbf{M}_1\mathbf{P}_1$  and  $\mathbf{P}^{(k)} = \mathbf{P}_k \cdots \mathbf{P}_1$ . It is sufficient to show that

$$\mathbf{P}^{(k)}(\mathbf{M}^{(k)})^{-1} = \mathbf{I} + (\mathbf{P}_k \cdots \mathbf{P}_2\mathbf{c}^{(1)})\mathbf{e}_1^T + (\mathbf{P}_k \cdots \mathbf{P}_3\mathbf{c}^{(2)})\mathbf{e}_2^T + \cdots + (\mathbf{P}_k\mathbf{c}^{(k-1)})\mathbf{e}_{k-1}^T + \mathbf{c}^{(k)}\mathbf{e}_k^T \quad (4.12)$$

for all  $k \in \{1, \dots, n-1\}$ . The statement is clear for  $k = 1$ , and we assume by induction that it

is true up to  $k - 1$ . Then notice that

$$\begin{aligned} \mathbf{P}^{(k)}(\mathbf{M}^{(k)})^{-1} &= \mathbf{P}_k \left( \mathbf{P}^{(k-1)}(\mathbf{M}^{(k-1)})^{-1} \right) \mathbf{P}_k^{-1} \mathbf{M}_k^{-1} \\ &= \mathbf{P}_k \left( \mathbf{I} + (\mathbf{P}_{k-1} \cdots \mathbf{P}_2 \mathbf{c}^{(1)}) \mathbf{e}_1^T + \cdots + (\mathbf{P}_{k-1} \mathbf{c}^{(k-2)}) \mathbf{e}_{k-2}^T + \mathbf{c}^{(k-1)} \mathbf{e}_{k-1}^T \right) \mathbf{P}_k^{-1} \mathbf{M}_k^{-1} \\ &= \left( \mathbf{I} + (\mathbf{P}_k \mathbf{P}_{k-1} \cdots \mathbf{P}_2 \mathbf{c}^{(1)}) \mathbf{e}_1^T + \cdots + (\mathbf{P}_k \mathbf{P}_{k-1} \mathbf{c}^{(k-2)}) \mathbf{e}_{k-2}^T + (\mathbf{P}_k \mathbf{c}^{(k-1)}) \mathbf{e}_{k-1}^T \right) \mathbf{M}_k^{-1}. \end{aligned}$$

In the last equality, we used that  $\mathbf{e}_i^T \mathbf{P}_k^{-1} = (\mathbf{P}_k \mathbf{e}_i)^T = \mathbf{e}_i^T$  for all  $i \in \{1, \dots, k-1\}$ , because the row permutation  $\mathbf{P}_k$  does not affect rows 1 to  $k-1$ . Using the expression  $\mathbf{M}_k^{-1} = \mathbf{I} + \mathbf{c}^{(k)} \mathbf{e}_k^T$ , expanding the product and noting that  $\mathbf{e}_j^T \mathbf{c}^{(k)} = 0$  if  $j \leq k$ , we obtain (4.12). The statement that the entries are bounded in absolute value from above by 1 follows from the choice of the pivot at each iteration.  $\square$

The expression of  $\mathbf{L}$  in Lemma 4.8 suggests the iterative procedure given in Algorithm 2 for performing the LU factorization with partial pivoting. A Julia implementation of this algorithm is presented in Listing 1.

---

**Algorithm 2** LU decomposition with partial pivoting

---

Assign  $\mathbf{A}^{(0)} \leftarrow \mathbf{A}$  and  $\mathbf{P} \leftarrow \mathbf{I}$   
**for**  $i \in \{1, \dots, n-1\}$  **do**  
    Find the row index  $k \geq i$  such that  $\mathbf{A}_{k,i}^{(i-1)}$  is maximum in absolute value.  
    Interchange the rows  $i$  and  $k$  of matrices  $\mathbf{A}^{(i-1)}$  and  $\mathbf{P}$ , and of vectors  $\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(i-1)}$ .  
    Construct  $\mathbf{M}_i$  with corresponding column vector  $\mathbf{c}^{(i)}$  as in Lemma 4.6.  
    Assign  $\mathbf{A}^{(i)} \leftarrow \mathbf{M}_i \mathbf{A}^{(i-1)}$   
**end for**  
Assign  $\mathbf{U} \leftarrow \mathbf{A}^{(n-1)}$ .  
Assign  $\mathbf{L} \leftarrow \mathbf{I} + \begin{pmatrix} \mathbf{c}^{(1)} & \dots & \mathbf{c}^{(n-1)} & \mathbf{0}_n \end{pmatrix}$ .

---

```
# Auxiliary function
function swap_rows!(i, j, matrices...)
    for M in matrices
        M_row_i = M[i, :]
        M[i, :] = M[j, :]
        M[j, :] = M_row_i
    end
end

n = size(A)[1]
L, U = zeros(n, 0), copy(A)
P = [i == j ? 1.0 : 0.0 for i in 1:n, j in 1:n]
for i in 1:n-1
    # Pivoting
    index_row_pivot = i - 1 + argmax(abs.(U[i:end, i]))
    swap_rows!(i, index_row_pivot, U, L, P)

    # Usual Gaussian transformation
    c = [zeros(i-1); 1.0; zeros(n-i)]
    for r in i+1:n
```

```

    ratio = U[r, i] / U[i, i]
    c[r] = ratio
    U[r, i:end] -= U[i, i:end] * ratio
end
L = [L c]
end
L = [L [zeros(n-1); 1.0]]
# It holds that P*A = L*U

```

Listing 1: LU factorization with partial pivoting.

*Remark 4.1.* It is possible to show that, if the matrix  $A$  is column diagonally dominant in the sense that

$$\forall j \in \{1, \dots, n\}, \quad |a_{jj}| \geq \sum_{i=1, i \neq j}^n |a_{ij}|,$$

then pivoting does not have an effect: at each iteration, the best pivot is already on the diagonal.

#### 4.2.4 Direct method for Hermitian positive definite matrices

The LU factorization with partial pivoting applies to any matrix  $A \in \mathbf{C}^{n \times n}$  that is invertible. If  $A$  is Hermitian positive definite, however, it is possible to compute a factorization into lower and upper triangular matrices at half the computational cost, using the so-called *Cholesky decomposition*.

**Lemma 4.9** (Cholesky decomposition). *If  $A$  is Hermitian positive definite, then there exists a lower-triangular matrix  $C \in \mathbf{C}^{n \times n}$  such that*

$$A = CC^*. \quad (4.13)$$

*Equation (4.13) is called the Cholesky factorization of  $A$ . The matrix  $C$  is unique if we require that all its diagonal entries are positive.*

*Proof.* Since  $A$  is positive definite, its LU decomposition exists and is unique by Propositions 4.4 and 4.7. Let  $D$  denote the diagonal matrix with the same diagonal as that of  $U$ . Then

$$A = LD(D^{-1}U).$$

Note that the matrix  $D^{-1}U$  is unit upper triangular. Since  $A$  is Hermitian, we have

$$A = A^* = (D^{-1}U)^*(LD)^*.$$

The first and second factors on the right-hand side are respectively unit lower triangular and upper triangular, and so we deduce, by uniqueness of the LU decomposition, that  $L = (D^{-1}U)^*$  and  $U = (LD)^*$ . But then

$$A = LU = LDL^* = (L\sqrt{D})(\sqrt{D}L)^*.$$

Here  $\sqrt{D}$  denotes the diagonal matrix whose diagonal entries are obtained by taking the square root of those of  $D$ , which are necessarily real and positive because  $A$  is positive definite. This implies the existence of a Cholesky factorization with  $C = L\sqrt{D}$ .  $\square$

### Calculation of the Choleski factor

The matrix  $C$  can be calculated from (4.13). For example, developing the matrix product gives that  $a_{1,1} = c_{1,1}^2$  and so  $c_{1,1} = \sqrt{a_{1,1}}$ . It is then possible to calculate  $c_{2,1}$  from the equation  $a_{2,1} = c_{2,1}c_{1,1}$ , and so on. Implementing the Cholesky factorization is the goal of Exercise 4.7.

### 4.2.5 Direct methods for banded matrices

In applications related to partial differential equations, the matrix  $A \in \mathbf{C}^{n \times n}$  very often has a bandwidth which is small in comparison with  $n$ .

**Definition 4.5.** The bandwidth of a matrix  $A \in \mathbf{C}^{n \times n}$  is the smallest number  $k \in \mathbf{N}$  such that  $a_{ij} = 0$  for all  $(i, j) \in \{1, \dots, n\}^2$  with  $|i - j| > k$ .

It is not difficult to show that, if  $A$  is a matrix with bandwidth  $k$ , then so are  $L$  and  $U$  in the absence of pivoting. This can be proved by equating the entries of the product  $LU$  with those of the matrix  $A$ . We emphasize, however, that the sparsity structure *within* the band of  $A$  may be destroyed in  $L$  and  $U$ ; this phenomenon is called *fill-in*.

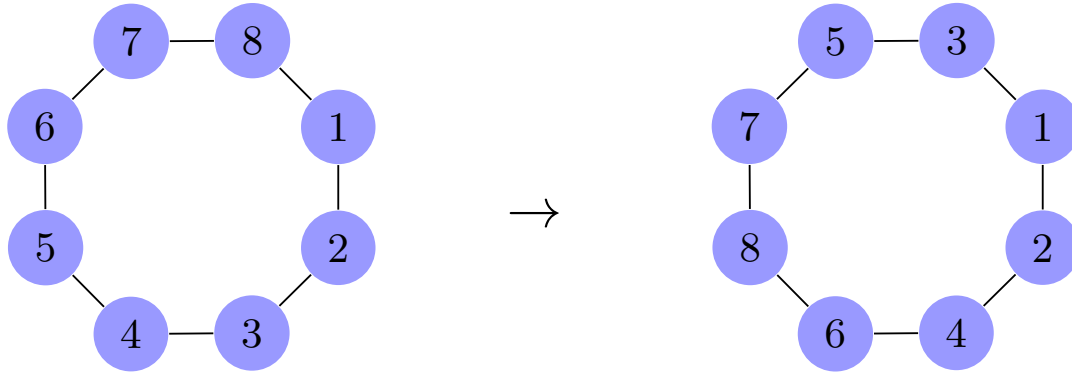
### Reducing the bandwidth: the Cuthill–McKee algorithm

The computational cost of calculating the LU or Cholesky decomposition of a matrix with bandwidth  $k$  scales as  $\mathcal{O}(nk^2)$ , which is much better than the general scaling  $\mathcal{O}(n^3)$  if  $k \ll n$ . In applications, the bandwidth  $k$  is often related to the matrix size  $n$ . For example, if  $A$  arises from the discretization of the Laplacian operator, then  $k = \mathcal{O}(\sqrt{n})$  provided that a good ordering of the vertices is employed. In this case, the computational cost scales as  $\mathcal{O}(n^2)$ .

Since a narrow band is associated with a lower computational cost of the LU decomposition, it is natural to wonder whether the bandwidth of a matrix  $A$  can be reduced. A possible strategy to this end is to use permutations. More precisely, is it possible to identify a row permutation matrix  $P$  such that  $PAP^T$  has minimal bandwidth? Given such a matrix, the solution to the linear system (4.1) can be obtained by first solving  $(PAP^T)\mathbf{y} = P\mathbf{b}$ , and then letting  $\mathbf{x} = P^T\mathbf{y}$ .

The Cuthill–McKee algorithm is a heuristic method for finding a good, but sometimes not optimal, permutation matrix  $P$  in the particular case where  $A$  is a *Hermitian* matrix. It is based on the fact that, to a Hermitian matrix  $A$ , we can associate a unique undirected graph whose adjacency matrix  $A_*$  has the same sparsity structure as that of  $A$ , i.e. zeros in the same places. For any row permutation matrix  $P_\sigma$  with corresponding permutation  $\sigma: \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  (see Definition 4.4), the matrices  $P_\sigma A P_\sigma^T$  and  $P_\sigma A_* P_\sigma^T$  also have the same sparsity structure. Therefore, minimizing the bandwidth of  $P_\sigma A P_\sigma^T$  is equivalent to minimizing the bandwidth of  $P_\sigma A_* P_\sigma^T$ . The key insight for understanding the Cuthill–McKee method is that  $P_\sigma A_* P_\sigma^T$  is the adjacency matrix of the graph obtained by renumbering the nodes according to the

permutation  $\sigma$ , i.e. by changing the number of the nodes from  $i$  to  $\sigma(i)$ . Consider, for example, the following graph and renumbering:



The associated adjacency matrices are given by:

$$\begin{pmatrix} 1 & 1 & & & & & & 1 \\ 1 & 1 & 1 & & & & & \\ & 1 & 1 & 1 & & & & \\ & & 1 & 1 & 1 & & & \\ & & & 1 & 1 & 1 & & \\ & & & & 1 & 1 & 1 & \\ & & & & & 1 & 1 & 1 \\ & & & & & & 1 & 1 \\ 1 & & & & & & & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 1 & 1 & & & & & \\ 1 & 1 & & 1 & & & & \\ 1 & & 1 & & 1 & & & \\ & 1 & & 1 & & 1 & & \\ & & 1 & & 1 & & 1 & \\ & & & 1 & & 1 & & 1 \\ & & & & 1 & & 1 & 1 \\ & & & & & 1 & 1 & 1 \\ & & & & & & 1 & 1 \end{pmatrix}$$

We assume that the nodes are all self-connected, although this is not depicted, and so the diagonal entries of the adjacency matrices are equal to 1. This renumbering corresponds to the permutation

$$\begin{pmatrix} i : & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ \sigma(i) : & 1 & 2 & 4 & 6 & 8 & 7 & 5 & 3 \end{pmatrix},$$

and we may verify that the adjacency matrix of the renumbered graph can be obtained from the associated row permutation matrix:

$$PA_*P^T = \begin{pmatrix} 1 & & & & & & & \\ & 1 & & & & & & \\ & & 1 & & & & & \\ & & & 1 & & & & \\ & & & & 1 & & & \\ & & & & & 1 & & \\ & & & & & & 1 & \\ & & & & & & & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & & & & & & 1 \\ 1 & 1 & 1 & & & & & \\ & 1 & 1 & 1 & & & & \\ & & 1 & 1 & 1 & & & \\ & & & 1 & 1 & 1 & & \\ & & & & 1 & 1 & 1 & \\ & & & & & 1 & 1 & 1 \\ & & & & & & 1 & 1 \\ 1 & & & & & & & \end{pmatrix} \begin{pmatrix} 1 & & & & & & & \\ & 1 & & & & & & \\ & & 1 & & & & & \\ & & & 1 & & & & \\ & & & & 1 & & & \\ & & & & & 1 & & \\ & & & & & & 1 & \\ & & & & & & & 1 \end{pmatrix}$$

In this example, renumbering the nodes of the graph enables a significant reduction of the bandwidth, from 7 to 2. The Cuthill–McKee algorithm, which was employed to calculate the

permutation, is an iterative method that produces an ordered  $n$ -tuple  $R$  containing the nodes in the new order; in other words, it returns  $(\sigma^{-1}(1), \dots, \sigma^{-1}(n))$ . The first step of the algorithm is to find the node  $i$  with the lowest *degree*, i.e. with the smallest number of connections to other nodes, and to initialize  $R = (i)$ . Then the following steps are repeated until  $R$  contains all the nodes of the graph:

- Define  $A_i$  as the set containing all the nodes which are adjacent to a node in  $R$  but not themselves in  $R$ ;
- Sort the nodes in  $A_i$  according to the following rules: a node  $i \in A_i$  comes before  $j \in A_i$  if  $i$  is connected to a node in  $R$  that comes before all the nodes in  $R$  to which  $j$  is connected. As a tiebreak, precedence is given to the node with highest degree.
- Append the nodes in  $A_i$  to  $R$ , in the order determined in the previous item.

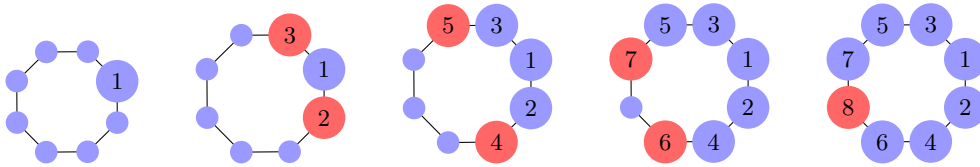


Figure 4.1: Illustration of the Cuthill–McKee algorithm. The new numbering of the nodes is illustrated. The first node was chosen randomly since all the nodes have the same degree. In this example, the ordered tuple  $R$  evolves as follows:  $(1) \rightarrow (1, 2, 8) \rightarrow (1, 2, 8, 3, 7) \rightarrow (1, 2, 8, 3, 7, 4, 6) \rightarrow (1, 2, 8, 3, 7, 4, 6, 5)$ .

The steps of the algorithm for the example above are depicted in Figure 4.1. Another example, taken from the original paper by Cuthill and McKee [1], is presented in Figure 4.2.

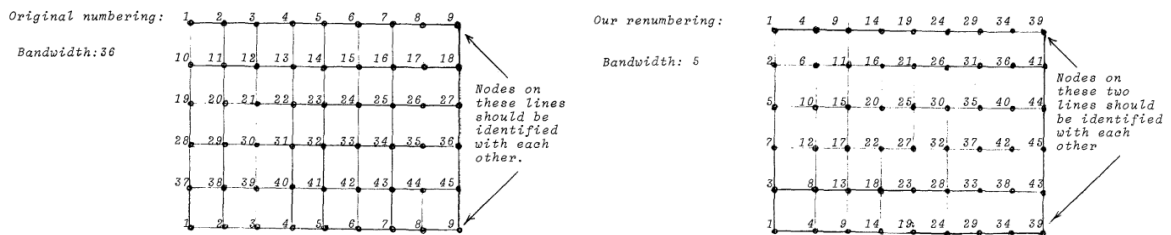


Figure 4.2: Example from the original Cuthill–McKee paper [1].

### 4.3 Iterative methods for linear systems

Iterative methods enjoy more flexibility than direct methods, because they can be stopped at any point if the residual is deemed sufficiently small. This generally enables to obtain a good solution at a computational cost that is significantly lower than that of direct methods. In this section, we present and study two classes of iterative methods: basic iterative methods based on a splitting of the matrix  $A$ , and the so-called Krylov subspace methods.



### 4.3.1 Basic iterative methods

The basic iterative methods are particular cases of a general *splitting method*. Given a splitting of the matrix of the linear system as  $A = M - N$ , for a nonsingular matrix  $M \in \mathbf{C}^{n \times n}$  and a matrix  $N \in \mathbf{C}^{n \times n}$ , together with an initial guess  $\mathbf{x}^{(0)}$  of the solution, one step of this general method reads

$$M\mathbf{x}^{(k+1)} = N\mathbf{x}^{(k)} + \mathbf{b}. \quad (4.14)$$

For any choice of splitting, the exact solution  $\mathbf{x}_*$  to the linear system is a fixed point of this iteration, in the sense that if  $\mathbf{x}^{(0)} = \mathbf{x}_*$ , then  $\mathbf{x}^{(k)} = \mathbf{x}_*$  for all  $k \geq 0$ . Equation (4.14) is a linear system with matrix  $M$ , unknown  $\mathbf{x}^{(k+1)}$ , and right-hand side  $N\mathbf{x}^{(k)} + \mathbf{b}$ . There is a compromise between the cost of a single step and the speed of convergence of the method. In the extreme case where  $M = A$  and  $N = 0$ , the method converges to the exact solution in one step, but performing this step amounts to solving the initial problem. In practice, in order for the method to be useful, the linear system (4.14) should be relatively simple to solve. Concretely, this means that the matrix  $M$  should be diagonal, triangular, block diagonal, or block triangular. The error  $\mathbf{e}^{(k)}$  and residual  $\mathbf{r}^{(k)}$  at iteration  $k$  are defined as follows:

$$\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}_*, \quad \mathbf{r}^{(k)} = A\mathbf{x}^{(k)} - \mathbf{b}.$$

#### Convergence of the splitting method

Before presenting concrete examples of splitting methods, we obtain a necessary and sufficient condition for the convergence of (4.14) for any initial guess  $\mathbf{x}^{(0)}$ .

**Proposition 4.10** (Convergence). *The splitting method (4.14) converges for any initial  $\mathbf{x}^{(0)}$  if and only if  $\rho(M^{-1}N) < 1$ . In addition, for any  $\varepsilon > 0$  there exists  $K > 0$  such that*

$$\forall k \geq K, \quad \|\mathbf{e}^{(k)}\| \leq (\rho(M^{-1}N) + \varepsilon)^k \|\mathbf{e}^{(0)}\|. \quad (4.15)$$

*Proof.* Let  $\mathbf{x}_*$  denote the solution to the linear system. Since  $M\mathbf{x}_* - N\mathbf{x}_* = \mathbf{b}$ , we have

$$M(\mathbf{x}^{(k+1)} - \mathbf{x}_*) = N(\mathbf{x}^{(k)} - \mathbf{x}_*).$$

Using the assumption that  $M$  is nonsingular, we obtain that the error satisfies the equation

$$\mathbf{e}^{(k+1)} = (M^{-1}N)\mathbf{e}^{(k)}.$$

Applying this equality repeatedly, we deduce that

$$\mathbf{e}^{(k)} = (M^{-1}N)\mathbf{e}^{(k-1)} = \dots = (M^{-1}N)^k \mathbf{e}^{(0)}. \quad (4.16)$$

**Proof that  $\rho(M^{-1}N) < 1$  is necessary for convergence.** We prove the equivalent claim that if  $\rho(M^{-1}N) \geq 1$ , then there exists  $\mathbf{x}^{(0)}$  such that the method is not convergent. Indeed, assume that  $\mathbf{x}^{(0)} = \mathbf{x}_* + \mathbf{v}^{(0)}$ , where  $\mathbf{v}^{(0)}$  is the eigenvector of  $M^{-1}N$  associated with the eigenvalue of largest modulus. Then  $\mathbf{e}^{(0)} = \mathbf{v}^{(0)}$  and the right-hand side of (4.16) does not converge to 0 in

the limit as  $k \rightarrow 0$ , because

$$\|(\mathbf{M}^{-1}\mathbf{N})^k \mathbf{e}^{(0)}\| = \rho(\mathbf{M}^{-1}\mathbf{N})^k \|\mathbf{v}^{(0)}\| \geq \|\mathbf{v}^{(0)}\|.$$

Thus, the condition  $\rho(\mathbf{M}^{-1}\mathbf{N}) < 1$  is necessary to ensure convergence for all initial guess  $\mathbf{x}^{(0)}$ .

**Proof that  $\rho(\mathbf{M}^{-1}\mathbf{N}) < 1$  is sufficient for convergence.** In order to show that the condition is also sufficient, note that by (4.16)

$$\forall k \geq 0, \quad \|\mathbf{e}^{(k)}\| \leq \|(\mathbf{M}^{-1}\mathbf{N})^k\| \|\mathbf{e}^{(0)}\|.$$

By Gelfand's formula, proved in Proposition A.10 of Appendix A, it holds that

$$\lim_{k \rightarrow \infty} \|(\mathbf{M}^{-1}\mathbf{N})^k\|^{\frac{1}{k}} = \rho(\mathbf{M}^{-1}\mathbf{N}). \quad (4.17)$$

Therefore, we deduce that if  $\rho(\mathbf{M}^{-1}\mathbf{N}) < 1$ , then  $\|(\mathbf{M}^{-1}\mathbf{N})^k\| \rightarrow 0$  and so  $\mathbf{e}^{(k)} \rightarrow \mathbf{0}$ . In addition, it follows from (4.17) that for all  $\varepsilon > 0$  there is  $K \in \mathbf{N}$  such that

$$\forall k \geq K, \quad \|(\mathbf{M}^{-1}\mathbf{N})^k\|^{\frac{1}{k}} \leq \rho(\mathbf{M}^{-1}\mathbf{N}) + \varepsilon.$$

Rearranging this inequality gives (4.15). □

At this point, it is natural to wonder whether there exist sufficient conditions on the matrix  $\mathbf{A}$  such that the inequality  $\rho(\mathbf{M}^{-1}\mathbf{N}) < 1$  is satisfied, which is best achieved on a case by case basis. In the next sections, we present four instances of splitting methods. For each of them, we obtain a sufficient condition for convergence. We are particularly interested in the case where the matrix  $\mathbf{A}$  is Hermitian and positive definite, which often arises in applications, and in the case where  $\mathbf{A}$  is strictly row or column diagonally dominant. We recall that a matrix  $\mathbf{A}$  is said to be row or column diagonally dominant if, respectively,

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}| \quad \forall i \quad \text{or} \quad |a_{jj}| \geq \sum_{i \neq j} |a_{ij}| \quad \forall j.$$

### Richardson's method

Arguably the simplest splitting of the matrix  $\mathbf{A}$  is given by  $\mathbf{A} = \frac{1}{\omega} \mathbf{I} - \left(\frac{1}{\omega} \mathbf{I} - \mathbf{A}\right)$ , for some parameter  $\omega \in \mathbf{R}$ , which leads to *Richardson's method*:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \omega(\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}). \quad (4.18)$$

In this case the spectral radius which enters in the asymptotic rate of convergence is given by

$$\rho(\mathbf{M}^{-1}\mathbf{N}) = \rho\left(\omega\left(\frac{1}{\omega}\mathbf{I} - \mathbf{A}\right)\right) = \rho(\mathbf{I} - \omega\mathbf{A})$$

The eigenvalues of the matrix  $I - \omega A$  are given by  $1 - \omega \lambda_i$ , where  $(\lambda_i)_{1 \leq i \leq L}$  are the eigenvalues of  $A$ . Therefore, the spectral radius is given by

$$\rho(M^{-1}N) = \max_{1 \leq i \leq L} |1 - \omega \lambda_i|.$$

If the eigenvalues of the matrix  $A$  do not either (i) all have a positive real part or (ii) all have a negative real part, then

$$\forall \omega \in \mathbf{R}, \quad \max_{1 \leq i \leq L} |1 - \omega \lambda_i| \geq 1.$$

In other words, by [Proposition 4.10](#), there is for any choice of  $\omega \in \mathbf{R}$  some  $\mathbf{x}^{(0)}$  such that Richardson's method is non-convergent. Therefore, in order for convergence to hold for all  $\mathbf{x}^{(0)}$ , it is necessary that the eigenvalues of  $A$  either all have a negative real part, or all have a positive real part. We focus in the next paragraph on the latter case and we also assume, for simplicity, that  $A$  is Hermitian.

**Case of symmetric positive definite  $A$ .** If the matrix  $A$  is Hermitian and positive definite, the eigenvalues of  $A$  are all real and positive, and it is possible to explicitly calculate the optimal value of  $\omega$  for convergence. In order for convergence to be as fast as possible, the spectral radius of  $M^{-1}N$  should be as small as possible, in view of [Proposition 4.10](#). Denoting by  $\lambda_{\min}$  and  $\lambda_{\max}$  the minimum and maximum eigenvalues of  $A$ , it is not difficult to show that

$$\rho(M^{-1}N) = \max_{1 \leq i \leq L} |1 - \omega \lambda_i| = \max\{|1 - \omega \lambda_{\min}|, |1 - \omega \lambda_{\max}|\}. \quad (4.19)$$

The right-hand side is minimized  $1 - \omega \lambda_{\min} = \omega \lambda_{\max} - 1$ , in which case the two arguments of the maximum in (4.19) are equal. From this we deduce the optimal value of  $\omega$  and the associated spectral radius:

$$\omega_{\text{opt}} = \frac{2}{\lambda_{\max} + \lambda_{\min}}, \quad \rho_{\text{opt}} = 1 - \frac{2\lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} = \frac{\kappa_2(A) - 1}{\kappa_2(A) + 1}.$$

We observe that the smaller the condition number of the matrix  $A$ , the better the asymptotic rate of convergence.

*Remark 4.2* (Link to optimization). In the case where  $A \in \mathbf{R}^{n \times n}$  is symmetric and positive definite, the Richardson update (4.18) may be viewed as a step of the steepest descent algorithm, which we study carefully in [Chapter 8](#), for the function  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{b}^T \mathbf{x}$ :

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \omega \nabla f(\mathbf{x}^{(k)}). \quad (4.20)$$

The gradient of this function is  $\nabla f(\mathbf{x}) = A\mathbf{x} - \mathbf{b}$ , and its Hessian matrix is  $A$ . Since the Hessian matrix is positive definite, the function is convex and attains its global minimum when  $\nabla f$  is zero, i.e. when  $A\mathbf{x} = \mathbf{b}$ .

### Jacobi's method

In Jacobi's method, the matrix  $M$  in the splitting is the diagonal matrix  $D$  with the same entries as those of  $A$  on the diagonal. We denote by  $L$  and  $U$  the lower and upper triangular parts of  $A$ , without the diagonal. One step of the method reads

$$D\mathbf{x}^{(k+1)} = (D - A)\mathbf{x}^{(k)} + \mathbf{b} = -(L + U)\mathbf{x}^{(k)} + \mathbf{b} \quad (4.21)$$

Since the matrix  $D$  on the left-hand side is diagonal, this linear system with unknown  $\mathbf{x}^{(k+1)}$  is very simple to solve. The equation (4.21) can be rewritten as

$$\begin{cases} a_{11}x_1^{(k+1)} + a_{12}x_2^{(k)} + \cdots + a_{1n}x_n^{(k)} = b_1 \\ a_{21}x_1^{(k)} + a_{22}x_2^{(k+1)} + \cdots + a_{2n}x_n^{(k)} = b_2 \\ \vdots \\ a_{n1}x_1^{(k)} + a_{n2}x_2^{(k)} + \cdots + a_{nn}x_n^{(k+1)} = b_n. \end{cases}$$

The updates for each of the entries of  $\mathbf{x}^{(k+1)}$  are independent, and so the Jacobi method lends itself well to parallel implementation. The computational cost of one iteration, measured in number of floating point operations required, scales as  $\mathcal{O}(n^2)$  if  $A$  is a full matrix, or  $\mathcal{O}(nk)$  if  $A$  is a sparse matrix with  $k$  nonzero elements per row on average. It is simple to prove the convergence of Jacobi's method is the case where  $A$  is diagonally dominant.

**Proposition 4.11.** *Assume that  $A$  is strictly (row or column) diagonally dominant. Then it holds that  $\rho(M^{-1}N) < 1$  for the Jacobi splitting.*

*Proof.* Assume that  $\lambda$  is an eigenvalue of  $M^{-1}N$  and  $\mathbf{v}$  is the associated unit eigenvector. Then

$$M^{-1}N\mathbf{v} = \lambda\mathbf{v} \quad \Leftrightarrow \quad N\mathbf{v} = \lambda M\mathbf{v} \quad \Leftrightarrow \quad (N - \lambda M)\mathbf{v} = 0.$$

In the case of Jacobi's splitting, this is equivalent to

$$-(L + \lambda D + U)\mathbf{v} = 0.$$

If  $|\lambda| > 1$ , then the matrix on the left-hand side of this equation is diagonally dominant and thus invertible (see [Exercise 4.9](#)). Therefore  $\mathbf{v} = 0$ , but this is a contradiction because  $\mathbf{v}$  is vector of unit norm. Consequently, all the eigenvalues are bounded from above strictly by 1 in modulus.  $\square$

### Gauss–Seidel's method

In Gauss Seidel's method, the matrix  $M$  in the splitting is the lower triangular part of  $A$ , including the diagonal. One step of the method then reads

$$(L + D)\mathbf{x}^{(k+1)} = -U\mathbf{x}^{(k)} + \mathbf{b} \quad (4.22)$$

The system is solved by forward substitution. The equation (4.22) can be rewritten equivalently as

$$\begin{cases} a_{11}x_1^{(k+1)} + a_{12}x_2^{(k)} + a_{13}x_3^{(k)} + \cdots + a_{1n}x_n^{(k)} = b_1 \\ a_{21}x_1^{(k+1)} + a_{22}x_2^{(k+1)} + a_{23}x_3^{(k)} + \cdots + a_{2n}x_n^{(k)} = b_2 \\ a_{32}x_1^{(k+1)} + a_{32}x_2^{(k+1)} + a_{33}x_3^{(k+1)} + \cdots + a_{3n}x_n^{(k)} = b_3 \\ \vdots \\ a_{n1}x_1^{(k+1)} + a_{n2}x_2^{(k+1)} + a_{n3}x_3^{(k+1)} + \cdots + a_{nn}x_n^{(k+1)} = b_n. \end{cases}$$

Given  $\mathbf{x}^{(k)}$ , the first entry of  $\mathbf{x}^{(k+1)}$  is obtained from the first equation. Then the value of the second entry is obtained from the second equation, etc. Unlike Jacobi's method, the Gauss–Seidel method is sequential and the entries of  $\mathbf{x}^{(k+1)}$  cannot be updated in parallel.

It is possible to prove the convergence of the Gauss–Seidel method in particular cases. For example, the method converges if  $\mathbf{A}$  is strictly diagonally dominant. Proving this, using an approach similar to that in the proof of Proposition 4.11, is the goal of Exercise 4.18. It is also possible to prove convergence when  $\mathbf{A}$  is Hermitian and positive definite. We show this in the next section for the relaxation method, which generalizes the Gauss–Seidel method.

### Relaxation method

The relaxation method generalizes the Gauss–Seidel method. It corresponds to the splitting

$$\mathbf{A} = \left( \frac{\mathbf{D}}{\omega} + \mathbf{L} \right) - \left( \frac{1-\omega}{\omega} \mathbf{D} - \mathbf{U} \right). \quad (4.23)$$

When  $\omega = 1$ , this is simply the Gauss–Seidel splitting. The idea of the relaxation method is that, by letting  $\omega$  be a parameter that can differ from 1, faster convergence can be achieved. This intuition will be verified later. The equation (4.14) for this splitting can be rewritten equivalently as

$$\begin{cases} a_{11}(x_1^{(k+1)} - x_1^{(k)}) = -\omega \left( a_{11}x_1^{(k)} + a_{12}x_2^{(k)} + \cdots + a_{1n}x_n^{(k)} - b_1 \right) \\ a_{22}(x_2^{(k+1)} - x_2^{(k)}) = -\omega \left( a_{21}x_1^{(k+1)} + a_{22}x_2^{(k)} + \cdots + a_{2n}x_n^{(k)} - b_2 \right) \\ \vdots \\ a_{nn}(x_n^{(k+1)} - x_n^{(k)}) = -\omega \left( a_{n1}x_1^{(k+1)} + a_{n2}x_2^{(k+1)} + \cdots + a_{nn}x_n^{(k)} - b_n \right). \end{cases}$$

The coefficient on the right-hand side is larger than in the Gauss–Seidel method if  $\omega > 1$ , and smaller if  $\omega < 1$ . These regimes are called *over-relaxation* and *under-relaxation*, respectively.

To conclude this section, we establish a sufficient condition for the convergence of the relaxation method, and also of the Gauss–Seidel method as a particular case when  $\omega = 1$ , when the matrix  $\mathbf{A}$  is Hermitian and positive definite. To this end, we begin by showing the following preparatory result, which concerns a general splitting  $\mathbf{A} = \mathbf{M} - \mathbf{N}$ .

**Proposition 4.12.** *Let  $\mathbf{A}$  be Hermitian and positive definite. If the Hermitian matrix  $\mathbf{M}^* + \mathbf{N}$*

is positive definite, then  $\rho(\mathbf{M}^{-1}\mathbf{N}) < 1$ .

*Proof.* First, notice that  $\mathbf{M}^* + \mathbf{N}$  is indeed Hermitian because

$$(\mathbf{M}^* + \mathbf{N})^* = \mathbf{M} + \mathbf{N}^* = \mathbf{A} + \mathbf{N} + \mathbf{N}^*.$$

We will show that  $\|\mathbf{M}^{-1}\mathbf{N}\|_{\mathbf{A}} < 1$ , where  $\|\bullet\|_{\mathbf{A}}$  is the matrix norm induced by the following norm on vectors:

$$\|\mathbf{x}\|_{\mathbf{A}} := \sqrt{\mathbf{x}^* \mathbf{A} \mathbf{x}}.$$

Showing that this indeed defines a vector norm is the goal of [Exercise 4.11](#). Since  $\mathbf{N} = \mathbf{M} - \mathbf{A}$ , it holds that  $\|\mathbf{M}^{-1}\mathbf{N}\|_{\mathbf{A}} = \|\mathbf{I} - \mathbf{M}^{-1}\mathbf{A}\|_{\mathbf{A}}$ , and so

$$\|\mathbf{M}^{-1}\mathbf{N}\|_{\mathbf{A}} = \sup\{\|\mathbf{x} - \mathbf{M}^{-1}\mathbf{A}\mathbf{x}\|_{\mathbf{A}} : \|\mathbf{x}\|_{\mathbf{A}} \leq 1\}.$$

Take  $\mathbf{x} \in \mathbf{C}^n$  with  $\|\mathbf{x}\|_{\mathbf{A}} \leq 1$  and let  $\mathbf{y} = \mathbf{M}^{-1}\mathbf{A}\mathbf{x}$ . We calculate

$$\begin{aligned} \|\mathbf{x} - \mathbf{M}^{-1}\mathbf{A}\mathbf{x}\|_{\mathbf{A}}^2 &= \mathbf{x}^* \mathbf{A} \mathbf{x} - \mathbf{y}^* \mathbf{A} \mathbf{x} - \mathbf{x}^* \mathbf{A} \mathbf{y} + \mathbf{y}^* \mathbf{A} \mathbf{y} \\ &= \mathbf{x}^* \mathbf{A} \mathbf{x} - \mathbf{y}^* \mathbf{M} \mathbf{M}^{-1} \mathbf{A} \mathbf{x} - (\mathbf{M}^{-1} \mathbf{A} \mathbf{x})^* \mathbf{M}^* \mathbf{y} + \mathbf{y}^* \mathbf{A} \mathbf{y} \\ &= \mathbf{x}^* \mathbf{A} \mathbf{x} - \mathbf{y}^* \mathbf{M} \mathbf{y} - \mathbf{y}^* \mathbf{M}^* \mathbf{y} + \mathbf{y}^* (\mathbf{M} - \mathbf{N}) \mathbf{y} \\ &= \mathbf{x}^* \mathbf{A} \mathbf{x} - \mathbf{y}^* (\mathbf{M}^* + \mathbf{N}) \mathbf{y} \leq 1 - \mathbf{y}^* (\mathbf{M}^* + \mathbf{N}) \mathbf{y} < 1, \end{aligned}$$

where we used in the last inequality the assumption that  $\mathbf{M}^* + \mathbf{N}$  is positive definite. This inequality holds true for all  $\mathbf{x} \in \mathbf{C}^n$  with  $\|\mathbf{x}\|_{\mathbf{A}} = 1$ , and so we conclude that  $\|\mathbf{M}^{-1}\mathbf{N}\|_{\mathbf{A}} < 1$ , which implies that  $\rho(\mathbf{M}^{-1}\mathbf{N}) < 1$ .  $\square$

As a corollary, we obtain a sufficient condition for the convergence of the relaxation method.

**Corollary 4.13.** *Assume that  $\mathbf{A}$  is Hermitian and positive definite. Then the relaxation method converges if  $\omega \in (0, 2)$ .*

*Proof.* For the relaxation method, we have

$$\mathbf{M} + \mathbf{N}^* = \left( \frac{\mathbf{D}}{\omega} + \mathbf{L} \right) + \left( \frac{1-\omega}{\omega} \mathbf{D} - \mathbf{U} \right)^*.$$

Since  $\mathbf{A}$  is Hermitian, it holds that  $\mathbf{D}^* = \mathbf{D}$  and  $\mathbf{U}^* = \mathbf{L}$ . Therefore,

$$\mathbf{M} + \mathbf{N}^* = \frac{2-\omega}{\omega} \mathbf{D}.$$

The diagonal elements of  $\mathbf{D}$  are all positive, because  $\mathbf{A}$  is positive definite. (Indeed, if there was an index  $i$  such that  $d_{ii} \leq 0$ , then it would hold that  $\mathbf{e}_i^T \mathbf{A} \mathbf{e}_i = d_{ii} \leq 0$ , contradicting the assumption that  $\mathbf{A}$  is positive definite.) We deduce that  $\mathbf{M} + \mathbf{N}^*$  is positive definite if and only if  $\omega \in (0, 2)$ . We can then conclude the proof by using [Proposition 4.12](#).  $\square$

Note that [Corollary 4.13](#) implies as a particular case the convergence of the Gauss–Seidel method when  $\mathbf{A}$  is Hermitian and positive definite. The condition  $\omega \in (0, 2)$  is in fact necessary for the

convergence of the relaxation method, not only in the case of a Hermitian positive definite matrix  $A$  but in general.

**Proposition 4.14** (Necessary condition for the convergence of the relaxation method). *Let  $A \in \mathbf{C}^{n \times n}$  be an invertible matrix, and let  $A = M_\omega - N_\omega$  denote the splitting of the relaxation method with parameter  $\omega$ . It holds that*

$$\forall \omega \neq 0, \quad \rho(M_\omega^{-1}N_\omega) \geq |\omega - 1|.$$

*Proof.* We recall the following facts:

- the determinant of a product of matrices is equal to the product of the determinants.
- the determinant of a triangular matrix is equal to the product of its diagonal entries;
- the determinant of a matrix is equal to the product of its eigenvalues, to the power of their algebraic multiplicity. This can be shown from the previous two items, by passing to the Jordan normal form.

Therefore, we have that

$$\det(M_\omega^{-1}N_\omega) = \det(M_\omega)^{-1} \det(N_\omega) = \frac{\det\left(\frac{1-\omega}{\omega}D - U\right)}{\det\left(\frac{D}{\omega} + L\right)} = (1 - \omega)^n.$$

Since the determinant on the left-hand side is the product of the eigenvalues of  $M_\omega^{-1}N_\omega$ , it is bounded from above in modulus by  $\rho(M_\omega^{-1}N_\omega)^n$ , and so we deduce  $\rho(M_\omega^{-1}N_\omega)^n \geq |1 - \omega|^n$ . The statement then follows by taking the  $n$ -th root.  $\square$

### Comparison between Jacobi and Gauss–Seidel for tridiagonal matrices

For tridiagonal matrices, the convergence rate of the Jacobi and Gauss–Seidel methods satisfy an explicit relation, which we prove in this section. We denote the Jacobi and Gauss–Seidel splittings by  $M_J - N_J$  and  $M_G - N_G$ , respectively, and use the following notation for the entries of the matrix  $A$ :

$$\begin{pmatrix} a_1 & b_1 & & \\ c_1 & \ddots & \ddots & \\ & \ddots & \ddots & b_{n-1} \\ & & c_{n-1} & a_n \end{pmatrix}.$$

Before presenting and proving the result, notice that for any  $\mu \neq 0$  it holds that

$$\begin{pmatrix} \mu & & & \\ & \mu^2 & & \\ & & \ddots & \\ & & & \mu^n \end{pmatrix} A \begin{pmatrix} \mu^{-1} & & & \\ & \mu^{-2} & & \\ & & \ddots & \\ & & & \mu^{-n} \end{pmatrix} = \begin{pmatrix} a_1 & \mu^{-1}b_1 & & \\ \mu c_1 & \ddots & \ddots & \\ & \ddots & \ddots & \mu^{-1}b_{n-1} \\ & & \mu c_{n-1} & a_n \end{pmatrix}. \quad (4.24)$$

**Proposition 4.15.** *Assume that  $\mathbf{A}$  is tridiagonal with nonzero diagonal elements, so that both  $\mathbf{M}_{\mathcal{G}} = \mathbf{D}$  and  $\mathbf{M}_{\mathcal{G}} = \mathbf{L} + \mathbf{D}$  are invertible. Then*

$$\rho(\mathbf{M}_{\mathcal{G}}^{-1}\mathbf{N}_{\mathcal{G}}) = \rho(\mathbf{M}_{\mathcal{J}}^{-1}\mathbf{N}_{\mathcal{J}})^2$$

*Proof.* If  $\lambda$  is an eigenvalue of  $\mathbf{M}_{\mathcal{G}}^{-1}\mathbf{N}_{\mathcal{G}}$  with associated unit eigenvector  $\mathbf{v}$ , then

$$\mathbf{M}_{\mathcal{G}}^{-1}\mathbf{N}_{\mathcal{G}}\mathbf{v} = \lambda\mathbf{v} \quad \Leftrightarrow \quad \mathbf{N}_{\mathcal{G}}\mathbf{v} = \lambda\mathbf{M}_{\mathcal{G}}\mathbf{v} \quad \Leftrightarrow \quad (\mathbf{N}_{\mathcal{G}} - \lambda\mathbf{M}_{\mathcal{G}})\mathbf{v} = 0.$$

For fixed  $\lambda$ , there exists a nontrivial solution  $\mathbf{v}$  to the last equation if and only if

$$p_{\mathcal{G}}(\lambda) := \det(\mathbf{N}_{\mathcal{G}} - \lambda\mathbf{M}_{\mathcal{G}}) = \det(-\lambda\mathbf{L} - \lambda\mathbf{D} - \mathbf{U}) = 0.$$

Likewise,  $\lambda$  is an eigenvalue of  $\mathbf{M}_{\mathcal{J}}^{-1}\mathbf{N}_{\mathcal{J}}$  if and only if

$$p_{\mathcal{J}}(\lambda) := \det(\mathbf{N}_{\mathcal{J}} - \lambda\mathbf{M}_{\mathcal{J}}) = \det(-\mathbf{L} - \lambda\mathbf{D} - \mathbf{U}) = 0.$$

Now notice that, for  $\lambda \neq 0$ ,

$$p_{\mathcal{G}}(\lambda^2) = \det(-\lambda^2\mathbf{L} - \lambda^2\mathbf{D} - \mathbf{U}) = \lambda^n \det(-\lambda\mathbf{L} - \lambda\mathbf{D} - \lambda^{-1}\mathbf{U}).$$

Applying (4.24) with  $\mu = \lambda \neq 0$ , we deduce

$$p_{\mathcal{G}}(\lambda^2) = \lambda^n \det(-\mathbf{L} - \lambda\mathbf{D} - \mathbf{U}) = \lambda^n p_{\mathcal{J}}(\lambda)$$

It is clear that this relation is true also if  $\lambda = 0$ . Consequently, it holds that if  $\lambda$  is an eigenvalue of the matrix  $\mathbf{M}_{\mathcal{J}}^{-1}\mathbf{N}_{\mathcal{J}}$  then  $\lambda^2$  is an eigenvalue of  $\mathbf{M}_{\mathcal{G}}^{-1}\mathbf{N}_{\mathcal{G}}$ . Conversely, if  $\lambda$  is a nonzero eigenvalue of  $\mathbf{M}_{\mathcal{G}}^{-1}\mathbf{N}_{\mathcal{G}}$ , then the two square roots of  $\lambda$  are eigenvalues of  $\mathbf{M}_{\mathcal{J}}^{-1}\mathbf{N}_{\mathcal{J}}$ .  $\square$

If a matrix  $\mathbf{A}$  is tridiagonal and Toeplitz, i.e. if it is of the form

$$\begin{pmatrix} a & b & & \\ c & \ddots & \ddots & \\ & \ddots & \ddots & b \\ & & c & a \end{pmatrix},$$

then it is possible to prove that the eigenvalues of  $\mathbf{A}$  are given by

$$\lambda_k = a + 2\sqrt{bc} \cos\left(\frac{k\pi}{n+1}\right), \quad k = 1, \dots, n. \quad (4.25)$$

In this case, the spectral radius of  $\mathbf{M}_{\mathcal{J}}^{-1}\mathbf{N}_{\mathcal{J}}$  can be determined explicitly.

### Monitoring the convergence

In practice, we have access to the residual  $\mathbf{r}^{(k)} = \mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}$  at each iteration, but not to the error  $\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}_*$ , as calculating the latter would require to know the exact solution of the



problem. Nevertheless, the two are related by the equation

$$\mathbf{r}^{(k)} = \mathbf{A}\mathbf{e}^{(k)} \quad \Leftrightarrow \quad \mathbf{e}^{(k)} = \mathbf{A}^{-1}\mathbf{r}^{(k)}.$$

Therefore, it holds that  $\|\mathbf{e}^{(k)}\| \leq \|\mathbf{A}^{-1}\|\|\mathbf{r}^{(k)}\|$ . Likewise, the relative error satisfies

$$\frac{\|\mathbf{e}^{(k)}\|}{\|\mathbf{x}_*\|} = \frac{\|\mathbf{A}^{-1}\mathbf{r}^{(k)}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|},$$

and since  $\|\mathbf{b}\| = \|\mathbf{A}\mathbf{A}^{-1}\mathbf{b}\| \leq \|\mathbf{A}\|\|\mathbf{A}^{-1}\mathbf{b}\|$ , we deduce

$$\frac{\|\mathbf{e}^{(k)}\|}{\|\mathbf{x}_*\|} \leq \kappa(\mathbf{A}) \frac{\|\mathbf{r}^{(k)}\|}{\|\mathbf{b}\|}.$$

The fraction on the right-hand side is the *relative residual*. If the system is well conditioned, that is if  $\kappa(\mathbf{A})$  is close to one, then controlling the relative residual enables a good control of the relative error.

### Stopping criterion

In practice, several criteria can be employed in order to decide when to stop iterating. Given a small number  $\varepsilon$  (unrelated to the machine epsilon in [Chapter 1](#)), the following alternatives are available:

- Stop when  $\|\mathbf{r}^{(k)}\| \leq \varepsilon$ . The downside of this approach is that it is not *scaling invariant*: when used for solving the following rescaled system

$$k\mathbf{A}\mathbf{x} = k\mathbf{b}, \quad k \neq 1,$$

a splitting method with rescaled initial guess  $k\mathbf{x}^{(0)}$  will require a number of iterations that depends on  $k$ : fewer if  $k \ll 1$  and more if  $k \gg 1$ . In practice, controlling the relative residual and the relative error is often preferable.

- Stop when  $\|\mathbf{r}^{(k)}\|/\|\mathbf{r}^{(0)}\| \leq \varepsilon$ . This criterion is scaling invariant, but the number of iterations is dependent on the quality of the initial guess  $\mathbf{x}^{(0)}$ .
- Stop when  $\|\mathbf{r}^{(k)}\|/\|\mathbf{b}\| \leq \varepsilon$ . This criterion is generally the best, because it is both scaling invariant and the quality of the final iterate is independent of that of the initial guess.

### 4.3.2 The conjugate gradient method

As already mentioned in [Remark 4.2](#), when the matrix  $\mathbf{A} \in \mathbf{C}^{n \times n}$  in the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  is symmetric and positive definite, the system can be interpreted as a minimization problem for the function

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} - \mathbf{b}^T \mathbf{x}. \quad (4.26)$$

The fact that the exact solution  $\mathbf{x}_*$  to the linear system is the unique minimizer of this function appears clearly when rewriting  $f$  as follows:

$$f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}_*)^T \mathbf{A}(\mathbf{x} - \mathbf{x}_*) - \frac{1}{2}\mathbf{x}_*^T \mathbf{A}\mathbf{x}_*. \quad (4.27)$$

The second term is constant with  $\mathbf{x}$ , and the first term is strictly positive if  $\mathbf{x} - \mathbf{x}_* \neq \mathbf{0}$ , because  $\mathbf{A}$  is positive definite. We saw that Richardson's method can be interpreted as a steepest descent with fixed step size,

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \omega \nabla f(\mathbf{x}^{(k)}).$$

In this section, we will present and study other methods for solving the linear system (4.1) which can be viewed as optimization methods. Since  $\mathbf{A}$  is symmetric, it is diagonalizable and the function  $f$  can be rewritten as

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{2}(\mathbf{x} - \mathbf{x}_*)^T \mathbf{Q}\mathbf{D}\mathbf{Q}^T(\mathbf{x} - \mathbf{x}_*) - \frac{1}{2}\mathbf{x}_*^T \mathbf{A}\mathbf{x}_* \\ &= \frac{1}{2}(\mathbf{Q}^T \mathbf{e})^T \mathbf{D}(\mathbf{Q}^T \mathbf{e}) - \frac{1}{2}\mathbf{x}_*^T \mathbf{A}\mathbf{x}_*, \quad \mathbf{e} = \mathbf{x} - \mathbf{x}_*. \end{aligned}$$

Therefore, we have that

$$f(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^n \lambda_i \eta_i^2 - \frac{1}{2}\mathbf{x}_*^T \mathbf{A}\mathbf{x}_*, \quad \boldsymbol{\eta} = \mathbf{Q}^T(\mathbf{x} - \mathbf{x}_*),$$

where  $(\lambda_i)_{1 \leq i \leq n}$  are the diagonal entries of  $\mathbf{D}$ . This shows that  $f$  is a paraboloid after a change of coordinates.

### Steepest descent method

The steepest descent method is more general than Richardson's method in the sense that the step size changes from iteration to iteration and the method is not restricted to quadratic functions of the form (4.26). Each iteration is of the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \omega_k \nabla f(\mathbf{x}^{(k)}).$$

It is natural to wonder whether the step size  $\omega_k$  can be fixed in such a way that  $f(\mathbf{x}^{(k+1)})$  is as small as possible. For the case of the quadratic function (4.26), this value of  $\omega_k$  can be calculated explicitly for a general search direction  $\mathbf{d}$ , and in particular also when  $\mathbf{d} = \nabla f(\mathbf{x}^{(k)})$ . We calculate that

$$\begin{aligned} f(\mathbf{x}^{(k+1)}) &= f(\mathbf{x}^{(k)} - \omega_k \mathbf{d}) = \frac{1}{2}(\mathbf{x}^{(k)} - \omega_k \mathbf{d})^T \mathbf{A}(\mathbf{x}^{(k)} - \omega_k \mathbf{d}) - (\mathbf{x}^{(k)} - \omega_k \mathbf{d})^T \mathbf{b} \\ &= f(\mathbf{x}^{(k)}) + \frac{\omega_k^2}{2} \mathbf{d}^T \mathbf{A} \mathbf{d} - \omega_k \mathbf{d}^T \mathbf{r}^{(k)}. \end{aligned} \quad (4.28)$$

When viewed as a function of the real parameter  $\omega_k$ , the right-hand side is a convex quadratic function. It is minimized when its derivative is equal to zero, i.e. when

$$\omega_k \mathbf{d}^T \mathbf{A} \mathbf{d} - \mathbf{d}^T (\mathbf{A} \mathbf{x}_k - \mathbf{b}) = 0 \quad \Rightarrow \quad \omega_k = \frac{\mathbf{d}^T \mathbf{r}^{(k)}}{\mathbf{d}^T \mathbf{A} \mathbf{d}}. \quad (4.29)$$

The steepest descent algorithm with step size obtained from this equation is summarized in [Algorithm 3](#) below. By construction, the function value  $f(\mathbf{x}^{(k)})$  is nonincreasing with  $k$ , which is equivalent to saying that the error  $\mathbf{x} - \mathbf{x}_*$  is nonincreasing in the norm  $\mathbf{x} \mapsto \sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}}$ . In order to quantify more precisely the decrease of the error in this norm, we introduce the notation

$$E_k = \|\mathbf{x} - \mathbf{x}_*\|_{\mathbf{A}}^2 := (\mathbf{x}^{(k)} - \mathbf{x}_*)^T \mathbf{A} (\mathbf{x}^{(k)} - \mathbf{x}_*) = (\mathbf{A} \mathbf{x}^{(k)} - \mathbf{b})^T \mathbf{A}^{-1} (\mathbf{A} \mathbf{x}^{(k)} - \mathbf{b}).$$

We begin by showing the following auxiliary lemma.

**Lemma 4.16** (Kantorovich inequality). *Let  $\mathbf{A} \in \mathbf{R}^{n \times n}$  be a symmetric and positive definite matrix, and let  $\lambda_0 \leq \dots \leq \lambda_n$  denote its eigenvalues. Then for all nonzero  $\mathbf{z} \in \mathbf{R}^n$  it holds that*

$$\frac{(\mathbf{z}^T \mathbf{z})^2}{(\mathbf{z}^T \mathbf{A} \mathbf{z})(\mathbf{z}^T \mathbf{A}^{-1} \mathbf{z})} \geq \frac{4\lambda_1 \lambda_n}{(\lambda_1 + \lambda_n)^2}.$$

*Proof.* By the AM-GM (arithmetic mean-geometric mean) inequality, it holds for all  $t > 0$  that

$$\begin{aligned} \sqrt{(\mathbf{z}^T \mathbf{A} \mathbf{z})(\mathbf{z}^T \mathbf{A}^{-1} \mathbf{z})} &= \sqrt{(t \mathbf{z}^T \mathbf{A} \mathbf{z})(t^{-1} \mathbf{z}^T \mathbf{A}^{-1} \mathbf{z})} \leq \frac{1}{2} \left( t \mathbf{z}^T \mathbf{A} \mathbf{z} + \frac{1}{t} \mathbf{z}^T \mathbf{A}^{-1} \mathbf{z} \right) \\ &= \frac{1}{2} \mathbf{z}^T \left( t \mathbf{A} + \frac{1}{t} \mathbf{A}^{-1} \right) \mathbf{z}. \end{aligned}$$

The matrix on the right-hand side is also symmetric and positive definite, with eigenvalues equal to  $t\lambda_i + (t\lambda_i)^{-1}$ . Therefore, we deduce

$$\forall t \geq 0, \quad \sqrt{(\mathbf{z}^T \mathbf{A} \mathbf{z})(\mathbf{z}^T \mathbf{A}^{-1} \mathbf{z})} \leq \frac{1}{2} \left( \max_{i \in \{1, \dots, n\}} t\lambda_i + (t\lambda_i)^{-1} \right) \mathbf{z}^T \mathbf{z}. \quad (4.30)$$

The function  $x \mapsto x + x^{-1}$  is convex, and so over any closed interval  $[x_{\min}, x_{\max}]$  it attains its maximum either at  $x_{\min}$  or at  $x_{\max}$ . Consequently, it holds that

$$\left( \max_{i \in \{1, \dots, n\}} t\lambda_i + (t\lambda_i)^{-1} \right) = \max \left\{ t\lambda_1 + \frac{1}{t\lambda_1}, t\lambda_n + \frac{1}{t\lambda_n} \right\}.$$

In order to obtain the best possible bound from (4.30), we should let  $t$  be such that the maximum is minimized, which occurs when the two arguments of the maximum are equal:

$$t\lambda_1 + \frac{1}{t\lambda_1} = t\lambda_n + \frac{1}{t\lambda_n} \quad \Rightarrow \quad t = \frac{1}{\sqrt{\lambda_1 \lambda_n}}.$$

For this value of  $t$ , the maximum in (4.30) is equal to

$$\sqrt{\frac{\lambda_1}{\lambda_n}} + \sqrt{\frac{\lambda_n}{\lambda_1}}.$$

By substituting this expression in (4.30) and rearranging, we obtain the statement.  $\square$

We are now able to prove the convergence of the steepest descent method.

**Theorem 4.17** (Convergence of the steepest descent method). *It holds that*

$$E_{k+1} \leq \left( \frac{\kappa_2(\mathbf{A}) - 1}{\kappa_2(\mathbf{A}) + 1} \right)^2 E_k.$$

*Proof.* Substituting  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \omega_k \mathbf{d}$  in the expression for  $E_{k+1}$ , we obtain

$$\begin{aligned} E_{k+1} &= (\mathbf{x}^{(k)} - \omega_k \mathbf{d} - \mathbf{x}_*)^T \mathbf{A} (\mathbf{x}^{(k)} - \omega_k \mathbf{d} - \mathbf{x}_*) \\ &= E_k - 2\omega_k \mathbf{d}^T \mathbf{r}^{(k)} + \omega_k^2 \mathbf{d}^T \mathbf{A} \mathbf{d} \\ &= E_k - \frac{(\mathbf{d}^T \mathbf{d})^2}{\mathbf{d}^T \mathbf{A} \mathbf{d}} = \left( 1 - \frac{(\mathbf{d}^T \mathbf{d})^2}{(\mathbf{d}^T \mathbf{A} \mathbf{d})(\mathbf{d}^T \mathbf{A}^{-1} \mathbf{d})} \right) E_k, \end{aligned}$$

Using the Kantorovich inequality, we have

$$E_{k+1} \leq \left( 1 - \frac{4\lambda_1\lambda_n}{(\lambda_1 + \lambda_n)^2} \right) E_k \leq \left( \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n} \right)^2 E_k = \left( \frac{\kappa_2(\mathbf{A}) - 1}{\kappa_2(\mathbf{A}) + 1} \right)^2 E_k.$$

We immediately deduce the statement from this inequality.  $\square$

---

**Algorithm 3** Steepest descent method

---

- 1: Pick  $\varepsilon$  and initial  $\mathbf{x}$
  - 2:  $\mathbf{r} \leftarrow \mathbf{A}\mathbf{x} - \mathbf{b}$
  - 3: **while**  $\|\mathbf{r}\| \geq \varepsilon\|\mathbf{b}\|$  **do**
  - 4:      $\mathbf{d} \leftarrow \mathbf{r}$
  - 5:      $\omega \leftarrow \mathbf{d}^T \mathbf{r} / \mathbf{d}^T \mathbf{A} \mathbf{d}$
  - 6:      $\mathbf{x} \leftarrow \mathbf{x} - \omega \mathbf{d}$
  - 7:      $\mathbf{r} \leftarrow \mathbf{A}\mathbf{x} - \mathbf{b}$
  - 8: **end while**
- 

*Remark 4.3.* The attentive reader will have noticed that the rate of convergence rate of the steepest descent method is not apparently better than that of Richardson's method with optimal  $\omega$ ; in both cases, some upper bound on the norm of the error is multiplied by

$$\frac{\kappa_2(\mathbf{A}) - 1}{\kappa_2(\mathbf{A}) + 1}$$

at each iteration. For the steepest descent method, this rate of convergence is always guaranteed, but for Richardson's method, this rate of convergence holds only for the optimal value of  $\omega$ , which itself depends on the condition number  $\kappa_2(\mathbf{A})$  and is computationally expensive to approximate. Indeed, as we shall see in [Chapter 6](#), the simplest method for calculating the smallest eigenvalue of a matrix, which is necessary for estimating the condition number, requires to solve linear systems with matrix  $\mathbf{A}$ .

### Preconditioned steepest descent

We observe from [Theorem 4.17](#) that the convergence of the steepest descent method is faster when the condition number of the matrix  $\mathbf{A}$  is low. This naturally leads to the following question: can we reformulate the minimization of  $f(\mathbf{x})$  in [\(4.26\)](#) as another optimization problem which is of the same form but involves a matrix with a lower condition number, thereby providing scope for faster convergence? In order to answer this question, we consider a linear change of coordinates  $\mathbf{y} = \mathbf{T}^{-1}\mathbf{x}$ , where  $\mathbf{T}$  is an invertible matrix, and we define

$$\tilde{f}(\mathbf{y}) = f(\mathbf{T}\mathbf{y}) = \frac{1}{2}\mathbf{y}^T(\mathbf{T}^T\mathbf{A}\mathbf{T})\mathbf{y} - (\mathbf{T}^T\mathbf{b})^T\mathbf{y}. \quad (4.31)$$

This function is of the same form as  $f$  in [\(4.26\)](#), with the matrix  $\tilde{\mathbf{A}} := \mathbf{T}^T\mathbf{A}\mathbf{T}$  instead of  $\mathbf{A}$  and the vector  $\tilde{\mathbf{b}} := \mathbf{T}^T\mathbf{b}$  instead of  $\mathbf{b}$ . Its minimizer is  $\mathbf{y}_* = \mathbf{T}^{-1}\mathbf{x}_*$ . The steepest descent algorithm can be applied to [\(4.31\)](#) and, from an approximation  $\mathbf{y}^{(k)}$  of the minimizer  $\mathbf{y}_*$ , an approximation  $\mathbf{x}^{(k)}$  of  $\mathbf{x}_*$  is obtained by the change of variable  $\mathbf{x}^{(k)} = \mathbf{T}\mathbf{y}^{(k)}$ . This approach is called *preconditioning*. By [Theorem 4.17](#), the steepest descent method satisfies the following error estimate when applied to the function [\(4.31\)](#):

$$E_{k+1} \leq \left( \frac{\kappa_2(\mathbf{T}^T\mathbf{A}\mathbf{T}) - 1}{\kappa_2(\mathbf{T}^T\mathbf{A}\mathbf{T}) + 1} \right)^2 E_k, \quad E_k = (\mathbf{y}^{(k)} - \mathbf{y}_*)^T \tilde{\mathbf{A}} (\mathbf{y}^{(k)} - \mathbf{y}_*), \\ = (\mathbf{x}^{(k)} - \mathbf{x}_*)^T \mathbf{A} (\mathbf{x}^{(k)} - \mathbf{x}_*).$$

Consequently, the convergence is faster than that of the usual steepest descent method if  $\kappa_2(\mathbf{T}^T\mathbf{A}\mathbf{T}) < \kappa_2(\mathbf{A})$ . The optimal change of coordinates is given by  $\mathbf{T} = \mathbf{C}^{-T}$ , where  $\mathbf{C}$  is the factor of the Cholesky factorization of  $\mathbf{A}$  as  $\mathbf{C}\mathbf{C}^T$ . Indeed, in this case

$$\mathbf{T}^T\mathbf{A}\mathbf{T} = \mathbf{C}^{-1}\mathbf{C}\mathbf{C}^T\mathbf{C}^{-T} = \mathbf{I} \quad \Rightarrow \quad \kappa_2(\mathbf{T}^T\mathbf{A}\mathbf{T}) = 1,$$

and the method converges in a single iteration! However, this iteration amounts to solving the linear system by direct Cholesky factorization of  $\mathbf{A}$ . In practice, it is usual to define  $\mathbf{T}$  from an approximation of the Cholesky factorization, such as the *incomplete Cholesky factorization*.

To conclude this section, we demonstrate that the change of variable from  $\mathbf{x}$  to  $\mathbf{y}$  need not be performed explicitly in practice. Indeed, one step of the steepest descent algorithm applied to function  $\tilde{f}$  reads

$$\mathbf{y}^{(k+1)} = \mathbf{y}^{(k)} - \tilde{\omega}_k (\tilde{\mathbf{A}}\mathbf{y}^{(k)} - \tilde{\mathbf{b}}), \quad \tilde{\omega}_k = \frac{(\tilde{\mathbf{A}}\mathbf{y}^{(k)} - \tilde{\mathbf{b}})^T (\tilde{\mathbf{A}}\mathbf{y}^{(k)} - \tilde{\mathbf{b}})}{(\tilde{\mathbf{A}}\mathbf{y}^{(k)} - \tilde{\mathbf{b}})^T \tilde{\mathbf{A}} (\tilde{\mathbf{A}}\mathbf{y}^{(k)} - \tilde{\mathbf{b}})}.$$

Letting  $\mathbf{x}^{(k)} = \mathbf{T}\mathbf{y}^{(k)}$ , this equation can be rewritten as the following iteration:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tilde{\omega}_k \mathbf{d}_k, \quad \tilde{\omega}_k = \frac{\mathbf{d}_k^T \mathbf{r}^{(k)}}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k}, \quad \mathbf{d}_k = \mathbf{T}\mathbf{T}^T (\mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}).$$

A comparison with [\(4.29\)](#) shows that the step size  $\tilde{\omega}_k$  is such that  $f(\mathbf{x}^{(k+1)})$  is minimized. This reasoning shows that the preconditioned conjugate gradient method amounts to choosing the

direction  $\mathbf{d}_k = \mathbb{T}\mathbb{T}^T \mathbf{r}^{(k)}$  at each iteration, instead of just  $\mathbf{r}^{(k)}$ , as is apparent in Algorithm 4. It is simple to check that  $-\mathbf{d}_k$  is a descent direction for  $f$ :

$$-\nabla f(\mathbf{x})^T (\mathbb{T}\mathbb{T}^T (\mathbf{A}\mathbf{x} - \mathbf{b})) = -(\mathbb{T}^T (\mathbf{A}\mathbf{x} - \mathbf{b}))^T (\mathbb{T}^T (\mathbf{A}\mathbf{x} - \mathbf{b})) \leq 0.$$

---

**Algorithm 4** Preconditioned steepest descent method
 

---

- 1: Pick  $\varepsilon$ , invertible  $\mathbb{T}$  and initial  $\mathbf{x}$
  - 2:  $\mathbf{r} \leftarrow \mathbf{A}\mathbf{x} - \mathbf{b}$
  - 3: **while**  $\|\mathbf{r}\| \geq \varepsilon\|\mathbf{b}\|$  **do**
  - 4:      $\mathbf{d} \leftarrow \mathbb{T}\mathbb{T}^T \mathbf{r}$
  - 5:      $\omega \leftarrow \mathbf{d}^T \mathbf{r} / \mathbf{d}^T \mathbf{A}\mathbf{d}$
  - 6:      $\mathbf{x} \leftarrow \mathbf{x} - \omega \mathbf{d}$
  - 7:      $\mathbf{r} \leftarrow \mathbf{A}\mathbf{x} - \mathbf{b}$
  - 8: **end while**
- 

**Conjugate directions method**

**Definition 4.6** (Conjugate directions). Let  $\mathbf{A}$  be a symmetric positive definite matrix. Two vectors  $\mathbf{d}_1$  and  $\mathbf{d}_2$  are called  $\mathbf{A}$ -orthogonal or conjugate with respect to  $\mathbf{A}$  if  $\mathbf{d}_1^T \mathbf{A}\mathbf{d}_2 = 0$ , i.e. if they are orthogonal for the inner product  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{A}} = \mathbf{x}^T \mathbf{A}\mathbf{y}$ .

Assume that  $\mathbf{d}_0, \dots, \mathbf{d}_{n-1}$  are  $n$  pairwise  $\mathbf{A}$ -orthogonal nonzero directions. By Exercise 4.19, these vectors are linearly independent, and so they form a basis of  $\mathbf{R}^n$ . Consequently, for any initial guess  $\mathbf{x}^{(0)}$ , the vector  $\mathbf{x}^{(0)} - \mathbf{x}_*$ , where  $\mathbf{x}_*$  is the solution to the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , can be decomposed as

$$\mathbf{x}^{(0)} - \mathbf{x}_* = \alpha_0 \mathbf{d}_0 + \dots + \alpha_{n-1} \mathbf{d}_{n-1}.$$

Taking the  $\langle \bullet, \bullet \rangle_{\mathbf{A}}$  inner product of both sides with  $\mathbf{d}_k$ , with  $k \in \{0, \dots, n-1\}$ , we obtain an expression for the scalar coefficient  $\alpha_k$ :

$$\alpha_k = \frac{\mathbf{d}_k^T \mathbf{A}(\mathbf{x}^{(0)} - \mathbf{x}_*)}{\mathbf{d}_k^T \mathbf{A}\mathbf{d}_k} = \frac{\mathbf{d}_k^T (\mathbf{A}\mathbf{x}^{(0)} - \mathbf{b})}{\mathbf{d}_k^T \mathbf{A}\mathbf{d}_k}.$$

Therefore, calculating the expression of the coefficient does not require to know the exact solution  $\mathbf{x}_*$ , but only the residual  $\mathbf{r}^{(0)}$ ! Given conjugate directions, the exact solution can be obtained as

$$\mathbf{x}_* = \mathbf{x}^{(0)} - \sum_{k=0}^{n-1} \alpha_k \mathbf{d}_k, \quad \alpha_k = \frac{\mathbf{d}_k^T \mathbf{r}^{(0)}}{\mathbf{d}_k^T \mathbf{A}\mathbf{d}_k}. \quad (4.32)$$

If  $\mathbf{x}^{(0)} = \mathbf{0}$ , then  $\mathbf{r}^{(0)} = -\mathbf{b}$  and this equations gives that

$$\mathbf{x}_* = \sum_{k=0}^{n-1} \frac{\mathbf{d}_k^T \mathbf{b}}{\mathbf{d}_k^T \mathbf{A}\mathbf{d}_k} \mathbf{d}_k = \left( \sum_{k=0}^{n-1} \frac{\mathbf{d}_k \mathbf{d}_k^T}{\mathbf{d}_k^T \mathbf{A}\mathbf{d}_k} \right) \mathbf{b},$$

which implies that the inverse of  $\mathbf{A}$  is given by

$$\mathbf{A}^{-1} = \sum_{k=0}^{n-1} \frac{\mathbf{d}_k \mathbf{d}_k^T}{\nu_k}, \quad \nu_k = \mathbf{d}_k^T \mathbf{A} \mathbf{d}_k.$$

The conjugate directions method is illustrated in [Algorithm 5](#). Its implementation is very similar to that of the steepest descent method, the only difference being that the descent direction at iteration  $k$  is given by  $\mathbf{d}_k$  instead of  $\mathbf{r}^{(k)}$ . In particular, the step size at each iteration is such that  $f(\mathbf{x}^{(k+1)})$  is minimized.

---

**Algorithm 5** Conjugate directions method

---

- 1: Assuming  $\mathbf{d}_0, \dots, \mathbf{d}_{n-1}$  are given.
  - 2: Pick initial  $\mathbf{x}^{(0)}$
  - 3: **for**  $k$  in  $\{0, \dots, n-1\}$  **do**
  - 4:      $\mathbf{r}^{(k)} = \mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}$
  - 5:      $\omega_k = \mathbf{d}_k^T \mathbf{r}^{(k)} / \mathbf{d}_k^T \mathbf{A} \mathbf{d}_k$
  - 6:      $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \omega_k \mathbf{d}_k$
  - 7: **end for**
- 

Let us now establish the connection between the [Algorithm 5](#) and (4.32), which may not be immediately apparent because (4.32) involves only the initial residual  $\mathbf{A}\mathbf{x}^{(0)} - \mathbf{b}$ , while the residual at the current iteration  $\mathbf{r}^{(k)}$  is used in the algorithm.

**Proposition 4.18** (Convergence of the conjugate directions method). *The vector  $\mathbf{x}^{(k)}$  obtained after  $k$  iterations of the conjugate directions method is given by*

$$\mathbf{x}^{(k)} = \mathbf{x}^{(0)} - \sum_{i=0}^{k-1} \alpha_i \mathbf{d}_i, \quad \alpha_i = \frac{\mathbf{d}_i^T \mathbf{r}^{(0)}}{\mathbf{d}_i^T \mathbf{A} \mathbf{d}_i}. \quad (4.33)$$

*In particular, the method converges in at most  $n$  iterations.*

*Proof.* Let us denote by  $\mathbf{y}^{(k)}$  the solution obtained after  $k$  steps of [Algorithm 5](#). Our goal is to show that  $\mathbf{y}^{(k)}$  coincides with  $\mathbf{x}^{(k)}$  defined in (4.33). The result is trivial for  $k = 0$ . Reasoning by induction, we assume that it is true up to  $k$ . Then performing step  $k + 1$  of the algorithm gives

$$\mathbf{y}^{(k+1)} = \mathbf{y}^{(k)} - \omega_k \mathbf{d}_k, \quad \omega_k = \frac{\mathbf{d}_k^T \mathbf{r}^{(k)}}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k}.$$

On the other hand, it holds from (4.33) that

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{d}_k, \quad \alpha_k = \frac{\mathbf{d}_k^T \mathbf{r}^{(0)}}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k}.$$

By the induction hypothesis, it holds that  $\mathbf{y}^{(k)} = \mathbf{x}^{(k)}$ , so in order to prove that  $\mathbf{y}^{(k+1)} = \mathbf{x}^{(k+1)}$ , it is sufficient to show that  $\omega_k = \alpha_k$ , i.e. that

$$\mathbf{d}_k^T \mathbf{r}^{(k)} = \mathbf{d}_k^T \mathbf{r}^{(0)} \quad \Leftrightarrow \quad \mathbf{d}_k^T (\mathbf{r}^{(k)} - \mathbf{r}^{(0)}) = 0 \quad \Leftrightarrow \quad \mathbf{d}_k^T \mathbf{A} (\mathbf{x}^{(k)} - \mathbf{x}^{(0)}) = 0.$$

The latter equality is obvious from the  $\mathbf{A}$ -orthonormality of the directions.  $\square$

Since  $\omega_k$  in Algorithm 5 coincides with the expression in (4.29), the conjugate directions algorithm satisfies the following “local optimization” property: the iterate  $\mathbf{x}^{(k+1)}$  minimizes  $f$  on the straight line  $\omega \mapsto \mathbf{x}^{(k)} - \omega \mathbf{d}_k$ . In contrast with the steepest descent method, however, the conjugate directions method also satisfies the following stronger property.

**Proposition 4.19** (Optimality of the conjugate directions method). *The iterate  $\mathbf{x}^{(k)}$  is the minimizer of  $f$  over the set  $\mathbf{x}^{(0)} + \mathcal{B}_k$ , where  $\mathcal{B}_k = \text{Span}\{\mathbf{d}_0, \dots, \mathbf{d}_{k-1}\}$ .*

*Proof.* By (4.32), it holds that

$$\mathbf{x}_* = \mathbf{x}^{(0)} - \sum_{i=0}^{n-1} \alpha_i \mathbf{d}_i, \quad \alpha_i = \frac{\mathbf{d}_i^T \mathbf{r}^{(0)}}{\mathbf{d}_i^T \mathbf{A} \mathbf{d}_i}$$

On the other hand, any vector  $\mathbf{y} \in \mathbf{x}^{(0)} + \mathcal{B}_k$  can be expanded as

$$\mathbf{y} = \mathbf{x}^{(0)} - \beta_0 \mathbf{d}_0 - \dots - \beta_{k-1} \mathbf{d}_{k-1}.$$

Employing these two expressions, the formula for  $f$  in (4.27), and the  $\mathbf{A}$ -orthogonality of the directions, we obtain

$$\begin{aligned} f(\mathbf{y}) &= \frac{1}{2}(\mathbf{y} - \mathbf{x}_*)^T \mathbf{A}(\mathbf{y} - \mathbf{x}_*) - \frac{1}{2} \mathbf{x}_*^T \mathbf{A} \mathbf{x}_* \\ &= \frac{1}{2} \sum_{i=0}^{k-1} (\beta_i - \alpha_i)^2 \mathbf{d}_i^T \mathbf{A} \mathbf{d}_i + \frac{1}{2} \sum_{i=k}^{n-1} \alpha_i^2 \mathbf{d}_i^T \mathbf{A} \mathbf{d}_i - \frac{1}{2} \mathbf{x}_*^T \mathbf{A} \mathbf{x}_* \end{aligned}$$

This is minimized when  $\beta_i = \alpha_i$  for all  $i \in \{0, \dots, k-1\}$ , in which case  $\mathbf{y}$  coincides with the  $k$ -th iterate  $\mathbf{x}^{(k)}$  of the conjugate directions method in view of Proposition 4.18.  $\square$

*Remark 4.4.* Let  $\|\bullet\|_{\mathbf{A}}$  denote the norm induced by the inner product  $\langle \bullet, \bullet \rangle_{\mathbf{A}}$ . Since

$$\|\mathbf{x}^{(k)} - \mathbf{x}_*\|_{\mathbf{A}} = \sqrt{2f(\mathbf{x}^{(k)}) + \mathbf{x}_*^T \mathbf{A} \mathbf{x}_*},$$

Proposition 4.19 shows that  $\mathbf{x}^{(k)}$  minimizes the norm  $\|\mathbf{x}^{(k)} - \mathbf{x}_*\|_{\mathbf{A}}$  over  $\mathbf{x}^{(0)} + \mathcal{B}_k$ . This is not surprising since, by construction, the vector  $\mathbf{x}^{(k)} - \mathbf{x}^{(0)}$  is the orthogonal projection of  $\mathbf{x}_* - \mathbf{x}^{(0)}$  onto  $\mathcal{B}_k$ , for the inner product  $\langle \bullet, \bullet \rangle_{\mathbf{A}}$ .

A corollary of (4.19) is that the gradient of  $f$  at  $\mathbf{x}^{(k)}$ , i.e. the residual  $\mathbf{r}^{(k)} = \mathbf{A} \mathbf{x}^{(k)} - \mathbf{b}$ , is orthogonal to any vector in  $\{\mathbf{d}_0, \dots, \mathbf{d}_{k-1}\}$  for the usual Euclidean inner product. This can also be checked directly from the formula

$$\mathbf{x}^{(k)} - \mathbf{x}_* = \sum_{i=k}^{n-1} \alpha_i \mathbf{d}_i, \quad \alpha_i = \frac{\mathbf{d}_i^T \mathbf{r}^{(0)}}{\mathbf{d}_i^T \mathbf{A} \mathbf{d}_i},$$



which follows directly from [Proposition 4.18](#). Indeed, it holds that

$$\forall j \in \{0, \dots, k-1\}, \quad \mathbf{d}_j^T \mathbf{r}^{(k)} = \mathbf{d}_j^T \mathbf{A}(\mathbf{x}^{(k)} - \mathbf{x}_*) = \sum_{i=k}^{n-1} \alpha_i \mathbf{d}_j^T \mathbf{A} \mathbf{d}_i = 0. \quad (4.34)$$

### The conjugate gradient method

In the previous section, we showed that, given  $n$  conjugate directions, the solution to the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be obtained in a finite number of iterations using [Algorithm 5](#). The conjugate gradient method can be viewed as a particular case of the conjugate directions method. Instead of assuming that the conjugate directions are given, they are constructed iteratively as part of the algorithm. Given an initial guess  $\mathbf{x}^{(0)}$ , the first direction is the residual  $\mathbf{r}^{(0)}$ , which coincides with the gradient of  $f$  at  $\mathbf{x}^{(0)}$ . The directions employed for the next iterations are obtained by applying the Gram-Schmidt process to the residuals. More precisely, given conjugate directions  $\mathbf{d}_0, \dots, \mathbf{d}_{k-1}$ , and letting  $\mathbf{x}^{(k)}$  denote the  $k$ -th iterate of the conjugate directions method, the direction  $\mathbf{d}_k$  is obtained by

$$\mathbf{d}_k = \mathbf{r}^{(k)} - \sum_{i=0}^{k-1} \frac{\mathbf{d}_i^T \mathbf{A} \mathbf{r}^{(k)}}{\mathbf{d}_i^T \mathbf{A} \mathbf{d}_i} \mathbf{d}_i, \quad \mathbf{r}^{(k)} = \mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}. \quad (4.35)$$

It is simple to check that  $\mathbf{d}_k$  is indeed  $\mathbf{A}$ -orthogonal to  $\mathbf{d}_i$  for  $i \in \{0, \dots, k-1\}$ , and that  $\mathbf{d}_k$  is nonzero if  $\mathbf{r}^{(k)}$  is nonzero. To prove the latter claim, we can take the Euclidean inner product of both sides with  $\mathbf{r}^{(k)}$  and use [Proposition 4.19](#) to deduce that

$$\mathbf{d}_k^T \mathbf{r}^{(k)} = (\mathbf{r}^{(k)})^T \mathbf{r}^{(k)} > 0. \quad (4.36)$$

Note also that since the directions are obtained by applying the Gram-Schmidt process to the residuals, it holds that

$$\forall k \in \{0, \dots, n-1\}, \quad \mathcal{B}_{k+1} := \text{Span}\{\mathbf{d}_0, \dots, \mathbf{d}_k\} = \text{Span}\{\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(k)}\}. \quad (4.37)$$

The following result characterizes precisely the subspace  $\mathcal{B}_{k+1}$ .

**Proposition 4.20.** *Assume that  $\|\mathbf{r}^{(k)}\| \neq 0$  for all  $k < m \leq n$ . Then it holds that*

$$\forall k \in \{0, \dots, m\}, \quad \text{Span}\{\mathbf{r}^{(0)}, \mathbf{r}^{(1)}, \dots, \mathbf{r}^{(k)}\} = \text{Span}\{\mathbf{r}^{(0)}, \mathbf{A}\mathbf{r}^{(0)}, \dots, \mathbf{A}^k \mathbf{r}^{(0)}\} \quad (4.38)$$

*The subspace on the right-hand side is called a Krylov subspace.*

*Proof.* The result is clear for  $k = 0$ . Reasoning by induction, we prove that if the result is true up to  $k < m$ , then it is also true for  $k + 1$ . A simple calculation gives that

$$\begin{aligned} \mathbf{r}^{(k+1)} &= \mathbf{A}(\mathbf{x}^{(k)} - \omega_k \mathbf{d}_k) - \mathbf{b} \\ &= \mathbf{r}^{(k)} - \omega_k \mathbf{A} \mathbf{d}_k. \end{aligned} \quad (4.39)$$

From (4.35), we deduce that

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \omega_k \mathbf{A} \left( \mathbf{r}^{(k)} - \sum_{i=0}^{k-1} \frac{\mathbf{d}_i^T \mathbf{A} \mathbf{r}^{(k)}}{\mathbf{d}_i^T \mathbf{A} \mathbf{d}_i} \mathbf{d}_i \right).$$

By (4.37) and the induction hypothesis, the bracketed expression on the right-hand side belongs to  $\mathcal{B}_{k+1}$ , so the inclusion  $\subset$  in (4.38) is clear. The inclusion  $\supset$  then follows from the fact the dimension of the subspace

$$\text{Span} \{\mathbf{d}_0, \dots, \mathbf{d}_k\} = \text{Span} \{\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(k)}\}$$

is equal to  $k + 1$ . □

It appears from (4.35) that the cost of calculating a new direction grows linearly with the iteration index. In fact, it turns out that only the last term in the sum is nonzero, and so the cost of calculating a new direction does not grow with the iteration index  $k$ . Indeed, notice that if  $i \leq k - 2$ , then

$$\mathbf{d}_i^T \mathbf{A} \mathbf{r}^{(k)} = (\mathbf{A} \mathbf{d}_i)^T \mathbf{A} \mathbf{r}^{(k)} = 0,$$

because  $\mathbf{A} \mathbf{d}_i \in \mathcal{B}_{i+2} \subset \mathcal{B}_k$  by Proposition 4.20, and  $\mathbf{r}^{(k)}$  orthogonal to  $\mathcal{B}_k$  for the Euclidean inner product by (4.34). This observation leads to Algorithm 6.

---

**Algorithm 6** Conjugate gradient method
 

---

- 1: Pick initial  $\mathbf{x}^{(0)}$
  - 2:  $\mathbf{d}_0 = \mathbf{r}^{(0)} = \mathbf{A} \mathbf{x}^{(0)} - \mathbf{b}$
  - 3: **for**  $k$  in  $\{0, \dots, n - 1\}$  **do**
  - 4:     **if**  $\|\mathbf{r}^{(k)}\| = 0$  **then**
  - 5:         Stop
  - 6:     **end if**
  - 7:      $\omega_k = \mathbf{d}_k^T \mathbf{r}^{(k)} / \mathbf{d}_k^T \mathbf{A} \mathbf{d}_k$
  - 8:      $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \omega_k \mathbf{d}_k$
  - 9:      $\mathbf{r}^{(k+1)} = \mathbf{A} \mathbf{x}^{(k+1)} - \mathbf{b}$
  - 10:      $\beta_k = \mathbf{d}_k^T \mathbf{A} \mathbf{r}^{(k+1)} / \mathbf{d}_k^T \mathbf{A} \mathbf{d}_k$ .
  - 11:      $\mathbf{d}_{k+1} = \mathbf{r}^{(k+1)} - \beta_k \mathbf{d}_k$ .
  - 12: **end for**
- 

Although the conjugate gradient method converges in a finite number of iterations, performing  $n$  iterations for very large systems would require an excessive computational cost, and so it is sometimes desirable to stop iterating when the residual is sufficiently small. To conclude this section, we study the convergence of the method.

**Theorem 4.21** (Convergence of the conjugate gradient method). *The error for the conjugate gradient method, measured as*

$$E_k := (\mathbf{x}^{(k)} - \mathbf{x}_*)^T \mathbf{A} (\mathbf{x}^{(k)} - \mathbf{x}_*),$$

satisfies the following inequality:

$$\forall q_k \in \mathbf{P}(k), \quad E_{k+1} \leq \max_{1 \leq i \leq n} (1 + \lambda_i q_k(\lambda_i))^2 E_0. \quad (4.40)$$

Here  $\mathbf{P}(k)$  is the vector space of polynomials of degree less than or equal to  $k$ .

*Proof.* In view of Proposition 4.20, the iterate  $\mathbf{x}^{(k+1)}$  can be written as

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(0)} + \sum_{i=0}^k \alpha_i \mathbf{A}^i \mathbf{r}^{(0)} = \mathbf{x}^{(0)} + p_k(\mathbf{A}) \mathbf{r}^{(0)},$$

where  $p_k$  is a polynomial of degree  $k$ . By Proposition 4.19,  $p_k$  is in fact the polynomial of degree  $k$  such that  $f(\mathbf{x}^{(k+1)})$  is minimized, and thus also  $E_{k+1}$  by (4.27). Noting that

$$\begin{aligned} \mathbf{x}^{(k+1)} - \mathbf{x}_* &= \mathbf{x}^{(0)} - \mathbf{x}_* + p_k(\mathbf{A}) \mathbf{r}^{(0)} = \mathbf{x}^{(0)} - \mathbf{x}_* + p_k(\mathbf{A}) \mathbf{A}(\mathbf{x}^{(0)} - \mathbf{x}_*) \\ &= (\mathbf{I} + \mathbf{A} p_k(\mathbf{A}))(\mathbf{x}^{(0)} - \mathbf{x}_*), \end{aligned}$$

we deduce that

$$\forall q_k \in \mathbf{P}(k), \quad E_{k+1} \leq (\mathbf{x}^{(0)} - \mathbf{x}_*)^T \mathbf{A} (\mathbf{I} + \mathbf{A} q_k(\mathbf{A}))^2 (\mathbf{x}^{(0)} - \mathbf{x}_*).$$

In order to exploit this inequality, it is useful to diagonalize  $\mathbf{A}$  as  $\mathbf{A} = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$ , for an orthogonal matrix  $\mathbf{Q}$  and a diagonal matrix  $\mathbf{D}$ . Since  $q_k(\mathbf{A}) = \mathbf{Q} q_k(\mathbf{D}) \mathbf{Q}^T$  for all  $q_k \in \mathbf{P}(k)$ , it holds that

$$\begin{aligned} \forall q_k \in \mathbf{P}(k), \quad E_{k+1} &= (\mathbf{Q}^T (\mathbf{x}^{(0)} - \mathbf{x}_*))^T \mathbf{D} (\mathbf{I} + \mathbf{D} q_k(\mathbf{D}))^2 (\mathbf{Q}^T (\mathbf{x}^{(0)} - \mathbf{x}_*)) \\ &\leq \max_{1 \leq i \leq n} (1 + \lambda_i q_k(\lambda_i))^2 \underbrace{(\mathbf{Q}^T (\mathbf{x}^{(0)} - \mathbf{x}_*))^T \mathbf{D} (\mathbf{Q}^T (\mathbf{x}^{(0)} - \mathbf{x}_*))}_{E_0}, \end{aligned}$$

which completes the proof.  $\square$

A corollary of Theorem 4.21 is that, if  $\mathbf{A}$  has  $m \leq n$  distinct eigenvalues, then the conjugate gradient method converges in at most  $m$  iterations. Indeed, in this case we can take

$$q_{m-1}(\lambda) = \frac{1}{\lambda} \left( \frac{(\lambda_1 - \lambda) \dots (\lambda_m - \lambda)}{\lambda_1 \dots \lambda_m} - 1 \right).$$

It is simple to check that the right-hand side is indeed a polynomial, and that  $1 + \lambda_i q_{m-1}(\lambda_i) = 0$  for all eigenvalues of  $\mathbf{A}$ .

In general, finding the polynomial that minimizes the right-hand side of (4.40) is not possible, because the eigenvalues of  $\mathbf{A}$  are unknown. However, it is possible to derive from this equation an error estimate with an explicit dependence on the condition number  $\kappa = \kappa_2(\mathbf{A})$ .

**Theorem 4.22.** *It holds that*

$$\forall k \geq 0, \quad E_k \leq 4 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2k} E_0,$$

*Proof.* Theorem 4.21 implies that

$$\forall q_k \in \mathbf{P}(k), \quad E_{k+1} \leq \max_{\lambda \in [\lambda_1, \lambda_n]} (1 + \lambda q_k(\lambda))^2 E_0,$$

where  $\lambda_1$  and  $\lambda_n$  are the minimum and maximum eigenvalues of  $\mathbf{A}$ . Notice that

$$\left\{ 1 + \lambda q_k : q_k \in \mathbf{P}(k) \right\} = \left\{ p_k : p_k \in \mathbf{P}(k+1) \text{ and } p_k(0) = 1 \right\}$$

Therefore, it follows from Exercise C.7 that the right-hand side is minimized when

$$1 + \lambda q_k(\lambda) = \frac{T_{k+1} \left( \frac{\lambda_n + \lambda_1 - 2\lambda}{\lambda_n - \lambda_1} \right)}{T_{k+1} \left( \frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} \right)}, \quad (4.41)$$

where  $T_{k+1}$  is the Chebyshev polynomial of degree  $k+1$ , see (C.1). We recall that  $|T_{k+1}(x)| \leq 1$  for all  $x \in [-1, 1]$ . Consequently, by the expression of Chebyshev polynomials given in Exercise C.3, the following inequality holds true for all  $\lambda \in [\lambda_1, \lambda_n]$ :

$$\begin{aligned} |1 + \lambda q_k(\lambda)| &\leq \frac{1}{T_{k+1} \left( \frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} \right)} = 2 \left( \left( r + \sqrt{r^2 - 1} \right)^{k+1} + \left( r - \sqrt{r^2 - 1} \right)^{k+1} \right)^{-1}, \\ &= 2 \left( \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^{k+1} + \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{k+1} \right)^{-1}. \end{aligned}$$

where  $r = \frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1}$ . Since the first term in the bracket converges to zero as  $k \rightarrow \infty$ , it is natural to bound this expression by keeping only the second term, which after simple algebraic manipulations leads to

$$\forall \lambda \in [\lambda_1, \lambda_n], \quad |1 + \lambda q_k(\lambda)| \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{k+1}.$$

From this inequality, the statement of the theorem follows immediately.  $\square$

## 4.4 Exercises

⚙️ **Exercise 4.1.** In the simple case where  $\mathbf{A}$  is symmetric, find values of  $\mathbf{x}$ ,  $\mathbf{b}$  and  $\Delta \mathbf{b}$  for which the inequality (4.2) is in fact an equality?

⚙️ **Exercise 4.2** (Inverse of Gaussian transformation). Prove the formula (4.8).

⚙️ **Exercise 4.3.** Prove that the product of two lower triangular matrices is lower triangular.

⚙️ **Exercise 4.4.** Assume that  $\mathbf{A} \in \mathbf{R}^{n \times n}$  is positive definite, i.e. that

$$\forall \mathbf{x} \in \mathbf{R}^n \setminus \{\mathbf{0}_n\}, \quad \mathbf{x}^T \mathbf{A} \mathbf{x} > 0.$$

Show that all the principal submatrices of  $\mathbf{A}$  are nonsingular.

□ **Exercise 4.5.** Implement the backward substitution algorithm for solving  $Ux = y$ . What is the computational cost of the algorithm?

□ **Exercise 4.6.** Compare the condition number of the matrices  $L$  and  $U$  with and without partial pivoting. For testing, use a matrix with pseudo-random entries generated as follows

```
import Random
# Set the seed so that the code is deterministic
Random.seed!(0)
n = 1000 # You can change this parameter
A = randn(n, n)
```

*Solution.* See the Jupyter notebook for this chapter. △

□ **Exercise 4.7.** Write a code for calculating the Cholesky factorization of a symmetric positive definite matrix  $A$  by comparing the entries of the product  $CC^T$  with those of the matrix  $A$ . What is the associated computational cost, and how does it compare with that of the LU factorization? **Extra credit:** ... if your code is able to exploit the potential banded structure of the matrix passed as argument for better efficiency. Specifically, your code will be tested with a matrix is of the type `BandedMatrix` defined in the `BandedMatrices.jl` package, which you will need to install. The following code can be useful for testing purposes.

```
import BandedMatrices
import LinearAlgebra

function cholesky(A)
    m, n = size(A)
    m != n && error("Matrix must be square")
    # Convert to banded matrix
    B = BandedMatrices.BandedMatrix(A)
    B.u != B.l && error("Matrix must be symmetric")
    # --> Your code comes here <--
end

n, u, l = 20000, 2, 2
A = BandedMatrices.brand(n, u, l)
A = A*A'
# so that A is symmetric and positive definite (with probability 1).
C = @time cholesky(A)
LinearAlgebra.norm(C*C' - A, Inf)
```

For information, my code takes about 1 second to run with the parameters given here.

⚙ **Exercise 4.8** (Matrix square root). Let  $A \in \mathbf{R}^{n \times n}$  be a symmetric positive definite matrix. Show that  $A$  has a positive definite square root, i.e. that there exists a symmetric matrix  $B$  such that  $BB = A$ .

*Solution.* Since  $A$  is symmetric, there exist a diagonal matrix  $D$  and an orthogonal matrix  $Q$  such that  $A = QDQ^T$ . Let  $D^{1/2}$  denote the diagonal matrix obtained by applying the square root function to the entries of  $D$ , and notice that  $D^{1/2}D^{1/2} = D$ . Then it holds that

$$A = (QD^{1/2}Q^T)(QD^{1/2}Q^T).$$

The matrix  $A^{1/2} := QD^{1/2}Q^T$  is a square root of the matrix  $A$ , in the sense that  $A^{1/2}A^{1/2} = A$ , and it is positive definite because the diagonal elements of  $D^{1/2}$  are strictly positive.  $\triangle$

⚙ **Exercise 4.9.** Show that if  $A$  is row or column diagonally dominant, then  $A$  is invertible.

⚙ **Exercise 4.10.** Let  $T$  be a nonsingular matrix. Show that

$$\|A\|_T := \|T^{-1}AT\|_2$$

defines a matrix norm induced by a vector norm.

⚙ **Exercise 4.11.** Let  $A \in \mathbf{R}^{n \times n}$  be a symmetric positive definite matrix. Show that the functional

$$\|\bullet\|_A : \mathbf{x} \mapsto \sqrt{\mathbf{x}^T A \mathbf{x}}$$

defines a norm on  $\mathbf{R}^n$ .

*Solution.* We need to prove that the three axioms of a norm are satisfied:

- **(Positivity)** Since  $A$  is positive definite, it holds that  $\|\mathbf{x}\|_A > 0$  for any  $\mathbf{x} \in \mathbf{R}^n \setminus \{\mathbf{0}\}$ .
- **(Homogeneity)** It is clear that  $\|c\mathbf{x}\|_A = |c|\|\mathbf{x}\|_A$  for any  $c \in \mathbf{R}$ .
- **(Triangle inequality)** Let  $A^{1/2}$  denote the positive definite square root of  $A$ , which exists by Exercise 4.8. Then

$$\|\mathbf{x}\|_A = \|A^{1/2}\mathbf{x}\|_2.$$

The triangle inequality for  $\|\bullet\|_A$  then follows from that for  $\|\bullet\|_2$ :

$$\|\mathbf{x} + \mathbf{y}\|_A = \|A^{1/2}\mathbf{x} + A^{1/2}\mathbf{y}\|_2 \leq \|A^{1/2}\mathbf{x}\|_2 + \|A^{1/2}\mathbf{y}\|_2 = \|\mathbf{x}\|_A + \|\mathbf{y}\|_A.$$

Another option for solving this exercise is to show that

$$\langle \mathbf{x}, \mathbf{y} \rangle_A := \mathbf{x}^T A \mathbf{y}$$

defines an inner product, with induced norm given by  $\|\bullet\|_A$ .  $\triangle$

⚙ **Exercise 4.12.** Show that the residual satisfies the equation

$$\mathbf{r}^{(k+1)} = N\mathbf{M}^{-1}\mathbf{r}^{(k)} = (I - \mathbf{A}\mathbf{M}^{-1})\mathbf{r}^{(k)}.$$

⚙ **Exercise 4.13.** Show that, if  $A$  and  $B$  are two square matrices, then  $\rho(AB) = \rho(BA)$ .

⚙ **Exercise 4.14.** Is  $\rho(\bullet)$  a norm? Prove or disprove.

⚙ **Exercise 4.15.** Prove that, if  $A$  is a diagonal matrix, then

$$\|A\|_1 = \|A\|_2 = \|A\|_\infty = \rho(A).$$

⚙️ **Exercise 4.16.** Show that, for any matrix norm  $\|\bullet\|$  induced by a vector norm,

$$\rho(\mathbf{A}) \leq \|\mathbf{A}\|.$$

⚙️ **Exercise 4.17.** Let  $\|\bullet\|$  denote the Euclidean vector norm on  $\mathbf{R}^n$ . We define in [Appendix A](#) the induced matrix norm as

$$\|\mathbf{A}\| = \sup\{\|\mathbf{A}\mathbf{x}\| : \|\mathbf{x}\| \leq 1\}.$$

Show from this definition that, if  $\mathbf{A}$  is symmetric and positive definite, then

$$\|\mathbf{A}\| = \|\mathbf{A}\|_* := \sup\{\mathbf{x}^T \mathbf{A} \mathbf{x} : \|\mathbf{x}\| \leq 1\}.$$

*Solution.* By the Cauchy–Schwarz inequality and the definition of  $\|\mathbf{A}\|$ , it holds that

$$\forall \mathbf{x} \in \mathbf{R}^n \text{ with } \|\mathbf{x}\| \leq 1, \quad |\mathbf{x}^T \mathbf{A} \mathbf{x}| \leq \|\mathbf{x}\| \|\mathbf{A} \mathbf{x}\| \leq \|\mathbf{x}\| \|\mathbf{A}\| \|\mathbf{x}\| \leq \|\mathbf{A}\|.$$

This shows that  $\|\mathbf{A}\|_* \leq \|\mathbf{A}\|$ . Conversely, letting  $\mathbf{B}$  denote a matrix square root of  $\mathbf{A}$  (see [Exercise 4.8](#)), we have

$$\begin{aligned} \forall \mathbf{x} \in \mathbf{R}^n \text{ with } \|\mathbf{x}\| \leq 1, \quad \|\mathbf{A} \mathbf{x}\| &= \sqrt{\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}} = \sqrt{(\mathbf{B} \mathbf{x})^T \mathbf{B} \mathbf{B} (\mathbf{B} \mathbf{x})} = \sqrt{(\mathbf{B} \mathbf{x})^T \mathbf{A} (\mathbf{B} \mathbf{x})} \\ &= \|\mathbf{B} \mathbf{x}\| \sqrt{\mathbf{y}^T \mathbf{A} \mathbf{y}}, \quad \mathbf{y} = \frac{\mathbf{B} \mathbf{x}}{\|\mathbf{B} \mathbf{x}\|}. \end{aligned}$$

It holds that  $\|\mathbf{B} \mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}} \leq \sqrt{\|\mathbf{A}\|_*}$ . In addition  $\|\mathbf{y}\| = 1$ , so the expression inside the square root is bounded from above by  $\|\mathbf{A}\|_*$ , which enables to conclude the proof.  $\triangle$

⚙️ **Exercise 4.18.** Prove that, if the matrix  $\mathbf{A}$  is strictly diagonally dominant (by rows or columns), then the Gauss–Seidel method converges, i.e.  $\rho(\mathbf{M}^{-1} \mathbf{N}) < 1$ . You can use the same approach as in the proof of [Proposition 4.11](#).

⚙️ **Exercise 4.19.** Let  $\mathbf{A} \in \mathbf{R}^{n \times n}$  denote a symmetric positive definite matrix, and assume that the vectors  $\mathbf{d}_1, \dots, \mathbf{d}_n$  are pairwise  $\mathbf{A}$ -orthogonal directions. Show that  $\mathbf{d}_1, \dots, \mathbf{d}_n$  are linearly independent.

📦 **Exercise 4.20** (Steepest descent algorithm). Consider the linear system

$$\mathbf{A} \mathbf{x} := \begin{pmatrix} 3 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} =: \mathbf{b}. \quad (4.42)$$

- Show that  $\mathbf{A}$  is positive definite.
- Draw the contour lines of the function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}.$$

- Plot the contour lines of  $f$  in Julia using the function `contourf` from the package `Plots`.
- Using [Theorem 4.17](#), estimate the number  $K$  of iterations of the steepest descent algorithm required in order to guarantee that  $E_K \leq 10^{-8}$ , when starting from the vector  $\mathbf{x}^{(0)} = (2 \ 3)^T$ .

- Implement the steepest descent method for finding the solution to (4.42), and plot the iterates as linked dots over the filled contour of  $f$ .
- Plot the error  $E_k$  as a function of the iteration index, using a linear scale for the  $x$  axis and a logarithmic scale for the  $y$  axis.

❁ **Exercise 4.21.** Compute the number of floating point operations required for performing one iteration of the conjugate gradient method, assuming that the matrix  $\mathbf{A}$  contains  $\alpha \ll n$  nonzero elements per row.

□ **Exercise 4.22** (Solving the Poisson equation over a rectangle). We consider in this exercise Poisson's equation in the domain  $\Omega = (0, 2) \times (0, 1)$ , equipped with homogeneous Dirichlet boundary conditions:

$$\begin{aligned} -\Delta f(x, y) &= b(x, y), & x \in \Omega, \\ f(x) &= 0, & x \in \partial\Omega. \end{aligned}$$

The right-hand side is

$$b(x, y) = \sin(4\pi x) + \sin(2\pi y).$$

A number of methods can be employed in order to discretize this partial differential equation. After discretization, a finite-dimensional linear system of the form  $\mathbf{A}\mathbf{x} = \mathbf{b}$  is obtained. A Julia function for calculating the matrix  $\mathbf{A}$  and the vector  $\mathbf{b}$  using the finite difference method is given to you on the course website, as well as a function to plot the solution. The goal of this exercise is to solve the linear system using the conjugate gradient method. Use the same stopping criterion as in [Exercise 4.25](#).

❁ **Exercise 4.23.** Show that if  $\mathbf{A} \in \mathbf{R}^{n \times n}$  is nonsingular, then the solution to the equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$  belongs to the Krylov subspace

$$\mathcal{K}_n(\mathbf{A}, \mathbf{b}) = \text{Span}\{\mathbf{b}, \mathbf{A}\mathbf{b}, \mathbf{A}^2\mathbf{b}, \dots, \mathbf{A}^{n-1}\mathbf{b}\}.$$

❁ **Exercise 4.24.** Write a function `lu(A)` for calculating the LU decomposition of a square matrix  $\mathbf{A} \in \mathbf{R}^{n \times n}$ , with  $\mathbf{L}$  unit lower triangular and  $\mathbf{U}$  upper triangular, not by Gaussian elimination but by comparing the entries of the product  $\mathbf{LU}$  with those of  $\mathbf{A}$ . To this end, one option is to compare the entries one by one in the order  $(1, 1)$ ,  $(1, 2)$ , ...,  $(1, n)$ ,  $(2, 1)$ ,  $(2, 2)$ , ..., i.e. row by row starting from the top. For example,

- Comparing the entry  $(1, k)$  with  $k \in \{1, \dots, n\}$  gives

$$\ell_{11}u_{1k} = a_{1k}.$$

Since  $\ell_{11} = 1$  as  $\mathbf{L}$  is unit lower triangular, this implies that  $u_{1k} = a_{1k}$ .

- Comparing the entry  $(2, 1)$  gives

$$\ell_{21}u_{11} = a_{21}$$

and so  $\ell_{21} = a_{21}/u_{11}$ .



- Comparing the entry  $(2, k)$  with  $k \in \{2, \dots, n\}$  gives

$$\ell_{21}u_{1k} + \ell_{22}u_{2k} = a_{2k}.$$

Given the previous items, the only unknown in this equation is  $u_{2k}$ .

- Comparing the entry  $(3, 1)$  gives

$$\ell_{31}u_{11} = a_{31},$$

and so  $\ell_{31} = a_{31}/u_{11}$ .

- Comparing the entry  $(3, 2)$  gives

$$\ell_{31}u_{12} + \ell_{32}u_{22} = a_{32},$$

Given the previous items, the only unknown in this equation is  $\ell_{32}$ .

- Comparing the entry  $(3, k)$  with  $k \in \{3, \dots, n\}$  gives

$$\ell_{31}u_{1k} + \ell_{32}u_{2k} + \ell_{33}u_{3k} = a_{3k},$$

Given the previous items, the only unknown in this equation is  $u_{3k}$ .

Notice that a pattern seems to be emerging: when going through the entries row by row starting from the top left corner of the matrix, comparing the entry  $(i, j)$  provides an equation for  $\ell_{ij}$  if  $j < i$ , and an equation for  $u_{ij}$  if  $j \geq i$ . Do not use any external package for this exercise.

**Extra credit:** ... if your code is able to exploit the potential banded structure of the matrix passed as argument for better efficiency. Specifically, your code will be tested with a matrix created as follows

```
b, n = 5, 10000
A = [abs(i-j) <= b ? rand() : 0.0 for i in 1:n, j in 1:n]
```

**□ Exercise 4.25.** Implement an iterative method based on a splitting for finding a solution to the following linear system on  $\mathbf{R}^n$ .

$$\frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & & & \\ -1 & 2 & -1 & & & & \\ & -1 & 2 & -1 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & -1 & 2 & -1 & \\ & & & & -1 & 2 & \\ & & & & & -1 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_{n-1} \\ x_n \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix}, \quad h = \frac{1}{n+1}.$$

Plot the norm of the residual as a function of the iteration index. Use as stopping criterion the condition

$$\|\mathbf{r}^{(k)}\| \leq \epsilon \|\mathbf{b}\|, \quad \epsilon = 10^{-8}.$$

As initial guess, use a vector of zeros. The code will be tested with  $n = 500$ . Do not use any library (except for plotting), and do not use the backslash operator.

**❁ Exercise 4.26.** Find a formula for the optimal value of  $\omega$  in the relaxation method given  $n$ , for the linear system in Exercise 4.25. The proof of Proposition 4.15, as well as the formula (4.25) for the eigenvalues of a tridiagonal matrix, are useful to this end.

*Solution.* Corollary 4.13 and Proposition 4.14 imply that a sufficient and necessary condition for convergence, when  $A$  is Hermitian and positive definite, is that  $\omega \in (0, 2)$ . Let  $M_\omega = \frac{1}{\omega}D + L$  and  $N_\omega = \frac{1-\omega}{\omega}D - U$ . A nonzero scalar  $\lambda \in \mathbf{C}$  is an eigenvalue of  $M_\omega^{-1}N_\omega$  if and only if

$$\det(M_\omega^{-1}N_\omega - \lambda I) = 0 \quad \Leftrightarrow \quad \det(M_\omega^{-1}) \det(N_\omega - \lambda M_\omega) = 0 \quad \Leftrightarrow \quad \det(\lambda M_\omega - N_\omega) = 0.$$

Substituting the expressions of  $M_\omega$  and  $N_\omega$ , we obtain that this condition can be equivalently rewritten as

$$\det\left(\lambda L + \left(\frac{\lambda + \omega - 1}{\omega}\right)D + U\right) = 0 \quad \Leftrightarrow \quad \det\left(\sqrt{\lambda}L + \left(\frac{\lambda + \omega - 1}{\omega}\right)D + \sqrt{\lambda}U\right) = 0$$

where we used (4.24) for the last equivalence. The equality of the determinants in these two equations is valid for  $\sqrt{\lambda}$  denoting either of the two complex square roots of  $\lambda$ . This condition is equivalent to

$$\det\left(L + \left(\frac{\lambda + \omega - 1}{\sqrt{\lambda}\omega}\right)D + U\right) = 0.$$

We recognize from the proof of Proposition 4.15 that this condition is equivalent to

$$\frac{\lambda + \omega - 1}{\sqrt{\lambda}\omega} \in \text{spectrum}(M_{\mathcal{J}}^{-1}N_{\mathcal{J}}).$$

In other words, for any  $(\lambda, \mu) \in \mathbf{C}^2$  such that

$$\frac{(\lambda + \omega - 1)^2}{\lambda\omega^2} = \mu^2, \tag{4.43}$$

it holds that  $\mu \in \text{spectrum}(M_{\mathcal{J}}^{-1}N_{\mathcal{J}})$  if and only if  $\lambda \in \text{spectrum}(M_\omega^{-1}N_\omega)$ . By (4.25), the eigenvalues of  $M_{\mathcal{J}}^{-1}N_{\mathcal{J}}$  are real and given by

$$\mu_j = \cos\left(\frac{j\pi}{n+1}\right), \quad 1 \leq j \leq n. \tag{4.44}$$

Rearranging (4.43), we find

$$\lambda^2 + \lambda(2(\omega - 1) - \omega^2\mu^2) + (\omega - 1)^2 = 0.$$

For given  $\omega \in (0, 2)$  and  $\mu \in \mathbf{R}$ , this is a quadratic equation for  $\lambda$  with solutions

$$\lambda_{\pm} = \left(\frac{\omega^2\mu^2}{2} + 1 - \omega\right) \pm \omega\mu\sqrt{\frac{\omega^2\mu^2}{4} + 1 - \omega},$$

Since the first bracket is positive when the argument of the square root is positive, it is clear that

$$\max\{|\lambda_-|, |\lambda_+|\} = \left|\frac{\omega^2\mu^2}{2} + 1 - \omega + \omega|\mu|\sqrt{\frac{\omega^2\mu^2}{4} + 1 - \omega}\right|.$$

Combining this with (4.44), we deduce that the spectral radius of  $M_\omega^{-1}N_\omega$  is given by

$$\rho(M_\omega^{-1}N_\omega) = \max_{j \in \{1, \dots, n\}} \left| \frac{\omega^2 \mu_j^2}{2} + 1 - \omega + \omega |\mu_j| \sqrt{\frac{\omega^2 \mu_j^2}{4} + 1 - \omega} \right|. \quad (4.45)$$

We wish to minimize this expression over the interval  $\omega \in (0, 2)$ . While this can be achieved by algebraic manipulations, we content ourselves here with graphical exploration. Figure 4.3 depicts the amplitude of the modulus in (4.45) for different values of  $\mu$ . It is apparent that, for given  $\omega$ , the modulus increases as  $\mu$  increases, which suggests that

$$\rho(M_\omega^{-1}N_\omega) = \left| \frac{\omega^2 \mu_*^2}{2} + 1 - \omega + \omega |\mu_*| \sqrt{\frac{\omega^2 \mu_*^2}{4} + 1 - \omega} \right|, \quad \mu_* = \rho(M_{\mathcal{J}}^{-1}N_{\mathcal{J}}). \quad (4.46)$$

The figure also suggests that for a given value of  $\mu$ , the modulus is minimized at the discontinuity of the first derivative, which occurs when the argument of the square root is zero. We conclude that the optimal  $\omega$  satisfies

$$\frac{\omega_{\text{opt}}^2 \mu_*^2}{4} + 1 - \omega_{\text{opt}} = 0 \quad \xrightarrow{\omega < 2} \quad \omega_{\text{opt}} = 2 \frac{1 - \sqrt{1 - \mu_*^2}}{\mu_*^2} = \frac{2}{1 + \sqrt{1 - \mu_*^2}} = \frac{2}{1 + \sin\left(\frac{\pi}{n+1}\right)}.$$

△

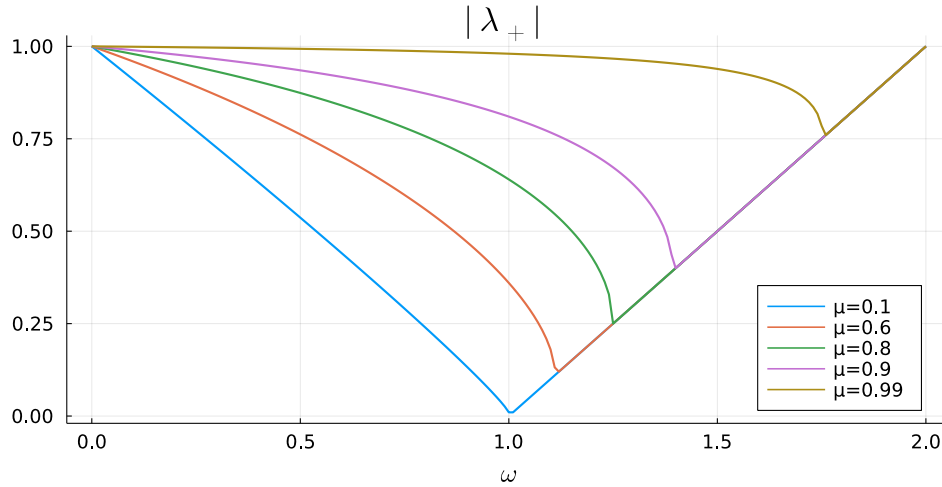


Figure 4.3: Modulus of  $|\lambda_+|$  as a function of  $\omega$ , for different eigenvalues of  $\mu$ .

⚙️ **Exercise 4.27** (Midterm 2022). Let  $A \in \mathbf{R}^{n \times n}$  be a symmetric positive definite matrix and let  $\mathbf{b} \in \mathbf{R}^n$ . The steepest descent algorithm for solving  $A\mathbf{x} = \mathbf{b}$  is given hereafter:

```

Pick  $\varepsilon > 0$  and initial  $\mathbf{x}$ 
 $\mathbf{r} \leftarrow A\mathbf{x} - \mathbf{b}$ 
while  $\|\mathbf{r}\| \geq \varepsilon \|\mathbf{b}\|$  do
     $\omega \leftarrow \mathbf{r}^T \mathbf{r} / \mathbf{r}^T A \mathbf{r}$ 
     $\mathbf{x} \leftarrow \mathbf{x} - \omega \mathbf{r}$ 
     $\mathbf{r} \leftarrow A\mathbf{x} - \mathbf{b}$ 
end while
    
```

- Why is this method called the steepest descent algorithm?
- How many floating point operations does an iteration of this algorithm require?
- Are the following statements true or false? (2 marks)

1. There exists a unique solution  $\mathbf{x}_*$  to the linear system  $\mathbf{Ax} = \mathbf{b}$ .
2. The iterates converge to  $\mathbf{x}_*$  in at most  $n$  iterations.
3. We consider the following modification of the algorithm:

Pick  $\varepsilon > 0$ ,  $\omega > 0$  and initial  $\mathbf{x}$

$\mathbf{r} \leftarrow \mathbf{Ax} - \mathbf{b}$

**while**  $\|\mathbf{r}\| \geq \varepsilon\|\mathbf{b}\|$  **do**

$\mathbf{x} \leftarrow \mathbf{x} - \omega\mathbf{r}$

$\mathbf{r} \leftarrow \mathbf{Ax} - \mathbf{b}$

**end while**

If  $\omega$  is sufficiently small, then this algorithm converges.

4. Here we no longer assume that  $\mathbf{A}$  is positive definite. Instead, we consider that

$$\mathbf{A} = \begin{pmatrix} -1 & 0 \\ 0 & -2 \end{pmatrix}.$$

In this case, the steepest descent algorithm is convergent for any initial  $\mathbf{x}$ .

⚙️ **Exercise 4.28** (Final exam Spring 2022). Assume that  $\mathbf{A} \in \mathbf{R}^{n \times n}$  is a nonsingular matrix and that  $\mathbf{b} \in \mathbf{R}^n$ . We wish to solve the linear system (4.1) using an iterative method where each iteration is of the form

$$\mathbf{M}\mathbf{x}_{k+1} = \mathbf{N}\mathbf{x}_k + \mathbf{b}. \quad (4.47)$$

Here  $\mathbf{A} = \mathbf{M} - \mathbf{N}$  is a splitting of  $\mathbf{A}$  such that  $\mathbf{M}$  is nonsingular, and  $\mathbf{x}_k \in \mathbf{R}^n$  denotes the  $k$ -th iterate of the numerical scheme.

1. Let  $\mathbf{e}_k := \mathbf{x}_k - \mathbf{x}_*$ , where  $\mathbf{x}_*$  is the exact solution to (4.1). Prove that

$$\mathbf{e}_{k+1} = \mathbf{M}^{-1}\mathbf{N}\mathbf{e}_k.$$

2. Let  $L = \|\mathbf{M}^{-1}\mathbf{N}\|_\infty$ . Prove that

$$\forall k \in \mathbf{N}, \quad \|\mathbf{e}_k\|_\infty \leq L^k \|\mathbf{e}_0\|_\infty. \quad (4.48)$$

3. Is the condition  $\|\mathbf{M}^{-1}\mathbf{N}\|_\infty < 1$  necessary for convergence when  $\mathbf{x}_0 \neq \mathbf{x}_*$ ?
4. Assume that  $\mathbf{A}$  is strictly row diagonally dominant, in the sense that

$$\forall i \in \{1, \dots, n\}, \quad |a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}|.$$

Show that, in this case, the inequality  $\|\mathbf{M}^{-1}\mathbf{N}\|_\infty < 1$  holds for the Jacobi method, i.e. when  $\mathbf{M}$  contains just the diagonal of  $\mathbf{A}$ . You may take for granted the following expression for the  $\infty$ -norm of a matrix  $\mathbf{X} \in \mathbf{R}^{n \times n}$ :

$$\|\mathbf{X}\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |x_{ij}|.$$

5. Write down a few iterations of the Jacobi method when

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \mathbf{x}_0 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$


Is the method convergent?

## 4.5 Discussion and bibliography

In this chapter, we presented direct methods and some of the standard iterative methods for solving linear systems. We focused particularly on linear systems with a symmetric positive definite matrix. Section 4.2 is based on [10, 18] and Section 4.3 roughly follows [15, Chapter 2]. The book [11] is a very detailed reference on iterative methods for solving sparse linear systems. The reference [13] is an excellent introduction to the conjugate gradient method.

# Chapter 5

## Solution of nonlinear systems

5.1	The bisection method . . . . .	128
5.2	Fixed point methods . . . . .	129
5.3	Convergence of fixed point methods . . . . .	130
5.4	Examples of fixed point methods . . . . .	134
5.4.1	The chord method . . . . .	134
5.4.2	The Newton–Raphson method . . . . .	135
5.4.3	The secant method  . . . . .	139
5.5	A numerical experiment . . . . .	142
5.6	Exercises . . . . .	144
5.7	Discussion and bibliography . . . . .	147

### Introduction

This chapter concerns the numerical solution of nonlinear equations of the general form

$$\mathbf{f}(\mathbf{x}) = 0, \quad \mathbf{f}: \mathbf{R}^n \rightarrow \mathbf{R}^n. \tag{5.1}$$

A solution to this equation is called a *zero* of the function  $f$ . Except in particular cases (for example linear systems), there does not exist a numerical method for solving (5.1) in a finite number of operations, so iterative methods are required.

In contrast with the previous chapter, it may not be the case that (5.1) admits one and only one solution. For example, the equation  $1 + x^2 = 0$  does not have a (real) solution, and the equation  $\cos(x) = 0$  has infinitely many. Therefore, convergence results usually contain assumptions on the function  $f$  that guarantee the existence and uniqueness of a solution in  $\mathbf{R}^n$  or a subset of  $\mathbf{R}^n$ .

For an iterative method generating approximations  $(\mathbf{x}_k)_{k \geq 0}$  of a root  $\mathbf{x}_*$ , we define the error

as  $\mathbf{e}_k = \mathbf{x}_k - \mathbf{x}_*$ . If the sequence  $(\mathbf{x}_k)_{k \geq 0}$  converges to  $\mathbf{x}_*$  in the limit as  $k \rightarrow \infty$  and if

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{e}_{k+1}\|}{\|\mathbf{e}_k\|^q} = r, \quad (5.2)$$

then we say that  $(\mathbf{x}_k)_{k \geq 0}$  converges with *order of convergence*  $q$  and *rate of convergence*  $r$ . In addition, we say that the convergence is linear  $q = 1$ , and quadratic if  $q = 2$ . The convergence is said to be superlinear if

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{e}_{k+1}\|}{\|\mathbf{e}_k\|} = 0. \quad (5.3)$$

In particular, the convergence is superlinear if the order of convergence is  $q > 1$ .

*Remark 5.1.* The notion of order of convergence may be defined also when the limit in (5.2) does not exist. A more general definition for the order of convergence of a sequence  $(\mathbf{x}_k)_{k \geq 0}$  converging to  $\mathbf{x}_*$  is the following:

$$q(\mathbf{x}_0) = \inf \left\{ p \in [1, \infty) : \limsup_{k \rightarrow \infty} \frac{\|\mathbf{e}_{k+1}\|}{\|\mathbf{e}_k\|^p} = \infty \right\},$$

or  $q(\mathbf{x}_0) = \infty$  if the numerator and denominator of the fraction are zero for sufficiently large  $k$ . It is possible to define similarly the order of convergence of an iterative method for an initial guess in a neighborhood  $V$  of  $\mathbf{x}_*$ :

$$q = \inf \left\{ p \in [1, \infty) : \sup_{\mathbf{x}_0 \in V} \left( \limsup_{k \rightarrow \infty} \frac{\|\mathbf{e}_{k+1}\|}{\|\mathbf{e}_k\|^p} \right) = \infty \right\},$$

where the fraction should be interpreted as 0 if the numerator and denominator are zero. A more detailed discussion of this subject is beyond the scope of this course.

The rest of chapter is organized as follows:

- In [Section 5.1](#), by way of introduction to the subject of numerical methods for nonlinear equations, we present and analyze the bisection method.
- In [Section 5.2](#), we present a general method based on a fixed point iteration for solving (5.1). The convergence of this method is analyzed in [Section 5.3](#).
- In [Section 5.4](#), two concrete examples of fixed point methods are studied: the chord method and the Newton–Raphson method.

## 5.1 The bisection method

As an introduction to numerical methods for solving nonlinear equations, we present the bisection method. This method applies only in the case of a real-valued function  $f: \mathbf{R} \rightarrow \mathbf{R}$ , and relies on the knowledge of two points  $a < b$  such that  $f(a)$  and  $f(b)$  have different signs. By the intermediate value theorem, there necessarily exists  $x_* \in (a, b)$  such that  $f(x_*) = 0$ . The idea of the bisection method is to successively divide the interval in two equal parts, and to retain, based on the sign of  $f$  at the midpoint  $x_{1/2}$ , the one that necessarily contains a root.

If  $f(x_{1/2})f(a) \geq 0$ , then  $f(x_{1/2})f(b) \leq 0$  and so there necessarily exists a root of  $f$  in the interval  $[x_{1/2}, b)$  by the intermediate value theorem. In contrast, if  $f(x_{1/2})f(a) < 0$ , then there necessarily is a root in the interval  $(a, x_{1/2})$ . The algorithm is presented in [Algorithm 7](#).

---

**Algorithm 7** Bisection method
 

---

Assume that  $f(a)f(b) < 0$  with  $a < b$ .

Pick  $\varepsilon > 0$ .

$x \leftarrow a/2 + b/2$

**while**  $|b - a| \geq \varepsilon$  **do**

**if**  $f(x)f(a) \geq 0$  **then**

$a \leftarrow x$

**else**

$b \leftarrow x$

**end if**

$x \leftarrow a/2 + b/2$

**end while**

---

The following result establishes the convergence of the method.

**Proposition 5.1.** *Assume that  $f: \mathbf{R} \rightarrow \mathbf{R}$  is a continuous function and  $f(a)f(b) < 0$ . Let  $[a_j, b_j]$  denote the interval obtained after  $j$  iterations of the bisection method, and let  $x_j = (a_j + b_j)/2$  denote the midpoint of the interval. Then there exists a root  $x_*$  of  $f$  such that*

$$|x_j - x_*| \leq (b_0 - a_0)2^{-(j+1)}. \quad (5.4)$$

*Proof.* By construction,  $f(a_j)f(b_j) \leq 0$  and  $f(b) \neq 0$ . Therefore, by the intermediate value theorem, there exists a root of  $f$  in the interval  $[a_j, b_j)$ , implying that

$$|x_j - x_*| \leq \frac{b_j - a_j}{2}.$$

Since  $b_j - a_j = 2^{-j}(b_0 - a_0)$ , the statement follows.  $\square$

Although the limit in (5.2) may not be well-defined (for example,  $x_1$  may be a root of  $f$ ), the error  $x_j - x_*$  is bounded in absolute value by the sequence  $(\tilde{e}_j)_{j \geq 0}$ , where  $\tilde{e}_j = (b_0 - a_0)2^{-(j+1)}$  by [Proposition 5.1](#). Since the latter sequence exhibits linear convergence to 0, the convergence of the bisection method is said to be linear, by a slight abuse of terminology.

## 5.2 Fixed point methods

Let  $\mathbf{x}_*$  denote a zero of the function  $\mathbf{f}$ . The idea of iterative methods for (5.1) is to construct, starting from an initial guess  $\mathbf{x}_0$ , a sequence  $(\mathbf{x}_k)_{k=0,1,\dots}$  approaching  $\mathbf{x}_*$ . A number of iterative methods for solving (5.1) are based on an iteration of the form

$$\mathbf{x}_{k+1} = \mathbf{F}(\mathbf{x}_k), \quad (5.5)$$



for an appropriate continuous function  $F$ . Assume that  $\mathbf{x}_k$  converges to some point  $\mathbf{x}_* \in \mathbf{R}^n$  in the limit as  $k \rightarrow \infty$ . Then, taking the limit  $k \rightarrow \infty$  in (5.5), we find that  $\mathbf{x}_*$  satisfies

$$\mathbf{F}(\mathbf{x}_*) = \mathbf{x}_*.$$

Such a point  $\mathbf{x}_*$  is called a *fixed point* of the function  $\mathbf{F}$ . Several definitions of the function  $\mathbf{F}$  can be employed in order to ensure that a fixed point of  $\mathbf{F}$  coincides with a zero of  $\mathbf{f}$ . One may, for example, define  $\mathbf{F}(\mathbf{x}) = \mathbf{x} - \alpha^{-1}\mathbf{f}(\mathbf{x})$ , for some nonzero scalar coefficient  $\alpha$ . Then  $\mathbf{F}(\mathbf{x}_*) = \mathbf{x}_*$  if and only if  $\mathbf{f}(\mathbf{x}_*) = 0$ . Later in this chapter, in Section 5.4, we study two instances of numerical methods which can be recast in the form (5.5). Before this, we study the convergence of the iteration (5.5) for a general function  $\mathbf{F}$ .

### 5.3 Convergence of fixed point methods

Equation (5.5) may be viewed as a *discrete-time* dynamical system. In order to study the behavior of the system as  $k \rightarrow \infty$ , it is important to understand the concept of stability of a fixed point. The concept of stability appears also in the field of ordinary differential equations, which are *continuous-time* dynamical systems. Before we define this concept, we introduce the following notation for the open ball of radius  $\delta$  around  $\mathbf{x} \in \mathbf{R}^n$ :

$$B_\delta(\mathbf{x}) := \{\mathbf{y} \in \mathbf{R}^n : \|\mathbf{y} - \mathbf{x}\| < \delta\}.$$

**Definition 5.1** (Stability of fixed points). Let  $(\mathbf{x}_k)_{k \geq 0}$  denote iterates obtained from (5.5) when starting from an initial vector  $\mathbf{x}_0$ . Then we say that a fixed point  $\mathbf{x}_*$  is

- an *attractor* if there exists a neighborhood  $\mathcal{V}$  of  $\mathbf{x}_*$  such that

$$\forall \mathbf{x}_0 \in \mathcal{V}, \quad \mathbf{x}_k \xrightarrow[k \rightarrow \infty]{} \mathbf{x}_*. \quad (5.6)$$

The largest neighborhood for which this is true, i.e. the set of values of  $\mathbf{x}_0$  such that (5.6) holds true, is called the basin of attraction of  $\mathbf{x}_*$ .

- stable (in the sense of Lyapunov) if for all  $\varepsilon > 0$ , there exists  $\delta > 0$  such that

$$\forall \mathbf{x}_0 \in B_\delta(\mathbf{x}_*), \quad \forall k \in \mathbf{N}, \quad \|\mathbf{x}_k - \mathbf{x}_*\| < \varepsilon.$$

- asymptotically stable if it is stable and an attractor.
- exponentially stable if there exists  $C > 0$ ,  $\alpha \in (0, 1)$ , and  $\delta > 0$  such that

$$\forall \mathbf{x}_0 \in B_\delta(\mathbf{x}_*), \quad \forall k \in \mathbf{N}, \quad \|\mathbf{x}_k - \mathbf{x}_*\| \leq C\alpha^k \|\mathbf{x}_0 - \mathbf{x}_*\|.$$

- globally exponentially stable if there exists  $C > 0$  and  $\alpha \in (0, 1)$  such that

$$\forall \mathbf{x}_0 \in \mathbf{R}^n, \quad \forall k \in \mathbf{N}, \quad \|\mathbf{x}_k - \mathbf{x}_*\| \leq C\alpha^k \|\mathbf{x}_0 - \mathbf{x}_*\|.$$

- unstable if it is not stable.

Clearly, global exponential stability implies exponential stability, which itself implies asymptotic stability and stability. If  $\mathbf{x}_*$  is globally exponentially stable, then  $\mathbf{x}_*$  is the unique fixed point of  $\mathbf{F}$ ; showing this is the aim of [Exercise 5.3](#). If  $\mathbf{x}_*$  is an attractor, then the dynamical system (5.5) is said to be locally convergent to  $\mathbf{x}_*$ . The larger the basin of attraction of  $\mathbf{x}_*$ , the less careful we need to be when picking the initial guess  $\mathbf{x}_0$ . Global exponential stability of a fixed point can sometimes be shown provided that  $\mathbf{F}$  satisfies a strong hypothesis.

**Definition 5.2** (Lipschitz continuity). A function  $\mathbf{F}: \mathbf{R}^n \rightarrow \mathbf{R}^n$  is said to be *Lipschitz continuous* with constant  $L$  if

$$\forall (\mathbf{x}, \mathbf{y}) \in \mathbf{R}^n \times \mathbf{R}^n, \quad \|\mathbf{F}(\mathbf{y}) - \mathbf{F}(\mathbf{x})\| \leq L\|\mathbf{y} - \mathbf{x}\|.$$

A function  $\mathbf{F}: \mathbf{R}^n \rightarrow \mathbf{R}^n$  that is Lipschitz continuous with a constant  $L < 1$  is called a *contraction*. For such a function, it is possible to prove that (5.5) has a unique globally exponentially stable fixed point.

**Theorem 5.2.** *Assume that  $\mathbf{F}$  is a contraction. Then there exists a unique fixed point of (5.5), and it holds that*

$$\forall \mathbf{x}_0 \in \mathbf{R}^n, \quad \forall k \in \mathbf{N}, \quad \|\mathbf{x}_k - \mathbf{x}_*\| \leq L^k \|\mathbf{x}_0 - \mathbf{x}_*\|. \quad (5.7)$$

*Proof.* Existence and uniqueness of the fixed point follows from the *Banach fixed point theorem*, see [Theorem A.3](#), so here we show only global exponential convergence. Since  $\mathbf{F}$  is a contraction, it holds that

$$\|\mathbf{x}_k - \mathbf{x}_*\| = \|\mathbf{F}(\mathbf{x}_{k-1}) - \mathbf{F}(\mathbf{x}_*)\| \leq L\|\mathbf{x}_{k-1} - \mathbf{x}_*\| \leq \dots \leq L^k \|\mathbf{x}_0 - \mathbf{x}_*\|, \quad (5.8)$$

which proves (5.7). □

It is possible to prove a weaker, local result under a less restrictive assumptions on the function  $\mathbf{F}$ .

**Theorem 5.3.** *Assume that  $\mathbf{x}_*$  is a fixed point of (5.5) and that  $\mathbf{F}: \mathbf{R}^n \rightarrow \mathbf{R}^n$  satisfies the local Lipschitz condition*

$$\forall \mathbf{x} \in B_\delta(\mathbf{x}_*), \quad \|\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}_*)\| \leq L\|\mathbf{x} - \mathbf{x}_*\|, \quad (5.9)$$

*with  $0 \leq L < 1$  and  $\delta > 0$ . Then  $\mathbf{x}_*$  is the unique fixed point of  $\mathbf{F}$  in  $B_\delta(\mathbf{x}_*)$  and, for all  $\mathbf{x}_0 \in B_\delta(\mathbf{x}_*)$ , it holds that*

- All the iterates  $(\mathbf{x}_k)_{k \in \mathbb{N}}$  belong to  $B_\delta(\mathbf{x}_*)$ .
- The sequence  $(\mathbf{x}_k)_{k \in \mathbb{N}}$  converges exponentially to  $\mathbf{x}_*$ .

*Proof.* See Exercise 5.4. □

It is possible to guarantee that condition (5.9) holds provided that we have sufficiently good control of the derivatives of the function  $\mathbf{F}$ . The function  $\mathbf{F}$  is said to be differentiable at  $\mathbf{x}$  (in the sense of Fréchet) if there exists a linear operator  $D\mathbf{F}_{\mathbf{x}}: \mathbf{R}^n \rightarrow \mathbf{R}^n$  such that

$$\lim_{\mathbf{h} \rightarrow 0} \frac{\|\mathbf{F}(\mathbf{x} + \mathbf{h}) - \mathbf{F}(\mathbf{x}) - D\mathbf{F}_{\mathbf{x}}(\mathbf{h})\|}{\|\mathbf{h}\|} = 0. \quad (5.10)$$

If  $\mathbf{F}$  is differentiable, then all its first partial derivatives  $\partial_j F_i$  exist and, in addition, it holds that  $D\mathbf{F}_{\mathbf{x}}(\mathbf{h}) = \mathbf{J}_F(\mathbf{x})\mathbf{h}$  where  $\mathbf{J}_F(\mathbf{x})$  is the Jacobian matrix of  $\mathbf{F}$  at  $\mathbf{x}$ :

$$\mathbf{J}_F(\mathbf{x}) = \begin{pmatrix} \partial_1 F_1(\mathbf{x}) & \dots & \partial_n F_1(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \partial_1 F_n(\mathbf{x}) & \dots & \partial_n F_n(\mathbf{x}) \end{pmatrix}.$$

**Proposition 5.4.** *Let  $\mathbf{x}_*$  be a fixed point of (5.5), and assume that there exists  $\delta$  and a subordinate matrix norm such that  $\mathbf{F}$  is differentiable everywhere in  $B_\delta(\mathbf{x}_*)$  and*

$$\forall \mathbf{x} \in B_\delta(\mathbf{x}_*), \quad \|\mathbf{J}_F(\mathbf{x})\| \leq L < 1.$$

*Then condition (5.9) is satisfied in the associated vector norm, and so the fixed point  $\mathbf{x}_*$  is locally exponentially stable.*

*Proof.* Let  $\mathbf{x} \in B_\delta(\mathbf{x}_*)$ . By the fundamental theorem of calculus and the chain rule, we have

$$\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}_*) = \int_0^1 \frac{d}{dt} \left( \mathbf{F}(\mathbf{x}_* + t(\mathbf{x} - \mathbf{x}_*)) \right) dt = \int_0^1 \mathbf{J}_F(\mathbf{x}_* + t(\mathbf{x} - \mathbf{x}_*)) (\mathbf{x} - \mathbf{x}_*) dt.$$

Therefore, it holds that

$$\|\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}_*)\| \leq \int_0^1 \|\mathbf{J}_F(\mathbf{x}_* + t(\mathbf{x} - \mathbf{x}_*))\| dt \|\mathbf{x} - \mathbf{x}_*\| \leq \int_0^1 L dt \|\mathbf{x} - \mathbf{x}_*\| = L \|\mathbf{x} - \mathbf{x}_*\|,$$

which is the statement. □

*Remark 5.2.* As a student observed during the lecture, in dimension  $n = 1$ , Proposition 5.4 can be proved by using the mean value theorem: since  $F$  is differentiable in  $(x_* - \delta, x_* + \delta)$ , there exists for all  $x$  in this interval a  $\xi = \xi(x)$  also in this interval such that

$$F(x) - F(x_*) = F'(\xi)(x - x_*).$$

It then follows immediately that

$$|F(x) - F(x_*)| = |F'(\xi)(x - x_*)| \leq L|x - x_*|.$$

This proof does not carry over to the multi-dimensional setting, however.

In fact, it is possible to prove that a fixed point  $\mathbf{x}_*$  is exponentially locally stable under an even weaker condition, involving only the derivative of  $\mathbf{F}$  at  $\mathbf{x}_*$ .

**Proposition 5.5.** *Let  $\mathbf{x}_*$  be a fixed point of (5.5) and that  $F$  is differentiable at  $\mathbf{x}_*$  with*

$$\|\mathbf{J}_F(\mathbf{x}_*)\| = L < 1,$$

*in a subordinate vector norm. Then the fixed point  $\mathbf{x}_*$  is locally exponentially stable.*

*Proof.* In this proof, the vector norm used is that associated with the matrix norm in the statement of the proposition. By the definition of differentiability (5.10), there exists for all  $\varepsilon > 0$  a  $\delta > 0$  such that

$$\forall \mathbf{x} \in B_\delta(\mathbf{x}_*) \setminus \{\mathbf{x}_*\}, \quad \frac{\|\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}_*) - \mathbf{J}_F(\mathbf{x}_*)(\mathbf{x} - \mathbf{x}_*)\|}{\|\mathbf{x} - \mathbf{x}_*\|} \leq \varepsilon.$$

By the triangle inequality, this implies that for all  $\mathbf{x} \in B_\delta(\mathbf{x}_*)$ ,

$$\begin{aligned} \|\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}_*)\| &\leq \|\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}_*) - \mathbf{J}_F(\mathbf{x}_*)(\mathbf{x} - \mathbf{x}_*)\| + \|\mathbf{J}_F(\mathbf{x}_*)(\mathbf{x} - \mathbf{x}_*)\| \\ &\leq \varepsilon\|\mathbf{x} - \mathbf{x}_*\| + \|\mathbf{J}_F(\mathbf{x}_*)\|\|\mathbf{x} - \mathbf{x}_*\| = (L + \varepsilon)\|\mathbf{x} - \mathbf{x}_*\|. \end{aligned}$$

We have thus shown that for all  $\varepsilon > 0$ , there exists  $\delta > 0$  such that condition (5.9) is satisfied with constant  $L + \varepsilon$ . By taking  $\varepsilon$  sufficiently small, we can ensure that  $L + \varepsilon < 1$ , and so the fixed point  $\mathbf{x}_*$  is locally exponentially stable by Theorem 5.3.  $\square$

The estimate in Theorem 5.2 suggests that when the fixed point iteration (5.5) converges, the convergence is linear. While this is usually the case, the convergence is superlinear if  $\mathbf{J}_F(\mathbf{x}_*) = 0$ .

**Proposition 5.6.** *Assume that  $\mathbf{x}_*$  is a fixed point of (5.5) and that  $\mathbf{J}_F(\mathbf{x}_*) = 0$ . Then the convergence to  $\mathbf{x}_*$  is superlinear, in the sense that if  $\mathbf{x}_k \rightarrow \mathbf{x}_*$  as  $k \rightarrow \infty$ , then*

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}_*\|}{\|\mathbf{x}_k - \mathbf{x}_*\|} = 0.$$

*Proof.* By Proposition 5.5, there exists  $\delta > 0$  such that  $(\mathbf{x}_k)_{k \geq 0}$  is a sequence converging to  $\mathbf{x}_*$  for all  $\mathbf{x}_0 \in B_\delta(\mathbf{x}_*)$ . It holds that

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}_*\|}{\|\mathbf{x}_k - \mathbf{x}_*\|} = \frac{\|\mathbf{F}(\mathbf{x}_k) - \mathbf{F}(\mathbf{x}_*)\|}{\|\mathbf{x}_k - \mathbf{x}_*\|} = \frac{\|\mathbf{F}(\mathbf{x}_k) - \mathbf{F}(\mathbf{x}_*) - \mathbf{J}_F(\mathbf{x}_*)(\mathbf{x}_k - \mathbf{x}_*)\|}{\|\mathbf{x}_k - \mathbf{x}_*\|}.$$

Since  $\mathbf{x}_k - \mathbf{x}_* \rightarrow \mathbf{0}$  as  $k \rightarrow \infty$ , the right-hand side converges to 0 by (5.10).  $\square$

Similarly, if there exist  $\delta > 0$ ,  $C > 0$  and  $q \in (1, \infty)$  such that

$$\forall \mathbf{x} \in B_\delta(\mathbf{x}_*), \quad \|\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{x}_*)\| \leq C\|\mathbf{x} - \mathbf{x}_*\|^q, \quad (5.11)$$

then assuming that  $(\mathbf{x}_k)_{k \geq 0}$  converges to  $\mathbf{x}_*$ , it holds for sufficiently large  $k$  that

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}_*\|}{\|\mathbf{x}_k - \mathbf{x}_*\|^q} = \frac{\|\mathbf{F}(\mathbf{x}_k) - \mathbf{F}(\mathbf{x}_*)\|}{\|\mathbf{x}_k - \mathbf{x}_*\|^q} \leq C.$$

In this case, the order of convergence is at least  $q$ .

## 5.4 Examples of fixed point methods

As we mentioned in Section 5.2, there are several choices for the function  $\mathbf{F}$  that guarantee the equivalence  $\mathbf{F}(\mathbf{x}) = \mathbf{x} \Leftrightarrow \mathbf{f}(\mathbf{x}) = \mathbf{0}$ .

### 5.4.1 The chord method

In the case where  $f$  is a function from  $\mathbf{R}$  to  $\mathbf{R}$ , the simplest approach, sometimes called the *chord method*, is to define

$$F(x) = x - \alpha^{-1}f(x).$$

The fixed point iteration (5.4) in this case admits a simple geometric interpretation: at each step, the function  $f$  is approximated by the affine function  $x \mapsto f(x_k) + \alpha(x - x_k)$ , and the new iterate is defined as the zero of this affine function, i.e.

$$x_{k+1} = x_k - \alpha^{-1}f(x_k) = F(x_k). \quad (5.12)$$

This is illustrated in Figure 5.1. By Proposition 5.5, a sufficient condition to ensure local convergence is that

$$|F'(x_*)| = |1 - \alpha^{-1}f'(x_*)| < 1. \quad (5.13)$$

In order for this condition to hold true, the slope  $\alpha$  must be of the same sign as  $f'(x_*)$  and the inequality  $|\alpha| \geq |f'(x_*)|/2$  must be satisfied. If  $f'(x_*) = 0$ , then the sufficient condition (5.13) is never satisfied; in this case, the convergence must be studied on a case-by-case basis. By Proposition 5.6, the convergence of the chord method is superlinear if  $\alpha = f'(x_*)$ . In practice, the solution  $x_*$  is unknown, and so this choice is not realistic. Nevertheless, the above reasoning suggests that, by letting the slope  $\alpha$  vary from iteration to iteration in such a manner that  $\alpha_k$  approaches  $f'(x_*)$  as  $k \rightarrow \infty$ , fast convergence can be obtained. This is precisely what the Newton–Raphson method aims to achieve; see Section 5.4.2

When  $\mathbf{f}$  is a function from  $\mathbf{R}^n$  to  $\mathbf{R}^n$ , the above approach generalizes to

$$x_{k+1} = \mathbf{F}(\mathbf{x}_k), \quad \mathbf{F}(\mathbf{x}) = \mathbf{x} - \mathbf{A}^{-1}\mathbf{f}(\mathbf{x}),$$

where  $\mathbf{A}$  is an invertible matrix. The geometric interpretation of the method in this case is the following: at each step, the function  $\mathbf{f}$  is approximated by the affine function  $\mathbf{x} \mapsto \mathbf{x}_k + \mathbf{A}(\mathbf{x} - \mathbf{x}_k)$ ,

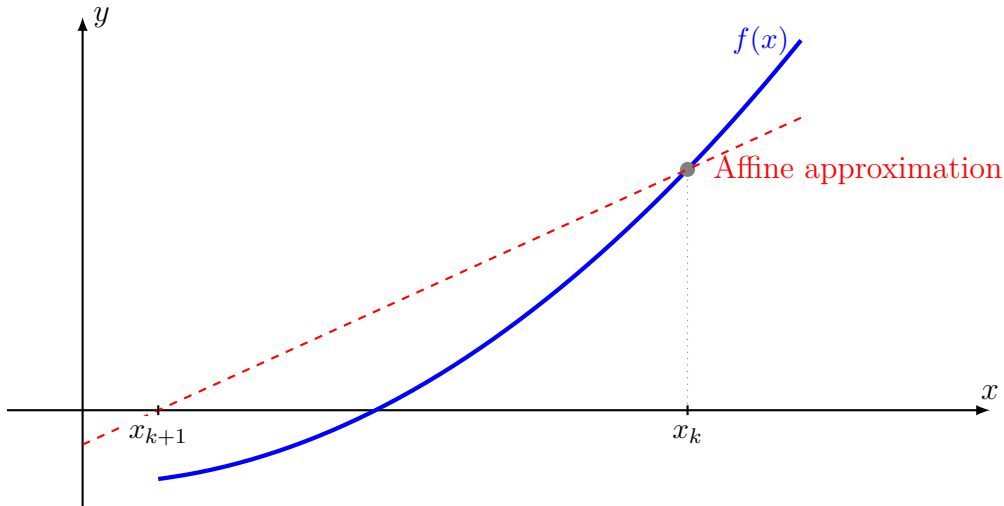


Figure 5.1: Graphical illustration of an iteration of the chord method.

and the next iterate is given by the unique zero of the latter function. Superlinear convergence is achieved when  $\mathbf{A} = \mathbf{J}_f(\mathbf{x}_*)$ . Notice that each iteration requires to calculate  $\mathbf{y} := \mathbf{A}^{-1}\mathbf{f}(\mathbf{x}_k)$ , which is generally achieved by solving the linear system  $\mathbf{A}\mathbf{y} = \mathbf{f}(\mathbf{x}_k)$ .

#### 5.4.2 The Newton–Raphson method

Let us first consider the case of a function from  $\mathbf{R}$  to  $\mathbf{R}$ . A necessary condition for the Newton–Raphson method to apply is that  $f$  is differentiable. At each step, the function  $f$  is approximated by the affine function  $x \mapsto f(x_k) + f'(x_k)(x - x_k)$  and the unique zero of this function is returned. In other words, one iteration of the Newton–Raphson method reads

$$x_{k+1} = x_k - f'(x_k)^{-1}f(x_k). \quad (5.14)$$

For this iteration to be well-defined, it is necessary that  $f'(x_k) \neq 0$ . The Newton–Raphson method may be viewed as a variation on (5.12) where the slope  $\alpha$  is adapted as the simulation progresses. If the method converges and  $f'$  is continuous, then  $f'(x_k) \rightarrow f'(x_*)$  in the limit as  $k \rightarrow \infty$ , which is an indication that superlinear convergence could occur in view of our discussion in the previous section. Equation (5.14) may be recast as a fixed point iteration of the form (5.4) with

$$F(x) = x - \frac{f(x)}{f'(x)}.$$

If  $x_*$  is a simple root of  $f$ , that is if  $f(x_*) = 0$  and  $f'(x_*) \neq 0$ , then  $x_*$  is a fixed point of the function  $F$ . If the function  $f$  is twice continuously differentiable, then the convergence of the Newton–Raphson method is superlinear by Proposition 5.6, because then

$$F'(x_*) = \frac{f(x_*)f''(x_*)}{f'(x_*)^2} = 0.$$

The geometric interpretation of the Newton–Raphson method in dimension 1 is the following: at each step, the function  $\mathbf{f}$  is approximated by the affine function  $x \mapsto x_k + f'(x_k)(x - x_k)$ ,

which is *the tangent line to  $f$  at  $x_k$* , and the next iterate is given by the unique zero of the latter function. This is illustrated in Figure 5.2.

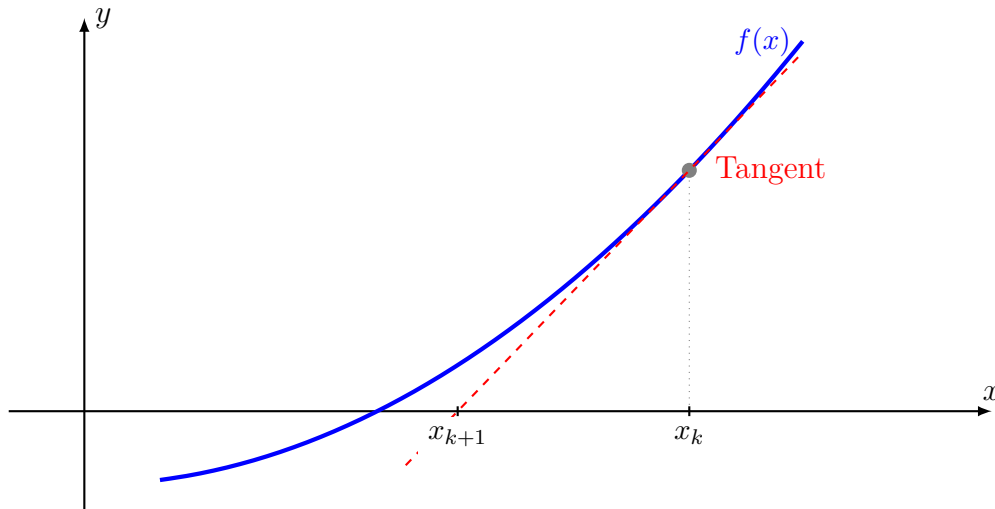


Figure 5.2: Graphical illustration of a Newton–Raphson iteration. The code used to create this figure is based on the answer <https://tex.stackexchange.com/a/551205/125558> on L<sup>A</sup>T<sub>E</sub>X stack exchange.

The Newton–Raphson method may be generalized to nonlinear equations in  $\mathbf{R}^n$  of the form (5.1). In this case  $\mathbf{F}(\mathbf{x}) = \mathbf{x} - \mathbf{J}_f(\mathbf{x})^{-1}\mathbf{f}(\mathbf{x})$ , and so an iteration of the method reads

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{J}_f(\mathbf{x}_k)^{-1}\mathbf{f}(\mathbf{x}_k). \quad (5.15)$$

In the rest of this section, we show that the iteration (5.15) is well-defined in a small neighborhood of a root of  $\mathbf{f}$  under appropriate assumptions, and we demonstrate the *second order* convergence of the method, first in dimension 1 under simplifying assumption involving the second derivative of  $f$ , and then in the multi-dimensional setting under more general assumptions.

### Convergence in the one-dimensional setting

We assume in this section that  $(x_k)_{k \geq 0}$  is generated from the Newton–Raphson method (5.14) and prove the following result.

**Theorem 5.7** (Quadratic convergence of Newton–Raphson). *Assume that  $f \in C^2(\mathbf{R})$  and that the following assumptions are satisfied:*

- *The first derivative of  $f$  is uniformly bounded away from zero:*

$$\inf_{x \in \mathbf{R}} |f'(x)| = m > 0.$$

- *The second derivative of  $f$  is uniformly bounded from above in absolute value:*

$$\sup_{x \in \mathbf{R}} |f''(x)| = M < \infty.$$

Then  $f(x)$  has a unique root  $x_*$  and it holds for all initial  $x_0 \in \mathbf{R}$  that

$$\forall k \in \mathbf{N}, \quad |x_{k+1} - x_*| \leq \frac{M}{2m} |x_k - x_*|^2. \quad (5.16)$$

*Proof.* By assumption, the function  $f$  is continuous and either strictly increasing everywhere or strictly decreasing everywhere. Therefore there exists a unique root  $x_* \in \mathbf{R}$  of  $f$ . In order to prove (5.16), we note that

$$x_{k+1} - x_* = x_k - \frac{f(x_k)}{f'(x_k)} - x_* = \frac{1}{f'(x_k)} \left( f'(x_k)(x_k - x_*) - f(x_k) \right). \quad (5.17)$$

By Taylor's theorem, there is  $\xi \in \mathbf{R}$  such that

$$f(x_*) = f(x_k) + f'(x_k)(x_* - x_k) + \frac{1}{2}f''(\xi)(x_* - x_k)^2.$$

Since  $x_*$  is a root of  $f$ , the left-hand side of this equation is zero. Combining this equation with (5.17), we deduce that

$$x_{k+1} - x_* = \frac{f''(\xi)(x_k - x_*)^2}{2f'(x_k)}.$$

Taking absolute values and using the assumptions gives

$$|x_{k+1} - x_*| \leq \frac{M}{2m} (x_k - x_*)^2,$$

which concludes the proof.  $\square$

*Remark 5.3.* As a corollary of [Theorem 5.7](#), we obtain that the Newton–Raphson method is convergent if

$$|x_k - x_*| \leq \frac{2m}{M}.$$

## Convergence in the multi-dimensional setting

As a first step towards a proof of quadratic convergence for the Newton–Raphson method in the multi-dimensional setting, we begin by proving the following preparatory lemma, which we will then employ in the particular case where the matrix-valued function  $\mathbf{A}$  is equal to  $\mathbf{J}_f$ .

**Lemma 5.8.** *Let  $\mathbf{A}: \mathbf{R}^n \rightarrow \mathbf{R}^{n \times n}$  denote a matrix-valued function on  $\mathbf{R}^n$  that is both continuous and nonsingular at  $\mathbf{x}_*$ , and let  $\mathbf{f}$  be a function that is differentiable at  $\mathbf{x}_*$  where  $\mathbf{f}(\mathbf{x}_*) = 0$ . Then the function*

$$\mathbf{G}(\mathbf{x}) = \mathbf{x} - \mathbf{A}(\mathbf{x})^{-1} \mathbf{f}(\mathbf{x})$$

*is well-defined in a neighborhood  $B_\delta(\mathbf{x}_*)$  of  $\mathbf{x}_*$ . In addition,  $\mathbf{G}$  is differentiable at  $\mathbf{x}_*$  with*

$$\mathbf{J}_G(\mathbf{x}_*) = \mathbf{I} - \mathbf{A}(\mathbf{x}_*)^{-1} \mathbf{J}_f(\mathbf{x}_*). \quad (5.18)$$



*Proof.* It holds that

$$\mathbf{A}(\mathbf{x}) = \left( \mathbf{A}(\mathbf{x}_*) - (\mathbf{A}(\mathbf{x}_*) - \mathbf{A}(\mathbf{x})) \right) = \mathbf{A}(\mathbf{x}_*) \left( \mathbf{I} - \mathbf{A}(\mathbf{x}_*)^{-1} (\mathbf{A}(\mathbf{x}_*) - \mathbf{A}(\mathbf{x})) \right). \quad (5.19)$$

Let  $\beta = \|\mathbf{A}(\mathbf{x}_*)^{-1}\|$  and  $\varepsilon = (2\beta)^{-1}$ . By continuity of the matrix-valued function  $\mathbf{A}$ , there exists  $\delta > 0$  such that

$$\forall \mathbf{x} \in B_\delta(\mathbf{x}_*), \quad \|\mathbf{A}(\mathbf{x}) - \mathbf{A}(\mathbf{x}_*)\| \leq \varepsilon.$$

For  $\mathbf{x} \in B_\delta(\mathbf{x}_*)$  we have  $\|\mathbf{A}(\mathbf{x}_*)^{-1}(\mathbf{A}(\mathbf{x}_*) - \mathbf{A}(\mathbf{x}))\| \leq \|\mathbf{A}(\mathbf{x}_*)^{-1}\| \|\mathbf{A}(\mathbf{x}_*) - \mathbf{A}(\mathbf{x})\| \leq \beta\varepsilon = \frac{1}{2}$ , and so Lemma 4.2 implies that the second factor on the right-hand side of (5.19) is invertible with a norm bounded from above by 2. Therefore, we deduce that  $\mathbf{A}(\mathbf{x})$  is invertible with

$$\forall \mathbf{x} \in B_\delta(\mathbf{x}_*), \quad \|\mathbf{A}(\mathbf{x})^{-1}\| \leq 2\|\mathbf{A}(\mathbf{x}_*)^{-1}\| = 2\beta, \quad (5.20)$$

which shows that  $\mathbf{G}$  is well-defined in  $B_\delta(\mathbf{x}_*)$ . In order to prove (5.18), we need to show that

$$\lim_{\|\mathbf{h}\| \rightarrow 0} \frac{\|\mathbf{G}(\mathbf{x}_* + \mathbf{h}) - \mathbf{G}(\mathbf{x}_*) - (\mathbf{I} - \mathbf{A}(\mathbf{x}_*)^{-1} \mathbf{J}_f(\mathbf{x}_*)) \mathbf{h}\|}{\|\mathbf{h}\|} = 0$$

By definition of  $\mathbf{G}$ , and using the fact that  $\mathbf{f}(\mathbf{x}_*) = \mathbf{0}$ , we obtain that the argument of the norm in the numerator is equal to

$$\begin{aligned} & \mathbf{A}(\mathbf{x}_*)^{-1} \mathbf{f}(\mathbf{x}_*) - \mathbf{A}(\mathbf{x}_* + \mathbf{h})^{-1} \mathbf{f}(\mathbf{x}_* + \mathbf{h}) + \mathbf{A}(\mathbf{x}_*)^{-1} \mathbf{J}_f(\mathbf{x}_*) \mathbf{h} \\ &= \underbrace{(\mathbf{A}^{-1}(\mathbf{x}_*) - \mathbf{A}(\mathbf{x}_* + \mathbf{h})^{-1}) \mathbf{J}_f(\mathbf{x}_*) \mathbf{h}}_{=: \mathbf{v}_1} - \underbrace{\mathbf{A}(\mathbf{x}_* + \mathbf{h})^{-1} (\mathbf{f}(\mathbf{x}_* + \mathbf{h}) - \mathbf{f}(\mathbf{x}_*) - \mathbf{J}_f(\mathbf{x}_*) \mathbf{h})}_{=: \mathbf{v}_2}. \end{aligned}$$

Noting that  $\mathbf{A}^{-1}(\mathbf{x}_*) - \mathbf{A}(\mathbf{x}_* + \mathbf{h})^{-1} = \mathbf{A}(\mathbf{x}_*)^{-1} (\mathbf{A}(\mathbf{x}_* + \mathbf{h}) - \mathbf{A}(\mathbf{x}_*)) \mathbf{A}(\mathbf{x}_* + \mathbf{h})^{-1}$ , we bound the norm of the first term on the right-hand side as follows:

$$\forall \mathbf{h} \in B_\delta(\mathbf{0}), \quad \|\mathbf{v}_1\| \leq 2\beta^2 \|\mathbf{A}(\mathbf{x}_* + \mathbf{h}) - \mathbf{A}(\mathbf{x}_*)\| \|\mathbf{J}_f(\mathbf{x}_*)\| \|\mathbf{h}\|.$$

Clearly  $\|\mathbf{v}_1\|/\|\mathbf{h}\| \rightarrow 0$  is the limit as  $\mathbf{h} \rightarrow \mathbf{0}$  by continuity of the matrix function  $\mathbf{A}$ . It also holds that  $\|\mathbf{v}_2\|/\|\mathbf{h}\| \rightarrow 0$  by differentiability of  $\mathbf{f}$  at  $\mathbf{x}_*$ , which concludes the proof.  $\square$

Using this lemma, we can show the following result on the convergence of the multi-dimensional Newton–Raphson method.

**Theorem 5.9** (Convergence of Newton–Raphson). *Let  $\mathbf{f}: \mathbf{R}^n \rightarrow \mathbf{R}^n$  denote a function that is differentiable in a neighborhood  $B_\delta(\mathbf{x}_*)$  of a point  $\mathbf{x}_*$  where  $\mathbf{f}(\mathbf{x}_*) = \mathbf{0}$ . Assume that the Jacobian matrix  $\mathbf{J}_f(\mathbf{x})$  is nonsingular and continuous at  $\mathbf{x}_*$ . Then  $\mathbf{x}_*$  is an attractor of the Newton–Raphson iteration (5.15) and the convergence is superlinear.*

*In addition, if there is  $\alpha > 0$  such that the Lipschitz condition*

$$\forall \mathbf{x} \in B_\delta(\mathbf{x}_*), \quad \|\mathbf{J}_f(\mathbf{x}) - \mathbf{J}_f(\mathbf{x}_*)\| \leq \alpha \|\mathbf{x} - \mathbf{x}_*\|$$

is satisfied, there exists  $d \in (0, \delta)$  and  $C > 0$  such that

$$\forall \mathbf{x}_k \in B_d(\mathbf{x}_*), \quad \|\mathbf{x}_{k+1} - \mathbf{x}_*\| \leq C \|\mathbf{x}_k - \mathbf{x}_*\|^2.$$

In other words, the convergence is at least quadratic in  $B_d(\mathbf{x}_*)$ .

*Proof.* Using Lemma 5.8, we obtain that the Newton–Raphson update

$$\mathbf{F}(\mathbf{x}) = \mathbf{x} - \mathbf{J}_f(\mathbf{x})^{-1} \mathbf{f}(\mathbf{x}),$$

is well-defined in a neighborhood  $B_\delta(\mathbf{x}_*)$  of  $\mathbf{x}_*$  for sufficiently small  $\delta$ . In addition, the second statement in Lemma 5.8 gives that  $\mathbf{J}_F(\mathbf{x}_*)^{-1} = \mathbf{I} - \mathbf{J}_F(\mathbf{x}_*)^{-1} \mathbf{J}_F(\mathbf{x}_*) = \mathbf{0}$ , which establishes the superlinear convergence by Proposition 5.6.

In order to show that the convergence is quadratic, we begin by noticing that, since

$$\mathbf{f}(\mathbf{x}_k) = \int_0^1 \frac{d}{dt} \mathbf{f}(\mathbf{x}_* + t(\mathbf{x}_k - \mathbf{x}_*)) dt = \int_0^1 \mathbf{J}_f(\mathbf{x}_* + t(\mathbf{x}_k - \mathbf{x}_*)) (\mathbf{x}_k - \mathbf{x}_*) dt,$$

it holds for all  $\mathbf{x}_k \in B_\delta(\mathbf{x}_*)$  that

$$\begin{aligned} \|\mathbf{f}(\mathbf{x}_k) - \mathbf{J}_f(\mathbf{x}_*)(\mathbf{x}_k - \mathbf{x}_*)\| &= \left\| \int_0^1 \left( \mathbf{J}_f(\mathbf{x}_* + t(\mathbf{x}_k - \mathbf{x}_*)) - \mathbf{J}_f(\mathbf{x}_*) \right) (\mathbf{x}_k - \mathbf{x}_*) dt \right\| \\ &\leq \int_0^1 \|\mathbf{J}_f(\mathbf{x}_* + t(\mathbf{x}_k - \mathbf{x}_*)) - \mathbf{J}_f(\mathbf{x}_*)\| \|\mathbf{x}_k - \mathbf{x}_*\| dt \\ &\leq \int_0^1 \alpha t \|\mathbf{x}_k - \mathbf{x}_*\|^2 dt \leq \frac{\alpha}{2} \|\mathbf{x}_k - \mathbf{x}_*\|^2. \end{aligned} \quad (5.21)$$

Let  $d \in (0, \delta)$  be sufficiently small to ensure that

$$\forall \mathbf{x} \in B_d(\mathbf{x}_*), \quad \|\mathbf{J}_f(\mathbf{x})^{-1}\| \leq 2\|\mathbf{J}_f(\mathbf{x}_*)^{-1}\|.$$

There exists such a  $d$  by (5.20). Using the inequality (5.21), we have that for all  $\mathbf{x}_k \in B_d(\mathbf{x}_*)$ ,

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_*\| &= \|\mathbf{F}(\mathbf{x}_k) - \mathbf{x}_*\| = \|\mathbf{x}_k - \mathbf{x}_* - \mathbf{J}_f(\mathbf{x}_k)^{-1} \mathbf{f}(\mathbf{x}_k)\| \\ &= \|\mathbf{J}_f(\mathbf{x}_k)^{-1} (\mathbf{f}(\mathbf{x}_k) - \mathbf{J}_f(\mathbf{x}_k)(\mathbf{x}_k - \mathbf{x}_*))\| \leq \|\mathbf{J}_f(\mathbf{x}_k)^{-1}\| \|\mathbf{f}(\mathbf{x}_k) - \mathbf{J}_f(\mathbf{x}_k)(\mathbf{x}_k - \mathbf{x}_*)\| \\ &\leq \|\mathbf{J}_f(\mathbf{x}_k)^{-1}\| \left( \|\mathbf{f}(\mathbf{x}_k) - \mathbf{J}_f(\mathbf{x}_*)(\mathbf{x}_k - \mathbf{x}_*)\| + \|\mathbf{J}_f(\mathbf{x}_*) - \mathbf{J}_f(\mathbf{x}_k)\| \|\mathbf{x}_k - \mathbf{x}_*\| \right) \\ &\leq \frac{3\alpha}{2} \|\mathbf{J}_f(\mathbf{x}_k)^{-1}\| \|\mathbf{x}_k - \mathbf{x}_*\|^2 \leq 3\alpha \|\mathbf{J}_f(\mathbf{x}_*)^{-1}\| \|\mathbf{x}_k - \mathbf{x}_*\|^2, \end{aligned}$$

which concludes the proof.  $\square$

### 5.4.3 The secant method

The Newton–Raphson method exhibits very fast convergence, but it requires the knowledge of the derivatives of the function  $\mathbf{f}$ . To conclude this chapter, we describe a root-finding algorithm, known as the secant method, that enjoys superlinear convergence but does not require the derivatives of  $\mathbf{f}$ . This method applies only when  $\mathbf{f}$  is a function from  $\mathbf{R}$  to  $\mathbf{R}$ , and so we drop

the vector notation in the rest of this section.

Unlike the other methods presented so far in Section 5.2, the secant method *can not* be recast as a fixed point iteration of the form  $x_{k+1} = F(x_k)$ . Instead, it is of the more general form  $x_{k+2} = F(x_k, x_{k+1})$ . The geometric intuition behind the method is the following: given  $x_k$  and  $x_{k+1}$ , the function  $f$  is approximated by the unique linear function that passes through  $(x_k, f(x_k))$  and  $(x_{k+1}, f(x_{k+1}))$ , and the iterate  $x_{k+2}$  is defined as the root of this linear function. In other words,  $f$  is approximated as follows:

$$\tilde{f}(x) = f(x_k) + \frac{f(x_{k+1}) - f(x_k)}{x_{k+1} - x_k}(x - x_k).$$

Solving  $\tilde{f}(x) = 0$  gives the following expression for  $x_{k+2}$ :

$$x_{k+2} = \frac{f(x_{k+1})x_k - f(x_k)x_{k+1}}{f(x_{k+1}) - f(x_k)}, \quad (5.22)$$

Showing the convergence of the secant method rigorously under general assumptions is tedious, so in this course we restrict our attention to the case where  $f$  is a quadratic function. Extending the proof of convergence to a more general smooth function can be achieved by using a quadratic Taylor approximation of  $f$  around the root  $x_*$ , which is accurate in a close neighborhood of  $x_*$ .

**Theorem 5.10** (Convergence of the secant method). *Assume that  $f$  is a convex quadratic polynomial with a simple root at  $x_*$  and that the secant method converges:  $\lim_{k \rightarrow \infty} x_k = x_*$ . Then the order of convergence is given by the golden ratio*

$$\varphi = \frac{1 + \sqrt{5}}{2}.$$

*More precisely, there exists a positive real number  $y_\infty$  such that*

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x_*|}{|x_k - x_*|^\varphi} = y_\infty. \quad (5.23)$$

*Proof.* Equation (5.22) implies that

$$x_{k+2} - x_* = \frac{f(x_{k+1})(x_k - x_*) - f(x_k)(x_{k+1} - x_*)}{f(x_{k+1}) - f(x_k)}.$$

By assumption, the function  $f$  may be expressed as

$$f(x) = \lambda(x - x_*) + \mu(x - x_*)^2, \quad \lambda \neq 0.$$

Substituting this expression in (5.4.3) and letting  $e_k = x_k - x_*$ , we obtain

$$e_{k+2} = \frac{\mu e_k e_{k+1} (e_{k+1} - e_k)}{\lambda(e_{k+1} - e_k) + \mu(e_{k+1}^2 - e_k^2)} = \frac{\mu e_k e_{k+1}}{\lambda + \mu(e_{k+1} + e_k)}.$$

Rearranging this equation, we have

$$\frac{e_{k+2}}{e_{k+1}} = \frac{\mu e_k}{\lambda + \mu(e_{k+1} + e_k)}. \quad (5.24)$$

By assumption, the right-hand side converges to zero, and so the left-hand side must also converge to zero; the convergence is superlinear.

To conclude the proof, we first reason formally in order to guess the order convergence, and then give a rigorous proof that our guess is correct. If  $e_k$  is small, then it holds approximately by (5.24) that

$$\frac{e_{k+2}}{e_{k+1}} \approx \mu e_k. \quad (5.25)$$

Assume that there exists  $q > 0$  such that the equation  $e_{k+1} = C e_k^q$  is valid for all  $k$ . Then it holds that  $e_{k+2} = C e_{k+1}^q = C(C e_k^q)^q$  and (5.25) enables to determine  $q$ :

$$\frac{C(C e_k^q)^q}{C e_k^q} = \frac{\mu}{\lambda} e_k \quad \Rightarrow \quad C^q e_k^{q^2 - q} = \frac{\mu}{\lambda} e_k \quad \Rightarrow \quad q^2 - q - 1 = 0. \quad \Rightarrow \quad q = \varphi.$$

Now comes the rigorous justification. Take absolute values in (5.24) to obtain, after rearranging,

$$\frac{|e_{k+2}|}{|e_{k+1}|^\varphi} = \left( \frac{|e_{k+1}|}{|e_k|^{\frac{1}{\varphi-1}}} \right)^{1-\varphi} \frac{\mu}{|\lambda + \mu(e_{k+1} + e_k)|} = \left( \frac{|e_{k+1}|}{|e_k|^\varphi} \right)^{1-\varphi} \frac{|\mu|}{|\lambda + \mu(e_{k+1} + e_k)|},$$

where we used that  $\varphi = \frac{1}{\varphi-1}$ , since  $\varphi$  is a root of the equation  $\varphi^2 - \varphi - 1 = 0$ . Thus, introducing the ratio  $y_k = |e_{k+1}|/|e_k|^\varphi$ , we have

$$y_{k+1} = y_k^{1-\varphi} \frac{|\mu|}{|\lambda + \mu(e_{k+1} + e_k)|}.$$

Taking logarithms in this equation, we deduce

$$\log(y_{k+1}) = (1 - \varphi) \log(y_k) + c_k, \quad c_k := \log \left( \frac{|\mu|}{|\lambda + \mu(e_{k+1} + e_k)|} \right).$$

This is a recurrence equation for  $\log(y_k)$ , whose explicit solution can be obtained from the variation-of-constants formula:

$$\log(y_k) = (1 - \varphi)^{k-1} \log(y_1) + \sum_{i=1}^{k-1} (1 - \varphi)^{k-1-i} c_i.$$

Since  $(c_k)_{k \geq 0}$  converges to the constant  $c_\infty = \log|\mu/\lambda|$  by the assumption that  $e_k \rightarrow 0$ , the sequence  $(\log(y_k))_{k \geq 0}$  converges to  $c_\infty/\varphi$  (prove this!). Therefore, by continuity of the exponential function, it holds that

$$y_k = \exp(\log(y_k)) \xrightarrow{k \rightarrow \infty} \exp \left( \frac{c_\infty}{\varphi} \right) = \left| \frac{\mu}{\lambda} \right|^{\frac{1}{\varphi}}$$

and so we deduce (5.23). □

## 5.5 A numerical experiment

To conclude this chapter, we present the results of a numerical experiment. Specifically, we consider four different fixed point methods for calculating the square root of 2, i.e. for solving the nonlinear equation

$$f(x) := x^2 - 2 = 0. \quad (5.26)$$

The unique positive solution to this equation is  $x_* = \sqrt{2}$ . The methods we consider are the following:

- The chord method with large  $\alpha = 10$ .
- The chord method with the optimal parameter  $\alpha$ , which is such that  $F'(x_*) = 0$ . The optimum value for  $\alpha$  for solving (5.26) is given by  $\alpha_* = 2\sqrt{2}$ .
- The Newton–Raphson method, where each iteration is given by

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} = x_k - \frac{x_k^2 - 2}{2x_k} = \frac{1}{2} \left( x_k + \frac{2}{x_k} \right) =: F(x_k),$$

with

$$F(x) = \frac{1}{2} \left( x + \frac{2}{x} \right).$$

Notice that  $F'(x_*) = 0$  and that  $F \in C^2((0, \infty))$ . Therefore, by Taylor's theorem it holds for all  $x \in (x_* - 1, x_* + 1)$  that

$$|F(x) - F(x_*)| = |F''(\xi(x))| \leq L(x - x_*)^2, \quad L := \sup_{|x-x_*| \leq 1} |F''(x)|.$$

We deduce that the convergence is at least quadratic by (5.11).

*Remark 5.4.* Note that the ancient Babylonian method coincides with the Newton–Raphson method applied to (5.26).

The following code implements these methods. Note that we use the arbitrary precision **BigFloat** format with a precision we manually set to 2000 bits, which enables using a very small  $\varepsilon$  in the stopping criterion.

```
function count_digits(x, y)
    xdigits = split(string(x), "")
    ydigits = split(string(y), "")
    len = min(length(xdigits), length(ydigits))
    for i in 1:len
        xdigits[i] != ydigits[i] && return i-2
    end
end

function my_sqrt(a)
```

```

exact = sqrt(a)
f(x) = x*x - a
fp(x) = 2x

# Uncomment desired line
F(x) = x - f(x)/10      # Chord method
# F(x) = x - f(x)/(2√a) # Chord method with optimal α
# F(x) = 1/2 * (x + a/x) # Babylonian / Newton Raphson

r, ε = 1, 1e-200
while abs(f(r)) > ε
    r = F(r)
    digits = ceil(Int, -log10(abs(r - exact)))
    println("Number of correct digits: $digits")
end
end

# Sets the precision of BigFloats to 1000 bits
setprecision(2000)
my_sqrt(BigFloat(2))

```

For each of the methods, the number of correct digits of the approximation as the iterations progress is illustrated in Table 5.1. Observe that for all the methods except the first one, the number of correct digits is approximately doubled at each iteration, which is consistent with quadratic convergence.

Method	Chord $\alpha = 10$	Chord $\alpha = 2\sqrt{2}$	Newton–Raphson
# Iterations	<b>1357</b>	<b>8</b>	<b>9</b>
# Correct digits $i = 1$	1	1	1
# Correct digits $i = 2$	1	3	3
# Correct digits $i = 3$	1	6	6
# Correct digits $i = 4$	1	12	12
# Correct digits $i = 5$	1	26	24
# Correct digits $i = 6$	1	53	48
# Correct digits $i = 7$	1	107	97
# Correct digits $i = 8$	1	214	196
# Correct digits $i = 9$	1	n/a	392

Table 5.1: Comparison of different fixed point methods for calculating  $\sqrt{2}$ . Here  $i$  denotes the iteration index.

## 5.6 Exercises

□ **Exercise 5.1.** Implement the bisection method for finding the solution(s) to the equation

$$x = \cos(x).$$

⚙️ **Exercise 5.2.** Find a discrete-time dynamical system over  $\mathbf{R}$  of the form

$$x_{k+1} = F(x_k)$$

for which 0 is an attractor but is not stable.

**Hint:** Use a function  $F$  that is discontinuous.

⚙️ **Exercise 5.3.** Show that if  $\mathbf{x}_*$  is a globally exponentially stable fixed point of  $F$ , then  $F$  does not have any other fixed point:  $\mathbf{x}_*$  is the unique fixed point.

⚙️ **Exercise 5.4.** Prove [Theorem 5.3](#).

⚙️ **Exercise 5.5.** Let  $\mathbf{x}_*$  be a fixed point of (5.5). Show that if

$$\rho(\mathbf{J}_F(\mathbf{x}_*)) < 1,$$

then  $\mathbf{x}_*$  is locally exponentially stable. It is sufficient by [Proposition 5.5](#) to find a subordinate matrix norm such that  $\|\mathbf{J}_F(\mathbf{x}_*)\| < 1$ . In other words, this exercise amounts to showing that for any matrix  $\mathbf{A} \in \mathbf{R}^{n \times n}$  with  $\rho(\mathbf{A}) < 1$ , there exists a matrix norm such that  $\|\mathbf{A}\| < 1$ .

**Hint:** One may employ a matrix norm of the form  $\|\mathbf{A}\|_{\mathbf{T}} := \|\mathbf{T}^{-1}\mathbf{A}\mathbf{T}\|_2$ , which is a subordinate norm by [Exercise 4.10](#). The Jordan normal form is useful for constructing the matrix  $\mathbf{T}$ , and equation (4.24) is also useful.

*Solution.* Let  $\mathbf{J} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}$  denote the Jordan normal form of  $\mathbf{A}$ , and let

$$\mathbf{E}_\varepsilon = \begin{pmatrix} \varepsilon & & & \\ & \varepsilon^2 & & \\ & & \ddots & \\ & & & \varepsilon^n \end{pmatrix}$$

By [Eq. \(4.24\)](#), the matrix  $\mathbf{J}_\varepsilon := \mathbf{E}_\varepsilon^{-1}\mathbf{J}\mathbf{E}_\varepsilon$  coincides with  $\mathbf{J}$ , except that the first superdiagonal is multiplied by  $\varepsilon$ . Let  $\mathbf{D}$  denote the diagonal part of  $\mathbf{J}_\varepsilon$ . We have that

$$\|\mathbf{J}_\varepsilon - \mathbf{D}\|_2 = \sqrt{\lambda_{\max}(\mathbf{E}_\varepsilon^T \mathbf{E}_\varepsilon)}.$$

The matrix  $\mathbf{E}_\varepsilon^T \mathbf{E}_\varepsilon$  is diagonal with entries equal to either 0 or  $\varepsilon^2$ , and so  $\|\mathbf{J}_\varepsilon - \mathbf{D}\|_2 < \varepsilon$ . By the triangle inequality, we have

$$\|\mathbf{J}_\varepsilon\| \leq \|\mathbf{D}\| + \|\mathbf{J}_\varepsilon - \mathbf{D}\|_2 \leq \rho(\mathbf{A}) + \varepsilon. \quad (5.27)$$

Let  $\|\mathbf{A}\|_\varepsilon := \|\mathbf{E}_\varepsilon^{-1}\mathbf{P}^{-1}\mathbf{A}\mathbf{P}\mathbf{E}_\varepsilon\|$ . By (4.10) with  $\mathbf{T} = \mathbf{P}\mathbf{E}_\varepsilon$ , this is indeed a subordinate matrix norm.

By (5.27) and the assumption that  $\rho(\mathbf{A}) < 1$ , it is clear that  $\|\mathbf{A}\|_\varepsilon < 1$  provided that  $\varepsilon$  is sufficiently small.  $\triangle$

*Remark 5.5.* A corollary of Exercise 4.10 is that the spectral radius of a matrix  $\mathbf{A}$  is the infimum of  $\|\mathbf{A}\|$  over all subordinate matrix norms.

⚙ **Exercise 5.6.** Calculate  $x = \sqrt[3]{3 + \sqrt[3]{3 + \sqrt[3]{3 + \sqrt{\dots}}}}$  using the bisection method.

⚙ **Exercise 5.7.** Solve the equation  $f(x) = e^x - 2 = 0$  using a fixed point iteration of the form

$$x_{k+1} = F(x_k), \quad F(x) = x - \alpha^{-1}f(x).$$

Using your knowledge of the exact solution  $x_* = \log 2$ , write a sufficient condition on  $\alpha$  to guarantee that  $x_*$  is locally exponentially stable. Verify your findings numerically and plot, using a logarithmic scale for the y axis, the error in absolute value as a function of  $k$ .

⚙ **Exercise 5.8.** Implement the Newton–Raphson method for solving  $f(x) = e^x - 2 = 0$ , and plot the error in absolute value as a function of the iteration index  $k$ .

⚙ **Exercise 5.9.** Find the point  $(x, y)$  on the parabola  $y = x^2$  that is closest to the point  $(3, 1)$ .

⚙ **Exercise 5.10.** Consider the linear system

$$\begin{cases} y = (x - 1)^2 \\ x^2 + y^2 = 4 \end{cases}$$

By drawing these two constraints in the  $xy$  plane, find an approximation of the solution(s). Then calculate the solution(s) using a fixed-point method.

⚙ **Exercise 5.11.** Find solutions  $(\psi, \lambda)$ , with  $\lambda > 0$ , to the following eigenvalue problem:

$$\psi'' = -\lambda^2\psi, \quad \psi(0) = 0, \quad \psi'(1) = \psi(1).$$

⚙ **Exercise 5.12.** Suppose that we have  $n$  data points  $(x_i, y_i)$  of an unknown function  $y = f(x)$ . We wish to approximate  $f$  by a function of the form

$$\tilde{f}(x) = \frac{a}{b + x}$$

by minimizing the sum of squares

$$\sum_{i=1}^n |\tilde{f}(x_i) - y_i|^2.$$

Write a system of nonlinear equations that the minimizer  $(a, b)$  must satisfy, and solve this system using the Newton–Raphson method starting from  $(1, 1)$ . The data is given below:

$\mathbf{x} = [0.0; 0.1; 0.2; 0.3; 0.4; 0.5; 0.6; 0.7; 0.8; 0.9; 1.0]$   
 $\mathbf{y} = [0.6761488864859304; 0.6345697680852508; 0.6396283580587062; 0.6132010027973919;$   
 $0.5906142598705267; 0.5718728461471725; 0.5524549902830562; 0.538938885654085;$   
 $0.5373495476994958; 0.514904589752926; 0.49243437874655027]$



Plot the data points together with the function  $\tilde{f}$  over the interval  $[0, 1]$ . Your plot should look like Figure 5.3.

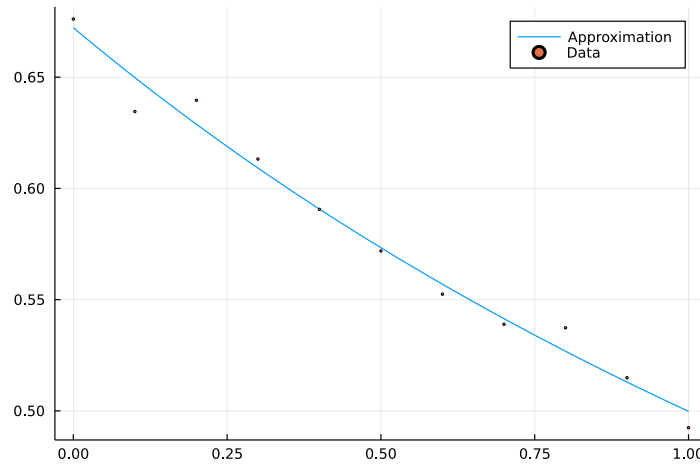


Figure 5.3: Solution to Exercise 5.12.

⚙️ **Exercise 5.13** (Nonlinear least-squares). Suppose that we are given  $n$  data points  $(x_i, y_i)$  of an unknown function  $y = f(x)$ . We wish to approximate  $f$  by a straight line

$$\tilde{f}(x) = ax + b$$

by minimizing the sum of squared Euclidean distances between the data points and the straight line  $\tilde{f}$ . Since the distance between a point  $(x_i, y_i)$  and the straight line is given by

$$\frac{|y_i - ax_i - b|}{\sqrt{1 + a^2}},$$

the objective function to minimize is given by

$$J(a, b) := \sum_{i=1}^n \frac{(y_i - ax_i - b)^2}{1 + a^2}.$$

This is a smooth function of  $a$  and  $b$ , and so a necessary condition for a pair  $(a_*, b_*) \in \mathbf{R}^2$  to be a minimizer is that

$$\nabla J(a_*, b_*) = \mathbf{0},$$

which is a nonlinear equation for the unknowns  $a_*$  and  $b_*$ . Solve this equation by using the Newton–Raphson method initialized at  $(1, 1)$ , and then plot the data points together with the function  $\tilde{f}$  over the interval  $[0, 1]$ . Your plot should look like Figure 5.4. The data is given hereafter:

```
x = [0.0; 0.1; 0.2; 0.3; 0.4; 0.5; 0.6; 0.7; 0.8; 0.9; 1.0]
y = [-0.9187980789440975; -0.6159791344678258; -0.25568734869121856;
     -0.14269370171581808; 0.3094396057228459; 0.6318327173549161;
     0.8370437988106428; 1.0970402798788812; 1.6057799131867696;
     1.869090784869698; 2.075369730726694]
```

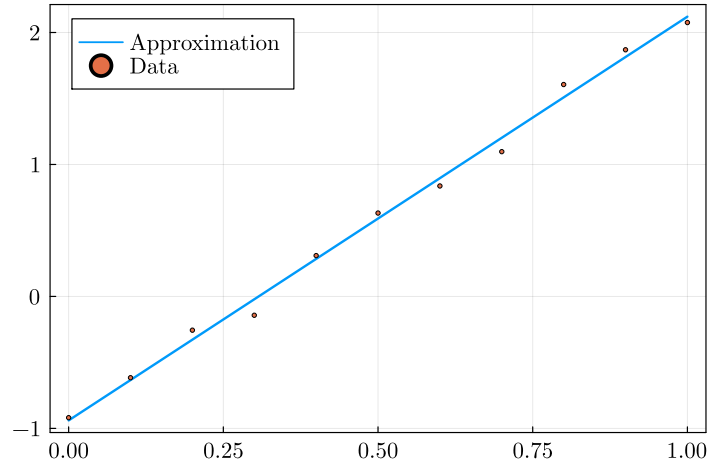


Figure 5.4: Solution to Exercise 5.13.

## 5.7 Discussion and bibliography

The content of this chapter is largely based on the lecture notes [15]. Several of the exercises are taken or inspired from [7]. The proof of convergence of the secant method is inspired from the general proof presented in the short paper [17]. For a detailed treatment of iterative methods for nonlinear equations, see the book [9].

# Chapter 6

## Numerical computation of eigenvalues

<b>6.1</b>	<b>Numerical methods for eigenvalue problems: general remarks . . .</b>	<b>149</b>
<b>6.2</b>	<b>Simple vector iterations . . . . .</b>	<b>149</b>
6.2.1	The power iteration . . . . .	150
6.2.2	Inverse iteration . . . . .	152
6.2.3	Rayleigh quotient iteration . . . . .	153
<b>6.3</b>	<b>Methods based on a subspace iteration . . . . .</b>	<b>153</b>
6.3.1	Simultaneous iteration . . . . .	153
6.3.2	The QR algorithm . . . . .	157
<b>6.4</b>	<b>Projection methods . . . . .</b>	<b>158</b>
6.4.1	Projection method in a Krylov subspace . . . . .	160
6.4.2	The Arnoldi iteration . . . . .	161
6.4.3	The Lanczos iteration . . . . .	163
<b>6.5</b>	<b>Exercises . . . . .</b>	<b>163</b>
<b>6.6</b>	<b>Discussion and bibliography . . . . .</b>	<b>168</b>

### Introduction

Calculating the eigenvalues and eigenvectors of a matrix is a task often encountered in scientific and engineering applications. Eigenvalue problems naturally arise in quantum physics, solid mechanics, structural engineering and molecular dynamics, to name just a few applications. The aim of this chapter is to present an overview of the standard methods for calculating eigenvalues and eigenvectors numerically. We focus predominantly on the case of a Hermitian matrix  $A \in \mathbf{C}^{n \times n}$ , which is technically simpler and arises in many applications. The reader is invited to go through the background material in [Appendix A.6](#) before reading this chapter. The rest of this chapter is organized as follows

- In [Section 6.1](#), we make general remarks concerning the calculation of eigenvalues.
- In [Section 6.2](#), we present standard methods based on a simple vector iteration.

- In [Section 6.3](#), we present a method for calculating several eigenvectors simultaneously, based on iterating a subspace.
- In [Section 6.4](#), we present method for constructing an approximation of the eigenvectors in a given subspace of  $\mathbf{C}^n$ .

## 6.1 Numerical methods for eigenvalue problems: general remarks

As mentioned in [Appendix A.6](#), a complex number  $\lambda \in \mathbf{C}$  is an eigenvalue of  $A \in \mathbf{C}^{n \times n}$  if and only if  $\lambda$  is a root of the characteristic polynomial  $p_A: \mathbf{C} \rightarrow \mathbf{C}$  of  $A$ , which is given by

$$p_A(\lambda) = \det(A - \lambda I).$$

One may, therefore, calculate the eigenvalues of  $A$  by calculating the roots of the polynomial  $p_A$  using, for example, one of the methods presented in [Chapter 5](#). While feasible for small matrices, this approach is not viable for large matrices, because the number of floating point operations required for calculating the coefficients of the characteristic polynomial scales as the factorial of  $n$ .

In view of the prohibitive computational cost required for calculating the characteristic polynomial, other methods are required for solving large eigenvalue problems numerically. All the methods that we study in this chapter are of iterative nature. While some of them are aimed at calculating all the eigenpairs of the matrix  $A$ , other methods enable to calculate only a small number of eigenpairs at a lower computational cost, which is often desirable. Indeed, calculating all the eigenvalues of a large matrix is computationally expensive; on a personal computer, the following Julia code takes well over a second to terminate:

```
import LinearAlgebra
A = rand(2000, 2000)
LinearAlgebra.eigen(A)
```

In many applications, the matrix  $A$  is sparse, and in this case it is important to use algorithms for eigenvalue problems that do not destroy the sparsity structure. Note that the eigenvectors of a sparse matrix are generally not sparse.

To conclude this section, we introduce some notation used throughout this chapter. For a diagonalizable matrix  $A$ , we denote the eigenvalues by  $\lambda_1, \dots, \lambda_n$ , with  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$ . The associated normalized eigenvectors are denoted by  $\mathbf{v}_1, \dots, \mathbf{v}_n$ . Therefore, it holds that

$$V^{-1}AV = D = \text{diag}(\lambda_1, \dots, \lambda_n), \quad \text{where } V = \begin{pmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_n \end{pmatrix}.$$

## 6.2 Simple vector iterations

In this section, we present simple iterative methods aimed at calculating just one eigenvector of the matrix  $A$ , which we assume to be diagonalizable for simplicity.

### 6.2.1 The power iteration

The power iteration is the simplest method for calculating the eigenpair associated with the eigenvalue of  $A$  with largest modulus. Since the eigenvectors of  $A$  span  $\mathbf{C}^n$ , any vector  $\mathbf{x}_0$  may be decomposed as

$$\mathbf{x}_0 = \alpha_1 \mathbf{v}_1 + \cdots + \alpha_n \mathbf{v}_n. \quad (6.1)$$

The idea of the power iteration is to repeatedly left-multiply this vector by the matrix  $A$ , in order to amplify the coefficient of  $\mathbf{v}_0$  relative to the other ones. Indeed, notice that

$$A^k \mathbf{x}_0 = \lambda_1^k \alpha_1 \mathbf{v}_1 + \cdots + \lambda_n^k \alpha_n \mathbf{v}_n.$$

If  $\lambda_1$  is strictly greater in modulus than the other eigenvalues, and if  $\alpha_1 \neq 0$ , then for large  $k$  the vector  $A^k \mathbf{x}_0$  is approximately aligned, in a sense made precise below, with the eigenvector  $\mathbf{v}_1$ . In order to avoid overflow errors at the numerical level, the iterates are normalized at each iteration. The power iteration is presented in [Algorithm 8](#).

---

**Algorithm 8** Power iteration

---

```

 $\mathbf{x} \leftarrow \mathbf{x}_0$ 
for  $i \in \{1, 2, \dots\}$  do
     $\mathbf{x} \leftarrow A\mathbf{x}$ 
     $\mathbf{x} \leftarrow \mathbf{x}/\|\mathbf{x}\|$ 
end for

```

---

To precisely quantify the convergence of the power method, we introduce the notion of *acute angle* between vectors of  $\mathbf{C}^n$ .

$$\begin{aligned} \angle(\mathbf{x}, \mathbf{y}) &= \arccos \left( \frac{|\mathbf{x}^* \mathbf{y}|}{\sqrt{\mathbf{x}^* \mathbf{x}} \sqrt{\mathbf{y}^* \mathbf{y}}} \right) \\ &= \arcsin \left( \frac{\|(I - P_{\mathbf{y}})\mathbf{x}\|}{\|\mathbf{x}\|} \right), \quad P_{\mathbf{y}} := \frac{\mathbf{y}\mathbf{y}^*}{\mathbf{y}^* \mathbf{y}}. \end{aligned}$$

This definition generalizes the familiar notion of angle for vectors in  $\mathbf{R}^2$  or  $\mathbf{R}^3$ , and we note that the angle function satisfies  $\angle(e^{i\theta_1} \mathbf{x}, e^{i\theta_2} \mathbf{y}) = \angle(\mathbf{x}, \mathbf{y})$  as well as  $\angle(\mathbf{x}, \mathbf{y}) \in [0, \pi/2]$ . We can then prove the following convergence result.

**Proposition 6.1** (Convergence of the power iteration). *Suppose that  $A$  is diagonalizable and that  $|\lambda_1| > |\lambda_2|$ . Then, for every initial guess with  $\alpha_1 \neq 0$ , the sequence  $(\mathbf{x}_k)_{k \geq 0}$  generated by the power iteration satisfies*

$$\lim_{k \rightarrow \infty} \angle(\mathbf{x}_k, \mathbf{v}_1) = 0.$$

*Proof.* By construction, it holds that

$$\mathbf{x}_k = \frac{\lambda_1^k \alpha_1 \mathbf{v}_1 + \cdots + \lambda_n^k \alpha_n \mathbf{v}_n}{\|\lambda_1^k \alpha_1 \mathbf{v}_1 + \cdots + \lambda_n^k \alpha_n \mathbf{v}_n\|} = e^{i\theta_k} \frac{\mathbf{v}_1 + \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \mathbf{v}_2 + \cdots + \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \mathbf{v}_n}{\left\| \mathbf{v}_1 + \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \mathbf{v}_2 + \cdots + \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \mathbf{v}_n \right\|}, \quad (6.2)$$

where

$$e^{i\theta_k} := \frac{\lambda_1^k \alpha_1}{|\lambda_1^k \alpha_1|}.$$

It follows from (6.2) that  $e^{-i\theta_k} \mathbf{x}_k \rightarrow \mathbf{v}_1 / \|\mathbf{v}_1\| = \mathbf{v}_1$  in the limit as  $k \rightarrow \infty$ , where we employed the fact that  $\|\mathbf{v}_1\| = 1$ . Using the definition of the angle between two vectors in  $\mathbf{C}^n$ , and the continuity with respect to either argument of the  $\mathbf{C}^n$  Euclidean inner product and of the arccos function, we obtain that

$$\begin{aligned} \angle(\mathbf{x}_k, \mathbf{v}_1) &= \arccos \left( \frac{|\mathbf{v}_1^* \mathbf{x}_k|}{\sqrt{\mathbf{v}_1^* \mathbf{v}_1} \sqrt{\mathbf{x}_k^* \mathbf{x}_k}} \right) = \arccos (|\mathbf{v}_1^* \mathbf{x}_k|) \\ &= \arccos \left( \left| \mathbf{v}_1^* \left( e^{-i\theta_k} \mathbf{x}_k \right) \right| \right) \xrightarrow[k \rightarrow \infty]{} \arccos(1) = 0, \end{aligned}$$

which concludes the proof.  $\square$

An inspection of the proof also reveals that the dominant term in the error, asymptotically in the limit as  $k \rightarrow \infty$ , is the one with coefficient  $\frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1}$ . Therefore, we deduce that

$$\angle(\mathbf{x}_k, \mathbf{v}_1) = \mathcal{O} \left( \left| \frac{\lambda_2}{\lambda_1} \right|^k \right).$$

The convergence is slow if  $|\lambda_2/\lambda_1|$  is close to one, and fast if  $|\lambda_2| \ll |\lambda_1|$ . Once an approximation of the eigenvector  $\mathbf{v}_1$  has been calculated, the corresponding eigenvalue  $\lambda_1$  can be estimated from the *Rayleigh quotient*:

$$\rho_{\mathbf{A}}: \mathbf{C}_*^n \rightarrow \mathbf{C}: \mathbf{x} \mapsto \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}}. \quad (6.3)$$

For any eigenvector  $\mathbf{v}$  of  $\mathbf{A}$ , the corresponding eigenvalue is equal to  $\rho_{\mathbf{A}}(\mathbf{v})$ . In order to study the error on the eigenvalue  $\lambda_1$  for the power iteration, we assume for simplicity that  $\mathbf{A}$  is Hermitian and that the eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$  are orthonormal. Substituting (6.2) in the Rayleigh quotient (6.3), we obtain

$$\rho_{\mathbf{A}}(\mathbf{x}_k) = \frac{\lambda_1 + \left| \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \right|^2 \lambda_2 + \dots + \left| \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \right|^2 \lambda_n}{1 + \left| \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \right|^2 + \dots + \left| \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \right|^2}.$$

Therefore, by reducing to a common denominator we deduce

$$\begin{aligned} |\rho_{\mathbf{A}}(\mathbf{x}_k) - \lambda_1| &= \left| \frac{\lambda_1 + \left| \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \right|^2 \lambda_2 + \dots + \left| \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \right|^2 \lambda_n}{1 + \left| \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \right|^2 + \dots + \left| \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \right|^2} - \lambda_1 \right| \\ &\leq \left| \frac{\lambda_2^k \alpha_2}{\lambda_1^k \alpha_1} \right|^2 |\lambda_2 - \lambda_1| + \dots + \left| \frac{\lambda_n^k \alpha_n}{\lambda_1^k \alpha_1} \right|^2 |\lambda_n - \lambda_1| = \mathcal{O} \left( \left| \frac{\lambda_2}{\lambda_1} \right|^{2k} \right). \end{aligned}$$

The convergence of the eigenvalue in the particular case of a Hermitian matrix is faster than for a general matrix in  $\mathbf{C}^{n \times n}$ . For general matrices, it is possible to show using a similar argument that the error is of order  $\mathcal{O}(|\lambda_2/\lambda_1|^k)$  in the limit as  $k \rightarrow \infty$ .

**Essential convergence.** It is useful at this point to introduce the concept of *essential convergence*. A sequence  $(\mathbf{x}_k)$  in  $\mathbf{C}^n$  is said to *converge essentially* to a vector  $\mathbf{x}_\infty$  if there exists a sequence of complex numbers  $(e^{i\phi_k})$  of modulus 1 such that the sequence  $(e^{i\phi_k}\mathbf{x}_k)$  converges to  $\mathbf{x}_\infty$ . For a sequence  $(\mathbf{x}_k)$  of normalized vectors, the essential convergence of  $(\mathbf{x}_k)$  to  $\mathbf{x}_\infty$  is equivalent to the convergence of  $\angle(\mathbf{x}_k, \mathbf{x}_\infty)$  to 0. Proving this equivalence is the goal of [Exercise 6.11](#). Reformulated in this new terminology, [Proposition 6.1](#) states that the sequence  $(\mathbf{x}_k)$  obtained from the power iteration converges essentially to  $\mathbf{v}_1$ .

### 6.2.2 Inverse iteration

The power iteration is simple but enables to calculate only the dominant eigenvalue of the matrix  $\mathbf{A}$ , i.e. the eigenvalue of largest modulus. In addition, the convergence of the method is slow when  $|\lambda_2| \approx |\lambda_1|$ .

The inverse iteration enables a more efficient calculation of not only the dominant eigenvalue but also the other eigenvalues of  $\mathbf{A}$ . It is based on applying the power iteration to  $(\mathbf{A} - \mu\mathbf{I})^{-1}$ , where  $\mu \in \mathbf{C}$  is a shift. The eigenvalues of  $(\mathbf{A} - \mu\mathbf{I})^{-1}$  are given by  $(\lambda_1 - \mu)^{-1}, \dots, (\lambda_n - \mu)^{-1}$ , with associated eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$ . If  $0 < |\lambda_J - \mu| < |\lambda_j - \mu|$  for all  $j \neq J$ , then the dominant eigenvalue of the matrix  $(\mathbf{A} - \mu\mathbf{I})^{-1}$  is  $(\lambda_J - \mu)^{-1}$ , and so the power iteration applied to this matrix yields an approximation of the eigenvector  $\mathbf{v}_J$ . In other words, the inverse iteration with shift  $\mu$  enables to calculate an approximation of the eigenvector of  $\mathbf{A}$  corresponding to the eigenvalue nearest  $\mu$ . The inverse iteration is presented in [Algorithm 9](#). In practice, the inverse matrix  $(\mathbf{A} - \mu\mathbf{I})^{-1}$  need not be calculated, and it is often preferable to solve a linear system at each iteration.

---

#### Algorithm 9 Inverse iteration

---

```

 $\mathbf{x} \leftarrow \mathbf{x}_0$ 
for  $i \in \{1, 2, \dots\}$  do
    Solve  $(\mathbf{A} - \mu\mathbf{I})\mathbf{y} = \mathbf{x}$ 
     $\mathbf{x} \leftarrow \mathbf{y}/\|\mathbf{y}\|$ 
end for
 $\lambda \leftarrow \mathbf{x}^*\mathbf{A}\mathbf{x}/\mathbf{x}^*\mathbf{x}$ 
return  $\mathbf{x}, \lambda$ 

```

---

An application of [Proposition 6.1](#) immediately gives the following convergence result for the inverse iteration.

**Proposition 6.2** (Convergence of the inverse iteration). *Assume that  $\mathbf{A} \in \mathbf{C}^n$  is diagonalizable and that there exist  $J$  and  $K$  such that*

$$0 < |\lambda_J - \mu| < |\lambda_K - \mu| \leq |\lambda_j - \mu| \quad \forall j \neq J.$$

*Assume also that  $\alpha_J \neq 0$ , where  $\alpha_J$  is the coefficient of  $\mathbf{v}_J$  in the expansion of  $\mathbf{x}_0$  given in (6.1). Then the iterates of the inverse iteration satisfy*

$$\lim_{k \rightarrow \infty} \angle(\mathbf{x}_k, \mathbf{v}_J) = 0.$$

More precisely,

$$\angle(\mathbf{x}_k, \mathbf{v}_J) = \mathcal{O}\left(\left|\frac{\lambda_J - \mu}{\lambda_K - \mu}\right|^k\right).$$

Proposition 6.2 states that  $\mathbf{x}_k$  converges essentially to  $\mathbf{v}_J$ . Notice that the closer  $\mu$  is to  $\lambda_J$ , the faster the inverse iteration converges. Note also that with  $\mu = 0$ , the inverse iteration enables to calculate the eigenvalue of  $\mathbf{A}$  of smallest modulus.

### 6.2.3 Rayleigh quotient iteration

Since the inverse iteration is fast when  $\mu$  is close to an eigenvalue  $\lambda_J$ , it is natural to wonder whether the method can be improved by progressively updating  $\mu$  as the simulation progresses. Specifically, an approximation of the eigenvalue associated with the current vector may be employed in place of  $\mu$ . This leads to the Rayleigh quotient iteration, presented in Algorithm 10.

---

#### Algorithm 10 Inverse iteration

---

```

 $\mathbf{x} \leftarrow \mathbf{x}_0$ 
for  $i \in \{1, 2, \dots\}$  do
   $\mu \leftarrow \mathbf{x}^* \mathbf{A} \mathbf{x} / \mathbf{x}^* \mathbf{x}$ 
  Solve  $(\mathbf{A} - \mu \mathbf{I}) \mathbf{y} = \mathbf{x}$ 
   $\mathbf{x} \leftarrow \mathbf{y} / \|\mathbf{y}\|$ 
end for
 $\lambda \leftarrow \mathbf{x}^* \mathbf{A} \mathbf{x} / \mathbf{x}^* \mathbf{x}$ 
return  $\mathbf{x}, \lambda$ 

```

---

It is possible to show that, when  $\mathbf{A}$  is Hermitian, the Rayleigh quotient iteration converges to an eigenvector for almost every initial guess  $\mathbf{x}_0$ . Furthermore, if convergence to an eigenvector occurs, then  $\mu$  converges cubically to the corresponding eigenvalue. See [12] and the references therein for more details.

## 6.3 Methods based on a subspace iteration

The subspace iteration resembles the power iteration but it is more general: not just one but several vectors are updated at each iteration.

### 6.3.1 Simultaneous iteration

Let  $\mathbf{X}_0 = (\mathbf{x}_1 \ \dots \ \mathbf{x}_p)$  denote an initial set of linearly independent vectors. Before we present the simultaneous iteration, we recall a statement concerning the QR decomposition of a matrix, which is related to the Gram–Schmidt orthonormalization process. We recall that the Gram–Schmidt method enables to construct, starting from an ordered set of vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_p\}$  in  $\mathbf{C}^n$ , a new set of vectors  $\{\mathbf{q}_1, \dots, \mathbf{q}_p\}$  which are *orthonormal* and span the same subspace of  $\mathbf{C}^n$  as the original vectors.



**Proposition 6.3** (Reduced QR decomposition). *Assume that  $\mathbf{X} \in \mathbf{C}^{n \times p}$  has linearly independent columns. Then there exist a matrix  $\mathbf{Q} \in \mathbf{C}^{n \times p}$  with orthonormal columns and an upper triangular matrix  $\mathbf{R} \in \mathbf{C}^{p \times p}$  such that the following factorization holds:*

$$\mathbf{X} = \mathbf{Q}\mathbf{R}. \quad (6.4)$$

*This decomposition is known as a reduced QR decomposition if  $p < n$ , or simply QR decomposition if  $p = n$ , in which case  $\mathbf{X}$  is a square matrix and  $\mathbf{Q}$  is a unitary matrix. The decomposition is unique if we require that the diagonal elements of  $\mathbf{R}$  are real and positive.*

*Proof.* The statement is clear when  $p = 1$ . Reasoning by induction, we assume that the result is true up to  $p - 1$ , and prove that it then also holds true for  $p$ . We wish to show that there is a unique matrix  $\mathbf{Q} \in \mathbf{C}^{n \times p}$  with orthonormal columns and a unique upper triangular matrix  $\mathbf{R} \in \mathbf{C}^{p \times p}$  with real and positive diagonal elements such that (6.4) is satisfied. To this end, we decompose the matrices  $\mathbf{Q}$  and  $\mathbf{R}$  as follows:

$$\mathbf{Q} = \begin{pmatrix} \mathbf{Q}_{p-1} & \mathbf{q} \end{pmatrix}, \quad \mathbf{R} = \begin{pmatrix} \mathbf{R}_{p-1} & \mathbf{r} \\ \mathbf{0}_{p-1}^T & r \end{pmatrix}. \quad (6.5)$$

Here  $\mathbf{Q}_{p-1} \in \mathbf{C}^{n \times (p-1)}$  is a matrix with orthonormal columns,  $\mathbf{R} \in \mathbf{C}^{(p-1) \times (p-1)}$  is an upper triangular matrix with positive real diagonal elements,  $\mathbf{q} \in \mathbf{C}^n$  is a normalized vector orthogonal to all the columns of  $\mathbf{Q}_{p-1}$ ,  $\mathbf{r} \in \mathbf{C}^{n-1}$  is a vector and  $r \in \mathbf{R}_{>0}$  is a scalar. Let us also denote by  $\mathbf{X}_{p-1} \in \mathbf{C}^{n \times (p-1)}$  the matrix containing the  $p - 1$  first columns of  $\mathbf{X}$ , and by  $\mathbf{x}_p \in \mathbf{C}^n$  the  $p$ -th column of  $\mathbf{X}$ . Substituting (6.5) into (6.4), we then obtain

$$\begin{pmatrix} \mathbf{X}_{p-1} & \mathbf{x}_p \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_{p-1}\mathbf{R}_{p-1} & \mathbf{Q}_{p-1}\mathbf{r} + \mathbf{q}r \end{pmatrix}, \quad (6.6)$$

By the induction hypothesis, there exist a unique choice of matrices  $\mathbf{Q}_{p-1}$  and  $\mathbf{R}_{p-1}$  with the required structure such that  $\mathbf{X}_{p-1} = \mathbf{Q}_{p-1}\mathbf{R}_{p-1}$ . Comparing the last column of both sides in (6.6), we obtain

$$\mathbf{x}_p = \mathbf{Q}_{p-1}\mathbf{r} + \mathbf{q}r. \quad (6.7)$$

Left-multiplying both sides by  $\mathbf{Q}_{p-1}^*$  and employing the orthogonality between  $\mathbf{q}$  and the columns of  $\mathbf{Q}_{p-1}$ , we deduce that necessarily  $\mathbf{r} = \mathbf{Q}_{p-1}^*\mathbf{x}_p$ . It then follows from (6.7) that

$$\mathbf{q} = \frac{1}{r} (\mathbf{x}_p - \mathbf{Q}_{p-1}\mathbf{Q}_{p-1}^*\mathbf{x}_p), \quad r = \|\mathbf{x}_p - \mathbf{Q}_{p-1}\mathbf{Q}_{p-1}^*\mathbf{x}_p\|.$$

It is simple to check that  $\mathbf{q}$  is indeed orthogonal to the columns of  $\mathbf{Q}$ , which concludes the proof. Note that  $\mathbf{Q}_{p-1}\mathbf{Q}_{p-1}^*\mathbf{x}_p$  is the orthogonal projection of  $\mathbf{x}_p$  onto the subspace spanned by the columns of  $\mathbf{Q}_{p-1}$ .  $\square$

Note that the columns of the matrix  $\mathbf{Q}$  of the decomposition coincide with the vectors that would be obtained by applying the Gram–Schmidt method to the columns of the matrix  $\mathbf{X}$ . In fact, the Gram–Schmidt process is one of several methods by which the QR decomposition can be calculated in practice.

**Algorithm 11** Simultaneous iteration

---

```

X ← X0
for k ∈ {1, 2, ...} do
  QkRk = AXk-1 (QR decomposition).
  Xk ← Qk.
end for

```

---

The simultaneous iteration method is presented in [Algorithm 11](#). Like the normalization in the power iteration [Algorithm 8](#), the QR decomposition performed at each step in [Algorithm 11](#) enables to avoid overflow errors. Notice that when  $p = 1$ , the simultaneous iteration reduces to the power iteration. We emphasize that the factorization step at each iteration does not influence the subspace spanned by the columns of  $X$ . Therefore, this subspace after  $k$  iterations coincides with that spanned by the columns of the matrix  $A^k X_0$ . In fact, in exact arithmetic, it would be equivalent to perform the QR decomposition only once as a final step, after the **for** loop. Indeed, denoting by  $Q_k R_k$  the QR decomposition of  $AX_{k-1}$ , we have

$$\begin{aligned} X_k &= AX_{k-1} R_k^{-1} = A^2 X_{k-2} R_{k-1}^{-1} R_k^{-1} = \dots = A^k X_0 R_1^{-1} \dots R_k^{-1} \\ &\Leftrightarrow X_k (R_k \dots R_1) = A^k X_0. \end{aligned}$$

Since  $X_k$  has orthonormal columns and  $R_k \dots R_1$  is an upper triangular matrix (see [Exercise 4.3](#)) with real positive elements on the diagonal (check this!), it follows that  $X_k$  can be obtained by QR factorization of  $A^k X_0$ . In order to show the convergence of the simultaneous iteration, we begin by proving the following preparatory lemma.

**Lemma 6.4** (Continuity of the reduced QR decomposition). *If  $Q_k R_k \rightarrow QR$ , where  $Q \in \mathbf{C}^{n \times p}$  has orthonormal columns and  $R \in \mathbf{C}^{p \times p}$  is upper triangular with positive real entries on the diagonal, then  $Q_k \rightarrow Q$ .*

*Proof.* We reason by contradiction and assume there is  $\varepsilon > 0$  and a subsequence  $(Q_{k_n})_{n \geq 0}$  such that  $\|Q_{k_n} - Q\| \geq \varepsilon$  for all  $n$ . Since the set of matrices with normalized columns is a compact subset of  $\mathbf{C}^{n \times p}$ , there exists a further subsequence  $(Q_{k_{n_m}})_{m \geq 0}$  that converges to a limit  $Q_\infty$  which also has orthonormal columns and is at least  $\varepsilon$  away in norm from  $Q$ . But then

$$R_{k_{n_m}} = Q_{k_{n_m}}^* (Q_{k_{n_m}} R_{k_{n_m}}) \xrightarrow{m \rightarrow \infty} Q_\infty^* (QR) =: R_\infty.$$

Since  $R_k$  is upper triangular with positive diagonal elements for all  $k$ , clearly  $R_\infty$  is also upper triangular with positive diagonal elements. But then  $Q_\infty R_\infty = QR$ , and by uniqueness of the decomposition we deduce that  $Q = Q_\infty$ , which is a contradiction.  $\square$

Before presenting the convergence theorem, we introduce the following terminology: we say that  $X_k \in \mathbf{C}^{n \times p}$  converges essentially to a matrix  $X_\infty$  if each column of  $X_k$  converges essentially to the corresponding column of  $X_\infty$ . We prove the convergence in the Hermitian case for simplicity. In the general case of  $A \in \mathbf{C}^{n \times n}$ , it cannot be expected that  $X_k$  converges essentially to  $V$ , because the columns of  $X_k$  are orthogonal but eigenvectors may not be orthogonal. In

this case, the columns of  $\mathbf{X}_k$  converge not to the eigenvectors but to the so-called Schur vectors of  $\mathbf{A}$ ; see [12] for more information.

**Theorem 6.5** (Convergence of the simultaneous iteration  $\textcircled{a}$ ). *Assume that  $\mathbf{A} \in \mathbf{C}^{n \times n}$  is Hermitian, that  $\mathbf{X}_0 \in \mathbf{C}^{n \times p}$  has linearly independent columns, and finally that the subspace spanned by the column of  $\mathbf{X}_0$  satisfies*

$$\text{col}(\mathbf{X}_0) \cap \text{Span}\{\mathbf{v}_{p+1}, \dots, \mathbf{v}_n\} = \emptyset. \quad (6.8)$$

If it holds that

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_p| > |\lambda_{p+1}| \geq |\lambda_{p+2}| \geq \dots \geq |\lambda_n|, \quad (6.9)$$

then  $\mathbf{X}_k$  converges essentially to  $\mathbf{V}_1 := (\mathbf{v}_1 \ \dots \ \mathbf{v}_p)$ .

*Proof.* Let  $\mathbf{B} = \mathbf{V}^{-1}\mathbf{X}_0 \in \mathbf{C}^{n \times p}$ , so that  $\mathbf{X}_0 = \mathbf{V}\mathbf{B}$ , and note that  $\mathbf{A}^k\mathbf{X}_0 = \mathbf{V}\mathbf{D}^k\mathbf{B}$ . We denote by  $\mathbf{B}_1 \in \mathbf{C}^{p \times p}$  and  $\mathbf{B}_2 \in \mathbf{C}^{(n-p) \times p}$  the upper  $p \times p$  and lower  $(n-p) \times p$  blocks of  $\mathbf{B}$ , respectively. The matrix  $\mathbf{B}_1$  is nonsingular, otherwise the assumption (6.8) would not hold. Indeed, if there was a nonzero vector  $\mathbf{z} \in \mathbf{C}^p$  such that  $\mathbf{B}_1\mathbf{z} = \mathbf{0}$ , then

$$\mathbf{X}_0\mathbf{z} = \mathbf{V} \begin{pmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{pmatrix} \mathbf{z} = \begin{pmatrix} \mathbf{V}_1 & \mathbf{V}_2 \end{pmatrix} \begin{pmatrix} \mathbf{0} \\ \mathbf{B}_2\mathbf{z} \end{pmatrix} = \mathbf{V}_2\mathbf{B}_2\mathbf{z}.$$

implying that  $\mathbf{X}_0\mathbf{z} \in \text{col}(\mathbf{X}_0)$  is a linear combination of the vectors in  $\mathbf{V}_2 = (\mathbf{v}_{p+1} \ \dots \ \mathbf{v}_n)$ , which contradicts the assumption. We also denote by  $\mathbf{D}_1$  and  $\mathbf{D}_2$  the  $p \times p$  upper-left and the  $(n-p) \times (n-p)$  lower-right blocks of  $\mathbf{D}$ , respectively. From the expression of  $\mathbf{A}^k\mathbf{X}_0$ , we have

$$\begin{aligned} \mathbf{A}^k\mathbf{X}_0 &= \begin{pmatrix} \mathbf{V}_1 & \mathbf{V}_2 \end{pmatrix} \begin{pmatrix} \mathbf{D}_1^k & \\ & \mathbf{D}_2^k \end{pmatrix} \begin{pmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{pmatrix} = \mathbf{V}_1\mathbf{D}_1^k\mathbf{B}_1 + \mathbf{V}_2\mathbf{D}_2^k\mathbf{B}_2, \\ &= \left( \mathbf{V}_1 + \mathbf{V}_2\mathbf{D}_2^k\mathbf{B}_2\mathbf{B}_1^{-1}\mathbf{D}_1^{-k} \right) \mathbf{D}_1^k\mathbf{B}_1. \end{aligned} \quad (6.10)$$

The second term in the bracket on the right-hand side converges to zero in the limit as  $k \rightarrow \infty$  by (6.9). Let  $\widetilde{\mathbf{Q}}_k\widetilde{\mathbf{R}}_k$  denote the reduced QR decomposition of the bracketed term. By Lemma 6.4, we deduce from  $\widetilde{\mathbf{Q}}_k\widetilde{\mathbf{R}}_k \rightarrow \mathbf{V}_1$  that  $\widetilde{\mathbf{Q}}_k \rightarrow \mathbf{V}_1$ . Rearranging (6.10), we have

$$\mathbf{A}^k\mathbf{X}_0 = \widetilde{\mathbf{Q}}_k(\widetilde{\mathbf{R}}_k\mathbf{D}_1^k\mathbf{B}_1).$$

Since the matrix between brackets is a  $p \times p$  square invertible matrix, this equation implies that  $\text{col}(\mathbf{A}^k\mathbf{X}_0) = \text{col}(\widetilde{\mathbf{Q}}_k)$ . Denoting by  $\mathbf{Q}_k\mathbf{R}_k$  the QR decomposition of  $\mathbf{A}_k\mathbf{X}_0$ , we therefore have  $\text{col}(\mathbf{Q}_k) = \text{col}(\widetilde{\mathbf{Q}}_k)$ , and so the projectors on these subspaces are equal. We recall that, for a set of orthonormal vectors  $\mathbf{y}_1, \dots, \mathbf{y}_p$  gathered in a matrix  $\mathbf{Y} = (\mathbf{y}_1 \ \dots \ \mathbf{y}_p)$ , the projector on  $\text{col}(\mathbf{Y}) = \text{Span}\{\mathbf{y}_1, \dots, \mathbf{y}_p\} \subset \mathbf{C}^n$  is the square  $n \times n$  matrix

$$\mathbf{Y}\mathbf{Y}^* = \mathbf{y}_1\mathbf{y}_1^* + \dots + \mathbf{y}_p\mathbf{y}_p^*.$$

Consequently, the equality of the projectors implies  $\mathbf{Q}_k\mathbf{Q}_k^* = \widetilde{\mathbf{Q}}_k\widetilde{\mathbf{Q}}_k^*$ , and so  $\mathbf{Q}_k\mathbf{Q}_k^* \rightarrow \mathbf{V}_1\mathbf{V}_1^*$ . Sim-

ilarly, noting that  $\mathbf{Q}_k[:, 1:i]\mathbf{R}[1:i, 1:i]$  is the QR decomposition of the matrix  $\mathbf{A}^k\mathbf{X}_0[:, 1:i]$  for all  $i \in \{1, \dots, p\}$ , we obtain the convergence

$$\forall i \in \{1, \dots, p\}, \quad \mathbf{Q}_k[:, 1:i]\mathbf{Q}_k[:, 1:i]^* \xrightarrow[k \rightarrow \infty]{} \mathbf{v}_1\mathbf{v}_1^* + \dots + \mathbf{v}_i\mathbf{v}_i^*. \quad (6.11)$$

This is not surprising given that the first  $k$  columns of  $\mathbf{X}_0$  undergo a simultaneous iteration independent of the other columns.

Next, we establish the essential convergence of  $\mathbf{Q}_k$  to  $\mathbf{V}_1$ . To this end, we denote the columns of  $\mathbf{Q}_k$  by  $\mathbf{q}_1^{(k)}, \dots, \mathbf{q}_p^{(k)}$  and first show by induction that  $\mathbf{q}_i^{(k)}\mathbf{q}_i^{(k)*} \rightarrow \mathbf{v}_i\mathbf{v}_i^*$  in the limit  $k \rightarrow \infty$ . For  $i = 1$  this follows from (6.11). Assume now that  $\mathbf{q}_\bullet^{(k)}\mathbf{q}_\bullet^{(k)*} \rightarrow \mathbf{v}_\bullet\mathbf{v}_\bullet^*$  for  $\bullet$  up to  $i - 1$ . Then

$$\begin{aligned} \mathbf{q}_i^{(k)}\mathbf{q}_i^{(k)*} &= \mathbf{Q}_k[:, 1:i]\mathbf{Q}_k[:, 1:i]^* - \mathbf{q}_1^{(k)}\mathbf{q}_1^{(k)*} - \dots - \mathbf{q}_{i-1}^{(k)}\mathbf{q}_{i-1}^{(k)*} \\ &\xrightarrow[k \rightarrow \infty]{} \mathbf{V}_1[:, 1:i]\mathbf{V}_1[:, 1:i]^* - \mathbf{v}_1\mathbf{v}_1^* - \dots - \mathbf{v}_{i-1}\mathbf{v}_{i-1}^* = \mathbf{v}_i\mathbf{v}_i^*. \end{aligned}$$

It remains to show that the convergence  $\mathbf{q}_i^{(k)}\mathbf{q}_i^{(k)*} \rightarrow \mathbf{v}_i\mathbf{v}_i^*$  implies the desired essential convergence. Noting that  $|a| = \sqrt{a\bar{a}}$  for every  $a \in \mathbf{C}$ , we have

$$|\mathbf{v}_i^*\mathbf{q}_i^{(k)}| = \sqrt{\mathbf{v}_i^*\mathbf{q}_i^{(k)}\mathbf{q}_i^{(k)*}\mathbf{v}_i} \xrightarrow[k \rightarrow \infty]{} \sqrt{\mathbf{v}_i^*\mathbf{v}_i\mathbf{v}_i^*\mathbf{v}_i} = 1,$$

Finally, observing that

$$\left\| e^{-i\theta_k}\mathbf{q}_i^{(k)} - \mathbf{v}_i \right\|^2 = 2 - 2|\mathbf{v}_i^*\mathbf{q}_i^{(k)}| \xrightarrow[k \rightarrow \infty]{} 0, \quad e^{i\theta_k} = \frac{\mathbf{v}_i^*\mathbf{q}_i^{(k)}}{|\mathbf{v}_i^*\mathbf{q}_i^{(k)}|},$$

we conclude that  $\mathbf{q}_i^{(k)}$  converges essentially to  $\mathbf{v}_i$ . □

In addition to this convergence result, it is possible to show that the error satisfies

$$\angle(\text{col}(\mathbf{X}_k), \text{col}(\mathbf{V}_1)) = \mathcal{O}\left(\left|\frac{\lambda_{p+1}}{\lambda_p}\right|^k\right).$$

Here, the angle between two subspaces  $\mathcal{A}$  and  $\mathcal{B}$  of  $\mathbf{C}^n$  is defined as

$$\angle(\mathcal{A}, \mathcal{B}) = \max_{\mathbf{a} \in \mathcal{A} \setminus \{0\}} \left( \min_{\mathbf{b} \in \mathcal{B} \setminus \{0\}} \angle(\mathbf{a}, \mathbf{b}) \right).$$

### 6.3.2 The QR algorithm

The QR algorithm, which is based on the QR decomposition, is one of the most famous algorithms for calculating *all* the eigenpairs of a matrix. We first present the algorithm and then relate it to the simultaneous iteration in [Section 6.3.1](#). The method is presented in [Algorithm 12](#).

Successive iterates of the QR algorithm are related by the equation

$$\mathbf{X}_k = \mathbf{Q}_k^{-1}\mathbf{X}_{k-1}\mathbf{Q}_k = \mathbf{Q}_k^*\mathbf{X}_{k-1}\mathbf{Q}_k = \dots = (\mathbf{Q}_1 \dots \mathbf{Q}_k)^*\mathbf{X}_0(\mathbf{Q}_1 \dots \mathbf{Q}_k) \quad (6.12)$$

**Algorithm 12** QR algorithm

---

```

 $X_0 = A$ 
for  $i \in \{1, 2, \dots\}$  do
     $Q_k R_k = X_{k-1}$  (QR decomposition)
     $X_k = R_k Q_k$ 
end for

```

---

Therefore, all the iterates are related by a unitary similarity transformation, and so they all have the same eigenvalues as  $X_0 = A$ . Rearranging (6.12), we have

$$(Q_1 \dots Q_k) X_k = A(Q_1 \dots Q_k),$$

and so, introducing  $\tilde{Q}_k = Q_1 \dots Q_k$  and noting that  $X_k = Q_{k+1} R_{k+1}$  by the algorithm, we deduce

$$\tilde{Q}_{k+1} R_{k+1} = A \tilde{Q}_k.$$

This reveals that the matrix sequence  $(\tilde{Q}_k)_{k \geq 1}$  undergoes a simultaneous iteration and so, assuming that  $A$  is Hermitian with  $n$  distinct nonzero eigenvalues, we deduce that  $\tilde{Q}_k \rightarrow V$  essentially in the limit as  $k \rightarrow \infty$ , by Theorem 6.5. As a consequence, by (6.12), it holds that  $X_k \rightarrow V^* X_0 V = D$ ; in other words, the matrix  $X_k$  converges to a diagonal matrix with the eigenvalues of  $A$  on the diagonal.

## 6.4 Projection methods

In this section, we begin by presenting a general method for constructing an approximation of the eigenvectors of  $A$  in a given subspace  $\mathcal{U}$  of  $\mathbb{C}^n$ . We then discuss a particular choice for the subspace  $\mathcal{U}$  as a Krylov subspace, which is very useful in practice.

Assume that  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  is an orthonormal basis of  $\mathcal{U}$ . Then for any vector  $\mathbf{v} \in \mathbb{C}^n$ , the vector of  $\mathcal{U}$  that is closest to  $\mathbf{v}$  in the Euclidean distance is given by the orthogonal projection

$$P_{\mathcal{U}} \mathbf{v} := U U^* \mathbf{v} = (\mathbf{u}_1 \mathbf{u}_1^* + \dots + \mathbf{u}_p \mathbf{u}_p^*) \mathbf{v}.$$

In practice, the eigenvectors of  $A$  are unknown, and so it is impossible to calculate approximations using this formula. The Rayleigh–Ritz method, which we present hereafter, is an alternative and practical method for constructing approximations of the eigenvectors and eigenvalues. In general, the subspace  $\mathcal{U}$  does not contain any eigenvector of  $A$ , and so the problem

$$A \mathbf{v} = \lambda \mathbf{v}, \quad \mathbf{v} \in \mathcal{U} \tag{6.13}$$

does not admit a solution. Let us denote by  $U$  the matrix with columns  $\mathbf{u}_1, \dots, \mathbf{u}_p$ . Since any vector  $\mathbf{v} \in \mathcal{U}$  is equal to  $U \mathbf{z}$  for some vector  $\mathbf{z} \in \mathbb{C}^p$ , equation (6.13) is equivalent to the problem

$$A U \mathbf{z} = \lambda U \mathbf{z},$$

which is a system of  $n$  equations with  $p < n$  unknowns. The Rayleigh–Ritz method is based on

the idea that, in order to obtain a problem with as many unknowns as there are equations, we can multiply this equation by  $U^*$ , which leads to the problem

$$Bz := (U^*AU)z = \lambda z. \quad (6.14)$$

This is standard eigenvalue problem for the matrix  $U^*AU \in \mathbf{C}^{p \times p}$ , which is much easier to solve than the original problem if  $p \ll n$ . Equivalently, equation (6.14) may be formulated as follows: find  $v \in \mathcal{U}$  such that

$$u^*(Av - \lambda v), \quad \forall u \in \mathcal{U}. \quad (6.15)$$

The solutions to (6.14) and (6.15) are related by the equation  $v = Uz$ . Of course, the eigenvalues of  $B$  in problem (6.14), which are called the Ritz values of  $A$  relative to  $\mathcal{U}$ , are in general different from those of  $A$ . Once an eigenvector  $y$  of  $B$  has been calculated, an approximate eigenvector of  $A$ , called a *Ritz vector* of  $A$  relative to  $\mathcal{U}$ , is obtained from the equation  $\hat{v} = Uy$ . The Rayleigh–Ritz algorithm is presented in full in [Algorithm 13](#).

---

**Algorithm 13** Rayleigh–Ritz

---

Choose  $\mathcal{U} \subset \mathbf{C}^n$

Construct a matrix  $U$  whose columns are orthonormal and span  $\mathcal{U}$

Find the eigenvalues  $\hat{\lambda}_i$  and eigenvectors  $y_i \in \mathbf{C}^p$  of  $B := U^*AU$

Calculate the corresponding Ritz vectors  $\hat{v}_i = Uy_i \in \mathbf{C}^n$ .

---

It is clear that if  $v_i \in \mathcal{U}$ , then  $\lambda_i$  is an eigenvalue of  $B$  in (6.14). In fact, we can show the following more general statement.

**Proposition 6.6.** *If  $\mathcal{U}$  is an invariant subspace of  $A$ , meaning that  $A\mathcal{U} \subset \mathcal{U}$ , then each Ritz vector of  $A$  relative to  $\mathcal{U}$  is an eigenvector of  $A$ .*

*Proof.* Let  $U \in \mathbf{C}^{n \times p}$  and  $W \in \mathbf{C}^{n \times (n-p)}$  be matrices whose columns form orthonormal bases of  $\mathcal{U}$  and  $\mathcal{U}^\perp$ , respectively. Here  $\mathcal{U}^\perp$  denotes the orthogonal complement of  $\mathcal{U}$  with respect to the Euclidean inner product. Then, since  $W^*AU = 0$  by assumption, it holds that

$$Q^*AQ = \begin{pmatrix} U^*AU & U^*AW \\ W^*AU & W^*AW \end{pmatrix} = \begin{pmatrix} U^*AU & U^*AW \\ 0 & W^*AW \end{pmatrix}, \quad Q = \begin{pmatrix} U & W \end{pmatrix}.$$

If  $(y, \hat{\lambda})$  is an eigenvector of  $U^*AU$ , then

$$Q^*AQ \begin{pmatrix} y \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} (U^*AU)y \\ \mathbf{0} \end{pmatrix} = \hat{\lambda} \begin{pmatrix} y \\ \mathbf{0} \end{pmatrix} =: \hat{\lambda}x,$$

and so  $(x, \hat{\lambda})$  is an eigenpair of  $Q^*AQ$ . But then  $(Qx, \hat{\lambda}) = (Uy, \hat{\lambda})$  is an eigenpair of  $A$ , which proves the statement.  $\square$

*Remark 6.1.* [Proposition 6.6](#) can be proved more directly from (6.15) by noting that, if  $\mathcal{U}$  is an invariant subspace of  $A$ , then this equation implies that  $Av - \lambda v$  belongs to both  $\mathcal{U}$  and

its orthogonal complement  $\mathcal{U}^\perp$ , and so this vector is  $\mathbf{0}$ .

If  $\mathcal{U}$  is close to being an invariant subspace of  $\mathbf{A}$ , then it is expected that the Ritz vectors and Ritz values of  $\mathbf{A}$  relative to  $\mathcal{U}$  will provide good approximations of some of the eigenpairs of  $\mathbf{A}$ . Quantifying this approximation is difficult, so we only present without proof the following error bound. See [11] for more information.

**Proposition 6.7.** *Let  $\mathbf{A}$  be a full rank Hermitian matrix and  $\mathcal{U}$  a  $p$ -dimensional subspace of  $\mathbb{C}^n$ . Then there exists eigenvalues  $\lambda_{i_1}, \dots, \lambda_{i_p}$  of  $\mathbf{A}$  which satisfy*

$$\forall j \in \{1, \dots, p\}, \quad |\lambda_{i_j} - \hat{\lambda}_j| \leq \|(\mathbf{I} - \mathbf{P}_{\mathcal{U}})\mathbf{A}\mathbf{P}_{\mathcal{U}}\|_2.$$

In the case where  $\mathbf{A}$  is Hermitian, it is possible to show that the Ritz values are bounded from above by the eigenvalues of  $\mathbf{A}$ . The proof of this result relies on the Courant–Fisher theorem for characterizing the eigenvalues of a Hermitian matrix, which is recalled in [Theorem A.6](#) in the appendix.

**Proposition 6.8.** *If  $\mathbf{A} \in \mathbb{C}^{n \times n}$  is Hermitian, then*

$$\forall i \in \{1, \dots, p\}, \quad \hat{\lambda}_i \leq \lambda_i$$

*Proof.* By the Courant–Fisher theorem, it holds that

$$\hat{\lambda}_i = \max_{\mathcal{S} \subset \mathbb{C}^p, \dim(\mathcal{S})=i} \left( \min_{\mathbf{x} \in \mathcal{S} \setminus \{0\}} \frac{\mathbf{x}^* \mathbf{B} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} \right)$$

Letting  $\mathbf{y} = \mathbf{U}\mathbf{x}$  and then  $\mathcal{R} = \mathbf{U}\mathcal{S}$ , we deduce that

$$\begin{aligned} \hat{\lambda}_i &= \max_{\mathcal{S} \subset \mathbb{C}^p, \dim(\mathcal{S})=i} \left( \min_{\mathbf{y} \in \mathbf{U}\mathcal{S} \setminus \{0\}} \frac{\mathbf{y}^* \mathbf{A} \mathbf{y}}{\mathbf{y}^* \mathbf{y}} \right) \\ &= \max_{\mathcal{R} \subset \mathcal{U}, \dim(\mathcal{R})=i} \left( \min_{\mathbf{y} \in \mathcal{R} \setminus \{0\}} \frac{\mathbf{y}^* \mathbf{A} \mathbf{y}}{\mathbf{y}^* \mathbf{y}} \right) \leq \max_{\mathcal{R} \subset \mathbb{C}^n, \dim(\mathcal{R})=i} \left( \min_{\mathbf{y} \in \mathcal{R} \setminus \{0\}} \frac{\mathbf{y}^* \mathbf{A} \mathbf{y}}{\mathbf{y}^* \mathbf{y}} \right) = \lambda_i, \end{aligned}$$

where we used the Courant–Fisher theorem for the matrix  $\mathbf{A}$  in the last equality.  $\square$

This projection approach is sometimes combined with a simultaneous subspace iteration: an approximation  $\mathbf{X}_k$  of the  $p$  first eigenvector is first calculated using [Algorithm 11](#), and then the matrix  $\mathbf{X}_k$  is used in place of  $\mathbf{U}$  in [Algorithm 13](#).

### 6.4.1 Projection method in a Krylov subspace

The power iteration constructs at iteration  $k$  an approximation of  $\mathbf{v}_1$  in the one-dimensional subspace spanned by the vector  $\mathbf{A}^k \mathbf{x}_0$ , and only the previous iteration  $\mathbf{x}_k$  is employed to construct  $\mathbf{x}^{k+1}$ . One may wonder whether, by employing all the previous iterates rather than only the previous one, a better approximation of  $\mathbf{v}_1$  can be constructed. More precisely, instead of looking for an approximation in the subspace  $\text{Span}\{\mathbf{A}^k \mathbf{x}_0\}$ , would it be useful to extend the

search area to the Krylov subspace

$$\mathcal{K}_{k+1}(\mathbf{A}, \mathbf{x}_0) := \text{Span}\{\mathbf{x}_0, \mathbf{A}\mathbf{x}_0, \dots, \mathbf{A}^k \mathbf{x}_0\}?$$

The answer to this question is positive, and the resulting method is often much faster than the power iteration. This is achieved by employing the Rayleigh–Ritz projection method [Algorithm 13](#) with the choice  $\mathcal{U} = \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{x}_0)$ . Applying this method requires to calculate an orthonormal basis of the Krylov subspace and to calculate the reduced matrix  $\mathbf{U}^* \mathbf{A} \mathbf{U}$ . The *Arnoldi method* enables to achieve these two goals simultaneously.

### 6.4.2 The Arnoldi iteration

This Arnoldi iteration is based on the Gram–Schmidt process and presented in [Algorithm 14](#). The iteration breaks down if  $h_{j+1,j} = 0$ , which indicates that  $\mathbf{A}\mathbf{u}_j$  belongs to the Krylov

---

**Algorithm 14** Arnoldi iteration for constructing an orthonormal basis of  $\mathcal{K}_p(\mathbf{A}, \mathbf{u}_1)$

---

Choose  $\mathbf{u}_1$  with unit norm.

**for**  $j \in \{1, \dots, p\}$  **do**

$\mathbf{u}_{j+1} \leftarrow \mathbf{A}\mathbf{u}_j$

**for**  $i \in \{1, \dots, j\}$  **do**

$h_{i,j} \leftarrow \mathbf{u}_i^* \mathbf{u}_{j+1}$

$\mathbf{u}_{j+1} \leftarrow \mathbf{u}_{j+1} - h_{i,j} \mathbf{u}_i$

**end for**

$h_{j+1,j} \leftarrow \|\mathbf{u}_{j+1}\|$

$\mathbf{u}_{j+1} \leftarrow \mathbf{u}_{j+1}/h_{j+1,j}$

**end for**

---

subspace  $\text{Span}\{\mathbf{u}_1, \dots, \mathbf{u}_j\} = \mathcal{K}_j(\mathbf{A}, \mathbf{u}_1)$ , implying that  $\mathcal{K}_{j+1}(\mathbf{A}, \mathbf{u}_1) = \mathcal{K}_j(\mathbf{A}, \mathbf{u}_1)$ . In this case, the subspace  $\mathcal{K}_j(\mathbf{A}, \mathbf{u}_1)$  is an invariant subspace of  $\mathbf{A}$  because, by [Exercise 6.2](#), we have

$$\mathbf{A}\mathcal{K}_j(\mathbf{A}, \mathbf{u}_1) \subset \mathcal{K}_{j+1}(\mathbf{A}, \mathbf{u}_1) = \mathcal{K}_j(\mathbf{A}, \mathbf{u}_1).$$

Therefore, applying the Rayleigh–Ritz with  $\mathcal{U} = \text{Span}\{\mathbf{u}_1, \dots, \mathbf{u}_j\}$  yields exact eigenpairs in view of [Proposition 6.6](#). If the iteration does not break down then, by construction, the vectors  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  at the end of the algorithm are orthonormal. It is also simple to show by induction that they form a basis of  $\mathcal{K}_p(\mathbf{A}, \mathbf{u}_1)$ . The scalar coefficients  $h_{i,j}$  can be collected in a matrix square  $p \times p$  matrix

$$\mathbf{H} = \begin{pmatrix} h_{1,1} & h_{1,2} & h_{1,3} & \cdots & h_{1,p} \\ h_{2,1} & h_{2,2} & h_{2,3} & \cdots & h_{2,p} \\ 0 & h_{3,2} & h_{3,3} & \cdots & h_{3,p} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & h_{p,p-1} & h_{p,p} \end{pmatrix}.$$

This matrix contains only zeros under the first subdiagonal; such a matrix is called a *Hessenberg* matrix. Inspecting the algorithm, we notice that the  $j$ -th column contains the coefficients of



the projection of the vector  $\mathbf{A}\mathbf{u}_j$  onto the basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ . In other words,

$$\mathbf{U}^*\mathbf{A}\mathbf{U} = \mathbf{H}, \quad (6.16)$$

We have thus shown that the Arnoldi algorithm enables to construct both an orthonormal basis of a Krylov subspace and the associated reduced matrix. In fact, we have the following equation

$$\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{H} + h_{p+1,p}(\mathbf{v}_{p+1}\mathbf{e}_p^*), \quad \mathbf{e}_p = \begin{pmatrix} 0 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{C}^p. \quad (6.17)$$

The Arnoldi algorithm, coupled with the Rayleigh–Ritz method, has very good convergence properties in the limit as  $p \rightarrow \infty$ , in particular for eigenvalues with a large modulus. The following result shows that the residual  $\mathbf{r} = \mathbf{A}\hat{\mathbf{v}} - \hat{\lambda}\hat{\mathbf{v}}$  associated with a Ritz vector can be estimated inexpensively. Specifically, the norm of the residual is equal to the last component of the associated eigenvector of  $\mathbf{H}$  multiplied by  $h_{p+1,p}$ .

**Proposition 6.9** (Formula for the residual <sup>ⓐ</sup>). *Let  $\mathbf{y}_i$  be an eigenvector of  $\mathbf{H}$  associated with the eigenvalues  $\hat{\lambda}_i$ , and let  $\hat{\mathbf{v}}_i = \mathbf{U}\mathbf{y}_i$  denote the corresponding eigenvector. Then*

$$\mathbf{A}\hat{\mathbf{v}}_i - \hat{\lambda}_i\hat{\mathbf{v}}_i = h_{p+1,p}(\mathbf{y}_i)_p\mathbf{v}_{p+1}.$$

Consequently, it holds that

$$\|\mathbf{A}\hat{\mathbf{v}}_i - \hat{\lambda}_i\hat{\mathbf{v}}_i\| = |h_{p+1,p}(\mathbf{y}_i)_p|.$$

*Proof.* Multiplying both sides of (6.17) by  $\mathbf{y}_i$ , we obtain

$$\mathbf{A}\mathbf{U}\mathbf{y}_i = \mathbf{U}\mathbf{H}\mathbf{y}_i + h_{p+1,p}(\mathbf{v}_{p+1}\mathbf{e}_p^*)\mathbf{y}_i.$$

Using the definition of  $\hat{\mathbf{v}}_i$  and rearranging the equation, we have

$$\mathbf{A}\hat{\mathbf{v}}_i - \hat{\lambda}_i\hat{\mathbf{v}}_i = h_{p+1,p}(\mathbf{v}_{p+1}\mathbf{e}_p^*)\mathbf{y}_i,$$

which immediately gives the result. □

In practice, the larger the dimension  $p$  of the subspace  $\mathcal{U}$  employed in the Rayleigh–Ritz method, the more memory is required for storing an orthonormal basis of  $\mathcal{U}$ . In addition, for large values of  $p$ , computing the reduced matrix (6.16) and its eigenpairs becomes computationally expensive; the computational cost of computing the matrix  $\mathbf{H}$  scales as  $\mathcal{O}(p^2)$ . To remedy these potential issues, the algorithm can be restarted periodically. For example, Algorithm 15 can be employed as an alternative to the power iteration in order to find the eigenvector associated with the eigenvalue with largest modulus.

**Algorithm 15** Restarted Arnoldi iteration

---

Choose  $\mathbf{u}_1 \in \mathbf{C}^n$  and  $p \ll n$   
**for**  $i \in \{1, 2, \dots\}$  **do**  
  Perform  $p$  iterations of the Arnoldi iteration and construct  $\mathcal{U}$ ;  
  Calculate the Ritz vector  $\hat{\mathbf{v}}_1$  associated with the largest Ritz value relative to  $\mathcal{U}$ ;  
  If this vector is sufficiently accurate, then stop. Otherwise, restart with  $\mathbf{u}_1 = \hat{\mathbf{v}}_1$ .  
**end for**

---

**6.4.3 The Lanczos iteration**

The Lanczos iteration may be viewed as a simplified version of the Arnoldi iteration in the case where the matrix  $A$  is Hermitian. Let us denote by  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  the orthonormal vectors generated by the Arnoldi iteration. When  $A$  is Hermitian, it holds that

$$h_{i,j} = \mathbf{u}_i^*(A\mathbf{u}_j) = (A\mathbf{u}_i)^*\mathbf{u}_j = \overline{h_{j,i}}.$$

Therefore, the matrix  $H$  is Hermitian. This is not surprising, since we showed that  $H = U^*AU$  and the matrix  $A$  is Hermitian. Since  $H$  is also of Hessenberg form, we deduce that  $H$  is tridiagonal. An inspection of [Algorithm 14](#) shows that the subdiagonal entries of  $H$  are real. Since  $A$  is Hermitian, the diagonal entries  $h_{i,i} = \mathbf{u}_i^*(A\mathbf{u}_j)$  are also real, and so we conclude that all the entries of the matrix  $H$  are in fact real. This matrix is of the form

$$H = \begin{pmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \alpha_2 & \beta_3 & & \\ & \beta_3 & \ddots & \ddots & \\ & & \ddots & \ddots & \beta_p \\ & & & \beta_p & \alpha_p \end{pmatrix}$$

Adapting the Arnoldi iteration to this setting leads to [Algorithm 16](#).

**Algorithm 16** Lanczos iteration for constructing an orthonormal basis of  $\mathcal{K}_p(A, \mathbf{u}_1)$ 


---

Choose  $\mathbf{u}_1$  with unit norm.  
 $\beta_1 \leftarrow 0, \mathbf{u}_0 \leftarrow \mathbf{0} \in \mathbf{C}^n$   
**for**  $j \in \{1, \dots, p\}$  **do**  
   $\mathbf{u}_{j+1} \leftarrow A\mathbf{u}_j - \beta_j\mathbf{u}_{j-1}$   
   $\alpha_j \leftarrow \mathbf{u}_j^*\mathbf{u}_{j+1}$   
   $\mathbf{u}_{j+1} \leftarrow \mathbf{u}_{j+1} - \alpha_j\mathbf{u}_j$   
   $\beta_{j+1} \leftarrow \|\mathbf{u}_{j+1}\|$   
   $\mathbf{u}_{j+1} \leftarrow \mathbf{u}_{j+1}/\beta_{j+1}$   
**end for**

---

**6.5 Exercises**

⚙ **Exercise 6.1.** *PageRank* is an algorithm for assigning a rank to the vertices of a directed graph. It is used by many search engines, notably Google, for sorting search results. In this

context, the directed graph encodes the links between pages of the World Wide Web: the vertices of the directed graph are webpages, and there is an edge going from page  $i$  to page  $j$  if page  $i$  contains a hyperlink to page  $j$ .

Let us consider a directed graph  $G(V, E)$  with vertices  $V = \{1, \dots, n\}$  and edges  $E$ . The graph can be represented by its adjacency matrix  $\mathbf{A} \in \{0, 1\}^{n \times n}$ , whose entries are given by

$$a_{ij} = \begin{cases} 1 & \text{if there is an edge from } i \text{ to } j, \\ 0 & \text{otherwise.} \end{cases}$$

Let  $r_i$  denote the “value” assigned to vertex  $i$ . The idea of PageRank, in its simplest form, is to assign values to the vertices by solving the following system of equations;

$$\forall i \in V, \quad r_i = \sum_{j \in \mathcal{N}(i)} \frac{r_j}{o_j}. \quad (6.18)$$

where  $o_j$  is the outdegree of vertex  $j$ , i.e. the number of edges leaving from  $j$ . Here the sum is over the set of nodes  $\mathcal{N}(i)$ , which denotes all the “incoming” neighbors of  $i$ , i.e. those that have an edge pointing towards node  $i$ .

- Read the Wikipedia page on PageRank to familiarize yourself with the algorithm.
- Let  $\mathbf{r} = (r_1 \ \dots \ r_n)^T$ . Show using (6.18) that  $\mathbf{r}$  satisfies

$$\mathbf{r} = \mathbf{A}^T \begin{pmatrix} \frac{1}{o_1} & & \\ & \ddots & \\ & & \frac{1}{o_n} \end{pmatrix} \mathbf{r} =: \mathbf{A}^T \mathbf{O}^{-1} \mathbf{r}.$$

In other words,  $\mathbf{r}$  is an eigenvector with eigenvalue 1 of the matrix  $\mathbf{M} = \mathbf{A}^T \mathbf{O}^{-1}$ .

- Show that  $\mathbf{M}$  is a left-stochastic matrix, i.e. that each column sums to 1.
- Prove that the eigenvalues of any matrix  $\mathbf{B} \in \mathbf{R}^{n \times n}$  coincide with those of  $\mathbf{B}^T$ . You may use the fact that  $\det(\mathbf{B}) = \det(\mathbf{B}^T)$ .
- Using the previous items, show that 1 is an eigenvalue and that  $\rho(\mathbf{M}) = 1$ . For the second part, find a subordinate matrix norm such that  $\|\mathbf{M}\| = 1$ .
- Implement PageRank in order to rank pages from a 2013 snapshot of English Wikipedia. You can use either the simplified version of the algorithm given in (6.18) or the improved version with a damping factor described on Wikipedia. In the former case, the following are both sensible stopping criteria:

$$\frac{\|\mathbf{M}\hat{\mathbf{r}} - \hat{\mathbf{r}}\|_1}{\|\hat{\mathbf{r}}\|_1} < 10^{-15} \quad \text{or} \quad \frac{\|\mathbf{M}\hat{\mathbf{r}} - \hat{\lambda}\hat{\mathbf{r}}\|_1}{\|\hat{\mathbf{r}}\|_1} < 10^{-15}, \quad \hat{\lambda} = \frac{\hat{\mathbf{r}}^T \mathbf{M} \hat{\mathbf{r}}}{\hat{\mathbf{r}}^T \hat{\mathbf{r}}},$$

where  $\hat{\mathbf{v}}$  is an approximation of the eigenvector corresponding to the dominant eigenvalue. A dataset is available on the course website to complete this part. This dataset contains

a subset of the data publicly available [here](#), and was generated from the full dataset by retaining only the 5% best rated articles. After decompressing the archive, you can load the dataset into Julia by using the following commands:

```
import CSV
import DataFrames

# Data (nodes and edges)
nodes = CSV.read("names.csv", DataFrames.DataFrame)
edges = CSV.read("edges.csv", DataFrames.DataFrame)

# Convert data to matrices
nodes = Matrix(nodes)
edges = Matrix(edges)
```

After you have assigned a rank to all the pages, print the 10 pages with the highest ranks. My code returns the following entries:

- |                   |                     |           |
|-------------------|---------------------|-----------|
| 1. United States  | 5. France           | 9. Canada |
| 2. United Kingdom | 6. Germany          | 10. India |
| 3. World War II   | 7. English language |           |
| 4. Latin          | 8. China            |           |

- **Extra credit:** Write a function `search(keyword)` that can be employed for searching the database. Here is an example of what it could return:

```
julia> search("New York")
481-element Vector{String}:
 "New York City"
 "New York"
 "The New York Times"
 "New York Stock Exchange"
 "New York University"
 ...
```

⚙️ **Exercise 6.2.** Show the following properties of the Krylov subspace  $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$ .

- $\mathcal{K}_p(\mathbf{A}, \mathbf{x}) \subset \mathcal{K}_{p+1}(\mathbf{A}, \mathbf{x})$ .
- $\mathbf{A}\mathcal{K}_p(\mathbf{A}, \mathbf{x}) \subset \mathcal{K}_{p+1}(\mathbf{A}, \mathbf{x})$ .
- The Krylov subspace  $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$  is invariant under rescaling: for all  $\alpha \in \mathbf{C}$ ,

$$\mathcal{K}_p(\mathbf{A}, \mathbf{x}) = \mathcal{K}_p(\alpha\mathbf{A}, \mathbf{x}) = \mathcal{K}_p(\mathbf{A}, \alpha\mathbf{x}).$$

- The Krylov subspace  $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$  is invariant under shift of the matrix  $\mathbf{A}$ : for all  $\alpha \in \mathbf{C}$ ,

$$\mathcal{K}_p(\mathbf{A}, \mathbf{x}) = \mathcal{K}_p(\mathbf{A} - \alpha \mathbf{I}, \mathbf{x}).$$

- Similarity transformation: If  $\mathbf{T} \in \mathbf{C}^{n \times n}$  is nonsingular, then

$$\mathcal{K}_p(\mathbf{T}^{-1} \mathbf{A} \mathbf{T}, \mathbf{T}^{-1} \mathbf{x}) = \mathbf{T}^{-1} \mathcal{K}_p(\mathbf{A}, \mathbf{x}).$$

⚙ **Exercise 6.3.** The minimal polynomial of a matrix  $\mathbf{A} \in \mathbf{C}^{n \times n}$  is the monic polynomial  $p$  of lowest degree such that  $p(\mathbf{A}) = \mathbf{0}$ . Prove that, if  $\mathbf{A}$  is Hermitian with  $m \leq n$  distinct eigenvalues, then the minimal polynomial is given by

$$p(t) = \prod_{i=1}^m (t - \lambda_i).$$

⚙ **Exercise 6.4.** The minimal polynomial for a general matrix  $\mathbf{A} \in \mathbf{C}^{n \times n}$  is given by

$$p(t) = \prod_{i=1}^m (t - \lambda_i)^{s_i}.$$

where  $s_i$  is the size of the largest Jordan block associated with the eigenvalue  $\lambda_i$  in the normal Jordan form of  $\mathbf{A}$ . Verify that  $p(\mathbf{A}) = \mathbf{0}$ .

⚙ **Exercise 6.5.** Let  $d$  denote the degree of the minimal polynomial of  $\mathbf{A}$ . Show that

$$\forall p \geq d, \quad \mathcal{K}_{p+1}(\mathbf{A}, \mathbf{x}) = \mathcal{K}_p(\mathbf{A}, \mathbf{x}).$$

Deduce that, for  $p \geq n$ , the subspace  $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$  is an invariant subspace of  $\mathbf{A}$ .

⚙ **Exercise 6.6.** Let  $\mathbf{A} \in \mathbf{C}^{n \times n}$ . Show that  $\mathcal{K}_n(\mathbf{A}, \mathbf{x})$  is the smallest invariant subspace of  $\mathbf{A}$  that contains  $\mathbf{x}$ .

□ **Exercise 6.7.** Consider the matrix

$$\mathbf{M} = \begin{pmatrix} 0 & 1 & 2 & 0 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 0 & 2 \\ 0 & 0 & 2 & 0 \end{pmatrix}$$

- Find the dominant eigenvalue of  $\mathbf{M}$  by using the power iteration.
- Find the eigenvalue of  $\mathbf{M}$  closest to 1 by using the inverse iteration.
- Find the other two eigenvalues of  $\mathbf{M}$  by using a method of your choice.

⚙ **Exercise 6.8** (A posteriori error bound). Assume that  $\mathbf{A} \in \mathbf{C}^{n \times n}$  is Hermitian, and that  $\widehat{\mathbf{v}}$  is a normalized approximation of an eigenvector which satisfies

$$\|\widehat{\mathbf{z}}\| := \|\mathbf{A}\widehat{\mathbf{v}} - \widehat{\lambda}\widehat{\mathbf{v}}\| = \delta, \quad \widehat{\lambda} = \frac{\widehat{\mathbf{v}}^* \mathbf{A} \widehat{\mathbf{v}}}{\widehat{\mathbf{v}}^* \widehat{\mathbf{v}}}.$$

Prove that there is an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$|\hat{\lambda} - \lambda| \leq \delta.$$

**Hint:** Consider first the case where  $\mathbf{A}$  is diagonal.

⚙️ **Exercise 6.9** (Bauer–Fike theorem). Assume that  $\mathbf{A} \in \mathbf{C}^{n \times n}$  is diagonalizable:  $\mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{D}$ . Show that, if  $\hat{\mathbf{v}}$  is a normalized approximation of an eigenvector which satisfies

$$\|\mathbf{r}\| := \|\mathbf{A}\hat{\mathbf{v}} - \hat{\lambda}\hat{\mathbf{v}}\| = \delta$$

for some  $\hat{\lambda} \in \mathbf{C}$ , then there is an eigenvalue  $\lambda$  of  $\mathbf{A}$  such that

$$|\hat{\lambda} - \lambda| \leq \kappa_2(\mathbf{V})\delta.$$

**Hint:** Rewrite

$$\|\hat{\mathbf{v}}\| = \|(\mathbf{A} - \hat{\lambda}\mathbf{I})^{-1}\mathbf{r}\| = \|\mathbf{V}(\mathbf{D} - \hat{\lambda}\mathbf{I})^{-1}\mathbf{V}^{-1}\mathbf{r}\|.$$

⚙️ **Exercise 6.10.** In *Exercise 6.8* and *Exercise 6.9*, we saw examples a posteriori error estimates which guarantee the existence of an eigenvalue of  $\mathbf{A}$  within a certain distance of the approximation  $\hat{\lambda}$ . In this exercise, our aim is to provide an answer to the following different question: given an approximate eigenpair  $(\hat{\mathbf{v}}, \hat{\lambda})$ , what is the smallest perturbation  $\mathbf{E}$  that we need to apply to  $\mathbf{A}$  in order to guarantee that  $(\hat{\mathbf{v}}, \hat{\lambda})$  is an exact eigenpair, i.e. that

$$(\mathbf{A} + \mathbf{E})\hat{\mathbf{v}} = \hat{\lambda}\hat{\mathbf{v}}?$$

Assume that  $\hat{\mathbf{v}}$  is normalized and let  $\mathcal{E} = \{\mathbf{E} \in \mathbf{C}^{n \times n} : (\mathbf{A} + \mathbf{E})\hat{\mathbf{v}} = \hat{\lambda}\hat{\mathbf{v}}\}$ . Prove that

$$\min_{\mathbf{E} \in \mathcal{E}} \|\mathbf{E}\|_2 = \|\mathbf{r}\|_2 := \|\mathbf{A}\hat{\mathbf{v}} - \hat{\lambda}\hat{\mathbf{v}}\|. \quad (6.19)$$

To this end, you may proceed as follows:

- Show first that any  $\mathbf{E} \in \mathcal{E}$  satisfies  $\mathbf{E}\hat{\mathbf{v}} = -\mathbf{r}$ .
- Deduce from the first item that

$$\inf_{\mathbf{E} \in \mathcal{E}} \|\mathbf{E}\|_2 \geq \|\mathbf{r}\|_2.$$

- Find a rank one matrix  $\mathbf{E}_*$  such that  $\|\mathbf{E}_*\|_2 = \|\mathbf{r}\|_2$ , and then conclude. Recall that any rank 1 matrix can be written in the form  $\mathbf{E}_* = \mathbf{u}\mathbf{w}^*$ , with norm  $\|\mathbf{u}\|_2\|\mathbf{w}\|_2$ .

Equation (6.19) is a simplified version of the Kahan–Parlett–Jiang theorem and is an example of a backward error estimate. Whereas forward error estimates quantify the distance between an approximation and the exact solution, backward error estimates give the size of the perturbation that must be applied to a problem so that an approximation is exact.

⚙️ **Exercise 6.11.** Assume that  $(\mathbf{x}_k)_{k \geq 0}$  is a sequence of normalized vectors in  $\mathbf{C}^n$ . Show that the following statements are equivalent:

- $(\mathbf{x}_k)_{k \geq 0}$  converges essentially to  $\mathbf{x}_\infty$  in the limit as  $k \rightarrow \infty$ .
- The angle  $\angle(\mathbf{x}_k, \mathbf{x}_\infty)$  converges to zero in the limit as  $k \rightarrow \infty$ .
- The projector  $\mathbf{P}_{\mathbf{x}_k}$  converges to  $\mathbf{P}_{\mathbf{x}_\infty}$  in the limit as  $k \rightarrow \infty$ .

⚙️ **Exercise 6.12.** Assume that  $\mathbf{A} \in \mathbf{C}^{n \times n}$  is skew-Hermitian. Derive a Lanczos-like algorithm for constructing an orthonormal basis of  $\mathcal{K}_p(\mathbf{A}, \mathbf{x})$  and calculating the reduced matrix

$$\mathbf{U}^* \mathbf{A} \mathbf{U},$$

where  $\mathbf{U} \in \mathbf{C}^{n \times p}$  contains the vectors of the basis as columns. Implement your algorithm with  $p = 20$  in order to approximate the dominant eigenvalue of the matrix  $\mathbf{A}$  constructed by the following piece of code:

```
n = 5000
A = rand(n, n) + im * randn(n, n)
A = A - A' # A is now skew-Hermitian
```

⚙️ **Exercise 6.13.** Assume that  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  and  $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$  are orthonormal bases of the same subspace  $\mathcal{U} \subset \mathbf{C}^n$ . Show that the projectors  $\mathbf{U}\mathbf{U}^*$  and  $\mathbf{W}\mathbf{W}^*$  are equal.

⚙️ **Exercise 6.14.** Assume that  $\mathbf{A} \in \mathbf{C}^{n \times n}$  is a Hermitian matrix with distinct eigenvalues, and suppose that we know the dominant eigenpair  $(\mathbf{v}_1, \lambda_1)$ , with  $\mathbf{v}_1$  normalized. Let

$$\mathbf{B} = \mathbf{A} - \lambda_1 \mathbf{v}_1 \mathbf{v}_1^*.$$

If we apply the power iteration to this matrix, what convergence can we expect?

⚙️ **Exercise 6.15.** Assume that  $\widehat{\mathbf{v}}_1$  and  $\widehat{\mathbf{v}}_2$  are two Ritz vectors of a Hermitian matrix  $\mathbf{A}$  relative to a subspace  $\mathcal{U} \subset \mathbf{C}^n$ . Show that, if the associated Ritz values are distinct, then  $\widehat{\mathbf{v}}_1^* \widehat{\mathbf{v}}_2 = 0$ .

## 6.6 Discussion and bibliography

The content of this chapter is inspired mainly from [15] and also from [12]. The latter volume contains a detailed coverage of the standard methods for eigenvalue problems. Some of the exercises are taken from [18], and others are inspired from [12].

# Chapter 7

## Numerical ordinary differential equations

<b>7.1</b>	<b>Analysis of the continuous problem</b>	<b>170</b>
<b>7.2</b>	<b>One-step methods</b>	<b>174</b>
7.2.1	Forward Euler method	175
7.2.2	Backward Euler method	176
7.2.3	Analysis of general one-step methods	177
7.2.4	Widely used one-step methods	180
<b>7.3</b>	<b>Multistep methods</b>	<b>183</b>
7.3.1	Adams–Bashforth methods	185
7.3.2	Adams–Moulton methods	187
<b>7.4</b>	<b>Absolute stability</b>	<b>187</b>
<b>7.5</b>	<b>Exercises</b>	<b>193</b>

### Introduction

This chapter concerns the numerical solution of ordinary differential equations (ODEs) of the following form:

$$\begin{cases} \mathbf{x}'(t) = \mathbf{f}(t, \mathbf{x}(t)), \\ \mathbf{x}(t_0) = \mathbf{x}_0. \end{cases} \quad (7.1)$$

Here  $\mathbf{f}: \mathbf{R} \times \mathbf{R}^n \rightarrow \mathbf{R}^n$  and  $\mathbf{x}_0$  is the initial condition. Equations of this type are the building blocks of a plethora of mathematical models in science and engineering. They have applications in celestial dynamics, molecular simulation and fluid mechanics, to mention just a few. Ordinary differential equations also arise after discretization of time-dependent partial differential equations, which are also ubiquitous in science. More often than not, it is not possible to find an explicit solution of (7.1), and so one has to resort to numerical simulation. The rest of the chapter is organized as follows:

- In Section 7.1, we define the concepts of local and global solutions for the continuous-time problem (7.1), and we recall fundamental results concerning the existence and uniqueness of a solution.



- In Section 7.2, we analyze the so-called *one-step* numerical methods to solve (7.1). We emphasize in particular the concepts of *consistency*, *stability* and *convergence*.
- In Section 7.3, we present *multistep* methods to solve (7.1), and discuss their drawbacks and advantages compared to one-step methods.
- Finally, in Section 7.4, we introduce the concept of *absolute stability* and discuss its relevance in the context of *stiff* differential equations.

## 7.1 Analysis of the continuous problem

A differentiable function  $\mathbf{x}: I \rightarrow \mathbf{R}^n$ , where  $I$  denotes an interval of  $\mathbf{R}$  containing  $t_0$ , is a solution of (7.1) if  $\mathbf{x}(t_0) = \mathbf{x}_0$  and the equation (7.1) is satisfied for all  $t \in I$ . The solution is called global if  $I = \mathbf{R}$ , and local otherwise.

**Integral formulation.** If  $\mathbf{x}$  is a solution to (7.1), then it holds that

$$\forall t \in I, \quad \mathbf{x}(t) = \mathbf{x}_0 + \int_{t_0}^t \mathbf{f}(s, \mathbf{x}(s)) \, ds. \quad (7.2)$$

The converse statement is not true in general, because a solution to (7.2) need not necessarily be differentiable everywhere. However, if the integral formulation (7.2) holds, then necessarily  $\mathbf{x}$  is absolutely continuous and (7.1) is satisfied for almost every  $t$ . Additionally, if (7.2) is satisfied and the function  $\mathbf{f}$  is continuous, then the function  $s \mapsto \mathbf{f}(s, \mathbf{x}(s))$  is continuous, and so (7.1) is satisfied for all  $t \in I$  by the fundamental theorem of analysis. We now focus on the integral formulation (7.2), and begin by establishing existence of a local solution.

**Theorem 7.1** (Existence of a solution). *Let  $\mathbf{x}_0 \in \mathbf{R}^n$  and let  $\Omega_{\mathcal{T}, \mathcal{R}}$  denote the set*

$$\{(t, \mathbf{x}) \in \mathbf{R} \times \mathbf{R}^n : |t - t_0| \leq \mathcal{T} \text{ and } \|\mathbf{x} - \mathbf{x}_0\| \leq \mathcal{R}\},$$

*Assume that the following conditions are satisfied for some  $\mathcal{T} > 0$  and  $\mathcal{R} > 0$ :*

- *The function  $\mathbf{f}$  is uniformly bounded on  $\Omega_{\mathcal{T}, \mathcal{R}}$ :*

$$\forall (t, \mathbf{x}) \in \Omega_{\mathcal{T}, \mathcal{R}}, \quad \|\mathbf{f}(t, \mathbf{x})\| \leq M. \quad (7.3)$$

- *The function  $\mathbf{f}$  satisfies the following Lipschitz condition: there is  $L > 0$  such that*

$$\forall ((t, \mathbf{x}_1), (t, \mathbf{x}_2)) \in \Omega_{\mathcal{T}, \mathcal{R}} \times \Omega_{\mathcal{T}, \mathcal{R}}, \quad \|\mathbf{f}(t, \mathbf{x}_1) - \mathbf{f}(t, \mathbf{x}_2)\| \leq L \|\mathbf{x}_1 - \mathbf{x}_2\|. \quad (7.4)$$

*Then there exists  $T \in (0, \mathcal{T}]$  depending on  $\mathcal{R}$ ,  $M$  and  $L$  such that the differential equation (7.2) has a local solution  $\mathbf{x}: [t_0 - T, t_0 + T] \rightarrow \mathbf{R}^n$ .*

*Proof.* Fix  $T \in (0, \mathcal{T}]$  and let  $I = [t_0 - T, t_0 + T]$ . Let also  $\mathcal{X}$  denote the following subset of

continuous functions defined from  $I$  to  $\mathbf{R}^n$ :

$$\mathcal{X} := \left\{ \mathbf{x} \in C(I, \mathbf{R}^n) : \sup_{t \in I} \|\mathbf{x}(t) - \mathbf{x}_0\| \leq \mathcal{R} \right\}$$

The set  $\mathcal{X}$  endowed with supremum metric is a closed subset of  $C(I, \mathbf{R}^n)$ . Since  $\mathcal{X}$  is a closed subset of a complete metric space, it is itself complete. Let  $\Phi: \mathcal{X} \rightarrow C(I, \mathbf{R}^n)$  denote the mapping

$$\Phi(\mathbf{x}): t \mapsto \mathbf{x}_0 + \int_0^t \mathbf{f}(s, \mathbf{x}(s)) \, ds.$$

The right-hand side, being the integral of a bounded function, is indeed a continuous function. We will show that, for  $T$  sufficiently small,

- the mapping  $\Phi$  maps  $\mathcal{X}$  into  $\mathcal{X}$ ;
- the mapping  $\Phi$  is a contraction.

From (7.3), it follows that

$$\forall \mathbf{x} \in \mathcal{X}, \quad \forall t \in I, \quad \|\Phi(\mathbf{x})(t) - \mathbf{x}_0\| = \left\| \int_{t_0}^t \mathbf{f}(s, \mathbf{x}(s)) \, ds \right\| \leq MT.$$

On the other hand, from the Lipschitz condition (7.4), it holds that

$$\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{X}, \quad \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\| \leq LT \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\|.$$

Therefore, it suffices to take  $T < \min \left\{ \mathcal{T}, \frac{\mathcal{R}}{M}, \frac{1}{L} \right\}$  to ensure that the above conditions are satisfied. For a value of  $T$  in this range, the Banach fixed point theorem, [Theorem A.3](#), gives the existence of a unique fixed point  $\mathbf{x}_* \in \mathcal{X}$  of  $\Phi$ . Since a fixed point of  $\Phi$  is a solution to (7.2) in view of the definition of  $\Phi$ , the statement is proved.  $\square$

It may seem at first glance that uniqueness of the solution to (7.2) follows from the uniqueness of the fixed point guaranteed by [Theorem A.3](#). However, this theorem implies uniqueness *only in the set  $\mathcal{X}$* , a property known as *conditional uniqueness*. In order to prove that the solution is unique over the full space  $C([t_0 - T, t_0 + T], \mathbf{R}^n)$ , additional assumptions and arguments are required. A simple approach is to rely on Grönwall's lemma.

**Lemma 7.2** (Grönwall's lemma, simplified integral form). *Suppose that  $u: [t_0 - T, t_0 + T] \rightarrow \mathbf{R}_{\geq 0}$  is continuous, nonnegative, and satisfies*

$$\forall t \in [t_0, t_0 + T], \quad u(t) \leq \alpha + \int_{t_0}^t \beta(s) u(s) \, ds, \quad (7.5)$$

*where  $\alpha \geq 0$  and  $\beta: [t_0, t_0 + T] \rightarrow \mathbf{R}_{\geq 0}$  is continuous and nonnegative. Then*

$$\forall t \in [t_0, t_0 + T], \quad u(t) \leq \alpha \exp \left( \int_{t_0}^t \beta(s) \, ds \right). \quad (7.6)$$

*Proof.* Assume first that  $\alpha > 0$ , so that the logarithm in (7.7) is well-defined. By the fundamental theorem of calculus and (7.5), it holds that

$$\frac{d}{dt} \left( \alpha + \int_{t_0}^t \beta(s) u(s) ds \right) \leq \beta(t) \left( \alpha + \int_{t_0}^t \beta(s) u(s) ds \right)$$

Therefore we have

$$\frac{d}{dt} \log \left( \alpha + \int_{t_0}^t \beta(s) u(s) ds \right) \leq \beta(t), \quad (7.7)$$

and after integrating and exponentiating, we obtain

$$\alpha + \int_{t_0}^t \beta(s) u(s) ds \leq \alpha \exp \left( \int_{t_0}^t \beta(s) ds \right)$$

The statement then follows by using (7.5) again. Assume next that  $\alpha = 0$ . If (7.5) is satisfied for  $\alpha = 0$ , then this condition is also satisfied for all  $\alpha > 0$ . Therefore the conclusion (7.6) holds for all  $\alpha > 0$ , and taking the limit  $\alpha \rightarrow 0$  in this equation, we obtain the statement.  $\square$

Note that the estimate (7.6) is sharp, since the function  $v: [t_0, t_0 + T]$  given by

$$v(t) = \alpha \exp \left( \int_{t_0}^t \beta(s) ds \right)$$

satisfies (7.5) with equality. We are now ready to prove uniqueness under an appropriate condition.

**Theorem 7.3** (Uniqueness of the solution). *Let  $\mathbf{x}_0 \in \mathbf{R}^n$  and let*

$$\Omega_{\mathcal{T}, \mathcal{R}} \{ (t, \mathbf{x}) \in \mathbf{R} \times \mathbf{R}^n : |t - t_0| \leq \mathcal{T} \text{ and } \|\mathbf{x} - \mathbf{x}_0\| \leq \mathcal{R} \},$$

*Assume that for all  $\mathcal{T} \in \mathbf{R}_{>0}$  and  $\mathcal{R} \in \mathbf{R}_{>0}$ , there is  $L_{\mathcal{T}, \mathcal{R}}$  such that*

$$\forall ((t, \mathbf{x}_1), (t, \mathbf{x}_2)) \in \Omega_{\mathcal{T}, \mathcal{R}} \times \Omega_{\mathcal{T}, \mathcal{R}}, \quad \|\mathbf{f}(t, \mathbf{x}_1) - \mathbf{f}(t, \mathbf{x}_2)\| \leq L_{\mathcal{T}, \mathcal{R}} \|\mathbf{x}_1 - \mathbf{x}_2\|. \quad (7.8)$$

*Then if  $\mathbf{x}_1$  and  $\mathbf{x}_2$  in  $C([t_0 - T, t_0 + T], \mathbf{R}^n)$  are local solutions to (7.2), it holds that  $\mathbf{x}_1 = \mathbf{x}_2$ .*

*Proof.* Suppose that  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are solutions to (7.2). Let  $I = [t_0 - T, t_0 + T]$  and

$$R := \max \left\{ \sup_{t \in I} \|\mathbf{x}_1(t) - \mathbf{x}_0\|, \sup_{t \in I} \|\mathbf{x}_2(t) - \mathbf{x}_0\| \right\} < \infty,$$

Since  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are solutions, it holds that

$$\forall t \in [t_0 - T, t_0 + T], \quad \mathbf{x}_1(t) - \mathbf{x}_2(t) = \int_{t_0}^t \left( \mathbf{f}(s, \mathbf{x}_1(s)) - \mathbf{f}(s, \mathbf{x}_2(s)) \right) ds.$$

Taking the norm and using (7.8), we obtain

$$\forall t \in [t_0, t_0 + T], \quad \|\mathbf{x}_1(t) - \mathbf{x}_2(t)\| \leq L_{T, R} \int_{t_0}^t \|\mathbf{x}_1(s) - \mathbf{x}_2(s)\| ds$$

Using Grönwall's lemma, we deduce that  $\mathbf{x}_1(t) = \mathbf{x}_2(t)$  for all  $t \in [t_0, t_0 + T]$ . A similar argument can be employed to show that  $\mathbf{x}_1 = \mathbf{x}_2$  on  $[t_0 - T, t_0]$ .  $\square$

**Corollary 7.4** (Maximal solutions). *Assume that  $\mathbf{f}$  is continuous in  $t$  and satisfies the local Lipschitz condition (7.8). Then there exists  $-\infty \leq T_- < T_+ \leq \infty$  such that  $t_0 \in (T_-, T_+)$  and the following properties are satisfied.*

- there exists a solution  $\mathbf{x}_*: (T_-, T_+) \rightarrow \mathbf{R}^n$  to (7.2);
- if  $\mathbf{x}: I \rightarrow \mathbf{R}^n$  is a local solution of (7.2), then  $I \subset (T_-, T_+)$  and  $\mathbf{x}(t) = \mathbf{x}_*(t)$  for all  $t \in I$ .
- If  $T_+$  is finite, then  $\lim_{t \rightarrow T_+} \|\mathbf{x}(t)\| = \infty$ , and if  $T_-$  is finite, then  $\lim_{t \rightarrow T_-} \|\mathbf{x}(t)\| = \infty$ .

The solution  $\mathbf{x}_*$  is called the maximal solution of (7.2).

*Proof.* Let  $\mathcal{I}$  denote the union of all the open intervals  $I$  such that there exists a solution in  $C(I, \mathbf{R}^n)$  to (7.2). The open set  $\mathcal{I}$  is connected and, by Theorem 7.1, it contains a neighborhood of  $t_0$ . Therefore  $\mathcal{I}$  is of the form  $(T_-, T_+)$ , where  $-\infty \leq T_- < t_0 < T_+ \leq \infty$ . In view of Theorem 7.3, all the local solutions coincide where they are defined, and so they can be patched together in order to construct a solution  $\mathbf{x}_*: (T_-, T_+) \rightarrow \mathbf{R}$ . It remains to prove the third item. To this end, suppose for contradiction that  $T_+$  was finite and that there was  $(t_n)_{n \in \mathbf{N}}$  such that  $t_n \rightarrow T_+$  in the limit  $n \rightarrow \infty$  and

$$K := \sup_{n \in \mathbf{N}} \|\mathbf{x}_*(t_n)\| < \infty.$$

Since  $\mathbf{f}$  is continuous, there is  $M$  such that  $|f(t, \mathbf{x})|$  is uniformly bounded from above by  $M$  for all  $(t, \mathbf{x}) \in [T_- - 1, T_+ + 1] \times B_{K+1}(\mathbf{0})$ . Furthermore, by the assumption (7.8), there is  $L$  such that for all  $t \in [T_- - 1, T_+ + 1]$ , the following Lipschitz condition holds:

$$\forall (\mathbf{x}_1, \mathbf{x}_2) \in B_{K+1}(\mathbf{0}) \times B_{K+1}(\mathbf{0}), \quad \|\mathbf{f}(t, \mathbf{x}_1) - \mathbf{f}(t, \mathbf{x}_2)\| \leq L \|\mathbf{x}_1 - \mathbf{x}_2\|.$$

Consequently, Theorem 7.1 with  $\mathcal{T} = \mathcal{R} = 1$  implies for all  $n$  the existence of a solution to

$$\begin{cases} \mathbf{x}'(t) = \mathbf{f}(t, \mathbf{x}(t)), \\ \mathbf{x}(t_n) = \mathbf{x}_*(t_n). \end{cases}$$

over the time interval  $[t_n - T, t_n + T]$ , where  $T > 0$  depends only on  $M$  and  $L$ , and not on  $n$ . But then, for  $n$  sufficiently large, this solution extends beyond  $T_+$ , which contradicts the maximality of  $\mathcal{I}$ . An analogous reasoning can be employed for  $T_-$ .  $\square$

*Example 7.1.* Consider the ODE

$$\begin{cases} x'(t) = x(t)^2, \\ x(0) = 1. \end{cases}$$

The maximal solution is  $x_*: (-\infty, 1) \rightarrow \mathbf{R}$  given by

$$\mathbf{x}_*(t) = \frac{1}{1-t}.$$

**Existence of a unique global solution.** In certain settings, it is possible to prove the maximal solution to (7.2) is globally defined for any initial condition. We discuss a few important examples.

- The first case is when  $\mathbf{f}: \mathbf{R} \times \mathbf{R}^n$  is globally Lipschitz in its second argument, with a Lipschitz constant that depends continuously on the first argument.
- The second case, generalizing the first, is when the growth of  $\mathbf{f}(t, \bullet)$  is at most affine:

$$\forall(t, \mathbf{x}) \in \mathbf{R} \times \mathbf{R}^n, \quad \|\mathbf{f}(t, \mathbf{x})\| \leq C(t) + L(t)\|\mathbf{x}\|,$$

with continuous constants  $C(t)$  and  $L(t)$ .

- The third case is when  $\mathbf{f}$  is independent of  $t$  and there is a function  $W \in C^1(\mathbf{R}^n)$  such that  $W(\mathbf{x}) \rightarrow \infty$  in the limit as  $\|\mathbf{x}\| \rightarrow \infty$  and

$$\forall \mathbf{x} \in \mathbf{R}^n, \quad \nabla W(\mathbf{x}) \cdot \mathbf{f}(\mathbf{x}) \leq c < \infty$$

Such a function is called a *Lyapunov function*.

The strategy of proof for global existence usually relies on an argument by contradiction. Consider for example the third setting. Since the assumptions of [Corollary 7.4](#) are satisfied, there exists a maximal solution  $\mathbf{x}_*: (T_-, T_+) \rightarrow \infty$ . Assume for contradiction that  $T_+$  is finite. Then the third item in [Corollary 7.4](#) implies that  $\lim_{t \rightarrow T_+} \|\mathbf{x}_*(t)\| \rightarrow \infty$ , and so  $W(\mathbf{x}_*(t))$  blows up as  $t$  approaches  $T_+$ . On the other hand, we have

$$\frac{d}{dt}W(\mathbf{x}_*(t)) = \nabla W(\mathbf{x}_*(t)) \cdot \mathbf{f}(\mathbf{x}_*(t)) \leq c.$$

Therefore  $\lim_{t \rightarrow T_+} W(\mathbf{x}_*(t)) \leq W(\mathbf{x}_*(t_0)) + |c|(T_+ - t_0)$ , which is a contradiction.

## 7.2 One-step methods

From now on, we assume for simplicity that  $t_0 = 0$  and that the initial value problem (7.1) admits a unique solution over the interval  $[0, T]$ . Most numerical methods for ODEs construct an approximation of the solution at discrete points:

$$\mathbf{x}_n \approx \mathbf{x}(t_n), \quad n = 0, 1, 2, \dots$$

The discretization points  $(t_n)_{n \in \mathbf{N}}$  are commonly equidistant, i.e.  $t_n = n\Delta$  where  $\Delta$  is the *discretization step*. Sometimes, it is useful to employ a variable time step, but we assume throughout this section that the time step is fixed, for simplicity. We begin in [Section 7.2.1](#) and [Section 7.2.2](#) by studying the simplest one-step methods, namely the forward and backward

Euler methods. Then, in Section 7.2.3, we present a general approach to the analysis of one-step methods. Finally, in Section 7.2.4, we present other widely used one-step methods in applications.

### 7.2.1 Forward Euler method

Assume that (7.1) has a unique solution  $\mathbf{x}(t)$  over the interval  $[0, T]$ . If  $\mathbf{x}(t)$  is twice continuously differentiable, then by Taylor's formula, we have

$$\mathbf{x}(t + \Delta) = \mathbf{x}(t) + \Delta \mathbf{f}(t, \mathbf{x}) + \frac{\Delta^2}{2} \mathbf{x}''(\tau), \quad \tau \in (t, t + \Delta). \quad (7.9)$$

This motivates a method known as the *forward* or *explicit* Euler method:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \mathbf{f}(t_n, \mathbf{x}_n),$$

with the same initial condition as for the continuous equation (7.1). The convergence of this method can be proved under a global Lipschitz assumption on the function  $\mathbf{f}$ .

**Theorem 7.5** (Convergence of the forward Euler method). *Assume that there is  $L \in \mathbf{R}_{>0}$  such that*

$$\forall (t, \mathbf{x}, \mathbf{y}) \in \mathbf{R} \times \mathbf{R}^n \times \mathbf{R}^n, \quad \|\mathbf{f}(t, \mathbf{x}) - \mathbf{f}(t, \mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|. \quad (7.10)$$

*Suppose in addition that there exists a unique, twice continuously differentiable of (7.1) over the interval  $[0, T]$ , and let*

$$M = \sup_{t \in [0, T]} \|\mathbf{x}''(t)\|$$

*Then the following error estimate holds:*

$$\forall n \in \left\{ 0, 1, \dots, \left\lfloor \frac{T}{\Delta} \right\rfloor \right\}, \quad \|\mathbf{x}(t_n) - \mathbf{x}_n\| \leq \frac{\Delta M}{2} \left( \frac{e^{Lt_n} - 1}{L} \right). \quad (7.11)$$

*Proof.* By Taylor's theorem, it holds that

$$\mathbf{x}(t_n) = \mathbf{x}(t_{n-1}) + \Delta \mathbf{f}(t_{n-1}, \mathbf{x}(t_{n-1})) + \frac{\Delta^2}{2} \boldsymbol{\alpha}_n, \quad \boldsymbol{\alpha}_n := 2 \int_0^1 (1-s) \mathbf{x}''(t_n + \Delta s) ds.$$

Notice that that  $\|\boldsymbol{\alpha}_n\| \leq M$ . Therefore, it holds that

$$\begin{aligned} \mathbf{x}(t_n) - \mathbf{x}_n &= \left( \mathbf{x}(t_{n-1}) + \Delta \mathbf{f}(t_{n-1}, \mathbf{x}(t_{n-1})) + \frac{\Delta^2}{2} \boldsymbol{\alpha}_n \right) - \left( \mathbf{x}_{n-1} + \Delta \mathbf{f}(t_{n-1}, \mathbf{x}_{n-1}) \right) \\ &= (\mathbf{x}(t_{n-1}) - \mathbf{x}_{n-1}) + \Delta \left( \mathbf{f}(t_{n-1}, \mathbf{x}(t_{n-1})) - \mathbf{f}(t_{n-1}, \mathbf{x}_{n-1}) \right) + \frac{\Delta^2}{2} \boldsymbol{\alpha}_n, \end{aligned}$$

Let  $\mathbf{e}_n = \mathbf{x}(t_n) - \mathbf{x}_n$  and  $\boldsymbol{\varepsilon}_n = \frac{\Delta^2}{2} \boldsymbol{\alpha}_n$ . The first term is the error at iteration  $n - 1$ , and the second may be bounded from (7.10), which gives

$$\|\mathbf{e}_n\| \leq (1 + \Delta L) \|\mathbf{e}_{n-1}\| + \|\boldsymbol{\varepsilon}_n\|.$$

The structure of this equation is important, as it appears in the analysis of all one-step methods for ODEs. The first term is an amplification of the error at the previous iteration, and the second term is an upper bound on the additional error introduced at step  $n$ . Applying this inequality to the previous time steps, we obtain

$$\begin{aligned} \|e_n\| &\leq (1 + \Delta L) \left( (1 + \Delta L) \|e_{n-2}\| + \|\varepsilon_{n-1}\| \right) + \|\varepsilon_n\| \\ &\leq \dots \leq (1 + \Delta L)^n \|e_0\| + \sum_{i=1}^n (1 + \Delta L)^{n-i} \|\varepsilon_i\|. \end{aligned} \quad (7.12)$$

Since  $\|\varepsilon_i\| \leq \Delta^2 M/2$ , we have by using the formula for geometric series that

$$\|e_n\| \leq (1 + \Delta L)^n \|e_0\| + \frac{(1 + \Delta L)^n - 1}{\Delta L} \left( \frac{\Delta^2 M}{2} \right).$$

The first term is zero because  $\|e_0\| = 0$ . Using the bound  $(1 + \Delta L)^n \leq (\exp(\Delta L))^n = e^{Lt_n}$  in the second term and rearranging, we finally obtain the statement (7.11).  $\square$

## 7.2.2 Backward Euler method

If we apply the Taylor expansion (7.9) backward around  $t + \Delta$ , instead of forward around  $t$ , then we obtain

$$\mathbf{x}(t) = \mathbf{x}(t + \Delta) - \Delta \mathbf{f}(t + \Delta, \mathbf{x}) + \mathcal{O}(\Delta^2).$$

This motivates the so-called *backward* or *implicit* Euler method:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \mathbf{f}(t_{n+1}, \mathbf{x}_{n+1}). \quad (7.13)$$

Observe that the right-hand side depends on  $\mathbf{x}_{n+1}$ . Therefore, given  $t_n$  and  $\mathbf{x}_n$ , this is a nonlinear equation for the unknown  $\mathbf{x}_{n+1}$ , which can be solved by using any of the methods studied in Chapter 5. Finding a solution to (7.13) amounts to finding a fixed point of the function

$$\mathbf{y} \mapsto \mathbf{F}(\mathbf{y}) := \mathbf{x}_n + \Delta \mathbf{f}(t_n, \mathbf{y}).$$

A priori, the existence and uniqueness of such a fixed point is not guaranteed. We proved in Theorem 5.2 that a sufficient condition for these two properties to hold is that  $\mathbf{F}$  is globally Lipschitz with a constant strictly less than 1, which holds if and only if the function  $\mathbf{y} \mapsto \mathbf{f}(t_n, \mathbf{y})$  is globally Lipschitz with a constant strictly less than  $1/\Delta$ . If the condition (7.10) holds, for example, then the backward Euler method (7.13) is guaranteed to be well defined for  $\Delta < \frac{1}{L}$ . Theorem 5.2 also ensures that, if  $\mathbf{F}$  is globally Lipschitz with a constant less than 1, then the fixed point can be approximated by using the iteration

$$\mathbf{y}_{k+1} = \mathbf{F}(\mathbf{y}_k). \quad (7.14)$$

and there is exponential convergence  $\mathbf{y}_k \rightarrow \mathbf{x}_{n+1}$  in the limit as  $k \rightarrow \infty$ . A natural starting point for (7.14) is  $\mathbf{y}_0 = \mathbf{x}_n$ . An alternative approach to the fixed point iteration (7.14) is to use the Newton–Raphson method for (7.13), which is faster in principle but must be initialized

sufficiently close to the fixed point.

Using a reasoning similar to that employed for proving [Theorem 7.5](#), we can prove the following result.

**Theorem 7.6** (Convergence of the backward Euler method). *If the assumptions of [Theorem 7.5](#) hold and  $\Delta < \frac{1}{L}$ , then the following error estimate holds:*

$$\forall n \in \left\{0, 1, \dots, \left\lfloor \frac{T}{\Delta} \right\rfloor\right\}, \quad \|\mathbf{x}(t_n) - \mathbf{x}_n\| \leq \frac{\Delta M}{2} \left( \frac{\left(\frac{1}{1-\Delta L}\right)^n - 1}{L} \right). \quad (7.15)$$

*Proof.* The proof is left as an exercise. □

*Remark 7.1.* Note that, if  $\Delta < \frac{1}{2L}$ , then

$$\begin{aligned} \frac{1}{1-\Delta L} &= 1 + \Delta L + (\Delta L)^2 + (\Delta L)^3 + (\Delta L)^4 \dots \\ &\leq 1 + \Delta L + (\Delta L)^2 + \frac{1}{2}(\Delta L)^2 + \frac{1}{4}(\Delta L)^2 + \dots \\ &\leq 1 + \Delta L + 2(\Delta L)^2 \leq \exp(\Delta L + (\Delta L)^2), \end{aligned}$$

and so the error estimate (7.15) gives

$$\|\mathbf{x}(t_n) - \mathbf{x}_n\| \leq \frac{\Delta M}{2} \left( \frac{\exp(Lt_n + \Delta L^2 t_n) - 1}{L} \right),$$

which makes it clear that the right-hand side of (7.15) is close, in absolute and relative terms, to that of (7.11) when  $\Delta \ll 1$ .

At this point, the reader may be wondering why one would use the backward Euler method instead of the forward Euler method, given that both methods have same order of convergence but iterations of the former are more computationally costly. The reason is that the backward Euler method, like many implicit methods, is more *stable* than its forward counterpart. Implicit methods are especially attractive in the context of *stiff* differential equations. We shall elaborate on this subject in [Section 7.4](#).

### 7.2.3 Analysis of general one-step methods

In general, one-step methods to solve differential equations are of the abstract form

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \Phi_{\Delta}(t_n, \mathbf{x}_n). \quad (7.16)$$

where  $\Phi_{\Delta}: \mathbf{R} \times \mathbf{R}^n \rightarrow \mathbf{R}^n$  is a function such that

$$\Phi_{\Delta}(t, \mathbf{x}) \approx \frac{1}{\Delta} \int_t^{t+\Delta} \mathbf{f}(s, \mathbf{x}^{t, \mathbf{x}}(s)) \, ds = \frac{\mathbf{x}^{t, \mathbf{x}}(t + \Delta) - \mathbf{x}}{\Delta}. \quad (7.17)$$



Here  $\mathbf{x}^{t,\mathbf{x}}$  denotes the solution to the differential equation (7.1) with initial condition  $\mathbf{x}(s) = \mathbf{x}$ . The main goal of this section is to establish general conditions, known as *consistency* and *stability*, under which the numerical scheme (7.16) is convergent. As we observed in the proof of Theorem 7.5 – specifically in equation (7.12) – the error at the final iteration for the forward Euler method is a sum of local errors, each amplified by a factor depending to the number of iterations left to reach the final time. *Consistency* of a numerical method enables to control the size of local errors when they arise, while *stability* enables to control their growth.

We emphasize that both the forward and the backward Euler methods can be recast in the form (7.16). For the forward Euler method  $\Phi_\Delta(t, \mathbf{x}) = \mathbf{f}(t, \mathbf{x})$ , while for the backward Euler method, the function  $\Phi_\Delta$  is defined implicitly as the function which to  $(t, \mathbf{x})$  associates the solution  $\phi \in \mathbf{R}^n$  to the equation

$$\phi = \mathbf{f}(t + \Delta, \mathbf{x} + \Delta\phi).$$

### Local truncation error and consistency

The local truncation error is the residual error obtained when substituting the exact solution of the differential equation in (7.16):

$$\boldsymbol{\eta}_{n+1} := \frac{\mathbf{x}(t_{n+1}) - \mathbf{x}(t_n)}{\Delta} - \Phi_\Delta(t_n, \mathbf{x}(t_n)).$$

Since there is a division by  $\Delta$ , the local truncation error has the same physical dimension as that of  $\mathbf{x}'$ , and so it should be viewed as an error *per time unit*.

**Definition 7.1** (Consistency). A numerical method is consistent if

$$\lim_{\Delta \rightarrow 0} \left( \max_{1 \leq n \leq N} \|\boldsymbol{\eta}_n\| \right) = 0, \quad N = \left\lfloor \frac{T}{\Delta} \right\rfloor.$$

It is consistent with order  $p$  if there exists  $C$  such that

$$\forall \Delta > 0, \quad \max_{1 \leq n \leq N} \|\boldsymbol{\eta}_n\| \leq C\Delta^p.$$

Proving the consistency of a numerical method is usually achieved on a case-by-case basis by application of Taylor's formula.

### Stability

The stability of a numerical method qualifies its sensitivity to perturbations. Roughly speaking, it expresses that small perturbation of the right-hand side of (7.16) lead to small perturbations of the numerical solution.

**Definition 7.2** (Stability). A numerical method of the form (7.16) is stable if there exists a

constant  $S(T) > 0$  independent of  $\Delta$  such that for all sequence  $(\mathbf{y}_n)_{1 \leq n \leq N}$  satisfying

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \Delta \Phi_\Delta(t_n, \mathbf{y}_n) + \Delta \boldsymbol{\delta}_{n+1}, \quad \mathbf{y}_0 = \mathbf{x}_0, \quad (7.18)$$

it holds that

$$\max_{1 \leq n \leq N} \|\mathbf{x}_n - \mathbf{y}_n\| \leq S(T) \Delta \sum_{n=1}^N \|\boldsymbol{\delta}_n\|. \quad (7.19)$$

It is convenient to introduce the following norms for sequences of vectors  $(\mathbf{u}_n)_{1 \leq n \leq N}$ :

$$\|\mathbf{u}_\bullet\|_{\ell_T^1} = \sum_{n=1}^N \|\mathbf{u}_n\|, \quad \|\mathbf{u}_\bullet\|_{\ell_T^\infty} = \max_{1 \leq n \leq N} \|\mathbf{u}_n\|,$$

With these notations, equation (7.19) may be rewritten compactly as follows:

$$\|\mathbf{x}_\bullet - \mathbf{y}_\bullet\|_{\ell_T^\infty} \leq S(T) \Delta \|\boldsymbol{\delta}_\bullet\|_{\ell_T^1}$$

One could argue that this equation is neater than (7.19); it bounds a norm of just one mathematical object, namely the sequence  $(\mathbf{x}_n - \mathbf{y}_n)_{1 \leq n \leq N}$ , by a norm of another object, namely the sequence  $(\boldsymbol{\delta}_n)_{1 \leq n \leq N}$ . Arguments for proving that a numerical scheme is stable often rely on some form of Lipschitz continuity. If the function  $\Phi_\Delta(t, \mathbf{y})$  is globally Lipschitz continuous with respect to  $\mathbf{y}$ , then stability is particularly simple to prove, as we now demonstrate.

**Proposition 7.7.** *Assume that there is  $L_\Phi > 0$  such that for all  $t \in [0, T]$  and  $\Delta > 0$ , the function  $\Phi_\Delta(t, \bullet)$  is globally Lipschitz continuous with constant  $L_\Phi$ . Then the one-step method (7.16) is stable.*

*Proof.* By (7.16) and (7.18), it holds that

$$\mathbf{x}_n - \mathbf{y}_n = \mathbf{x}_{n-1} - \mathbf{y}_{n-1} + \Delta \left( \Phi_\Delta(t_{n-1}, \mathbf{x}_{n-1}) - \Phi_\Delta(t_{n-1}, \mathbf{y}_{n-1}) \right) - \Delta \boldsymbol{\delta}_n.$$

Taking the Euclidean norm and employing the Lipschitz continuity assumption, we obtain

$$\|\mathbf{x}_n - \mathbf{y}_n\| \leq (1 + \Delta L_\Phi) \|\mathbf{x}_{n-1} - \mathbf{y}_{n-1}\| + \Delta \|\boldsymbol{\delta}_n\|.$$

By a reasoning similar to that in the proof of Theorem 7.5, we then obtain

$$\|\mathbf{x}_n - \mathbf{y}_n\| \leq (1 + \Delta L_\Phi)^n \|\mathbf{x}_0 - \mathbf{y}_0\| + \sum_{i=1}^n (1 + \Delta L_\Phi)^{n-i} \Delta \|\boldsymbol{\delta}_i\| \leq 0 + e^{L_\Phi t_n} \Delta \sum_{i=1}^n \|\boldsymbol{\delta}_i\|.$$

We conclude that (7.19) is satisfied with  $S(T) = e^{L_\Phi T}$ . □

## Convergence

We are now ready to prove that consistency and stability of the numerical (7.16) together imply convergence, in the sense that

$$\lim_{\Delta \rightarrow 0} \left( \max_{1 \leq n \leq N} \|\mathbf{x}(t_n) - \mathbf{x}_n\| \right) = 0, \quad N = \left\lfloor \frac{T}{\Delta} \right\rfloor.$$

This result is an instance of the *Lax equivalence theorem*, a pillar of numerical analysis with far-reaching applications.

**Theorem 7.8** (Consistence and stability imply convergence). *Assume that the one-step numerical method (7.16) is consistent and stable. Then the method is also convergent.*

*Proof.* By definition of the local truncation error, it holds that

$$\mathbf{x}(t_{n+1}) = \mathbf{x}(t_n) + \Delta \Phi_{\Delta}(t_n, \mathbf{x}(t_n)) + \Delta \boldsymbol{\eta}_{n+1}.$$

Therefore, the sequence  $(\mathbf{x}(t_n))_{1 \leq n \leq N}$  satisfies (7.18) with  $\boldsymbol{\delta}_n = \boldsymbol{\eta}_n$ , and so the stability estimate (7.19) implies that

$$\max_{1 \leq n \leq N} \|\mathbf{x}(t_n) - \mathbf{x}_n\| \leq S(T) \Delta \sum_{n=1}^N \|\boldsymbol{\eta}_n\|.$$

By consistency, the right-hand side converges to zero in the limit as  $\Delta \rightarrow 0$ , which concludes the proof.  $\square$

*Remark 7.2.* If we assume in [Theorem 7.8](#) that the method is consistent with order  $p$ , then by adapting the proof, we find that the error satisfies

$$\max_{1 \leq n \leq N} \|\mathbf{x}(t_n) - \mathbf{x}_n\| \leq CS(T) \Delta^p.$$

In this setting, the numerical scheme is said to be *convergent with order  $p$* .

### 7.2.4 Widely used one-step methods

In this section, we motivate and describe some of the other widely-used one-step methods, namely methods of Taylor and Runge–Kutta type. We assume in this section that the equation (7.1) admits a unique smooth solution over the interval  $[0, T]$ .

#### Taylor methods

In order to construct a method with a smaller local truncation error than that of the forward Euler method, a Taylor expansion of higher order than (7.9) can be employed:

$$\mathbf{x}(t + \Delta) = \mathbf{x}(t) + \Delta \mathbf{x}'(t) + \cdots + \frac{\Delta^p}{p!} \mathbf{x}^{(p)}(t) + \mathcal{O}(\Delta^{p+1}). \quad (7.20)$$

Since  $\mathbf{x}: [0, T] \rightarrow \mathbf{R}$  is a smooth solution to (7.1) by assumption, the time derivatives of  $\mathbf{x}$  can be obtained by differentiation of (7.1):

$$\mathbf{x}'(t) = \mathbf{f}(t, \mathbf{x}(t)), \quad \mathbf{x}''(t) = \partial_t \mathbf{f}(t, \mathbf{x}(t)) + \left( \mathbf{f}(t, \mathbf{x}(t)) \cdot \nabla_{\mathbf{x}} \right) \mathbf{f}(t, \mathbf{x}(t)), \quad \dots$$

In general, it is immediate to show inductively that  $\mathbf{x}^{(p)}(t) = \mathbf{f}^{(p-1)}(t, \mathbf{x}(t))$ , where the functions  $\mathbf{f}^{(p)}: \mathbf{R} \times \mathbf{R}^n \rightarrow \mathbf{R}$  are defined recursively from the following equation:

$$\mathbf{f}^{(p+1)} = \partial_t \mathbf{f}^{(p)}(t, \mathbf{x}(t)) + \left( \mathbf{f}(t, \mathbf{x}(t)) \cdot \nabla_{\mathbf{x}} \right) \mathbf{f}^{(p)}(t, \mathbf{x}(t)).$$

The Taylor expansion (7.20) motivates the so-called Taylor methods for integrating (7.1) numerically, which are based on the following iteration:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \mathbf{T}_{\Delta}^p(t_n, \mathbf{x}_n), \quad (7.21)$$

where

$$\mathbf{T}_{\Delta}^p(t, \mathbf{x}) := \mathbf{f}(t, \mathbf{x}) + \frac{\Delta}{2!} \mathbf{f}^{(1)}(t, \mathbf{x}) + \dots + \frac{\Delta^{p-1}}{p!} \mathbf{f}^{(p-1)}(t, \mathbf{x}).$$

Note that, for any  $p$ , the Taylor scheme (7.21) may be rewritten as

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \frac{d\mathbf{x}^{t_n, \mathbf{x}_n}}{dt}(t_n) + \dots + \frac{\Delta^p}{p!} \frac{d^p \mathbf{x}^{t_n, \mathbf{x}_n}}{dt^p}(t_n).$$

For  $p = 1$ , this scheme coincides with the forward Euler scheme.

### Runge–Kutta methods

Runge–Kutta methods resemble Taylor methods, but they do not require to calculate the derivatives of the function  $\mathbf{f}$ . This is achieved by approximating the derivatives in Taylor methods by difference quotients. Consider for example the Taylor method of order  $p = 2$ :

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \frac{d\mathbf{x}^{t_n, \mathbf{x}_n}}{dt}(t_n) + \frac{\Delta^2}{2!} \frac{d^2 \mathbf{x}^{t_n, \mathbf{x}_n}}{dt^2}(t_n). \quad (7.22)$$

Substituting the approximation

$$\begin{aligned} \frac{d^2 \mathbf{x}^{t_n, \mathbf{x}_n}}{dt^2}(t_n) &\approx \frac{1}{\Delta} \left( \frac{d\mathbf{x}^{t_n, \mathbf{x}_n}}{dt}(t_n + \Delta) - \frac{d\mathbf{x}^{t_n, \mathbf{x}_n}}{dt}(t_n) \right) \\ &= \frac{1}{\Delta} \left( \mathbf{f}(t_n + \Delta, \mathbf{x}^{t_n, \mathbf{x}_n}(t_n + \Delta)) - \mathbf{f}(t_n, \mathbf{x}_n) \right) \\ &\approx \frac{1}{\Delta} \left( \mathbf{f}(t_n + \Delta, \mathbf{x}_n + \Delta \mathbf{f}(t_n, \mathbf{x}_n)) - \mathbf{f}(t_n, \mathbf{x}_n) \right) \end{aligned} \quad (7.23)$$

in (7.22), we obtain an explicit method known as *Heun's method*:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \frac{\Delta}{2} \mathbf{f}(t_n, \mathbf{x}_n) + \frac{\Delta}{2} \mathbf{f}(t_n + \Delta, \mathbf{x}_n + \Delta \mathbf{f}(t_n, \mathbf{x}_n)).$$

It is possible to show that the local truncation error for this method also scales as  $\Delta^2$ . Heun's method is a particular instance of a Runge–Kutta method. In general, an explicit Runge–Kutta

method with  $s$  stages is of the form

$$\begin{aligned}\mathbf{x}_{n+1} &= \mathbf{x}_n + \Delta \sum_{i=1}^s b_i \mathbf{k}_i \\ \mathbf{k}_1 &= \mathbf{f}(t_n, \mathbf{x}_n), \\ \mathbf{k}_2 &= \mathbf{f}(t_n + c_2 \Delta, \mathbf{x}_n + \Delta(a_{21} \mathbf{k}_1)), \\ \mathbf{k}_3 &= \mathbf{f}(t_n + c_3 \Delta, \mathbf{x}_n + \Delta(a_{31} \mathbf{k}_1 + a_{32} \mathbf{k}_2)), \\ &\vdots \\ \mathbf{k}_s &= \mathbf{f}\left(t_n + c_s \Delta, \mathbf{x}_n + \Delta \sum_{j=1}^{s-1} a_{sj} \mathbf{k}_j\right),\end{aligned}$$

with appropriate coefficients  $c_i$  and  $a_{ij}$ . Heun's iteration can be recast in this form as follows:

$$\begin{aligned}\mathbf{x}_{n+1} &= \mathbf{x}_n + \frac{\Delta}{2}(\mathbf{k}_1 + \mathbf{k}_2) \\ \mathbf{k}_1 &= \mathbf{f}(t_n, \mathbf{x}_n) \\ \mathbf{k}_2 &= \mathbf{f}(t_n + \Delta, \mathbf{x}_n + \Delta \mathbf{k}_1).\end{aligned}$$

The approach we employed to construct Heun's method may be generalized to higher orders. For example, the most widely known Runge–Kutta method approximates the Taylor method of order  $p = 4$  with the following iteration:

$$\begin{aligned}\mathbf{x}_{n+1} &= \mathbf{x}_n + \frac{\Delta}{6}(\mathbf{k}_1 + 2\mathbf{k}_2 + 2\mathbf{k}_3 + \mathbf{k}_4), \\ \mathbf{k}_1 &= \mathbf{f}(t_n, \mathbf{x}_n), & \mathbf{k}_2 &= \mathbf{f}\left(t_n + \frac{\Delta}{2}, \mathbf{x}_n + \Delta \frac{\mathbf{k}_1}{2}\right), \\ \mathbf{k}_3 &= \mathbf{f}\left(t_n + \frac{\Delta}{2}, \mathbf{x}_n + \Delta \frac{\mathbf{k}_2}{2}\right), & \mathbf{k}_4 &= \mathbf{f}(t_n + \Delta, \mathbf{x}_n + \Delta \mathbf{k}_3).\end{aligned}$$

The local truncation error for this method scales as  $\Delta^4$  and, when  $\mathbf{f}(t, \mathbf{x}) = \mathbf{f}(t)$ , this method coincides with Simpson's formula (3.6) for the approximation of the integral in (7.17). The systematic derivation of Runge–Kutta methods is cumbersome, and so we do not address this issue in this course.

*Remark 7.3.* Explicit Runge–Kutta methods of a given order are not uniquely defined. For example, if we employ instead of (7.23) the approximation

$$\frac{d^2 \mathbf{x}^{t_n, \mathbf{x}_n}}{dt^2}(t_n) \approx \frac{2}{\Delta} \left( \mathbf{f}\left(t_n + \frac{\Delta}{2}, \mathbf{x}_n + \frac{\Delta}{2} \mathbf{f}(t_n, \mathbf{x}_n)\right) - \mathbf{f}(t_n, \mathbf{x}_n) \right),$$

then we obtain by substitution in (7.22) the so-called *explicit midpoint method*, which is also a Runge–Kutta method of order 2:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \mathbf{f}\left(t_n + \frac{1}{2} \Delta, \mathbf{x}_n + \frac{\Delta}{2} \mathbf{f}(t_n, \mathbf{x}_n)\right).$$

### Implicit methods

To conclude this section, we mention two common implicit methods with a better order of convergence than that of the backward Euler method.

- The Crank–Nicolson method:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \frac{\Delta}{2} (f(t_n, \mathbf{x}_n) + f(t_n + \Delta, \mathbf{x}_{n+1})). \quad (7.24)$$

When  $\mathbf{f}$  is independent of  $\mathbf{x}$  and depends only on  $t$ , this method coincides with the trapezoidal rule for numerical integration.

- The implicit midpoint method:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta f \left( t_n + \frac{\Delta}{2}, \frac{\mathbf{x}_n + \mathbf{x}_{n+1}}{2} \right).$$

Similarly to the backward Euler method, each iteration of these methods requires the resolution of a nonlinear equation. Implicit methods often enjoy better stability than their explicit counterparts. This subject is further discussed in [Section 7.4](#).

## 7.3 Multistep methods

The idea of multistep methods is to use, in the construction of a new iterate, information from not only the current but also previous iterations. This degree of freedom enables to construct more economical numerical methods than one-step methods for the same order of convergence, at the cost of a more difficult initialization. In this section we focus on *linear* multistep methods of the form

$$\begin{aligned} \mathbf{x}_{n+1} = & a_0 \mathbf{x}_n + a_1 \mathbf{x}_{n-1} + \cdots + a_k \mathbf{x}_{n-k} \\ & + \Delta \left( b_{-1} \mathbf{f}(t_{n+1}, \mathbf{x}_{n+1}) + b_0 \mathbf{f}(t_n, \mathbf{x}_n) + \cdots + b_k \mathbf{f}(t_{n-k}, \mathbf{x}_{n-k}) \right). \end{aligned} \quad (7.25)$$

This equation defines an explicit method if  $b_{-1} = 0$ , and an implicit method if  $b_{-1} \neq 0$ . Note that explicit methods of the form (7.25) require only one additional evaluation  $f(t_n, \mathbf{x}_n)$  per iteration, in contrast with Runge–Kutta methods. When  $b_{-1} \neq 0$ , the iteration (7.25) is a nonlinear equation for the unknown  $\mathbf{x}_{n+1}$ , which must itself be solved by resorting to a numerical method.

**Initialization.** In order to initiate the numerical method (7.25), the values  $\mathbf{x}_0, \dots, \mathbf{x}_k$  are required. These can be calculated by using a one-step method with an order of convergence matching that of the multistep method.

**Local truncation error.** Consistently with the setting of one-step methods, the local truncation error for (7.25) is defined as the residual error left when the exact solution is substituted in the

numerical scheme:

$$\begin{aligned} \Delta \boldsymbol{\eta}_{n+1} := & \boldsymbol{x}(t_n + \Delta) - a_0 \boldsymbol{x}(t_n) - a_1 \boldsymbol{x}(t_n - \Delta) - \cdots - a_k \boldsymbol{x}(t_n - k\Delta) \\ & - \Delta \left( b_{-1} \boldsymbol{x}'(t_n + \Delta) + b_0 \boldsymbol{x}'(t_n) + \cdots + b_k \boldsymbol{x}'(t_n - k\Delta) \right). \end{aligned} \quad (7.26)$$

The multistep method (7.25) is said to be of order  $p$  if the maximum local truncation error over all the discretization points, in norm, scales as  $\mathcal{O}(\Delta^p)$ . The following result is useful for estimating the order of consistency of a linear multistep method.

**Proposition 7.9.** *The linear multistep method (7.25) is consistent with order  $p$  for any smooth  $\boldsymbol{x}: [0, T] \rightarrow \mathbf{R}^n$  if and only if the local truncation error (7.26) is everywhere zero when  $\boldsymbol{x}(t)$  is of the scalar form*

$$x(t) = t^q, \quad q \in \{0, \dots, p\}. \quad (7.27)$$

*Proof.* Assume that the method is consistent with order  $p$ , fix  $q \in \{1, \dots, p\}$ , and let  $x(t) = t^q$ . Fix also  $t \in [0, T]$  and consider the function  $\xi: \{\Delta: t/\Delta \in \mathbf{N}_{>0}\} \rightarrow \mathbf{R}$  given by

$$\begin{aligned} \Delta \xi(\Delta) = \Delta \eta_{(t/\Delta)+1} = & x(t + \Delta) - a_1 x(t) - a_2 x(t - \Delta) - \cdots - a_k x(t - (k-1)\Delta) \\ & - \Delta \left( b_0 x'(t + \Delta) + b_1 x'(t) + \cdots + b_k x'(t - (k-1)\Delta) \right). \end{aligned}$$

The quantity  $\xi(\Delta)$  should be understood as the local truncation error at  $t$  for time step  $\Delta$ . It is a polynomial in  $\Delta$  of degree at most  $p$  and scaling as  $\mathcal{O}(\Delta^{p+1})$ . Therefore, it holds necessarily that  $\xi(\Delta) = 0$ .

Conversely, assume that the right-hand side of (7.26) is equal to zero for any function of the form (7.27). If  $\boldsymbol{x}(t)$  denotes a smooth solution of (7.1), then by Taylor's theorem there is  $C > 0$  independent of  $t_n$  such that

$$\forall t \in [0, T], \quad \begin{cases} \|\boldsymbol{x}(t) - \boldsymbol{y}(t)\| \leq C|t - t_n|^{p+1} \\ \|\boldsymbol{x}'(t) - \boldsymbol{y}'(t)\| \leq C|t - t_n|^p \end{cases}, \quad \boldsymbol{y}(t) := \boldsymbol{x}(t_n) + \sum_{i=1}^p \boldsymbol{e}_i (t - t_n)^i,$$

for appropriate vectors  $\boldsymbol{e}_i \in \mathbf{R}^n$  depending on  $t_n$ . Substituting  $\boldsymbol{x}(t) = \boldsymbol{y}(t) + (\boldsymbol{x}(t) - \boldsymbol{y}(t))$  in the right-hand side of (7.26), we obtain

$$\Delta \|\boldsymbol{\eta}_{n+1}\| = \mathcal{O}(\Delta^{p+1}) + \Delta \mathcal{O}(\Delta^p) = \mathcal{O}(\Delta^{p+1}),$$

with the constant implicit in the big  $\mathcal{O}$  notation independent of  $n$ . This concludes the proof.  $\square$

*Example 7.2.* In the one-dimensional setting, we wish to find parameters  $a_0$ ,  $a_1$  and  $b_1$  such that the order of consistency of the following multistep scheme is as high as possible:

$$x_{n+1} = a_0 x_n + a_2 x_{n-1} + b_0 \Delta f(t_n, x_n).$$

Substituting  $x(t) = 1$  in the formula (7.26) for the local truncation error, we obtain

$$\eta_{n+1} = x(t_n + \Delta) - a_0x(t_n) - a_1x(t_n - \Delta) - b_0\Delta x'(t_n) = 1 - a_0 - a_1.$$

Therefore  $a_1 = (1 - a_0)$ . Next, substituting  $x(t) = t - t_n$  in (7.26), we obtain

$$\eta_{n+1} = \Delta(2 - a_0 - b_0),$$

which gives  $b_0 = 2 - a_0$ . Finally, substituting  $x(t) = (t - t_n)^2$ , we obtain

$$\eta_{n+1} = \Delta^2 a_0,$$

and so  $a_0 = 1$ . We conclude that the best parameters, leading to a local truncation error scaling as  $\mathcal{O}(\Delta^2)$ , are given by  $a_0 = 0$ ,  $a_1 = 1$  and  $b_0 = 2$ . The resulting method reads

$$x_{n+1} = x_{n-1} + 2\Delta f(t_n, x_n),$$

and is known as the *multistep midpoint method*.

We now present two widely used systematic approaches for constructing multistep methods, known as the Adams–Bashforth and Adams–Moulton approaches.

### 7.3.1 Adams–Bashforth methods

Let  $\mathbf{x}: [0, T] \rightarrow \mathbf{R}^n$  denote the exact solution to the differential equation (7.1). Integrating this equation between  $t_n$  and  $t_{n+1}$ , we obtain

$$\mathbf{x}(t_{n+1}) = \mathbf{x}(t_n) + \int_{t_n}^{t_{n+1}} \mathbf{f}(t, \mathbf{x}(t)) dt. \quad (7.28)$$

The key idea of the Adams–Bashforth method is to approximate the function  $t \mapsto \mathbf{f}(t, \mathbf{x}(t))$  by the interpolating polynomial  $\widehat{\mathbf{f}}$  of degree  $k$  at the nodes  $t_{n-k}, \dots, t_n$ :

$$\widehat{\mathbf{f}}(t) = \sum_{i=0}^k \mathbf{f}(t_{n-i}, \mathbf{x}(t_{n-i})) L_i(t), \quad L_i(t) := \prod_{\substack{j=0 \\ j \neq i}}^k \frac{t - t_{n-j}}{t_{n-i} - t_{n-j}}. \quad (7.29)$$

Substituting this approximation in (7.28), we obtain

$$\mathbf{x}(t_{n+1}) \approx \mathbf{x}(t_n) + \sum_{i=0}^k \mathbf{f}(t_{n-i}, \mathbf{x}(t_{n-i})) \int_{t_n}^{t_{n+1}} L_i(t) dt.$$

This motivates the following *Adams–Bashforth* numerical scheme:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \sum_{i=0}^k b_i \mathbf{f}(t_{n-i}, \mathbf{x}_{n-i}), \quad b_i := \int_0^1 \prod_{\substack{j=0 \\ j \neq i}}^k \frac{s+j}{-i+j} ds. \quad (7.30)$$



Since the Lagrange polynomials  $(L_i)_{0 \leq i \leq k}$  depend on  $k$ , so do the coefficients  $b_i$ . However, these are independent of  $\Delta$ , and so they can be tabulated. The value of these coefficients for the first few Adams–Bashforth methods are collated in Table 7.1.

$i$	0	1	2	3
$k = 0$	1			
$k = 1$	$\frac{3}{2}$	$-\frac{1}{2}$		
$k = 2$	$\frac{23}{12}$	$-\frac{16}{12}$	$\frac{5}{12}$	
$k = 3$	$\frac{55}{24}$	$-\frac{59}{24}$	$\frac{37}{24}$	$-\frac{9}{24}$

Table 7.1: Coefficients  $(b_i)_{0 \leq i \leq k}$  of the Adams–Bashforth methods.

**Local truncation error.** Assuming  $\mathbf{x} \in C^{k+2}([0, T], \mathbf{R}^n)$  and applying Theorem 2.3 for the interpolation error component-wise, we obtain

$$\forall t \in [0, T], \quad \left\| \mathbf{x}'(t) - \widehat{\mathbf{f}}(t) \right\|_{\infty} \leq \frac{|t - t_{n-k}| \cdots |t - t_n|}{(k+1)!} \sup_{t \in [0, T]} \left\| \mathbf{x}^{(k+2)}(t) \right\|_{\infty},$$

where  $\widehat{\mathbf{f}}$  is the function defined in (7.29). Since

$$\Delta \boldsymbol{\eta}_{n+1} = \mathbf{x}(t_{n+1}) - \mathbf{x}(t_n) - \Delta \sum_{i=0}^k b_i \mathbf{f}(t_{n-i}, \mathbf{x}(t_{n-i})) = \int_{t_n}^{t_{n+1}} (\mathbf{x}'(t) - \widehat{\mathbf{f}}(t)) dt,$$

we deduce that

$$\|\boldsymbol{\eta}_{n+1}\| \leq C_k M_{k+2} \Delta^{k+1}, \quad M_{k+2} := \sup_{t \in [0, T]} \left\| \mathbf{x}^{(k+2)}(t) \right\|_{\infty}, \quad (7.31)$$

for an appropriate numerical constant  $C_k$  independent of  $n$  and of the problem data. Therefore the Adams–Bashforth method (7.30) is consistent with order  $k+1$ .

**Convergence.** By using a reasoning similar to that in the proof of Theorem 7.5, we can prove a convergence result of the Adams–Bashforth method.

**Theorem 7.10.** *Assume that the solution  $\mathbf{x}: [0, T] \rightarrow \mathbf{R}^n$  to (7.1) is  $k+2$  times continuously differentiable and that the global Lipschitz condition (7.10) is satisfied. Suppose also that*

$$\forall i \in \{0, \dots, k\}, \quad \|\mathbf{x}(t_i) - \mathbf{x}_i\| \leq \delta.$$

*Then the following error estimate holds for the Adams–Bashforth method (7.30):*

$$\forall n \in \left\{ 0, 1, \dots, \left\lfloor \frac{T}{\Delta} \right\rfloor \right\}, \quad \|\mathbf{x}(t_n) - \mathbf{x}_n\| \leq \delta e^{LB} + C_k M_{k+2} \Delta^{k+1} \left( \frac{e^{LBt_n} - 1}{LB} \right),$$

*where  $C_k$  and  $M_{k+2}$  are the constants from (7.31), and with  $B := |b_0| + \dots + |b_k|$ .*

*Sketch of proof.* Let  $\mathbf{e}_n := \mathbf{x}(t_{n+1}) - \mathbf{x}_{n+1}$ . From the equation

$$\mathbf{x}(t_{n+1}) - \mathbf{x}_{n+1} = \mathbf{x}(t_n) - \mathbf{x}_n + \Delta \sum_{i=0}^k b_i \left( \mathbf{f}(t_{n-i}, \mathbf{x}(t_{n-i})) - \mathbf{f}(t_{n-i}, \mathbf{x}_{n-i}) \right) + \Delta \boldsymbol{\eta}_{n+1},$$

which is valid for  $n \geq k$ , we deduce that

$$\max \left\{ \|\mathbf{e}_0\|, \dots, \|\mathbf{e}_{n+1}\| \right\} \leq (1 + \Delta LB) \max \left\{ \|\mathbf{e}_0\|, \dots, \|\mathbf{e}_n\| \right\} + C_k M_{k+2} \Delta^{k+2}.$$

Since  $\max \left\{ \|\mathbf{e}_0\|, \dots, \|\mathbf{e}_k\| \right\} \leq \delta$  by assumption, the statement easily follows.  $\square$

### 7.3.2 Adams–Moulton methods

The Adams–Moulton methods are very similar to their Adams–Bashforth cousins. The only difference is that the former are obtained by interpolating the function  $t \mapsto \mathbf{f}(t, \mathbf{x}(t))$  in (7.28) at nodes shifted forward by  $\Delta$ , i.e. at the nodes  $t_{n-k+1}, \dots, t_{n+1}$ . This leads to the method

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \sum_{i=-1}^{k-1} b_i \mathbf{f}(t_{n-i}, \mathbf{x}_{n-i}), \quad b_i := \int_0^1 \prod_{\substack{j=-1 \\ j \neq i}}^{k-1} \frac{s+j}{-i+j} ds. \quad (7.32)$$

Unlike the Adams–Bashforth methods, which are *explicit*, the Adams–Moulton methods are *implicit*. The value of the coefficients for the first few Adams–Moulton methods are collated in Table 7.2. Notice that, for  $k = 0$ , the Adams–Moulton method coincides with the backward Euler method, and for  $k = 1$  it coincides with the Crank–Nicolson method.

$i$	$-1$	$0$	$1$	$2$
$k = 0$	$1$			
$k = 1$	$\frac{1}{2}$	$\frac{1}{2}$		
$k = 2$	$\frac{5}{12}$	$\frac{8}{12}$	$-\frac{1}{12}$	
$k = 3$	$\frac{9}{24}$	$\frac{19}{24}$	$-\frac{5}{24}$	$\frac{1}{24}$

Table 7.2: Coefficients  $(b_i)_{0 \leq i \leq k}$  of the Adams–Moulton methods.

## 7.4 Absolute stability

To conclude this chapter, we introduce the notion of *absolute stability* and explain its relevance. Absolute stability is a property of a numerical method in relation to the model equation

$$\begin{cases} x'(t) = \lambda x(t), \\ x(0) = 1. \end{cases} \quad (7.33)$$

A numerical scheme for approximating (7.33) is called *absolutely stable* if

$$|x_n| \rightarrow 0 \quad \text{in the limit as } n \rightarrow \infty. \quad (7.34)$$

where  $(x_n)_{n=0,1,\dots}$  denotes the numerical solution to (7.33). Whether a numerical method is absolutely stable or not depends on the parameters  $\lambda$  and  $\Delta$ .

*Example 7.3.* The forward Euler method for (7.33) reads

$$x_{n+1} = x_n + \Delta\lambda x_n = (1 + \Delta\lambda)x_n.$$

Therefore  $x_n \rightarrow 0$  if and only if  $|1 + \Delta\lambda| \leq 1$ .

As [Example 7.3](#) illustrates, whether absolute stability holds for the forward Euler methods depends only the value of the product  $\Delta\lambda \in \mathbf{C}$ . This dependence on  $\lambda$  and  $\Delta$  only through the product  $\Delta\lambda$  holds in fact generally. Indeed, all the numerical schemes we considered in this chapter are invariant under linear time rescaling of the ordinary differential equation: the numerical solution of the rescaled equation, when the time step is rescaled by the same factor, coincides with the discrete function obtained by linear rescaling of the numerical solution to the original equation. This motivates the definition of *absolute stability region* as

$$\mathcal{A} := \{z \in \mathbf{C} : (7.34) \text{ holds when } \Delta\lambda = z\} \subset \mathbf{C}.$$

The exact solution to the model equation (7.33) diverges to  $\infty$  as  $t \rightarrow \infty$  if  $\Re(\lambda) > 0$ , and it converges to 0 if  $\Re(\lambda) < 0$ . Numerical schemes which exhibit a similar property at the discrete level are called *A-stable*. More precisely, a numeric method is A-stable if the absolute stability region  $\mathcal{A}$  contains the left half-plane, i.e. if

$$\{z \in \mathbf{C} : \Re(z) < 0\} \subset \mathcal{A}.$$

Before investigating whether the numerical schemes introduced previously in this chapter are absolutely stable, we address the following natural question: why focus on the simple model equation (7.33)? We provide a couple of motivations:

- First, note that equations of the form (7.33) are more relevant in science than might appear at first glance. Indeed, when discretizing in space a linear parabolic partial differential equation, one often obtains a linear differential equation of the form

$$\mathbf{x}'(t) = \mathbf{A}\mathbf{x},$$

where  $\mathbf{A} \in \mathbf{C}^{n \times n}$ . If the matrix  $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*$  is diagonalizable, then the vector  $\mathbf{z}(t) := \mathbf{Q}^*\mathbf{x}(t)$  satisfies the differential equation

$$\mathbf{z}'(t) = \mathbf{D}\mathbf{z}.$$

In other words, each component of  $\mathbf{z}$  satisfies an ordinary differential equation of the same form as the model equation (7.33). In applications, the components of  $\mathbf{z}$  often encode the

amplitudes of Fourier modes of the solution to the partial differential equation, and for dissipative equations all the eigenvalues of  $\mathbf{A}$  have a negative real part. However, the spectral radius of  $\mathbf{A}$  usually diverges as the number of discretization points increases. In this context, A-stability is particularly attractive, as it ensures that the numerical approximation remains well-behaved as the number of discretization points increases.

- Second, the model equation (7.33) may be viewed as a linearized approximation of a more interesting equation. Consider, for example, the following one-dimensional autonomous differential equation:

$$\begin{cases} x'(t) = f(x(t)), \\ x(0) = x_0. \end{cases} \quad (7.35)$$

Assume that  $f(x_*) = 0$  for some  $x_* \in \mathbf{R}$ . Such a point is called a *critical point* of the differential equation. If  $f'(x_*) < 0$ , then  $x_*$  is an attractor of the equation, in the sense that  $x(t) \rightarrow x_*$  provided that  $x(0)$  is sufficiently close to  $x_*$ . This result, which is the counterpart of Proposition 5.5 for differential equations, is a particular case of a theorem due to Poincaré and Lyapunov; see [16, Theorem 7.1]. If  $|x_0 - x_*|$  is sufficiently small, then the solution to (7.35) is expected to be close to that of the linearized equation

$$\begin{cases} y'(t) = f'(x_*)(y(t) - x_*), \\ y(0) = x_0, \end{cases} \quad (7.36)$$

which is of the form (7.33). Often, studying the linearized equation (7.36) enables to gain insight into the behavior of the original equation (7.35), and analyzing the performance of a numerical method for the linearized equation (7.36) is useful to inform the choice of a numerical scheme for (7.35).

- More generally, if  $x(t)$  and  $x_\varepsilon(t)$  are respectively the solutions to

$$\begin{cases} x'(t) = f(t, x(t)), \\ x(0) = x_0 + \varepsilon, \end{cases} \quad \text{and} \quad \begin{cases} x'_\varepsilon(t) = f(t, x_\varepsilon(t)), \\ x_\varepsilon(0) = x_0 + \varepsilon, \end{cases} \quad (7.37)$$

then the difference  $e(t) := x_\varepsilon(t) - x(t)$  satisfies the equation

$$\begin{aligned} e'(t) &= f(t, x_\varepsilon(t)) - f(t, x(t)) \approx \partial_x f(t, x(t))e(t), \\ e(0) &= \varepsilon, \end{aligned} \quad (7.38)$$

which looks similar to (7.33) with  $\partial_x f(t, x(t))$  in place of  $\lambda$ . At a given time  $t$ , the solutions tend to converge to each other as time increases if  $\partial_x f(t, x(t)) < 0$ , and diverge from each other if  $\partial_x f(t, x(t)) > 0$ . Testing absolute stability with  $\lambda = \partial_x f(t, x(t))$  enables to determine whether this property holds true also at the discrete level. Although the latter statement is difficult to state precisely and prove generally, we illustrate its validity for the forward Euler method in Example 7.4.

*Example 7.4.* Let  $(x_n)$  and  $(x_n^\varepsilon)$  denote the numerical solutions obtained by applying the forward Euler method to the differential equations in (7.37). If  $\varepsilon \ll 1$ , then

$$\begin{aligned} x_{n+1}^\varepsilon - x_{n+1} &= x_n^\varepsilon - x_n + \Delta f(t_n, x_n^\varepsilon) - \Delta f(t_n, x_n) \\ &\approx x_n^\varepsilon - x_n + \Delta \partial_x f(t_n, x_n)(x_n^\varepsilon - x_n) = (1 + \Delta \partial_x f(t_n, x_n))(x_n^\varepsilon - x_n). \end{aligned}$$

Therefore, the numerical solutions  $(x_n^\varepsilon)$  and  $(x_n)$  tend to become closer as  $n$  increases if

$$\Delta \partial_x f(t_n, x_n) \in \mathcal{A}. \quad (7.39)$$

The absolute stability regions of the forward and backward Euler methods are illustrated in green in Figure 7.1. For the forward Euler method, absolute stability holds if and only if  $|1 + \Delta\lambda| < 1$ , as we proved in Example 7.3. A similar reasoning gives that the absolute stability region for backward Euler method is given by  $\{z \in \mathbf{C} : |1 - z|^{-1} < 1\}$ . The backward Euler method is A-stable but the forward Euler method is not. Notice that, if the time step is sufficiently large, then the backward Euler method is absolutely stable even for values of  $\lambda$  with a positive real part, for which exact solutions to the model equation are divergent.

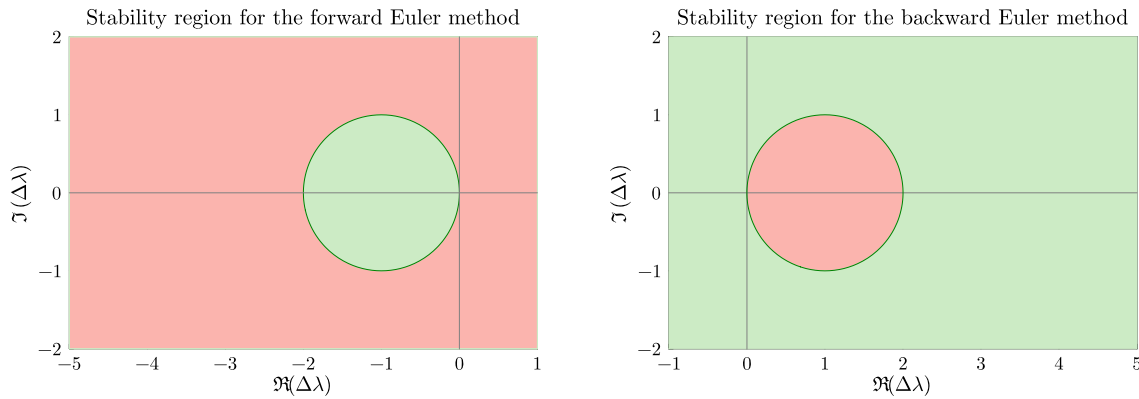


Figure 7.1: Absolute stability regions for the forward (left) and backward (right) Euler methods.

*Example 7.5* (Absolute stability region of the Taylor methods). When applied to (7.33), the Taylor method of order  $p$  given in (7.21) reads

$$x_{n+1} = \left( 1 + \Delta\lambda + \frac{\Delta^2 \lambda^2}{2} + \cdots + \frac{\Delta^p \lambda^p}{p!} \right) x_n.$$

Thus, the absolute stability region is given by

$$\left\{ z \in \mathbf{C} : \left| 1 + z + \frac{z^2}{2} + \cdots + \frac{z^p}{p!} \right| < 1 \right\}.$$

This region is illustrated for various values of  $p$  in Figure 7.2. We observe that the absolute stability region grows as  $p$  increases.

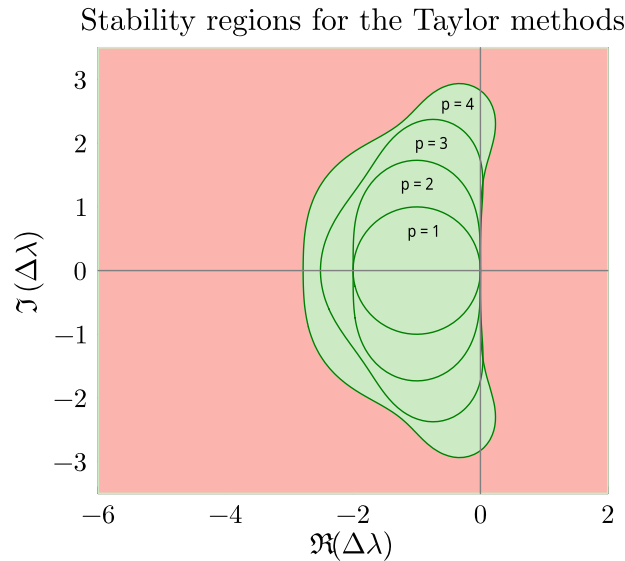


Figure 7.2: Stability regions for the first few Taylor methods.

### Stiff differential equations

In the context of ordinary differential equations, stiffness is not a precisely defined concept, but rather a generic term employed to describe equations with widely separated time scales. Roughly speaking, a differential equation of the form (7.1) is called stiff if the Jacobian matrix of  $f$ , with respect to the variable  $\mathbf{x}$ , has at least one eigenvalue with a large negative real part. In the one-dimensional setting, the solutions to stiff differential equations which are close at the initial time tend to converge quickly to each other, in view of (7.38). This is illustrated in Example 7.6.

*Example 7.6* (Stiff differential equation). Consider the following equation [7, Chapter 4]:

$$\begin{cases} x'(t) = -\alpha(x(t) - \sin(t)) + \cos(t) \\ x(0) = x_0 \end{cases} \quad (7.40)$$

The exact solution to this equation is given by

$$x(t) = \sin(t) + x_0 e^{-\alpha t}.$$

When  $\alpha \in \mathbf{R}$  is large, the distance between the solution and the function  $t \mapsto \sin(t)$  converges to zero very quickly, regardless of the initial condition. This behavior is illustrated in Figure 7.3.

In the rest of this section, we use the differential equation (7.40) as a guiding example. For this problem, we have  $\partial_x f(t, x) = \alpha$ . Therefore, in view of (7.39), we expect that the forward Euler scheme is non-divergent if  $|1 - \alpha\Delta| < 1$ , i.e. if

$$\Delta < \Delta_* = \frac{2}{\alpha}.$$

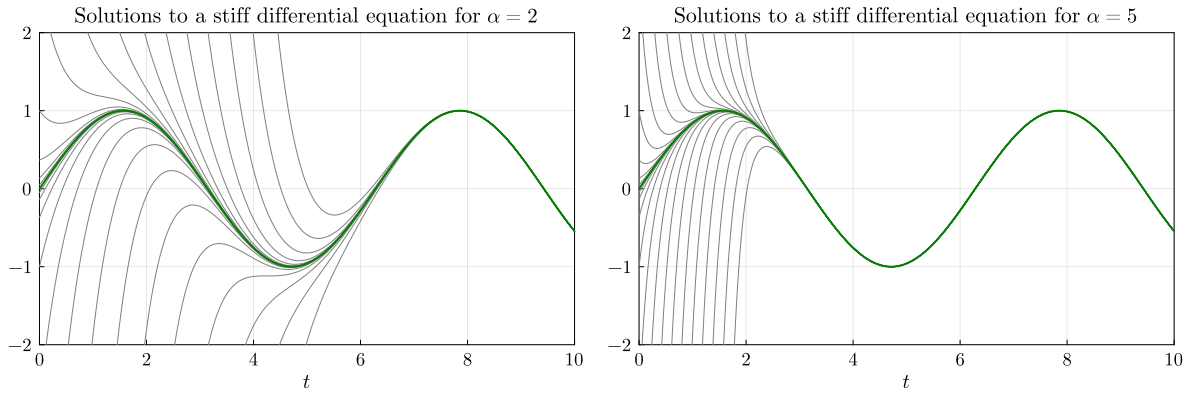


Figure 7.3: Solutions to (7.40) for various initial conditions when  $\alpha = 2$  (left) and  $\alpha = 5$  (right).

It turns out that this prediction is precise, as depicted in Figure 7.4. Note that if the equation is very stiff, that is to say if  $\alpha \gg 1$ , then a very small time step is required to ensure stability.

In contrast with the forward Euler scheme, the backward Euler scheme is stable regardless of the time step. Since the right-hand side of (7.40) is linear in  $x$ , the value of the iterate  $x_{n+1}$  can be calculated explicitly from  $x_n$  for the backward scheme:

$$x_{n+1} = \frac{x_n + \Delta\alpha \sin(t_{n+1}) + \Delta \cos(t_{n+1})}{1 + \Delta\alpha}.$$

Numerical approximations obtained using this scheme are illustrated in Figure 7.5. We observe that the method is stable even for the large time step  $\Delta = 2\Delta_*$ .

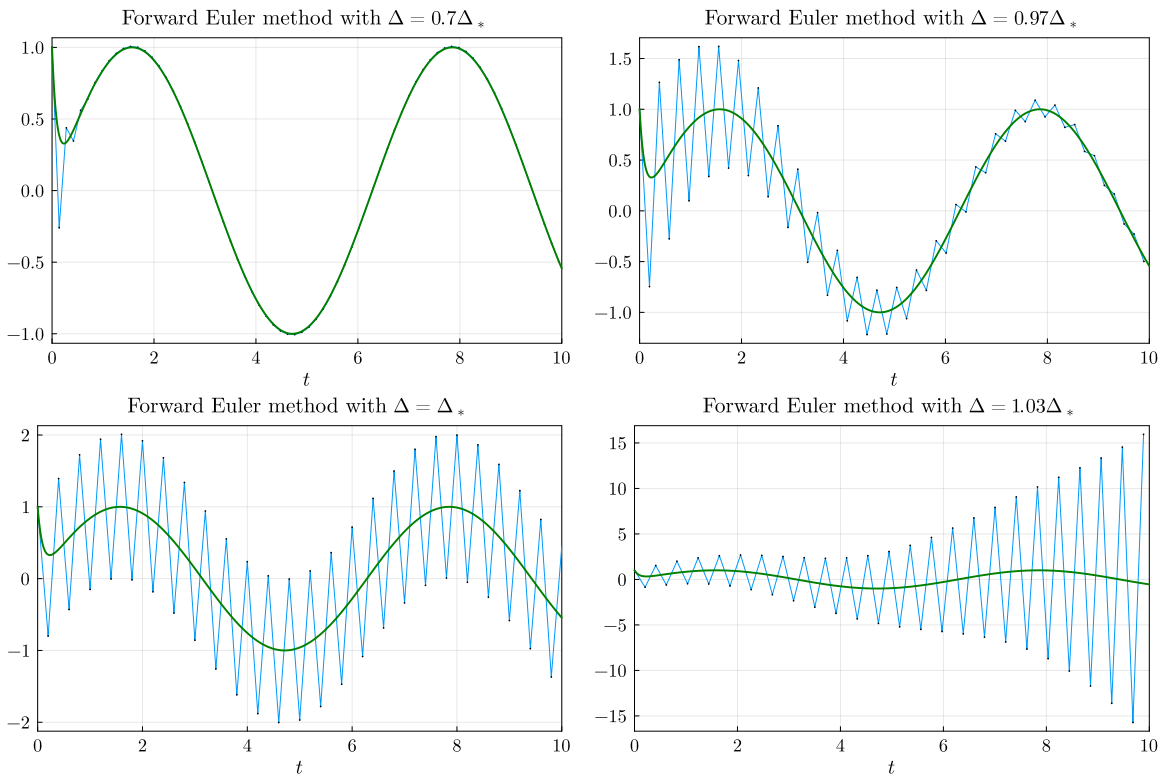


Figure 7.4: Numerical approximations of the solution to (7.40) with  $\alpha = 10$  obtained with the forward Euler method, for four different values of  $\Delta$ .

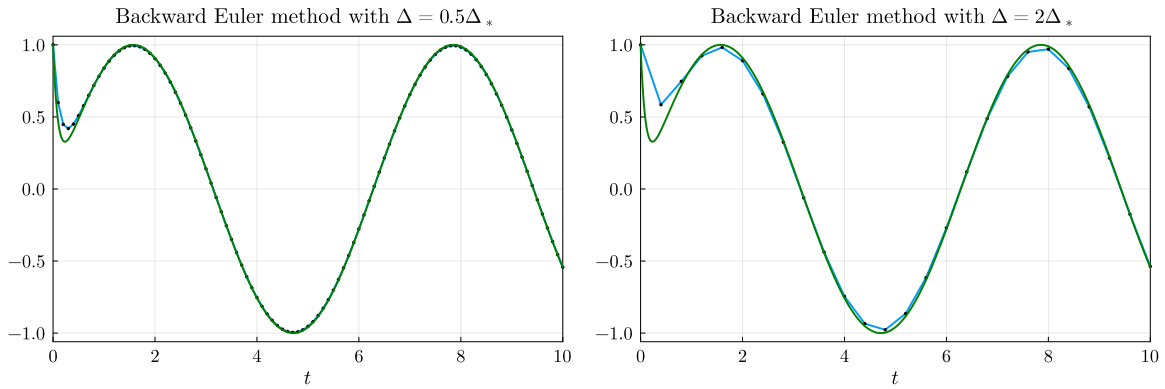


Figure 7.5: Numerical approximations of the solution to (7.40) with  $\alpha = 10$  obtained with the backward Euler method, for two different values of  $\Delta$ .

## 7.5 Exercises

⚙️ **Exercise 7.1.** Show that the absolute stability region of the Crank–Nicolson method (7.24) is given by the left half-plane; see Figure 7.6.

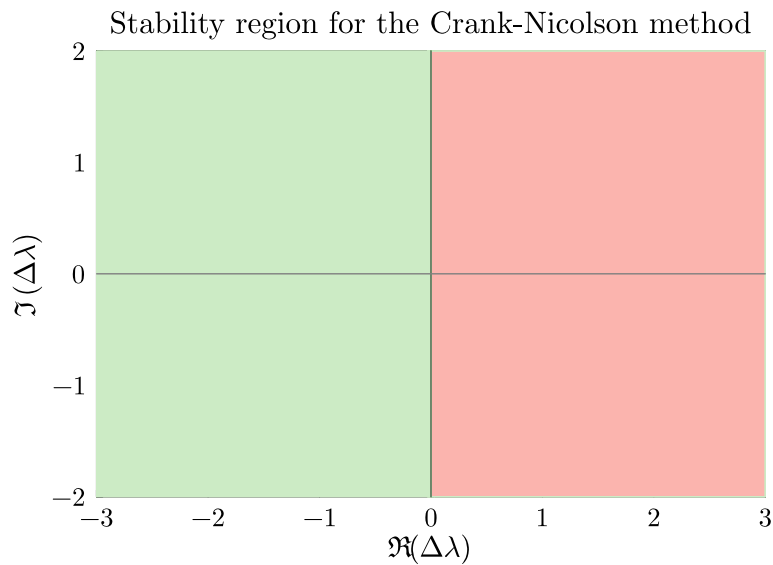


Figure 7.6: Absolute stability regions for the Crank Nicolson method.

⚙️ **Exercise 7.2.** Calculate the absolute stability region for Gear’s method.



# Chapter 8

## Optimization

8.1 Definition and characterization of convexity . . . . .	195
8.2 Unconstrained optimization . . . . .	197
8.3 Constrained optimization . . . . .	199

In this chapter, we focus on optimization problems of the following form:

$$\text{Find } \mathbf{x}_* \in \arg \min_{\mathbf{x} \in \mathcal{K}} J(\mathbf{x}), \tag{8.1}$$

where  $\mathcal{K}$  is a given subset of  $\mathbf{R}^n$  and  $J: \mathcal{K} \rightarrow \mathbf{R}$  is a given *objective function*. We came across several examples of such problems earlier in these notes:

- In [Chapter 2](#), in the context of least-squares approximation, we considered the problem of minimizing

$$J(\boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{A}\boldsymbol{\alpha} - \mathbf{b}\|^2.$$

- In [Chapter 4](#), we observed that, if  $\mathbf{A}$  is a symmetric and positive definite matrix, then solving the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  amounts to finding the minimizer of the functional

$$J(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}.$$

When  $\mathcal{K} = \mathbf{R}^n$ , equation (8.1) is an *unconstrained* optimization problem, and when  $\mathcal{K} \subsetneq \mathbf{R}^n$ , equation (8.1) is a *constrained* optimization problem. In practice, the set  $\mathcal{K}$  is often an intersection of sets of the form

$$\{\mathbf{x} \in \mathbf{R}^n : \phi(\mathbf{x}) \leq 0\}, \quad \text{or} \quad \{\mathbf{x} \in \mathbf{R}^n : \phi(\mathbf{x}) = 0\},$$

for appropriate  $\phi: \mathbf{R}^n \rightarrow \mathbf{R}$ . Constraints of the former form are called *inequality constraints*, while constraints of the latter form are called *equality constraints*. Our aim in this chapter is to give a brief introduction to numerical optimization. We focus on the simplest method, namely the *steepest descent method* with fixed step. The rest of this chapter is organized as follows:

- We begin in Section 8.1 by defining the notions of *convexity*, *strict convexity* and *strong convexity*, which play an important role in optimization.
- Then, in Section 8.2, we analyze the steepest descent method with fixed step in the setting of unconstrained optimization. To this end, we first establish conditions under which (8.1) is well posed.
- Finally, in Section 8.3, we extend the steepest descent method to the case of optimization with constraints.

*Remark 8.1.* For generality, we could consider the setting where the set  $\mathcal{K}$  in (8.1) is a subset of some finite dimensional or infinite dimensional vector space  $V$ . An optimization problem over (a subset of) a finite dimensional vector space of dimension  $n$  can always be recast as an optimization problem over (a subset of)  $\mathbf{R}^n$  – the type we study in this chapter – by fixing a basis. The case of an infinite dimensional vector space, however, is more delicate, and we do not address it here.

## 8.1 Definition and characterization of convexity

**Definition 8.1** (Convexity). Assume that  $J: \mathcal{K} \rightarrow \mathbf{R}$ .

- The function  $J$  is said to be *convex* if

$$\forall(\mathbf{x}, \mathbf{y}) \in \mathcal{K} \times \mathcal{K}, \quad \forall\theta \in [0, 1], \quad J(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta J(\mathbf{x}) + (1 - \theta)J(\mathbf{y}). \quad (8.2)$$

- The function  $J$  is called *strictly convex* if (8.2) holds with strict inequality if  $\mathbf{x} \neq \mathbf{y}$  and  $\theta \in (0, 1)$ .
- The function  $J$  is called *strongly convex* with parameter  $\alpha > 0$  if for all  $(\mathbf{x}, \mathbf{y}) \in \mathcal{K} \times \mathcal{K}$  and for all  $\theta \in [0, 1]$ ,

$$J(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta J(\mathbf{x}) + (1 - \theta)J(\mathbf{y}) - \frac{\alpha}{2}\theta(1 - \theta)\|\mathbf{x} - \mathbf{y}\|^2. \quad (8.3)$$

If the function  $J$  is differentiable, then convexity, strict convexity and strong convexity can be characterized in terms of the gradient  $\nabla J$ . We illustrate this for strong convexity, noting that a characterization of convexity is obtained by substituting  $\alpha = 0$  in the following result.

**Proposition 8.1.** *A differentiable function  $J: \mathbf{R}^n \rightarrow \mathbf{R}$  is strongly convex with parameter  $\alpha$  if and only if*

$$\forall(\mathbf{x}, \mathbf{y}) \in \mathbf{R}^n \times \mathbf{R}^n, \quad J(\mathbf{x}) \geq J(\mathbf{y}) + \langle \nabla J(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\alpha}{2}\|\mathbf{x} - \mathbf{y}\|^2, \quad (8.4)$$

or, equivalently,

$$\forall(\mathbf{x}, \mathbf{y}) \in \mathbf{R}^n \times \mathbf{R}^n, \quad \langle \nabla J(\mathbf{x}) - \nabla J(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \alpha \|\mathbf{x} - \mathbf{y}\|^2. \quad (8.5)$$

*Proof.* For clarity, we divide the proof into items and prove one implication per item.

- (8.3)  $\Rightarrow$  (8.4). Rearranging (8.3), we have

$$\frac{J(\mathbf{y} + \theta(\mathbf{x} - \mathbf{y})) - J(\mathbf{y})}{\theta} \leq J(\mathbf{x}) - J(\mathbf{y}) - \frac{\alpha}{2}(1 - \theta)\|\mathbf{x} - \mathbf{y}\|^2.$$

Taking the limit  $\theta \rightarrow 0$ , we deduce that

$$\langle \nabla J(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq J(\mathbf{x}) - J(\mathbf{y}) - \frac{\alpha}{2}\|\mathbf{x} - \mathbf{y}\|^2.$$

This gives (8.4) after rearranging.

- (8.4)  $\Rightarrow$  (8.3). To prove this implication, suppose that (8.4) holds, take  $(\mathbf{x}, \mathbf{y}) \in \mathbf{R}^n \times \mathbf{R}^n$  and let  $\mathbf{z} = \theta\mathbf{x} + (1 - \theta)\mathbf{y}$ . Using (8.4) successively with  $(\mathbf{x}, \mathbf{z})$  and  $(\mathbf{y}, \mathbf{z})$ , we deduce

$$\begin{aligned} J(\mathbf{x}) &\geq J(\mathbf{z}) + \langle \nabla J(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle + \frac{\alpha}{2}\|\mathbf{x} - \mathbf{z}\|^2, \\ J(\mathbf{y}) &\geq J(\mathbf{z}) + \langle \nabla J(\mathbf{z}), \mathbf{y} - \mathbf{z} \rangle + \frac{\alpha}{2}\|\mathbf{y} - \mathbf{z}\|^2. \end{aligned}$$

Combining these inequalities, we deduce that

$$\begin{aligned} \theta J(\mathbf{x}) + (1 - \theta)J(\mathbf{y}) &\geq J(\mathbf{z}) + \langle \nabla J(\mathbf{z}), \theta\mathbf{x} + (1 - \theta)\mathbf{y} - \mathbf{z} \rangle \\ &\quad + \frac{\alpha\theta}{2}\|\mathbf{x} - \mathbf{z}\|^2 + \frac{\alpha(1 - \theta)}{2}\|\mathbf{y} - \mathbf{z}\|^2 \\ &= J(\mathbf{z}) + 0 + \frac{\alpha}{2}\theta(1 - \theta)\|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

Rearranging gives (8.3).

- (8.4)  $\Rightarrow$  (8.5). Assuming that (8.4) holds and applying this inequality first to  $(\mathbf{x}, \mathbf{y})$  and then to  $(\mathbf{y}, \mathbf{x})$ , we obtain

$$\begin{aligned} J(\mathbf{x}) &\geq J(\mathbf{y}) + \langle \nabla J(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\alpha}{2}\|\mathbf{x} - \mathbf{y}\|^2 \\ J(\mathbf{y}) &\geq J(\mathbf{x}) + \langle \nabla J(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\alpha}{2}\|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

Adding these equations and rearranging, we deduce (8.5).

- (8.5)  $\Rightarrow$  (8.4). Suppose that (8.5) holds and take  $(\mathbf{x}, \mathbf{y}) \in \mathbf{R}^n \times \mathbf{R}^n$ . Using the fundamental

theorem of analysis and (8.5), we have

$$\begin{aligned} J(\mathbf{x}) &= J(\mathbf{y}) + \int_0^1 \langle \nabla J(\mathbf{y} + \theta(\mathbf{x} - \mathbf{y})), \mathbf{x} - \mathbf{y} \rangle d\theta \\ &\geq J(\mathbf{y}) + \int_0^1 \langle \nabla J(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \alpha\theta \|\mathbf{x} - \mathbf{y}\|^2 d\theta \\ &= J(\mathbf{y}) + \langle \nabla J(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|^2, \end{aligned}$$

which gives (8.4).

We have proved all the implications required to conclude the proof.  $\square$

## 8.2 Unconstrained optimization

Throughout this section  $\mathcal{K} = \mathbf{R}^n$ . We begin by establishing conditions under which the optimization problem (8.1) admits a unique solution in this setting. We first prove existence of a global minimizer under appropriate conditions.

**Proposition 8.2** (Existence of a global minimizer). *Suppose that  $J: \mathbf{R}^n \rightarrow \mathbf{R}$  is continuous and coercive, the latter meaning that  $J(\mathbf{x}) \rightarrow \infty$  when  $\|\mathbf{x}\| \rightarrow \infty$ . Then there exists a global minimizer of  $J$  in  $\mathbf{R}^n$ .*

*Proof.* Let  $(\mathbf{x}_n)_{n \in \mathbf{N}}$  be a minimizing sequence of  $J$ , i.e. a sequence in  $\mathbf{R}^n$  such that

$$J(\mathbf{x}_n) \rightarrow \inf_{\mathbf{x} \in \mathbf{R}^n} J(\mathbf{x}) \quad \text{as } n \rightarrow \infty.$$

The sequence  $(\mathbf{x}_n)$  is bounded, because otherwise it would hold that  $J(\mathbf{x}_n) \rightarrow \infty$  by coercivity. Therefore, since closed bounded sets in  $\mathbf{R}^n$  are compact, there is a subsequence  $(\mathbf{x}_{n_k})_{k \in \mathbf{N}}$  converging to some  $\mathbf{x}_* \in \mathbf{R}^n$ . Since  $J$  is continuous, we have that

$$J(\mathbf{x}_*) = \lim_{k \rightarrow \infty} J(\mathbf{x}_{n_k}) = \inf_{\mathbf{x} \in \mathbf{R}^n} J(\mathbf{x}).$$

We conclude that  $\mathbf{x}_*$  is a minimizer of  $J$ .  $\square$

*Remark 8.2.* We relied crucially in the proof of Proposition 8.2 on the fact that closed bounded sets in  $\mathbf{R}^n$  are compact. In the infinite-dimensional setting, coercivity and continuity alone are not sufficient to guarantee the existence of a minimizer.

Uniqueness of the minimizer can be established under a strict convexity assumption.

**Proposition 8.3** (Uniqueness of the minimizer). *If  $J$  is strictly convex, then there exists at most one global minimizer.*

*Proof.* Suppose for contradiction that there were two minimizers  $\mathbf{x}_*$  and  $\mathbf{y}_*$ . Then by strict convexity we have

$$J\left(\frac{\mathbf{x}_* + \mathbf{y}_*}{2}\right) < \frac{1}{2}(J(\mathbf{x}_*) + J(\mathbf{y}_*)) = J(\mathbf{x}_*),$$

which contradicts the minimality of  $J(\mathbf{x}_*)$ .  $\square$

Finally, before introducing the steepest descent algorithm, we recall the following standard result from analysis, the proof of which is left as an exercise.

**Theorem 8.4** (Euler condition). *Suppose that  $J: \mathbf{R}^n \rightarrow \mathbf{R}$  is differentiable.*

- *If  $\mathbf{x}_*$  is a local minimizer of  $J$ , then  $\nabla J(\mathbf{x}_*) = 0$ .*
- *If  $J$  is convex, then  $\nabla J(\mathbf{x}_*) = 0$  if and only if  $\mathbf{x}_*$  is a global minimizer.*

**Steepest descent method.** In this section, we study the more general version of the steepest descent with *fixed step* given in [Algorithm 17](#).

---

**Algorithm 17** Steepest descent method

---

- 1: Pick  $\lambda$ , and initial  $\mathbf{x}_0$ .
  - 2: **for**  $k \in \{0, 1, \dots\}$  **do**
  - 3:      $\mathbf{d}_k \leftarrow \nabla J(\mathbf{x}_k)$
  - 4:      $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \lambda \mathbf{d}_k$
  - 5: **end for**
- 

*Remark 8.3.* We encountered the steepest descent with fixed step for a quadratic objective function when we analyzed Richardson's method for solving linear equations in [Chapter 4](#).

In practice, [Algorithm 17](#) must be supplemented with an appropriate stopping criterion. This could be, for example, a criterion of the form  $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| \leq \varepsilon$ , or  $|J(\mathbf{x}_{k+1}) - J(\mathbf{x}_k)| \leq \varepsilon$ . It is sometimes also useful to use a normalized criterion of the form  $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| \leq \varepsilon \|\mathbf{x}_0\|$ . The steepest descent method may be viewed as a fixed point iteration for the function

$$F_\lambda(\mathbf{x}) = \mathbf{x} - \lambda \nabla J(\mathbf{x}). \quad (8.6)$$

A point  $\mathbf{x}_* \in \mathbf{R}^n$  is a fixed point of this function if and only if  $\mathbf{x}_*$  is a solution to the nonlinear equation  $\nabla J(\mathbf{x}_*) = 0$ . We shall now prove the convergence of the steepest descent under appropriate assumptions on the function  $J$ .

**Theorem 8.5** (Convergence of the steepest descent method). *Suppose that  $J$  is differentiable, strongly convex with parameter  $\alpha$ , and that its gradient  $\nabla J: \mathbf{R}^n \rightarrow \mathbf{R}^n$  is Lipschitz continuous with parameter  $L$ :*

$$\forall(\mathbf{x}, \mathbf{y}) \in \mathbf{R}^n \times \mathbf{R}^n, \quad \|\nabla J(\mathbf{x}) - \nabla J(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|. \quad (8.7)$$

*Then provided that*

$$0 < \lambda < \frac{2\alpha}{L}, \quad (8.8)$$

*the steepest descent method with fixed step is convergent. More precisely, there exists  $\rho \in (0, 1)$*

such that for all  $k \geq 0$

$$\|\mathbf{x}_k - \mathbf{x}_*\| \leq \rho^k \|\mathbf{x}_0 - \mathbf{x}_*\|. \quad (8.9)$$

*Proof.* Under the assumptions of the theorem, there exists a unique global minimizer of  $J$ , which is the unique fixed point of  $\mathbf{F}_\lambda$ . We begin by proving that  $\mathbf{F}_\lambda$  defined in (8.6) is globally Lipschitz continuous. We have

$$\begin{aligned} \|\mathbf{F}_\lambda(\mathbf{x}) - \mathbf{F}_\lambda(\mathbf{y})\|^2 &= \|\mathbf{x} - \mathbf{y} - \lambda(\nabla J(\mathbf{x}) - \nabla J(\mathbf{y}))\|^2 \\ &= \|\mathbf{x} - \mathbf{y}\|^2 - 2\lambda\langle \mathbf{x} - \mathbf{y}, \nabla J(\mathbf{x}) - \nabla J(\mathbf{y}) \rangle + \lambda^2 \|\nabla J(\mathbf{x}) - \nabla J(\mathbf{y})\|^2 \\ &\leq (1 - 2\alpha\lambda + \lambda^2 L)\|\mathbf{x} - \mathbf{y}\|^2, \end{aligned}$$

where we employed (8.5) for the second term and (8.7) for the third term. Thus,  $\mathbf{F}_\lambda$  is globally Lipschitz continuous with constant  $\rho = \sqrt{1 - 2\alpha\lambda + \lambda^2 L}$ , which is less than 1 if and only if (8.8) is satisfied. The bound (8.9) then follows by noting that

$$\|\mathbf{x}_k - \mathbf{x}_*\| = \|\mathbf{F}_\lambda(\mathbf{x}_{k-1}) - \mathbf{F}_\lambda(\mathbf{x}_*)\| \leq \rho \|\mathbf{x}_{k-1} - \mathbf{x}_*\| \leq \dots \leq \rho^k \|\mathbf{x}_0 - \mathbf{x}_*\|,$$

which concludes the proof. (Note that (8.9) also follows from Theorem 5.2.)  $\square$

*Remark 8.4* (Convergence speed). The choice of  $\lambda$  minimizing the Lipschitz constant  $\rho$  is given by  $\lambda_* = \frac{\alpha}{L^2}$ , which corresponds to  $\rho_* = 1 - \left(\frac{\alpha}{L}\right)^2$ . Often, in practice, it holds that  $\alpha \ll L$ , in which case the convergence of the steepest descent with fixed step is slow.

### 8.3 Constrained optimization

In this section, we assume that  $\mathcal{K} \subset \mathbf{R}^n$ . We begin by establishing well-posedness of the optimization problem (8.1) in this setting.

**Proposition 8.6** (Well posedness of (8.1) in the constrained setting). *The two items below concern existence and uniqueness, respectively.*

- Suppose that  $\mathcal{K} \subset \mathbf{R}^n$  is closed and that  $J: \mathcal{K} \rightarrow \mathbf{R}$  is continuous and coercive. Then there exists a global minimizer of  $J$  in  $\mathcal{K}$ .
- Suppose that  $\mathcal{K} \subset \mathbf{R}^n$  is convex and that  $J: \mathcal{K} \rightarrow \mathbf{R}$  is strictly convex. Then there exists at most one global minimizer.

*Proof.* The proof is very similar to those of Proposition 8.2 and Proposition 8.3, and so we leave it to the reader. Note that the set  $\mathcal{K}$  must be closed to ensure existence, and convex to guarantee uniqueness. These assumptions are clearly satisfied when  $\mathcal{K} = \mathbf{R}^n$ , so Proposition 8.6 indeed generalizes Propositions 8.2 and 8.3.  $\square$

The following theorem, which generalizes (8.4), establishes a characterization of the minimizer when  $J$  is differentiable.

**Theorem 8.7** (Euler–Lagrange conditions). *Suppose that  $J: \mathcal{K} \rightarrow \mathbf{R}$  is differentiable and that  $\mathcal{K} \subset \mathbf{R}^n$  is closed and convex. Then the following statements hold.*

- If  $\mathbf{x}_*$  is a local minimizer of  $J$ , then

$$\forall \mathbf{x} \in \mathcal{K}, \quad \langle \nabla J(\mathbf{x}_*), \mathbf{x} - \mathbf{x}_* \rangle \geq 0. \quad (8.10)$$

- Conversely, if (8.10) is satisfied and  $J$  is convex, then  $\mathbf{x}_*$  is a global minimizer of  $J$ .

*Proof.* Suppose that  $\mathbf{x}_*$  is a local minimizer of  $J$ . This means that there exists  $\delta > 0$  such that

$$\forall \mathbf{x} \in B_\delta(\mathbf{x}_*) \cap \mathcal{K}, \quad J(\mathbf{x}_*) \leq J(\mathbf{x}).$$

Therefore  $J(\mathbf{x}_*) \leq J((1-t)\mathbf{x}_* + t\mathbf{x})$  for all  $t \in [0, 1]$  sufficiently small. But then

$$\langle \nabla J(\mathbf{x}_*), \mathbf{x} - \mathbf{x}_* \rangle = \lim_{t \rightarrow 0} \frac{J((1-t)\mathbf{x}_* + t\mathbf{x}) - J(\mathbf{x}_*)}{t} \geq 0.$$

Conversely, suppose that (8.10) is satisfied and that  $J$  is convex. Since  $J$  is convex, equation (8.4) holds with  $\alpha = 0$ , and applying this equation with  $\mathbf{y} = \mathbf{x}_*$ , we deduce that  $\mathbf{x}_*$  is a global minimizer.  $\square$

The steepest descent [Algorithm 17](#) can be extended to optimization problems with constraints by introducing an additional projection step. In order to precisely formulate the algorithm, we begin by introducing the projection operator  $\Pi_{\mathcal{K}}$ .

**Proposition 8.8** (Projection on a closed convex set). *Suppose that  $\mathcal{K}$  is a closed convex subset of  $\mathbf{R}^n$ . Then for all  $\mathbf{x} \in \mathbf{R}^n$  there is a unique  $\Pi_{\mathcal{K}}\mathbf{x} \in \mathcal{K}$ , called the orthogonal projection of  $\mathbf{x}$  onto  $\mathcal{K}$ , such that*

$$\|\Pi_{\mathcal{K}}\mathbf{x} - \mathbf{x}\| = \inf_{\mathbf{y} \in \mathcal{K}} \|\mathbf{y} - \mathbf{x}\|.$$

*Proof.* The functional  $J_{\mathbf{x}}(\mathbf{y}) = \|\mathbf{y} - \mathbf{x}\|^2$  is strongly convex, and so [Proposition 8.6](#) immediately implies the existence and uniqueness of  $\Pi_{\mathcal{K}}\mathbf{x}$ .  $\square$

*Remark 8.5.* In view of [Theorem 8.7](#), the projection  $\Pi_{\mathcal{K}}\mathbf{x}$  is the unique element of  $\mathcal{K}$  which satisfies

$$\forall \mathbf{y} \in \mathcal{K}, \quad \langle \Pi_{\mathcal{K}}\mathbf{x} - \mathbf{x}, \mathbf{y} - \Pi_{\mathcal{K}}\mathbf{x} \rangle \geq 0. \quad (8.11)$$

We are now ready to present the steepest descent method with projection: see [Algorithm 18](#). Like [Algorithm 17](#), the steepest descent with projection may be viewed as a fixed point iteration, this time for the function

$$\mathbf{F}_\lambda(\mathbf{x}) := \Pi_{\mathcal{K}}(\mathbf{x} - \lambda \nabla J(\mathbf{x})). \quad (8.12)$$

We now prove the convergence of the method.

**Algorithm 18** Steepest descent with projection

- 
- 1: Pick  $\lambda$ , and initial  $\mathbf{x}_0$ .
  - 2: **for**  $k \in \{0, 1, \dots\}$  **do**
  - 3:      $\mathbf{d}_k \leftarrow \nabla J(\mathbf{x}_k)$
  - 4:      $\mathbf{x}_{k+1} \leftarrow \Pi_{\mathcal{K}}(\mathbf{x}_k - \lambda \mathbf{d}_k)$
  - 5: **end for**
- 

**Theorem 8.9** (Convergence of steepest descent with projection). *Suppose that  $J$  is differentiable, strongly convex with parameter  $\alpha$ , and that its gradient  $\nabla J: \mathbf{R}^n \rightarrow \mathbf{R}^n$  is Lipschitz continuous with parameter  $L$ . Assume also that  $\mathcal{K} \subset \mathbf{R}^n$  is closed and convex. Then provided that*

$$0 < \lambda < \frac{2\alpha}{L}, \quad (8.13)$$

*the steepest descent method with fixed step is convergent. More precisely, there exists  $\rho \in (0, 1)$  such that for all  $k \geq 0$*

$$\|\mathbf{x}_k - \mathbf{x}_*\| \leq \rho^k \|\mathbf{x}_0 - \mathbf{x}_*\|.$$

*Proof.* Under the assumptions, there exists a unique global minimizer  $\mathbf{x}_* \in \mathcal{K}$ . We already showed in the proof of [Theorem 8.5](#) that the mapping  $\mathbf{x} \mapsto \mathbf{x} - \lambda \nabla J(\mathbf{x})$  is a contraction if and only if  $\lambda$  satisfies (8.13). In order to prove that  $\mathbf{F}_\lambda$  given in (8.12) is a contraction under the same condition, it is sufficient to prove that  $\Pi_{\mathcal{K}}: \mathbf{R}^n \rightarrow \mathcal{K}$  satisfies the following estimate:

$$\forall (\mathbf{x}, \mathbf{y}) \in \mathbf{R}^n \times \mathbf{R}^n, \quad \|\Pi_{\mathcal{K}}\mathbf{x} - \Pi_{\mathcal{K}}\mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}\|.$$

To this end, take  $(\mathbf{x}, \mathbf{y}) \in \mathbf{R}^n \times \mathbf{R}^n$  and let  $\boldsymbol{\delta} = \Pi_{\mathcal{K}}\mathbf{x} - \Pi_{\mathcal{K}}\mathbf{y}$ . By (8.11), it holds that

$$\begin{aligned} \|\boldsymbol{\delta}\|^2 &= \langle \boldsymbol{\delta}, \Pi_{\mathcal{K}}\mathbf{x} - \mathbf{x} \rangle + \langle \boldsymbol{\delta}, \mathbf{x} - \mathbf{y} \rangle + \langle \boldsymbol{\delta}, \mathbf{y} - \Pi_{\mathcal{K}}\mathbf{y} \rangle \\ &\leq 0 + \langle \boldsymbol{\delta}, \mathbf{x} - \mathbf{y} \rangle + 0 \leq \|\boldsymbol{\delta}\| \|\mathbf{x} - \mathbf{y}\|, \end{aligned}$$

which yields the required inequality. Therefore  $\mathbf{F}_\lambda$  in (8.12) is a contraction and so, by the Banach fixed point theorem, it admits a unique fixed point  $\mathbf{y}_* \in \mathcal{K}$ . To show that  $\mathbf{y}_* = \mathbf{x}_*$ , note that if  $\mathbf{F}_\lambda(\mathbf{y}_*) = \mathbf{y}_*$ , then by (8.11) it holds that

$$\forall \mathbf{y} \in \mathcal{K}, \quad \langle \lambda \nabla J(\mathbf{y}_*), \mathbf{y} - \mathbf{y}_* \rangle \geq 0.$$

Therefore, using [Theorem 8.7](#), we obtain that  $\mathbf{y}_*$  is a global minimizer of  $J$ , so  $\mathbf{y}_* = \mathbf{x}_*$ .  $\square$

*Remark 8.6.* The applicability of [Algorithm 18](#) is limited in practice, as computing  $\Pi_{\mathcal{K}}(\mathbf{x})$  analytically is possible only in simple settings.



# Appendix A

## Background material

A.1 Inner products and norms . . . . .	202
A.2 Completeness . . . . .	205
A.3 Contraction mappings and the Banach fixed point theorem . . . . .	206
A.4 Vector norms . . . . .	207
A.5 Matrix norms . . . . .	207
A.6 Diagonalization and spectral theorem . . . . .	209
A.7 Similarity transformation and Jordan normal form . . . . .	212
A.8 Oldenburger's theorem and Gelfand's formula . . . . .	213

In this chapter, we collect basic results that are useful for this course.

### A.1 Inner products and norms

We begin by recalling the definitions of the fundamental concepts of *norm* and *inner product*. For generality, we consider the case of a *complex* vector space, i.e. a vector space for which the scalar field is  $\mathbf{C}$ .

**Definition A.1.** A norm on a complex vector space  $\mathcal{X}$  is a function  $\|\bullet\| : \mathcal{X} \rightarrow \mathbf{R}$  satisfying the following axioms:

- **Positivity:**  $\forall \mathbf{x} \in \mathcal{X} \setminus \{\mathbf{0}\}, \quad \|\mathbf{x}\| > 0.$
- **Homogeneity:**  $\forall (c, \mathbf{x}) \in \mathbf{C} \times \mathcal{X}, \quad \|c\mathbf{x}\| = |c| \|\mathbf{x}\|.$
- **Triangular inequality:**  $\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{X}, \quad \|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|.$

For example, the Euclidean norm on  $\mathbf{C}^n$  is given by

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

**Definition A.2.** An inner product on a *complex* vector space  $\mathcal{X}$  is a function

$$\langle \bullet, \bullet \rangle : \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{C}$$

satisfying the following axioms:

- **Conjugate symmetry:** Here  $\bar{\phantom{x}}$  denotes the complex conjugate.

$$\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{X}, \quad \langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle}.$$

- **Linearity:** For all  $(\alpha, \beta) \in \mathbf{C}^2$  and all  $(\mathbf{x}, \mathbf{y}, \mathbf{z}) \in \mathcal{X}^3$ , it holds that

$$\langle \alpha \mathbf{x} + \beta \mathbf{y}, \mathbf{z} \rangle = \alpha \langle \mathbf{x}, \mathbf{z} \rangle + \beta \langle \mathbf{y}, \mathbf{z} \rangle.$$

- **Positive-definiteness:**

$$\forall \mathbf{x} \in \mathcal{X} \setminus \{0\}, \quad \langle \mathbf{x}, \mathbf{x} \rangle > 0.$$

For example, the familiar Euclidean inner product on  $\mathbf{C}^n$  is given by

$$\langle \mathbf{x}, \mathbf{y} \rangle := \sum_{i=1}^n x_i \bar{y}_i.$$

A vector space with an inner product is called an *inner product space*. Any inner product on  $\mathcal{X}$  induces a norm via the formula

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}. \quad (\text{A.1})$$

The Cauchy–Schwarz inequality enables to bound inner products using norms. It is also useful for showing that the functional defined in (A.1) satisfies the triangle inequality, which is the goal of [Exercise A.2](#).

**Proposition A.1** (Cauchy–Schwarz inequality). *Let  $\mathcal{X}$  be an inner product space. Then*

$$\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{X}, \quad |\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|. \quad (\text{A.2})$$

*Proof.* The statement is obvious if  $\mathbf{y} = \mathbf{0}$ , so we assume in the rest of the proof that  $\mathbf{y} \neq \mathbf{0}$ . Let us define  $p: \mathbf{R} \ni t \mapsto \|\mathbf{x} + t\mathbf{y}\|^2$ . Using the bilinearity of the inner product, we have

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

This shows that  $p$  is a convex second-order polynomial with a minimum at  $t_* = -\langle \mathbf{x}, \mathbf{y} \rangle / \|\mathbf{y}\|^2$ . Substituting this value in the expression of  $p$ , we obtain

$$p(t_*) = \|\mathbf{x}\|^2 - 2 \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^2} + \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^2} = \|\mathbf{x}\|^2 - \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^2}.$$

Since  $p(t_*) \geq 0$  by definition of  $p$ , we obtain (A.2). □

Several norms can be defined on the same vector space  $\mathcal{X}$ . Two norms  $\|\bullet\|_\alpha$  and  $\|\bullet\|_\beta$  on  $\mathcal{X}$  are said to be equivalent if there exist positive real numbers  $c_\ell$  and  $c_u$  such that

$$\forall \mathbf{x} \in \mathcal{X}, \quad c_\ell \|\mathbf{x}\|_\alpha \leq \|\mathbf{x}\|_\beta \leq c_u \|\mathbf{x}\|_\alpha. \quad (\text{A.3})$$

As the terminology indicates, norm equivalence is an *equivalence relation*. When working with norms on finite-dimensional vector spaces, it is important to keep in mind the following result. The proof is provided for information purposes only.

**Proposition A.2.** *Assume that  $\mathcal{X}$  is a finite-dimensional vector space. Then all the norms defined on  $\mathcal{X}$  are pairwise equivalent.*

*Proof.* Let  $(\mathbf{e}_1, \dots, \mathbf{e}_n)$  be a basis of  $\mathcal{X}$ . Any  $\mathbf{x} \in \mathcal{X}$  admits a unique representation in this basis as  $\mathbf{x} = \lambda_1 \mathbf{e}_1 + \dots + \lambda_n \mathbf{e}_n$ . We will show that any norm  $\|\bullet\|$  on  $\mathcal{X}$  is equivalent to the norm  $\|\bullet\|_*$  given by

$$\|\mathbf{x}\|_* = |\lambda_1| + \dots + |\lambda_n|. \quad (\text{A.4})$$

By the triangle inequality, it holds that

$$\begin{aligned} \|\mathbf{x}\| &\leq |\lambda_1| \|\mathbf{e}_1\| + \dots + |\lambda_n| \|\mathbf{e}_n\| \leq \left( |\lambda_1| + \dots + |\lambda_n| \right) \max \left\{ \|\mathbf{e}_1\|, \dots, \|\mathbf{e}_n\| \right\} \\ &= \|\mathbf{x}\|_* \max \left\{ \|\mathbf{e}_1\|, \dots, \|\mathbf{e}_n\| \right\}. \end{aligned} \quad (\text{A.5})$$

It remains to show that there exists a positive constant  $\ell$  such that

$$\forall \mathbf{x} \in \mathcal{X}, \quad \|\mathbf{x}\| \geq \ell \left( |\lambda_1| + \dots + |\lambda_n| \right). \quad (\text{A.6})$$

To this end, we reason by contradiction. If this inequality were not true, then there would exist a sequence  $(\mathbf{x}^{(i)})_{i \in \mathbf{N}}$  such that  $\|\mathbf{x}^{(i)}\| \rightarrow 0$  as  $i \rightarrow \infty$  and  $\|\mathbf{x}^{(i)}\|_* = 1$  for all  $i \in \mathbf{N}$ . Since  $\lambda_1^{(i)} \in [-1, 1]$  for all  $i \in \mathbf{N}$ , we can extract a subsequence, still denoted by  $(\mathbf{x}^{(i)})_{i \in \mathbf{N}}$  for simplicity, such that the corresponding coefficient  $\lambda_1^{(i)}$  satisfies  $\lambda_1^{(i)} \rightarrow \lambda_1^* \in [-1, 1]$ , by compactness of the interval  $[-1, 1]$ . Repeating this procedure for  $\lambda_2, \lambda_3, \dots$ , taking a new subsequence every time, we obtain a subsequence  $(\mathbf{x}^{(i)})_{i \in \mathbf{N}}$  such that  $\lambda_j^{(i)} \rightarrow \lambda_j^*$  in the limit as  $i \rightarrow \infty$ , for all  $j \in \{1, \dots, n\}$ . Therefore, it holds that  $\mathbf{x}^{(i)} \rightarrow \mathbf{x}^* := \lambda_1^* \mathbf{e}_1 + \dots + \lambda_n^* \mathbf{e}_n$  in the  $\|\bullet\|_*$  norm, and thus also in the  $\|\bullet\|$  norm by (A.5). Since  $\mathbf{x}^{(i)} \rightarrow \mathbf{0}$  in the latter norm by assumption, we deduce that  $\mathbf{x}^* = \mathbf{0}$ . But the vectors  $\mathbf{e}_1, \dots, \mathbf{e}_n$  are linearly independent, and so this implies that  $\lambda_1^* = \dots = \lambda_n^* = 0$ , which is a contradiction because we also have that

$$|\lambda_1^*| + \dots + |\lambda_n^*| = \lim_{i \rightarrow \infty} |\lambda_1^{(i)}| + \dots + |\lambda_n^{(i)}| = 1.$$

This concludes the proof of (A.6). □

⚙️ **Exercise A.1.** *Show that  $\|\bullet\|_*: \mathcal{X} \rightarrow \mathbf{R}$  defined in (A.4) is indeed a norm.*

⚙️ **Exercise A.2.** *Using Proposition A.1, show that the function  $\|\bullet\|$  defined by (A.1) satisfies the triangle inequality.*

## A.2 Completeness

Assume that  $\mathcal{X}$  is a vector space with a norm  $\|\bullet\|$ . Together,  $(\mathcal{X}, \|\bullet\|)$  form a *normed vector space*. A sequence  $(\mathbf{x}_n)_{n \geq 0}$  in  $\mathcal{X}$  is convergent in this space if there exists  $\mathbf{x}_* \in \mathcal{X}$  such that

$$\|\mathbf{x}_n - \mathbf{x}_*\| \rightarrow 0 \quad \text{in the limit } n \rightarrow \infty.$$

In this case, we write  $\lim_{n \rightarrow \infty} \mathbf{x}_n = \mathbf{x}_*$  or  $\mathbf{x}_n \rightarrow \mathbf{x}_*$ .

**Definition A.3** (Cauchy sequence). A sequence  $(\mathbf{x}_n)_{n \geq 0}$  in  $\mathcal{X}$  is said to be Cauchy if

$$\lim_{m, n \rightarrow +\infty} \|\mathbf{x}_n - \mathbf{x}_m\| = 0.$$

The normed vector space  $(\mathcal{X}, \|\bullet\|)$  is called *complete* if every Cauchy sequence is convergent.

Every convergent sequence is Cauchy, but the converse is not always true.

*Example A.1.* Consider the case where  $\mathcal{X} = C([-1, 1])$ , the space of continuous functions from  $[-1, 1]$  to  $\mathbf{R}$ , endowed with the norm

$$\|f\| = \int_{-1}^1 |f(x)| \, dx.$$

The sequence of functions  $(f_n)_{n \geq 0}$  in  $\mathcal{X}$  given by

$$f_n(x) = x^{\frac{1}{2n+1}} \tag{A.7}$$

is Cauchy but not convergent. Indeed, assume for contradiction that there existed  $f_* \in \mathcal{X}$  such that

$$\|f_n - f_*\| \xrightarrow[n \rightarrow \infty]{} 0. \tag{A.8}$$

Then  $f_*$  necessarily coincides with sign function:

$$\text{sgn}(x) := \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ 1 & \text{if } x > 0. \end{cases}$$

However, this function is discontinuous, which contradicts the statement that  $f_* \in \mathcal{X}$ .

In this course, all the vector spaces encountered are complete. For example,

- $\mathbf{R}^n$  with any vector norm is complete;
- $\mathbf{C}^{m \times n}$  with any matrix norm is complete.

In order to show that a sequence is convergent in a complete normed vector space, it is sufficient to show that the sequence is Cauchy. This approach is used in [Lemma 4.2](#) and [Theorem A.3](#).

⚙️ **Exercise A.3.** Prove that every convergent sequence is Cauchy.

### A.3 Contraction mappings and the Banach fixed point theorem

Let  $(\mathcal{X}, \|\bullet\|)$  denote a normed vector space. A map  $\phi: \mathcal{X} \rightarrow \mathcal{X}$  is called a contraction mapping if there is a constant  $L \in (0, 1)$  such that

$$\forall (x, y) \in \mathcal{X} \times \mathcal{X}, \quad \|\phi(x) - \phi(y)\| \leq L\|x - y\|.$$

The importance of contraction mappings in this course stems from the following theorem.

**Theorem A.3** (Banach fixed point theorem). *Let  $(\mathcal{X}, \|\bullet\|)$  be a complete normed space, and let  $\phi: \mathcal{X} \rightarrow \mathcal{X}$  be a contraction mapping. Then  $\phi$  has a unique fixed point in  $\mathcal{X}$ .*

*Proof.* We prove first existence and then uniqueness.

**Existence.** Take  $x_0 \in \mathcal{X}$ , and define the sequence  $(x_k)_{k \in \mathbb{N}}$  inductively by

$$x_{k+1} = \phi(x_k). \tag{A.9}$$

It holds that

$$\|x_{k+1} - x_k\| = \|\phi(x_k) - \phi(x_{k-1})\| \leq L\|x_k - x_{k-1}\| \leq \dots \leq L^k\|x_1 - x_0\|.$$

Therefore, for any  $n \geq m$ , we have by the triangle inequality

$$\begin{aligned} \|x_n - x_m\| &\leq \|x_n - x_{n-1}\| + \dots + \|x_{m+1} - x_m\| \\ &\leq (L^{n-1} + \dots + L^m)\|x_1 - x_0\| \\ &\leq L^m(1 + L + L^2 + \dots)\|x_1 - x_0\| = \frac{L^m}{1-L}\|x_1 - x_0\|. \end{aligned}$$

It follows that the sequence  $(x_k)_{k \in \mathbb{N}}$  is Cauchy in  $\mathcal{X}$ , implying by completeness that  $x_k \rightarrow x_*$  in the limit as  $k \rightarrow \infty$ , for some limit  $x_* \in \mathcal{X}$ . Being a contraction, the mapping  $\phi$  is continuous, and so taking the limit  $k \rightarrow \infty$  in (A.9), we obtain that

$$x_* = \lim_{k \rightarrow \infty} x_{k+1} = \lim_{k \rightarrow \infty} \phi(x_k) = \phi\left(\lim_{k \rightarrow \infty} x_k\right) = \phi(x_*).$$

In other words,  $x_*$  is a fixed point of  $\phi$ .

**Uniqueness.** Assume that  $y_* \in \mathcal{X}$  is a fixed point. Then,

$$\|y_* - x_*\| = \|\phi(y_*) - \phi(x_*)\| \leq L\|y_* - x_*\|,$$

which implies that  $y_* = x_*$  since  $L < 1$ . □

*Remark A.1.* The Banach fixed point theorem holds also in complete metric spaces.

## A.4 Vector norms

In the vector space  $\mathbf{C}^n$ , the most commonly used norms are particular cases of the  $p$ -norm, also called Hölder norm.

**Definition A.4.** Given  $p \in [1, \infty]$ , the  $p$ -norm of a vector  $\mathbf{x} \in \mathbf{C}^n$  is defined as follows:

$$\|\mathbf{x}\|_p := \begin{cases} (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}} & \text{if } p < \infty, \\ \max\{|x_1|, \dots, |x_n|\} & \text{if } p = \infty. \end{cases}$$

The values of  $p$  most commonly encountered in applications are 1, 2 and  $\infty$ . The 1-norm is sometimes called the *taxicab* or *Manhattan* norm, and the 2-norm is usually called the *Euclidean norm*. The explicit expressions of these norms are

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|, \quad \|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}.$$

Notice that the infinity norm  $\|\bullet\|_\infty$  may be defined as the limit of the  $p$ -norm as  $p \rightarrow \infty$ :

$$\|\mathbf{x}\|_\infty := \lim_{p \rightarrow \infty} \|\mathbf{x}\|_p.$$

In the rest of this chapter, the notations  $\langle \bullet, \bullet \rangle$  and  $\|\bullet\|$  without subscript always refer to the Euclidean inner product (A.1) and induced norm, unless specified otherwise.

## A.5 Matrix norms

Given two norms  $\|\bullet\|_\alpha$  and  $\|\bullet\|_\beta$  on  $\mathbf{C}^m$  and  $\mathbf{C}^n$ , respectively, we define the *operator norm* induced by  $\|\bullet\|_\alpha$  and  $\|\bullet\|_\beta$  of the matrix  $\mathbf{A}$  as

$$\|\mathbf{A}\|_{\alpha,\beta} = \sup\{\|\mathbf{A}\mathbf{x}\|_\alpha : \mathbf{x} \in \mathbf{C}^n, \|\mathbf{x}\|_\beta \leq 1\}. \quad (\text{A.10})$$

The term *operator norm* is motivated by the fact that, to any matrix  $\mathbf{A} \in \mathbf{C}^{m \times n}$ , there naturally corresponds the linear operator from  $\mathbf{C}^n$  to  $\mathbf{C}^m$  with action  $\mathbf{x} \mapsto \mathbf{A}\mathbf{x}$ . Matrix norms of the type (A.10) are also called *subordinate* matrix norms. An immediate corollary of the definition (A.10) is that, for all  $\mathbf{x} \in \mathbf{C}^n$ ,

$$\|\mathbf{A}\mathbf{x}\|_\alpha = \|\mathbf{A}\hat{\mathbf{x}}\|_\alpha \|\mathbf{x}\|_\beta \leq \sup\{\|\mathbf{A}\mathbf{y}\|_\alpha : \|\mathbf{y}\|_\beta \leq 1\} \|\mathbf{x}\|_\beta = \|\mathbf{A}\|_{\alpha,\beta} \|\mathbf{x}\|_\beta, \quad \hat{\mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|_\beta}. \quad (\text{A.11})$$

⚙️ **Exercise A.4.** Show that equation (A.10) defines a norm on  $\mathbf{C}^{m \times n}$ .

The matrix  $p$ -norm is defined as the operator norm (A.10) in the particular case where  $\|\bullet\|_\alpha$  and  $\|\bullet\|_\beta$  are both Hölder norms with the same value of  $p$ .

**Definition A.5.** Given  $p \in [1, \infty]$ , the  $p$ -norm of a matrix  $\mathbf{A} \in \mathbf{C}^{m \times n}$  is given by

$$\|\mathbf{A}\|_p := \sup\{\|\mathbf{A}\mathbf{x}\|_p : \mathbf{x} \in \mathbf{C}^n, \|\mathbf{x}\|_p \leq 1\}. \quad (\text{A.12})$$

Not all matrix norms are induced by vector norms. For example, the Frobenius norm, which is widely used in applications, is not induced by a vector norm. It is, however, induced by an inner product on  $\mathbf{C}^{m \times n}$ .

**Definition A.6.** The Frobenius norm of  $\mathbf{A} \in \mathbf{C}^{m \times n}$  is given by

$$\|\mathbf{A}\|_{\text{F}} = \left( \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{\frac{1}{2}}. \quad (\text{A.13})$$

A matrix norm  $\|\bullet\|$  is said to be submultiplicative if, for any two matrices  $\mathbf{A} \in \mathbf{C}^{m \times n}$  and  $\mathbf{B} \in \mathbf{C}^{n \times \ell}$ , it holds that

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|.$$

All subordinate matrix norms, for example the  $p$ -norms, are submultiplicative, and so is the Frobenius norm.

⚙️ **Exercise A.5.** Write down the inner product on  $\mathbf{C}^{m \times n}$  corresponding to (A.13).

⚙️ **Exercise A.6.** Show that the matrix  $p$ -norm is submultiplicative.

## A.6 Diagonalization and spectral theorem

**Definition A.7.** A square matrix  $A \in \mathbf{C}^{n \times n}$  is said to be diagonalizable if there exists an invertible matrix  $P \in \mathbf{C}^{n \times n}$  and a diagonal matrix  $D \in \mathbf{C}^{n \times n}$  such that

$$AP = PD. \quad (\text{A.14})$$

In this case, the diagonal elements of  $D$  are called the eigenvalues of  $A$ , and the columns of  $P$  are called the eigenvectors of  $A$ .

Denoting by  $\mathbf{e}_i$  the  $i$ -th column of  $P$  and by  $\lambda_i$  the  $i$ -th diagonal element of  $D$ , we have by (A.14) that  $A\mathbf{e}_i = \lambda_i\mathbf{e}_i$  or, equivalently,  $(A - \lambda_i I_n)\mathbf{e}_i = \mathbf{0}$ . Here  $I_n$  is the  $\mathbf{C}^{n \times n}$  identity matrix. Therefore, a complex number  $\lambda$  is an eigenvalue of  $A$  if and only if  $\det(A - \lambda I_n) = 0$ . In other words, the eigenvalues of  $A$  are the roots of  $\det(A - \lambda I_n)$ , which is called the *characteristic polynomial*.

### Symmetric matrices and spectral theorem

The transpose of a matrix  $A \in \mathbf{C}^{m \times n}$  is denoted by  $A^T \in \mathbf{C}^{n \times m}$  and defined as the matrix with entries  $a_{ij}^T = a_{ji}$ . The conjugate transpose of  $A$ , denoted by  $A^*$ , is the matrix obtained by taking the transpose and taking the complex conjugate of all the entries. A real matrix that is equal to its transpose is necessarily square and called *symmetric*, and a complex matrix that is equal to its conjugate transpose is called *Hermitian*. Hermitian matrices, of which real symmetric matrices are a subset, enjoy many nice properties, the main one being that they are diagonalizable with a matrix  $Q$  that is unitary, i.e. such that  $Q^{-1} = Q^*$ . This is the content of the *spectral theorem*, a pillar of linear algebra with important generalizations to infinite-dimensional operators.

**Theorem A.4** (Spectral theorem for Hermitian matrices). *If  $A \in \mathbf{C}^{n \times n}$  is Hermitian, then there exists a unitary matrix  $Q \in \mathbf{C}^{n \times n}$  and a diagonal matrix  $D \in \mathbf{R}^{n \times n}$  such that*

$$AQ = QD.$$

*Sketch of the proof.* The result is trivial for  $n = 1$ . Reasoning by induction, we assume that the result is true for Hermitian matrices in  $\mathbf{C}^{(n-1) \times (n-1)}$  and prove that it then also holds for  $A \in \mathbf{C}^{n \times n}$ .

**Step 1. Existence of a real eigenvalue.** By the fundamental theorem of algebra, there exists at least one solution  $\lambda_1 \in \mathbf{C}$  to the equation  $\det(A - \lambda I_n) = 0$ , to which there corresponds at least one solution  $\mathbf{q}_1 \in \mathbf{C}^n$  of norm 1 to the equation  $(A - \lambda_1 I_n)\mathbf{q}_1 = \mathbf{0}$ . The eigenvalue  $\lambda_1$  is necessarily real because

$$\lambda_1 \langle \mathbf{q}_1, \mathbf{q}_1 \rangle = \langle \lambda_1 \mathbf{q}_1, \mathbf{q}_1 \rangle = \langle A\mathbf{q}_1, \mathbf{q}_1 \rangle = \langle \mathbf{q}_1, A\mathbf{q}_1 \rangle = \langle \mathbf{q}_1, \lambda_1 \mathbf{q}_1 \rangle = \bar{\lambda}_1 \langle \mathbf{q}_1, \mathbf{q}_1 \rangle.$$

**Step 2. Using the induction hypothesis.** Next, take an orthonormal basis  $(\mathbf{e}_2, \dots, \mathbf{e}_n)$  of the orthogonal complement  $\text{Span}\{\mathbf{q}_1\}^\perp$  and construct the unitary matrix

$$V = \begin{pmatrix} \mathbf{q}_1 & \mathbf{e}_2 & \dots & \mathbf{e}_n \end{pmatrix},$$



i.e. the matrix with columns  $\mathbf{q}_1, \mathbf{e}_2$ , etc. A calculation gives,

$$\mathbf{V}^* \mathbf{A} \mathbf{V} = \begin{pmatrix} \langle \mathbf{q}_1, \mathbf{A} \mathbf{q}_1 \rangle & \langle \mathbf{q}_1, \mathbf{A} \mathbf{e}_2 \rangle & \dots & \langle \mathbf{q}_1, \mathbf{A} \mathbf{e}_n \rangle \\ \langle \mathbf{e}_2, \mathbf{A} \mathbf{q}_1 \rangle & \langle \mathbf{e}_2, \mathbf{A} \mathbf{e}_2 \rangle & \dots & \langle \mathbf{e}_2, \mathbf{A} \mathbf{e}_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \mathbf{e}_n, \mathbf{A} \mathbf{q}_1 \rangle & \langle \mathbf{e}_n, \mathbf{A} \mathbf{e}_2 \rangle & \dots & \langle \mathbf{e}_n, \mathbf{A} \mathbf{e}_n \rangle \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \langle \mathbf{e}_2, \mathbf{A} \mathbf{e}_2 \rangle & \dots & \langle \mathbf{e}_2, \mathbf{A} \mathbf{e}_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \langle \mathbf{e}_n, \mathbf{A} \mathbf{e}_2 \rangle & \dots & \langle \mathbf{e}_n, \mathbf{A} \mathbf{e}_n \rangle \end{pmatrix}.$$

Let us denote the  $(n-1) \times (n-1)$  lower right block of this matrix by  $\mathbf{V}_{n-1}$ . This is a Hermitian matrix of size  $n-1$  so, using the induction hypothesis, we deduce that  $\mathbf{V}_{n-1} = \mathbf{Q}_{n-1} \mathbf{D}_{n-1} \mathbf{Q}_{n-1}^*$  for appropriate matrices  $\mathbf{Q}_{n-1} \in \mathbf{C}^{(n-1) \times (n-1)}$  and  $\mathbf{D}_{n-1} \in \mathbf{R}^{(n-1) \times (n-1)}$  which are unitary and diagonal, respectively.

**Step 3. Constructing Q and D.** Define now

$$\mathbf{Q} = \mathbf{V} \begin{pmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{Q}_{n-1} \end{pmatrix}.$$

It is not difficult to verify that  $\mathbf{Q}$  is a unitary matrix, and we have

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{Q}_{n-1}^* \end{pmatrix} \mathbf{V}^* \mathbf{A} \mathbf{V} \begin{pmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{Q}_{n-1} \end{pmatrix} = \begin{pmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{Q}_{n-1}^* \end{pmatrix} \begin{pmatrix} \lambda_1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{V}_{n-1} \end{pmatrix} \begin{pmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{Q}_{n-1} \end{pmatrix}.$$

Developing the last expression, we obtain

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} \lambda_1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{D}_{n-1} \end{pmatrix},$$

which concludes the proof.  $\square$

We deduce, as a corollary of the spectral theorem, that if  $\mathbf{e}_1$  and  $\mathbf{e}_2$  are eigenvectors of a Hermitian matrix associated with different eigenvalues, then they are necessarily orthogonal for the Euclidean inner product. Indeed, since  $\mathbf{A} = \mathbf{A}^*$  and the eigenvalues are real, it holds that

$$\begin{aligned} (\lambda_1 - \lambda_2) \langle \mathbf{e}_1, \mathbf{e}_2 \rangle &= \langle \lambda_1 \mathbf{e}_1, \mathbf{e}_2 \rangle - \langle \mathbf{e}_1, \bar{\lambda}_2 \mathbf{e}_2 \rangle \\ &= \langle \mathbf{A} \mathbf{e}_1, \mathbf{e}_2 \rangle - \langle \mathbf{e}_1, \mathbf{A} \mathbf{e}_2 \rangle = \langle \mathbf{A} \mathbf{e}_1, \mathbf{e}_2 \rangle - \langle \mathbf{A}^* \mathbf{e}_1, \mathbf{e}_2 \rangle = 0. \end{aligned}$$

The largest eigenvalue of a matrix, in modulus, is called the *spectral radius* and denoted by  $\rho$ . The following result relates the 2-norm of a matrix to the spectral radius of  $\mathbf{A} \mathbf{A}^*$ .

**Proposition A.5.** *It holds that  $\|\mathbf{A}\|_2 = \sqrt{\rho(\mathbf{A}^* \mathbf{A})}$ .*

*Proof.* Since  $\mathbf{A}^* \mathbf{A}$  is Hermitian, it holds by the spectral theorem that  $\mathbf{A}^* \mathbf{A} = \mathbf{Q} \mathbf{D} \mathbf{Q}^*$  for some unitary matrix  $\mathbf{Q}$  and real diagonal matrix  $\mathbf{D}$ . Therefore, denoting by  $(\mu_i)_{1 \leq i \leq n}$  the (positive)

diagonal elements of  $\mathbf{D}$  and introducing  $\mathbf{y} := \mathbf{Q}^* \mathbf{x}$ , we have

$$\begin{aligned} \|\mathbf{Ax}\| &= \sqrt{\mathbf{x}^* \mathbf{A}^* \mathbf{A} \mathbf{x}} = \sqrt{\mathbf{x}^* \mathbf{Q} \mathbf{D} \mathbf{Q}^* \mathbf{x}} \\ &= \sqrt{\sum_{i=1}^n \mu_i y_i^2} \leq \sqrt{\rho(\mathbf{A}^* \mathbf{A})} \sqrt{\sum_{i=1}^n y_i^2} = \sqrt{\rho(\mathbf{A}^* \mathbf{A})} \|\mathbf{x}\|, \end{aligned} \quad (\text{A.15})$$

where we used in the last equality the fact that  $\mathbf{y}$  has the same norm as  $\mathbf{x}$ , because  $\mathbf{Q}$  is unitary. It follows from (A.15) that  $\|\mathbf{A}\| \leq \sqrt{\rho(\mathbf{A}^* \mathbf{A})}$ , and the converse inequality also holds true since  $\|\mathbf{Ax}\| = \sqrt{\rho(\mathbf{A}^* \mathbf{A})} \|\mathbf{x}\|$  if  $\mathbf{x}$  is the eigenvector of  $\mathbf{A}^* \mathbf{A}$  corresponding to an eigenvalue of modulus  $\rho(\mathbf{A}^* \mathbf{A})$ .  $\square$

To conclude this section, we recall and prove the Courant–Fisher theorem.

**Theorem A.6** (Courant–Fisher Min-Max theorem). *The eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  of a Hermitian matrix are characterized by the relation*

$$\lambda_k = \max_{\mathcal{S}, \dim(\mathcal{S})=k} \left( \min_{\mathbf{x} \in \mathcal{S} \setminus \{0\}} \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} \right). \quad (\text{A.16})$$

*Proof.* Let  $\mathbf{v}_1, \dots, \mathbf{v}_n$  be normalized and pairwise orthogonal eigenvectors associated with the eigenvalues  $\lambda_1, \dots, \lambda_n$ , and let  $\mathcal{S}_k = \text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ . Any  $\mathbf{x} \in \mathcal{S}_k$  may be expressed as a linear combination  $\mathbf{x} = \alpha_1 \mathbf{v}_1 + \dots + \alpha_k \mathbf{v}_k$ , and so

$$\forall \mathbf{x} \in \mathcal{S}_k, \quad \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} = \frac{\sum_{i=1}^k \lambda_i |\alpha_i|^2}{\sum_{i=1}^k |\alpha_i|^2} \geq \lambda_k.$$

Therefore, it holds that

$$\min_{\mathbf{x} \in \mathcal{S}_k \setminus \{0\}} \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} \geq \lambda_k,$$

which proves the  $\geq$  direction of (A.16). For the  $\leq$  direction, let  $\mathcal{U}_k = \text{Span}\{\mathbf{v}_k, \dots, \mathbf{v}_n\}$ . Using a well-known result from linear algebra, we calculate that, for any subspace  $\mathcal{S} \subset \mathbf{C}^n$  of dimension  $k$ ,

$$\begin{aligned} \dim(\mathcal{S} \cap \mathcal{U}_k) &= \dim(\mathcal{S}) + \dim(\mathcal{U}_k) - \dim(\mathcal{S} + \mathcal{U}_k) \\ &\geq k + (n - k + 1) - n = 1. \end{aligned}$$

Therefore, any  $\mathcal{S} \subset \mathbf{C}^n$  of dimension  $k$  has a nonzero intersection with  $\mathcal{U}_k$ . But since any vector in  $\mathcal{U}_k$  can be expanded as  $\beta_1 \mathbf{v}_k + \dots + \beta_n \mathbf{v}_n$ , we have

$$\forall \mathbf{x} \in \mathcal{U}_k, \quad \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} = \frac{\sum_{i=k}^n \lambda_i |\alpha_i|^2}{\sum_{i=k}^n |\alpha_i|^2} \leq \lambda_k.$$

This shows that

$$\forall \mathcal{S} \subset \mathbf{C}^n \text{ with } \dim(\mathcal{S}) = k, \quad \min_{\mathbf{x} \in \mathcal{S} \setminus \{0\}} \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{x}} \leq \lambda_k,$$

which enables to conclude the proof.  $\square$

**⚙️ Exercise A.7.** *Prove that if  $\mathbf{A} \in \mathbf{R}^{n \times n}$  is diagonalizable as in (A.14), then  $\mathbf{A}^n = \mathbf{P} \mathbf{D}^n \mathbf{P}^{-1}$ .*

## A.7 Similarity transformation and Jordan normal form

In this section, we work with matrices in  $\mathbf{C}^{n \times n}$ . A *similarity transformation* is a mapping of the type  $\mathbf{C}^{n \times n} \ni \mathbf{A} \mapsto \mathbf{P}^{-1}\mathbf{A}\mathbf{P} \in \mathbf{C}^{n \times n}$ , where  $\mathbf{P} \in \mathbf{C}^{n \times n}$  is a nonsingular matrix. If two matrices are related by a similarity transformation, then they are called *similar*, because they may be viewed as two representations of the same linear mapping in different bases.

**Definition A.8** (Jordan block). A Jordan block with dimension  $n$  is a matrix of the form

$$\mathbf{J}_n(\lambda) = \begin{pmatrix} \lambda & 1 & & & \\ & \lambda & 1 & & \\ & & \ddots & \ddots & \\ & & & \lambda & 1 \\ & & & & \lambda \end{pmatrix}$$

The parameter  $\lambda \in \mathbf{C}$  is called the eigenvalue of the Jordan block.

A Jordan block is diagonalizable if and only if it is of dimension 1. The only eigenvector of a Jordan block is  $(1 \ 0 \ \dots \ 0)^T$ . The power of a Jordan block admits an explicit expression.

**Lemma A.7.** *It holds that*

$$\mathbf{J}_n(\lambda)^k = \begin{pmatrix} \lambda^k & \binom{k}{1}\lambda^{k-1} & \binom{k}{2}\lambda^{k-2} & \dots & \dots & \binom{k}{n-1}\lambda^{k-n+1} \\ & \lambda^k & \binom{k}{1}\lambda^{k-1} & \dots & \dots & \binom{k}{n-2}\lambda^{k-n+2} \\ & & \ddots & \ddots & & \vdots \\ & & & \ddots & \ddots & \vdots \\ & & & & \lambda^k & \binom{k}{1}\lambda^{k-1} \\ & & & & & \lambda^k \end{pmatrix}. \quad (\text{A.17})$$

*Proof.* The explicit expression of the Jordan block can be obtained by decomposing the block as  $\mathbf{J}_n(\lambda) = \lambda\mathbf{I} + \mathbf{N}$  and using the binomial formula:

$$(\lambda\mathbf{I} + \mathbf{N})^k = \sum_{i=0}^k \binom{k}{i} (\lambda\mathbf{I})^{k-i} \mathbf{N}^i.$$

To conclude the proof, we use the fact that  $\mathbf{N}^i$  is a matrix with zeros everywhere except for  $i$ -th super-diagonal, which contains only ones. Moreover  $\mathbf{N}^i = \mathbf{0}_{n \times n}$  if  $i \geq n$ .  $\square$

A matrix is said to be of *Jordan normal form* if it is block-diagonal with Jordan blocks on

the diagonal. In other words, a matrix  $J \in \mathbf{C}^{n \times n}$  is of Jordan normal form if

$$J = \begin{pmatrix} J_{n_1}(\lambda_1) & & & & \\ & J_{n_2}(\lambda_2) & & & \\ & & \ddots & & \\ & & & J_{n_{k-1}}(\lambda_{k-1}) & \\ & & & & J_{n_k}(\lambda_k) \end{pmatrix}$$

with  $n_1 + \dots + n_k = n$ . Note that  $\lambda_1, \dots, \lambda_k$  are the eigenvalues of  $A$ . We state without proof the following important result.

**Proposition A.8** (Jordan normal form). *Any matrix  $A \in \mathbf{C}^{n \times n}$  is similar to a matrix in Jordan normal form. In other words, there exists an invertible matrix  $P \in \mathbf{C}^{n \times n}$  and a matrix in normal Jordan form  $J \in \mathbf{C}^{n \times n}$  such that*

$$A = PJP^{-1}$$

## A.8 Oldenburger's theorem and Gelfand's formula

The following result establishes a necessary and sufficient condition for the convergence of  $\|A^k\|$  to 0 in terms of the spectral radius of  $A$ , and for any matrix norm  $\|\bullet\|$ .

**Proposition A.9** (Oldenburger). *Let  $\rho(A)$  denote the spectral radius of  $A \in \mathbf{C}^{n \times n}$  and  $\|\bullet\|$  be a matrix norm. Then*

- $\|A^k\| \rightarrow 0$  in the limit as  $k \rightarrow \infty$  if and only if  $\rho(A) < 1$ .
- $\|A^k\| \rightarrow \infty$  in the limit as  $k \rightarrow \infty$  if and only if  $\rho(A) > 1$ .

*Proof.* Since all matrix norms are equivalent, we can assume without loss of generality that  $\|\bullet\|$  is the 2-norm. We prove only the equivalence  $\|A^k\| \rightarrow 0 \Leftrightarrow \rho(A) < 1$ . The other statement can be proved similarly. By [Proposition A.8](#), there exists a nonsingular matrix  $P$  such that  $A = PJP^{-1}$ , for a matrix  $J \in \mathbf{C}^{n \times n}$  which is in normal Jordan form. Since  $\rho(A) = \rho(J)$  and  $\|A^k\| \rightarrow 0$  if and only if  $\|J^k\| \rightarrow 0$ , it is sufficient to show that  $\|J^k\| \rightarrow 0 \Leftrightarrow \rho(J) < 1$ . The latter statement follows from the expression of the power of a Jordan block given in [Lemma A.7](#).  $\square$

With this result, we can prove Gelfand's formula, which relates the spectral radius to the asymptotic growth of  $\|A^k\|$ , and is used in [Chapter 4](#).

**Proposition A.10** (Gelfand's formula). *Let  $A \in \mathbf{C}^{n \times n}$ . It holds for any norm that*

$$\lim_{k \rightarrow \infty} \|A^k\|^{1/k} = \rho(A)$$

*Proof.* Let  $0 < \varepsilon < \rho(A)$  and define  $A^+ = \frac{A}{\rho(A) + \varepsilon}$  and  $A^- = \frac{A}{\rho(A) - \varepsilon}$ . It holds by construction

that  $\rho(\mathbf{A}^+) < 1$  and  $\rho(\mathbf{A}^-) > 1$ . Using Proposition A.9, we deduce that

$$\lim_{k \rightarrow \infty} \|(\mathbf{A}^+)^k\| = 0, \quad \lim_{k \rightarrow \infty} \|(\mathbf{A}^-)^k\| = \infty.$$

Therefore, it holds that

$$\limsup_{k \rightarrow \infty} \|(\mathbf{A}^+)^k\|^{\frac{1}{k}} \leq 1, \quad \liminf_{k \rightarrow \infty} \|(\mathbf{A}^-)^k\|^{\frac{1}{k}} \geq 1.$$

Substituting the expressions of  $\mathbf{A}^+$  and  $\mathbf{A}^-$ , we deduce that

$$\limsup_{k \rightarrow \infty} \|\mathbf{A}^k\|^{\frac{1}{k}} \leq \rho(\mathbf{A}) + \varepsilon, \quad \liminf_{k \rightarrow \infty} \|\mathbf{A}^k\|^{\frac{1}{k}} \geq \rho(\mathbf{A}) - \varepsilon.$$

Since  $\varepsilon$  was arbitrary, we obtain that

$$\rho(\mathbf{A}) \leq \liminf_{k \rightarrow \infty} \|\mathbf{A}^k\|^{\frac{1}{k}} \leq \limsup_{k \rightarrow \infty} \|\mathbf{A}^k\|^{\frac{1}{k}} \leq \rho(\mathbf{A}),$$

which implies the statement. □

# Appendix B

## Brief introduction to Julia

In this chapter, we very briefly present some of the basic features and functions of Julia. Most of the information contained in this chapter can be found in the online manual, to which we provide pointers in each section.

### Installing Julia

The suggested programming environment for this course is the open-source text editor Visual Studio Code. You may also use *Vim* or *Emacs*, if you are familiar with any of these.

☐ **Task 1.** *Install Visual Studio Code. Install also the Julia and Jupyter Notebook extensions.*

### Obtaining documentation

To find documentation on a function from the Julia console, type “?” to access “help mode”, and then the name of the function. Tab completion is helpful for listing available function names.

☐ **Task 2.** *Read the help pages for **if**, **while** and **for**. More information on these keywords is available in the [online documentation](#).*

*Remark B.1* (Shorthand **if** notation). If there is no **elseif** clause, it is sometimes convenient to use the following shorthand notations instead of an **if** block.

```
condition = true

# Assign x = 0 if `condition` is true, else assign x = 2
x = condition ? 0 : 2

# Print "true" if `condition` is true
condition && println("true")
```

```
# Print "false" if `condition` is false
condition || println("false")
```

## Installing and using a package [\[link to relevant manual section\]](#)

To install a package from the Julia REPL (Read Evaluate Print Loop, also more simply called the Julia console), first type `]` to enter the package REPL, and then type `add` followed by the name of the package to install. After it has been added, a package can be used with the `import` keyword. A function `fun` defined in a package `pack` can be accessed as `pack.fun`. For example, to plot the cosine function from the Julia console or in a script, write

```
import Plots
Plots.plot(cos)
```

Alternatively, a package may be imported with the `using` keyword, and then functions can be accessed without specifying the package name. While convenient, this approach is less descriptive; it does not explicitly show what package a function comes from. For this reason, it is often recommended to use `import`, especially in a large codebase.

**Task 3.** *Install the `Plots` package, read the documentation of the `Plots.plot` function, and plot the function  $f(x) = \exp(x)$ . The tutorial on plotting available at [this link](#) may be useful for this exercise.*

*Remark B.2.* We have seen that `?` and `]` enable to access “help mode” and “package mode”, respectively. Another mode which is occasionally useful is “shell mode”, which is accessed with the character `;` and allows to type `bash` commands, such as `cd` to change directory. See [this part](#) of the manual for additional documentation on Julia modes.

## Printing output

The functions `println` and `print` enable to display output. The former adds a new line at the end and the latter does not. The symbol `$`, followed by a variable name or an expression within brackets, can be employed to perform *string interpolation*. For instance, the following code prints `a = 2`, `a2 = 4`.

```
a = 2
println("a = $a, a^2 = $(a*a)")
```

To print a matrix in an easily readable format, the `display` function is very useful.

## Defining functions [\[link to relevant manual section\]](#)

Functions can be defined using a `function` block. For example, the following code block defines a function that prints “Hello, NAME!”, where `NAME` is the string passed as argument.

```
function hello(name)
    # Here * is the string concatenation operator
```

```

println("Hello, " * name)
end

# Call the function
hello("Bob")

```

If the function definition is short, it is convenient to use the following more compact syntax:

```
hello(name) = println("Hello, " * name)
```

Sometimes, it is useful to define a function without giving it a name, called an *anonymous function*. This can be achieved in Julia using the arrow notation `->`. For example, the following expressions calculate the squares and cubes of the first 5 natural numbers. Here, the function `map` enables to transform the collection passed as second argument by applying the function passed as first argument to each element.

```
squares = map(x -> x^2, [1, 2, 3, 4, 5])
cubes = map(x -> x^3, [1, 2, 3, 4, 5])
```

The `return` keyword can be used for returning a value to the function caller. Several values, separated by commas, can be returned at once. For instance, the following function takes a number  $x$  and returns a tuple  $(x, x^2, x^3)$ .

```

function powers(x)
    return x, x^2, x^3
end

# This is an equivalent definition in short notation
short_powers(x) = x, x^2, x^3

# This assigns a = 2, b = 4, c = 8
a, b, c = powers(2)

```

Like many other languages, including Python and Scheme, Julia follows a convention for argument-passing called “pass-by-sharing”: values passed as arguments to a function are not copied, and the arguments act as new bindings within the function body. It is possible, therefore, to modify a value passed as argument, provided this value is of mutable type. Functions that modify some of their arguments usually end with an exclamation mark `!`. For example, the following code prints first `[4, 3, 2, 1]`, because the function `sort` does not modify its argument, and then it prints `[1, 2, 3, 4]`, because the function `sort!` does.

```

x = [4, 3, 2, 1]
y = sort(x) # y is sorted
println(x); sort!(x); println(x)

```

Similarly, when displaying several curves in a figure, we first start with the function `plot`, and then we use `plot!` to modify the existing figure.



```
import Plots
Plots.plot(cos)
Plots.plot!(sin)
```

As a final example to illustrate argument-passing, consider the following code. Here two arguments are passed to the function `test`: an array, which is a mutable value, and an integer, which is immutable. The instruction `arg1[1] = 0` modifies the array to which both `a` and `arg1` are bindings. The instruction `arg2 = 2`, on the other hand, just causes the variable `arg2` to point to a new immutable value (3), but it does not change the destination of the binding `b`, which remains the immutable value 2. Therefore, the code prints `[0, 2, 3]` and 3.

```
function test(arg1, arg2)
    arg1[1] = 0
    arg2 = 2
end
a = [1, 2, 3]
b = 3
test(a, b)
println(a, b)
```

□ **Task 4** (Euler–Mascheroni constant for the harmonic series). *Euler showed that*

$$\lim_{N \rightarrow \infty} \left( -\ln(N) + \sum_{n=1}^N \frac{1}{n} \right) = \gamma := 0.577\dots$$

Write a function that returns an approximation of the Euler–Mascheroni constant  $\gamma$  by evaluating the expression between brackets at a finite value of  $N$ .

```
function euler_constant(N)
    # Your code comes here
end
```

□ **Task 5** (Ancient algorithms). *The goal of this exercise is to explore three of the oldest algorithms ever invented.*

- Circa 1600 BC, the Babylonians invented an iterative method for calculating the square root of a number. Read the relevant information on the associated [Wikipedia page](#) and write a function that calculates the square root of the argument using this algorithm.

```
function babylonian_square_root(n)
    # Your code comes here
end
# The function should return the square root of n
```

- Circa 300 BC, the Greek mathematician Euclid of Alexandria published the *Elements*, his famous mathematical treatise. In one of the books, he proposes an algorithm for calculating

the greatest common divisor of two numbers. This algorithm, which is still in common use today, is based on the observation that if  $a > b \geq 0$  are natural numbers, then

$$\gcd(a, b) = \gcd(b, r), \quad (\text{B.1})$$

where  $r$  is the remainder of the division of  $a$  by  $b$ . Indeed, in view of the equation

$$a = qb + r,$$

the common divisors of  $\{a, b\}$  coincide with those of  $\{b, r\}$ . Using (B.1), write a function to calculate the greatest common divisor of two numbers.

```
function euclid_gcd(a, b)
    # Your code comes here
end
```

- Circa 200 BC, the Greek mathematician Eratosthenes of Cyrene invented a method for efficiently calculating the prime numbers, which is now known as the sieve of Eratosthenes. Read the associated [Wikipedia page](#) and write a function implementing this algorithm.

```
function eratosthenes_sieve(n)
    # Your code comes here
end
# The function should return an array containing all the prime
# numbers less than or equal to n.
```

**□ Task 6** (Tower of Hanoi). We consider a variation on the classic Tower of Hanoi problem, in which the number  $r$  of pegs is allowed to be larger than 3. We denote the pegs by  $p_1, \dots, p_r$ , and assume that the problem includes  $n$  disks with radii 1 to  $n$ . The tower is initially constructed in  $p_1$ , with the disks arranged in order of decreasing radius, the largest at the bottom. The goal of the problem is to reconstruct the tower at  $p_r$  by moving the disks one at the time, with the constraint that a disk may be placed on top of another only if its radius is smaller.

It has been conjectured that the optimal solution, which requires the minimum number of moves, can always be decomposed into the following three steps, for some  $k \in \{1, n-1\}$ :

- First move the top  $k$  disks of the tower to peg  $p_2$ ;
- Then move the bottom  $n - k$  disks of the tower to  $p_r$  without using  $p_2$ ;
- Finally, move the top of the tower from  $p_2$  to  $p_r$ .

This suggests a recursive procedure for solving the problem, known as the Frame-Stewart algorithm. Write a Julia function  $T(n, r)$  returning the minimal number of moves necessary.

### Local and global scopes [\[link to relevant manual section\]](#)

Some constructs in Julia introduce scope blocks, notably **for** and **while** loops, as well as **function** blocks. The variables defined within these structures are not available outside them. For example

```

if true
    a = 1
end
println(a)

```

prints 1, because **if** does not introduce a scope block, but

```

for i in [1, 2, 3]
    a = 1
end
println(a)

```

produces **ERROR: LoadError: UndefVarError: a not defined.** The variable `a` defined within the **for** loop is said to be in the *local scope* of the loop, whereas a variable defined outside of it is in the *global scope*. In order to modify a global variable from a local scope, the **global** keyword must be used. For instance, the following code

```

a = 1
for i in [1, 2, 3]
    global a += 1
end
println(a)

```

modifies the global variable `a` and prints 4.

### Multi-dimensional arrays [\[link to relevant manual section\]](#)

A working knowledge of multi-dimensional arrays is important for this course, because vectors and matrices are ubiquitous in numerical algorithms. In Julia, a two-dimensional array can be created by writing its lines one by one, separating them with a semicolon `;`. Within a line, elements are separated by a space. For example, the instruction

```
M = [1 2 3; 4 5 6]
```

creates the matrix

$$M = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$$

More generally, the semicolon enables vertical concatenation while space concatenates horizontally. For example, `[M M]` defines the matrix

$$\begin{pmatrix} 1 & 2 & 3 & 1 & 2 & 3 \\ 4 & 5 & 6 & 4 & 5 & 6 \end{pmatrix}$$

The expression `M[r, c]` gives the  $(r, c)$  matrix element of  $M$ , located at row  $r$  and column  $c$ . The special entry **end** can be used to access the last row or column. For instance, `M[end-1, end]` gives the matrix entry in the second to last row and the last column. From the matrix  $M$  above, the submatrix `[2 3; 5 6]` can be obtained with `M[:, 2:3]`. Here the row index `:` means “select all lines” and the column index `2:3` means “select columns 2 to 3”. Likewise, the submatrix `[1 3; 4 6]` may be extracted with `M[:, [1; 3]]`.

*Remark B.3* (One-dimensional arrays). The comma `,` can also be employed for creating one-dimensional arrays, but its behavior differs slightly from that of the vertical concatenation operator `;`. For example, `x = [1, [2; 3]]` creates a **Vector** object with two elements, the first one being 1 and the second one being `[2; 3]`, which is itself a **Vector**. In contrast, the instruction `x = [1; [1; 2]]` creates the same **Vector** as `[1; 2; 3]` would.

We also mention that the expression `x = [1 2 3]` produces not a one-dimensional **Vector** but a two-dimensional **Matrix**, with one row and three columns. This can be checked using the `size` function, which for `x = [1 2 3]` returns the tuple `(1, 3)`.

There are many built-in functions for quickly creating commonly used arrays. For example,

- `transpose(M)` gives the transpose of  $M$ , and `adjoint(M)` or `M'` gives the transpose conjugate. For a matrix with real-valued entries, both functions deliver the same result.
- `zeros(Int, 4, 5)` creates a  $4 \times 5$  matrix of zeros of type **Int**;
- `ones(2, 2)` creates a  $2 \times 2$  matrix of ones of type **Float64**;
- `range(0, 1, length=101)`, or `LinRange(0, 1, 101)`, creates an array of size 101 with elements evenly spaced between 0 and 1 included. More precisely, `range` returns an array-like object, which can be converted to a vector using the `collect` function.
- `collect(reshape(1:9, 3, 3))` creates a  $3 \times 3$  matrix with elements

$$\begin{pmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{pmatrix}$$

Let us also mention the following shorthand notation, called *array comprehension*, for creating vectors and matrices:

- `[i^2 for i in 1:5]` creates the vector `[1, 4, 9, 16, 25]`.
- `[i + 10*j for i in 1:4, j in 1:4]` creates the matrix

$$\begin{pmatrix} 11 & 21 & 31 & 41 \\ 12 & 22 & 32 & 42 \\ 13 & 23 & 33 & 43 \\ 14 & 24 & 34 & 44 \end{pmatrix}.$$

- `[i for i in 1:10 if ispow2(i)]` creates the vector `[1, 2, 4, 8]`. The same result can be achieved with the `filter` function: `filter(ispow2, 1:10)`.

In contrast with Matlab, array assignment in Julia does not perform a copy. For example the following code prints `[1, 2, 3, 4]`, because the instruction `b = a` defines a new binding to the array `a`.

```

a = [2; 2; 3]
b = a
b[1] = 1
append!(b, 4)
println(a)

```

A similar behavior applies when passing an array as argument to a function, as we saw previously. The `copy` function can be used to perform a copy.

□ **Task 7.** Create a 10 by 10 diagonal matrix with the  $i$ -th entry on the diagonal equal to  $i$ .

## Broadcasting

To conclude this chapter, we briefly discuss *broadcasting*, which enables to apply functions to array elements and to perform operations on arrays of different sizes. Julia really shines in this area, with syntax that is both explicit and concise. Rather than providing a detailed definition of broadcasting, which is available in [this part](#) of the official documentation, we illustrate the concept using examples. Consider first the following code block:

```

function welcome(name)
    return "Hello, " * name * "!"
end
result = broadcast(welcome, ["Alice", "Bob"])

```

Here `broadcast` returns an array with elements `"Hello, Alice!"` and `"Hello, Bob!"`, as would the `map` function. Broadcasting, however, is much more flexible because it can handle arrays with different sizes. For instance, `broadcast(gcd, 24, [10, 20, 30])` returns an array of size 3 containing the greatest common divisors of the pairs (24,10), (24,20) and (24,30). Similarly, the instruction `broadcast(+, 1, [1, 2, 3])` returns `[2, 3, 4]`. To understand the latter example, note that `+` (as well as `*`, `-` and `/`) can be called like any other Julia functions; the notation `a + b` is just syntactic sugar for `+(a, b)`.

Since broadcasting is so often useful in numerical mathematics, Julia provides a shorthand notation for it: the instruction `broadcast(welcome, ["Alice", "Bob"])` can be written compactly as `welcome.(["Alice", "Bob"])`. Likewise, the line `broadcast(+, 1, [1, 2, 3])` can be shortened to `(+).(1, [1, 2, 3])`, or to the more readable expression `1 .+ [1, 2, 3]`.

□ **Task 8.** Explain in words what the following instructions do.

```

reshape(1:9, 3, 3) .* [1 2 3]
reshape(1:9, 3, 3) .* [1; 2; 3]
reshape(1:9, 3, 3) * [1; 2; 3]

```

## Appendix C

# Chebyshev polynomials

The Chebyshev polynomials  $(T_n)_{n \in \mathbf{N}}$  are given on  $[-1, 1]$  by the formula

$$\forall x \in [-1, 1], \quad T_n(x) = \cos(n \arccos(x)). \quad (\text{C.1})$$

Although this formula makes sense only if  $x \in [-1, 1]$ , the polynomials are defined for all  $x \in \mathbf{R}$ . Equivalently, the Chebyshev polynomials can be defined from the equation

$$\forall x \in [1, \infty), \quad T_n(x) = \cosh(n \operatorname{arccosh}(x)), \quad (\text{C.2})$$

where  $\cosh(\theta) = \frac{1}{2}(e^\theta + e^{-\theta})$  and  $\operatorname{arccosh}: [1, \infty) \rightarrow [0, \infty)$  is the inverse function of  $\cosh$ . The first few Chebyshev polynomials are illustrated in [Figure C.1](#). It is immediate to show the following properties from (C.1):

- The roots of  $T_n$  are given by

$$z_k = \cos\left(\frac{\pi}{2n} + \frac{k\pi}{n}\right), \quad k = 0, \dots, n-1.$$

These are illustrated in [Figure C.2](#).

- The polynomial  $T_n$  takes the value 1 or -1 when evaluated at

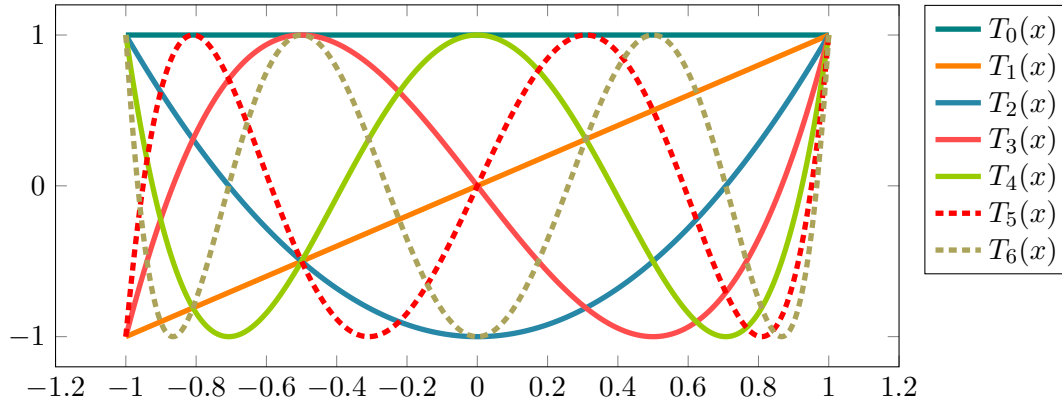
$$x_k = \cos\left(\frac{k\pi}{n}\right), \quad k = 0, \dots, n. \quad (\text{C.3})$$

More precisely, it holds that  $T_n(x_k) = (-1)^k$ .

**⚙️ Exercise C.1.** Show that (C.1) defines a polynomial of degree  $n$ , and find its expression in the usual polynomial notation.

*Solution.* The key idea is to rewrite the cosine function in terms of the complex exponential:

$$\cos(n\theta) = \frac{1}{2}(e^{in\theta} + e^{-in\theta}) = \frac{1}{2}\left((\cos(\theta) + i\sin(\theta))^n + (\cos(\theta) - i\sin(\theta))^n\right).$$


 Figure C.1: Illustration of the first few Chebyshev polynomials over the interval  $[-1, 1]$ .

By expanding the powers on the right-hand side, we obtain

$$\begin{aligned} (\cos(\theta) + i \sin(\theta))^n &= \sum_{j=0}^n \binom{n}{j} \cos(\theta)^{n-j} i^j \sin(\theta)^j \\ (\cos(\theta) - i \sin(\theta))^n &= \sum_{j=0}^n \binom{n}{j} \cos(\theta)^{n-j} (-i)^j \sin(\theta)^j. \end{aligned}$$

The terms corresponding to odd values of  $j$  cancel out in the expression of  $\cos(n\theta)$ , and so we obtain the following expression for  $\cos(n\theta)$  in terms of  $\cos(\theta)$  and  $\sin(\theta)$ :

$$\begin{aligned} \cos(n\theta) &= \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j} \cos(\theta)^{n-2j} i^{2j} \sin(\theta)^{2j} \\ &= \sum_{j=0}^{\lfloor n/2 \rfloor} (-1)^j \binom{n}{2j} \cos(\theta)^{n-2j} (1 - \cos(\theta)^2)^j. \end{aligned}$$

Therefore, we conclude that

$$T_n(x) = \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j} x^{n-2j} (x^2 - 1)^j. \quad (\text{C.4})$$

△

⚙️ **Exercise C.2.** Show that the same polynomials are obtained from (C.2).

*Solution.* Notice that

$$\begin{aligned} \cosh(n\xi) &= \frac{1}{2} (e^{n\xi} + e^{-n\xi}) \\ &= \frac{1}{2} \left( (\cosh(\xi) + \sinh(\xi))^n + (\cosh(\xi) - \sinh(\xi))^n \right). \end{aligned}$$

Using the binomial formula, we obtain

$$\begin{aligned} \cosh(n\xi) &= \frac{1}{2} \sum_{j=0}^n \binom{n}{j} (\cosh(\xi)^{n-j} \sinh(\xi)^j + \cosh(\xi)^{n-j} (-1)^j \sinh(\xi)^j) \\ &= \frac{1}{2} \sum_{j=0}^n \binom{n}{j} \cosh(\xi)^{n-j} (\sinh(\xi)^j + (-1)^j \sinh(\xi)^j). \end{aligned}$$

The contributions of the odd values of  $j$  cancel out, and so we obtain

$$\cosh(n\xi) = \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j} \cosh(\xi)^{n-2j} \sinh(\xi)^{2j}.$$

Since  $\cosh(\xi)^2 - \sinh(\xi)^2 = 1$ , we deduce that

$$\cosh(n\xi) = \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{j} \cosh(\xi)^{n-2j} (\cosh(\xi)^2 - 1)^j,$$

which after the substitution of  $\xi = \operatorname{arccosh}(x)$  leads to (C.4).  $\triangle$

⚙️ **Exercise C.3** (Yet another expression for the Chebyshev polynomials). *Show that  $T_n(x)$  may be defined from the formula*

$$T_n(x) = \frac{1}{2} \left( x + \sqrt{x^2 - 1} \right)^n + \frac{1}{2} \left( x - \sqrt{x^2 - 1} \right)^n \quad \text{for } |x| \geq 1. \quad (\text{C.5})$$

*Solution.* We showed in the solution of Exercise C.2 that

$$\cosh(n\xi) = \frac{1}{2} \left( (\cosh(\xi) + \sinh(\xi))^n + (\cosh(\xi) - \sinh(\xi))^n \right).$$

Letting  $\xi = \operatorname{arccosh}(x)$  in this equation and using that  $\cosh(\xi)^2 - \sinh(\xi)^2 = 1$ , we obtain

$$T_n(x) = \frac{1}{2} \left( (x + \sqrt{x^2 - 1})^n + (x - \sqrt{x^2 - 1})^n \right),$$

which is the required formula.  $\triangle$

⚙️ **Exercise C.4** (Recursion relation). *Show that the Chebyshev polynomials satisfy the relation*

$$\forall n \in \{1, 2, \dots\}, \quad T_{n+1} = 2xT_n - T_{n-1}. \quad (\text{C.6})$$

*Solution.* It is sufficient to show the identity for  $x \in [-1, 1]$ , where the formula (C.1) applies. Using well-known trigonometric identities, we have

$$\begin{aligned} \cos((n+1)\theta) &= \cos(n\theta)\cos(\theta) - \sin(n\theta)\sin(\theta) \\ \cos((n-1)\theta) &= \cos(n\theta)\cos(\theta) + \sin(n\theta)\sin(\theta). \end{aligned}$$

Adding both equations and rearranging, we obtain

$$\cos((n+1)\theta) = 2\cos(n\theta)\cos(\theta) - \cos((n-1)\theta).$$



Therefore, using this equation with  $\theta = \arccos(x)$ , we obtain the statement.  $\triangle$

*Remark C.1.* The recursion relation in Exercise C.4 can be employed to show by recursion that  $T_n(x)$  is indeed a polynomial of degree  $n$ .

⚙️ **Exercise C.5.** Since  $T_n: \mathbf{R} \rightarrow \mathbf{R}$  is a polynomial, it may be written in the standard form

$$T_n(x) = \alpha_n^{(n)}x^n + \dots + \alpha_1^{(n)}x + \alpha_0^{(n)}.$$

Prove that  $\alpha_n^{(n)} = 2^{(n-1)}$  provided that  $n \geq 1$ .

*Solution.* From the definition (C.1), the Chebyshev polynomials of degrees 0 and 1 are given by  $T_0(x) = 1$  and  $T_1(x) = x$ . The statement then follows by recursion, using Exercise C.4.  $\triangle$

⚙️ **Exercise C.6.** Let  $\xi \in \mathbf{R} \setminus (-1, 1)$ . Show that, among all the polynomials in  $\mathbf{P}(n)$  that are bounded from above by 1 in absolute value uniformly over the interval  $(-1, 1)$ , the Chebyshev polynomial  $T_n$  achieves the largest absolute value when evaluated at  $\xi$ .

*Solution.* Reasoning by contradiction, we assume that there exists  $p \in \mathbf{P}(n)$  that satisfies

$$\sup_{x \in (-1, 1)} |p(x)| \leq 1 \quad \text{and} \quad |p(\xi)| > |T_n(\xi)|.$$

Let  $q(x) = p(x)T_n(\xi)/p(\xi)$ . Then by construction  $q(\xi) = T_n(\xi)$  and

$$\sup_{x \in (-1, 1)} |q(x)| < 1.$$

Consequently, denoting by  $x_k$  the points defined in (C.3), we have that

$$\forall k \in \{0, \dots, n\}, \quad (-1)^k (T_n - q)(x_k) > 0.$$

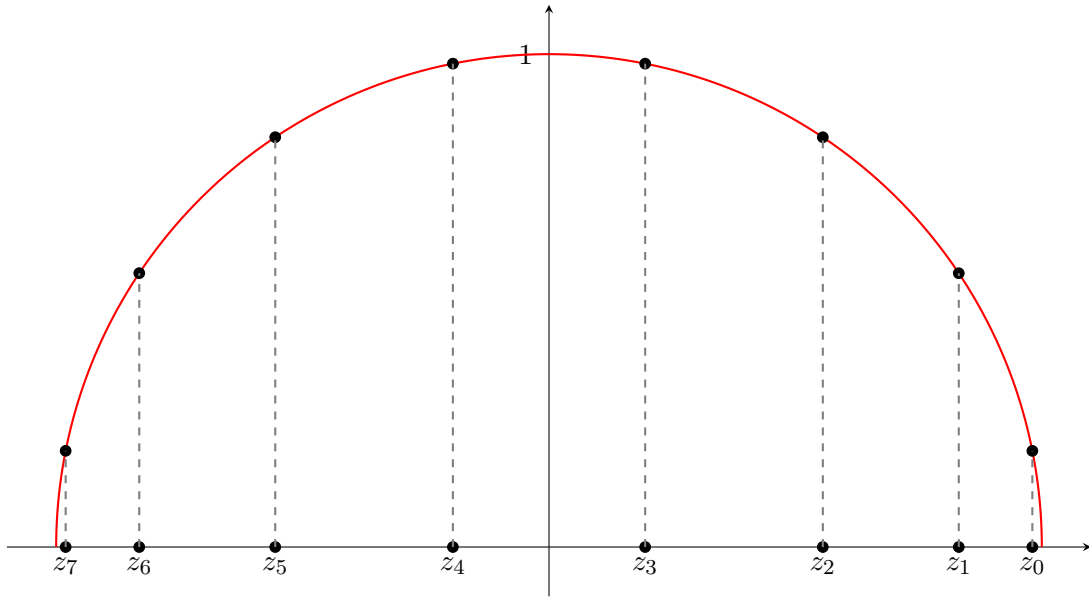
In other words, the polynomial  $T_n - q$  takes positive values at  $\{x_0, x_2, x_4, \dots\}$  and negative values at  $\{x_1, x_3, x_5, \dots\}$ . Consequently, by the intermediate value theorem,  $T_n - q$  possesses  $n$  distinct roots in the open interval  $(-1, 1)$ . Since, in addition,  $(T_n - q)(\xi) = 0$ , we deduce that  $T_n - q$  has  $n + 1$  distinct roots, which is a contradiction given that  $T_n - q$  is a nonzero polynomial of degree at most  $n$ .  $\triangle$

⚙️ **Exercise C.7.** Assume that  $0 < \lambda_1 < \lambda_2$ . Prove that for any polynomial  $p \in \mathbf{P}(n)$  that satisfies  $p(0) = 1$ , it holds that

$$\sup_{\lambda \in (\lambda_1, \lambda_2)} |p(\lambda)| \geq \frac{1}{T_n(\xi)}, \quad \xi := \frac{\lambda_2 + \lambda_1}{\lambda_2 - \lambda_1},$$

with equality for

$$p_*(\lambda) = \frac{T_n\left(\frac{\lambda_1 + \lambda_2 - 2\lambda}{\lambda_2 - \lambda_1}\right)}{T_n\left(\frac{\lambda_1 + \lambda_2}{\lambda_2 - \lambda_1}\right)}. \quad (\text{C.7})$$


 Figure C.2: Roots of the Chebyshev polynomial  $T_8$ .

*Solution.* Assume that  $p \in \mathbf{P}(n)$  is such that  $p(0) = 1$ , and let  $q \in \mathbf{P}(n)$  be given by

$$q(\mu) = p\left(\frac{\lambda_1 + \lambda_2 - (\lambda_2 - \lambda_1)\mu}{2}\right) \quad \Leftrightarrow \quad p(\lambda) = q\left(\frac{\lambda_1 + \lambda_2 - 2\lambda}{\lambda_2 - \lambda_1}\right).$$

Since  $\xi > 1$ , it holds from (C.5) that  $T_n(\xi) > 0$  and it follows from Exercise C.6 that

$$p(0) = q(\xi) \leq T_n(\xi) \sup_{\mu \in (-1,1)} |q(\mu)| = T_n(\xi) \sup_{\lambda \in (\lambda_1, \lambda_2)} |p(\lambda)|,$$

with equality when  $q \propto T_n$ , i.e. when

$$p(\lambda) \propto T_n\left(\frac{\lambda_1 + \lambda_2 - 2\lambda}{\lambda_2 - \lambda_1}\right).$$

The expression (C.7) then follows from the fact that  $p_*(0) = 1$ . △

# Bibliography

- [1] E. CUTHILL and J. MCKEE. Reducing the bandwidth of sparse symmetric matrices. In *Proceedings of the 1969 24th national conference*, pages 157–172, 1969.
- [2] A. ERNST and G. STOLTZ. *Calcul Scientifique*. Lecture notes at École des Ponts, 2014.  
URL: <https://cermics.enpc.fr/cours/CS/poly.pdf>.
- [3] D. GOLDBERG. What every computer scientist should know about floating-point arithmetic. *ACM computing surveys (CSUR)*, **23**(1):5–48, 1991.
- [4] IEEE Standard for Binary Floating-Point Arithmetic. *ANSI/IEEE Std 754-1985*:1–20, 1985.  
DOI: [10.1109/IEEESTD.1985.82928](https://doi.org/10.1109/IEEESTD.1985.82928).
- [5] IEEE Standard for Floating-Point Arithmetic. *IEEE Std 754-2008*:1–70, 2008.  
DOI: [10.1109/IEEESTD.2008.4610935](https://doi.org/10.1109/IEEESTD.2008.4610935).
- [6] IEEE Standard for Floating-Point Arithmetic. *IEEE Std 754-2019*:1–84, 2019.  
DOI: [10.1109/IEEESTD.2019.8766229](https://doi.org/10.1109/IEEESTD.2019.8766229).
- [7] V. LEGAT. *Mathématiques et méthodes numériques*. Lecture notes for the course EPL1104 at École polytechnique de Louvain, 2009.  
URL: <https://perso.uclouvain.be/vincent.legat/documents/epl1104/epl1104-notes-v8-2.pdf>.
- [8] A. MAGNUS. *Analyse numérique: approximation, interpolation, intégration*. Lecture notes for the course INMA2171 at École polytechnique de Louvain, 2010.  
URL: <https://perso.uclouvain.be/alphonse.magnus/num1a/m217111.pdf>.
- [9] J. M. ORTEGA and W. C. RHEINBOLDT. *Iterative solution of nonlinear equations in several variables*, volume **30** of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2000.  
DOI: [10.1137/1.9780898719468](https://doi.org/10.1137/1.9780898719468).  
URL: <https://doi-org.extranet.enpc.fr/10.1137/1.9780898719468>.
- [10] A. QUARTERONI, R. SACCO, and F. SALERI. *Numerical mathematics*, volume **37** of *Texts in Applied Mathematics*. Springer-Verlag, Berlin, second edition, 2007.  
DOI: [10.1007/b98885](https://doi.org/10.1007/b98885).
- [11] Y. SAAD. *Iterative methods for sparse linear systems*. Society for Industrial and Applied Mathematics, Philadelphia, PA, second edition, 2003.  
DOI: [10.1137/1.9780898718003](https://doi.org/10.1137/1.9780898718003).  
URL: <https://doi-org.extranet.enpc.fr/10.1137/1.9780898718003>.
- [12] Y. SAAD. *Numerical methods for large eigenvalue problems*, volume **66** of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2011.  
DOI: [10.1137/1.9781611970739.ch1](https://doi.org/10.1137/1.9781611970739.ch1).  
URL: <https://doi-org.extranet.enpc.fr/10.1137/1.9781611970739.ch1>.

- [13] J. R. SHEWCHUK et al. An introduction to the conjugate gradient method without the agonizing pain, 1994.  
URL: <https://www.cs.cmu.edu/~quake-papers/painless-conjugate-gradient.pdf>.
- [14] L. N. TREFETHEN. The definition of numerical analysis. Technical report, Cornell University, 1992.  
URL: <https://cims.nyu.edu/~oneil/courses/sp18-math252/trefethen-def-na.pdf>.
- [15] P. VAN DOOREN. *Analyse numérique*. Lecture notes for the course INMA1170 at École polytechnique de Louvain, 2012.
- [16] F. VERHULST. *Nonlinear differential equations and dynamical systems*. Universitext. Springer-Verlag, Berlin, second edition, 1996.  
DOI: 10.1007/978-3-642-61453-8.  
URL: <https://doi-org.extranet.enpc.fr/10.1007/978-3-642-61453-8>.
- [17] M. VIANELLO and R. ZANOVELLO. On the superlinear convergence of the secant method. *Amer. Math. Monthly*, **99**(8):758–761, 1992.  
DOI: 10.2307/2324244.  
URL: <https://doi-org.extranet.enpc.fr/10.2307/2324244>.
- [18] C VUIK and D. J. P. LAHAYE. *Scientific Computing*. Lecture notes for the course wi4201 at Delft University of Technology, 2019.  
URL: [http://ta.twi.tudelft.nl/users/vuik/wi4201/wi4201\\_notes.pdf](http://ta.twi.tudelft.nl/users/vuik/wi4201/wi4201_notes.pdf).