

**Refresher course in maths
and
a project on numerical modeling done in twos**

Course number: MAP 502

Master of Science and Technology STEEM
École Polytechnique, Autumn 2022

Samuel Amstutz¹ and Thomas Wick^{1,2}

¹ École Polytechnique, Centre de Mathématiques Appliquées,
91128 Palaiseau, France

<https://samuelamstutz.wordpress.com/>

² Leibniz University Hannover, Institute of Applied Mathematics,
Welfengarten 1, 30167 Hannover, Germany

<https://thomaswick.org>



Organization of this class

Class number	Date	Topic	Chapter	lecturer
1	23/09	parts I, II	1,2,3,4,5,6, 9.1-9.17	W+A
2	30/09	part III	10, 12, 15.1-15.11	W+A
3	07/10	part II,III	8, 11, 9.20-9.23	W+A
4	14/10	part III	13, 14, 15.12-15.21	W+A
5	21/10	part IV	16, 17, 18.1-18.5	W+A
6	28/10	mid-term exam		W+A
7	18/11	projects		W+A
8	25/11	projects		W+A
9	09/12	projects		W+A
10	16/12	presentations (exam)		W+A

Some chapters and sections are indicated as **complement**. These are not relevant for the mid-term exam, but may need to be consulted (not all, but some) for the project work.

Preface

These lecture notes have been written since September 2016 for various purposes. Originally, they served as handout for a numerical modeling project within the STEEM (Energy Environment: Science Technology and Management Master) program. Later in 2019, this class was redesigned as a refresher class in maths plus the numerical modeling project. Third, Part I-III of these notes serve currently as preparation for admission in Numerics 1 for foreign Master students at the Leibniz University Hannover, Germany. Consequently, by nature, these notes address an audience that is international and interdisciplinary with participants not necessarily having a previous Bachelor degree in mathematics.

Content-wise, we strive the motivation for scientific computing, floating point numbers, linear algebra (pure and numerics), analysis (some pure elements and numerics), up to the Fourier transform. Moreover, for the sake of a state-of-the-art preparation of the students for their further studies and career, we dedicate our emphasis to differential equations in the final part. Parts that are not relevant for the mid-term evaluation in the STEEM program are indicated as *complement*. However, some of these parts might become relevant for the later numerical projects. In these projects, the tasks are fivefold: mathematical modeling, algorithmic design, implementation, presenting scientific findings, and learning how to interpret those results in order to draw scientific conclusions. Therefore, we provide measures how to build confidence into numerical findings such as intuition, error analysis, convergence analysis, and comparison to manufactured solutions. These projects are partially purely academic, but modern directions such as regression problems and neural networks can be chosen as well.

Specifically, with respect to the lastly mentioned applications and the impact of mathematics reaching out in science, technical developments, economics, and the society, it is our sincere aim to not only carry over the beauty of mathematics, but also to show its usefulness in other fields.

In case of questions, please do not hesitate to contact us!

Paris and Hannover
September 2022

Samuel Amstutz
Thomas Wick

Contents

I	Introduction	11
1	Guiding questions in numerical modeling	13
1.1	What is numerical modeling?	13
1.2	Concepts in numerical mathematics	14
1.2.1	Definitions	14
1.2.2	Examples	15
1.3	First related mathematical notions	17
1.3.1	Well-posedness	17
1.3.2	Measuring errors	17
1.3.3	Illustration of three errors (experimental, theory, numerical) with the clothesline problem	18
1.3.4	Measuring distances: metrics and norms	18
1.3.5	Application in numerical mathematics	20
1.3.6	Notation for the complexity and convergence orders: Landau symbols	20
1.4	Differential equations as an example and guiding questions	21
1.4.1	Definitions of differential equations	21
1.4.2	Some important questions and tasks with regard to numerical concepts	22
2	Basic ingredients: sets and numbers	23
2.1	Sets	23
2.2	Classical sets of numbers	24
2.3	Real numbers and their usage in computers: floating-point number system	24
2.3.1	Round-off errors and machine precision	24
2.3.2	Influence of machine precision in numerical mathematics	25
2.3.3	Stability and condition number	26
2.4	Complex numbers: a brief reminder	28
3	Exercises	31
II	Linear algebra and related numerical notions	33
4	Linear systems and matrices	35
4.1	Gaussian elimination for solving linear systems	35
4.1.1	General definitions	35
4.1.2	Elementary row operations, row echelon systems	36
4.1.3	The Gaussian elimination method	36
4.2	Vectors and matrices	37
4.2.1	Basic definitions	37
4.2.2	Operations on matrices	38
4.2.3	Matrix of a linear system	41
4.2.4	Invertible matrices	41

4.3	Linear subspaces of \mathbb{K}^n	42
4.3.1	Definitions and first properties	42
4.3.2	Bases	43
4.3.3	Dimension	43
4.4	Determinant	43
4.4.1	Expansion with respect to a row or a column	44
4.4.2	Properties	45
4.4.3	Determinant of a product, characterization of invertible matrices	45
4.5	LU decomposition (complement)	46
4.5.1	Matrix representation of Gaussian elimination	46
4.5.2	LU decomposition for diagonal-dominant matrices	55
4.5.3	Case of Hermitian positive definite matrices: Cholesky decomposition	55
4.6	Canonical inner product, Euclidean norm, matrix norm	55
4.6.1	Canonical inner product and Euclidean norm	55
4.6.2	Matrix norm	56
5	Diagonalization of square matrices and applications	57
5.1	Eigenvalues, eigenvectors, diagonalization	57
5.1.1	Eigenvalues and eigenvectors	57
5.1.2	Diagonalization	58
5.1.3	Case of Hermitian matrices	58
5.2	Applications	58
5.3	The power method (complement)	59
6	Vector spaces and elements of topology	61
6.1	Vector spaces	61
6.1.1	General definition	61
6.1.2	Vector spaces of finite dimension	62
6.2	Basic topology	62
6.2.1	Norms and balls	62
6.2.2	Open sets, closed sets (complement)	63
6.2.3	Closure, interior (complement)	64
6.2.4	Converging sequences (complement)	64
6.2.5	Equivalent norms	64
6.3	Special normed spaces (complement)	64
6.3.1	Inner product spaces	64
6.3.2	Banach and Hilbert spaces	65
6.3.3	Complement on matrix norms	66
7	Linear equation systems and iterative solvers (complement)	67
7.1	Stability analysis of linear equation systems	67
7.1.1	Motivation	67
7.1.2	Stability analysis	68
7.2	Basic solvers	70
7.2.1	Fixed-point solvers: Richardson, Jacobi, Gauss-Seidel	70
7.2.2	Gradient descent	72
7.3	Conjugate gradient method	73
7.3.1	Formulation of the CG scheme.	74
7.3.2	Convergence analysis of the CG scheme	76
7.4	Preconditioning	78
7.5	GMRES - generalized minimal residual method	80

8	Convex sets and systems of linear inequalities	85
8.1	Convex sets: definition and first properties	85
8.2	Convex combinations and convex hull (complement)	86
8.2.1	Convex combinations	86
8.2.2	Convex hull	86
8.3	Projection	88
8.4	Cones	89
8.4.1	Definitions	89
8.4.2	Normal cone (complement)	90
8.4.3	Polar cones (complement)	90
8.5	Systems of linear inequalities, introduction to linear programming	91
8.5.1	Linear inequality systems: Fourier-Motzkin elimination	91
8.5.2	Conical hull and Farkas lemma (complement)	92
8.5.3	Introduction to linear programming (complement)	93
9	Exercises	95
III	Functions and related numerical notions	101
10	Functions of one or several variables	103
10.1	Basic concepts	103
10.2	Differentiation	103
10.2.1	Partial derivatives	103
10.2.2	Fréchet derivative, Jacobian matrix	104
10.2.3	Chain rule	104
10.2.4	Schwarz' theorem	105
10.2.5	Multiindex notation (complement)	105
10.3	Classical differential operators	105
10.4	Taylor expansions, Hessian matrix	106
10.5	Convex functions	107
10.5.1	Definitions	107
10.5.2	Characterization	108
10.6	Integration by parts	108
10.6.1	Substitution rule	108
10.6.2	Integration by parts and Green's formulae	109
11	Introduction to nonlinear optimization	111
11.1	General concepts	111
11.1.1	Brief classification of optimization problems	111
11.1.2	Problem setting	111
11.1.3	Global and local minimizers	111
11.1.4	Optimality conditions	112
11.2	Convex case (complement)	113
11.2.1	Specific aspects of convex optimization	113
11.2.2	Towards constrained optimization	114
11.3	Steepest descent method	114
11.3.1	Descent methods	114
11.3.2	Steepest descent	114
11.3.3	Line search	115
11.4	Gradient descent with projection (complement)	116
11.5	Application of Newton-type methods (complement)	118
11.5.1	Newton's method for optimization	118

11.5.2	Quasi-Newton methods	119
12	Solving nonlinear equations	121
12.1	Introduction to iterative methods	121
12.2	The bisection method	122
12.3	Fixed points	123
12.3.1	Reminder: the (Banach) fixed point theorem	123
12.3.2	Attractive and repulsive fixed points	124
12.3.3	Calculating fixed points	124
12.4	The Newton method	125
12.5	Generalization to higher dimensions (complement)	126
12.5.1	Newton's method: going from \mathbb{R} to higher dimensions	126
12.5.2	A basic algorithm for a residual-based Newton method	127
12.5.3	Example of the basic Newton method	127
12.5.4	Example using a Newton defect-correction scheme including line search	128
12.5.5	Newton's method in higher dimensions and the Newton-Kantorovich theorem	129
12.5.6	Globalization of Newton's method	133
12.6	Newton's method for a coupled, nonlinear system of ODEs (complement)	136
12.7	Iteration schemes in nonlinear optimization (complement)	137
12.7.1	Linear regression	137
12.7.2	Neural networks	138
13	Interpolation and approximation	141
13.1	Polynomial interpolation	141
13.1.1	Introduction	141
13.1.2	Existence, uniqueness, expression	141
13.1.3	Interpolation error (complement)	142
13.2	Numerical differentiation	143
13.2.1	Approximation of the first derivative	143
13.2.2	Approximation of the second derivative (complement)	144
13.3	Numerical integration (complement)	144
13.3.1	Goal	144
13.3.2	General principle	145
13.3.3	Elementary quadrature formula	145
13.3.4	Simplest rules: box, mid-point, trapezoidal, Simpson	145
13.3.5	Composite quadrature formula	146
13.3.6	Order of a quadrature formula	146
13.3.7	Interpolatory quadrature	146
13.3.8	Gauss quadrature	147
14	Trigonometric interpolation, Fourier series and Fourier transform	149
14.1	Trigonometric interpolation: discrete Fourier transform	149
14.2	Fourier series	151
14.2.1	Series	151
14.2.2	Convergence of Fourier series	151
14.2.3	Parseval's equality	152
14.2.4	Applications	152
14.3	Fourier transform	152
14.3.1	General concept	152
14.3.2	A glimpse on the Dirac distribution	153
15	Exercises	155

IV	Differential equations	161
16	Introduction to ODEs (ordinary differential equations)	163
16.1	An introductory example	163
16.2	The model ODE	164
16.3	Well-posedness	165
16.4	Second order linear ODE's with constant coefficients	165
17	Finite differences for ODEs	167
17.1	Problem statement of an IVP (initial value problem)	167
17.2	Stiff problems	167
17.3	One-step schemes	168
17.3.1	The Euler method	168
17.3.2	Implicit schemes	168
17.4	Numerical analysis	169
17.4.1	Stability	171
17.4.2	Consistency / local discretization error - convergence order	173
17.4.3	Convergence	173
17.5	Detailed numerical tests	174
17.5.1	Problem statement	174
17.5.2	Discussion of the results for test 1 with $a=0.25$	175
17.5.3	Investigating the instability of the forward Euler method: test 3 with $a=-10$	176
17.6	A BVP: boundary value problem	176
17.6.1	A 1D model problem: Poisson	176
17.6.2	Well-posedness of the continuous problem	176
17.6.3	Spatial discretization	179
17.6.4	Solving the linear equation system	180
17.6.5	Well-posedness of the discrete problem	181
17.6.6	Numerical analysis: consistency, stability, and convergence	182
17.6.7	Numerical test: 1D Poisson	186
17.7	Computational convergence analysis (complement)	189
17.7.1	Discretization error	189
17.7.2	Computationally-obtained convergence order	191
17.7.3	Spatial discretization error	192
17.7.4	Temporal discretization error for fixed spatial numerical solution	193
17.7.5	Extrapolation to the limit	193
17.7.6	Iteration error	194
18	Exercises	195
V	Projects on numerical modeling in teams	197
19	Project work	199
19.1	Idea and formal aspects	199
19.1.1	Choice of your project and contact email addresses	199
19.1.2	Final exam mid December	200
19.2	List of projects (to be confirmed)	200
19.2.1	Direct and iterative solution of linear equation systems	200
19.2.2	Numerical methods for eigenvalues	201
19.2.3	Image compression with SVD (singular valued decomposition) - related to eigenvalues	201
19.2.4	Numerical optimization and application to regression problems	202

19.2.5	Neural network for image classification	204
19.2.6	Numerical methods for root-finding problems (I)	204
19.2.7	Numerical methods for root-finding problems (II)	205
19.2.8	Shape optimization for fluids using FreeFem++	205
19.2.9	Discrete Fourier transform	206
19.2.10	Numerical solution of ODEs with finite differences	206
19.2.11	Predator-prey systems	207
19.2.12	Neural network approximation of ODEs	208
19.3	Practical hints during the project work and final presentation	208
19.3.1	Typical guideline questions	208
19.3.2	Hints for presenting results in a talk (final exam)	209
19.3.3	Some example pages of a presentation	210
	Bibliography	213
	Index	216

Part I

Introduction

Chapter 1

Guiding questions in numerical modeling

In this class, we refresh basic notions of numerical linear algebra and numerical analysis in combination with numerical modeling up to differential equations. The content comprises in brief classical notions of analytical methods and linear algebra. In view of technological applications, our goal is the design and analysis of algorithms that can be used for computer simulations of prototype and complex problems.

1.1 What is numerical modeling?

Numerical modeling is a part of **scientific computing**. The latter comprises three main fields:

1. **Mathematical modeling** and analysis of physical, biological, chemical, economical, financial processes, and so forth;
2. Development of reliable and efficient **numerical methods** and algorithms and their analysis;
3. **Research software development**: Implementation of the derived algorithms

All these steps work in a **feed-back manner** and the different subtasks interact with each other. It is in fact the third above aspect, namely *software and computers*, who helped to establish this third category of science. Thus, a new branch of mathematics, *computational science*, has been established. This kind of mathematics may become experimental like experiments in physics/chemistry/biology.

A key task is the design and analysis of algorithms:

Definition 1.1 (Algorithm). *An algorithm is an instruction for a schematic solution of a mathematical problem statement. The main purpose of an algorithm is to formulate a scheme that can be implemented into a computer to carry out so-called numerical simulations. Direct schemes solve the given problem up to round-off errors (for instance Gaussian elimination). Iterative schemes approximate the solution up to a certain accuracy (for instance Richardson iteration for solving linear equation systems, or fixed-point iterations). Algorithms differ in terms of **accuracy, robustness, and efficiency**.*

An important feedback task is to analyse (rigorously or computationally) these algorithms in order to detect shortcomings and suggest improvements. These can be nowadays on the algorithmic side (classical numerical analysis with convergence proofs) or the computational side (analysis of different discretization levels, parameter variations by just running the code again) or coding improvements (re-organization of, e.g., for-loops in a finite element program), or hardware-specific aspects (e.g., CPU computing, GPU computing).

Computational science allows us to investigate research fields that have partially not been addressable in the past. Why? On the one hand experiments are often too expensive, too far away (Mars, Moon, astronomy in general), the scales are too small (nano-scale for example); or experiments are simply too dangerous. On the other hand, mathematical theory or the explicit solution of an (ambitious) engineering problem in an analytical manner is often impossible!

Mathematics helps to formalise and structure given systems from physics, engineering, economics and other disciplines. Hereupon, established algorithms for the numerical discretization and numerical solution for known parts of the given system shall be further extended for new problems and new applications.

1.2 Concepts in numerical mathematics

1.2.1 Definitions

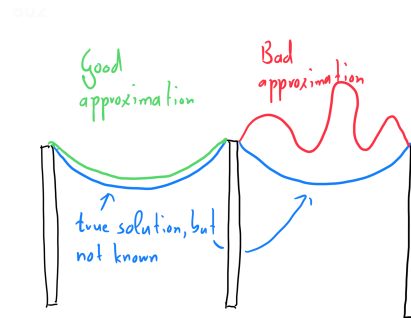
In [15] seven concepts that are very characteristic of numerical modeling were identified. They will be frequently encountered in the forthcoming chapters.

1. **Approximation:** For most mathematical problems, specifically those of practical relevance, the derivation of analytical solutions is difficult or even impossible. Consequently, we must be content with *approximations* (in appropriate function spaces and corresponding topologies), which are often obtained via numerical procedures and computer simulations.
2. **Convergence:** Whether approximations have something to do with a (unique) limit value is studied with the notion of convergence. It is qualitative in the sense that it tells us when a sequence $(a_n)_{n \in \mathbb{N}}$ admits a limit a . In numerical mathematics this limit is often the (unknown) solution we are aiming for, and the a_n are approximate solutions.
3. **Order of convergence:** The speed (or rate) of convergence tells us how fast a numerical approximation tends to limit values. While in analysis we are often merely interested in the convergence itself, in numerical mathematics we must pay attention to how long it takes until an approximate solution has sufficient accuracy. The longer a simulation takes, the more time and more energy (electricity to run the computer, air conditioning of servers, etc.) are consumed. Therefore, we are heavily interested in developing fast algorithms.
4. **Errors:** Numerical mathematics can be considered as the branch of ‘mathematics of errors’. What does this mean? Numerical modeling is not wrong, inexact or non-precise! Since we cut sequences after a finite number of steps or accept sufficiently accurate solutions obtained from a software, we need to say how well this numerical solution approximates the (unknown) exact solution. In other words, we need to determine the errors, which can arise in various forms. Due to its importance, we also provide the extra Section 1.3.2
5. **Error estimation:** This is one of the biggest branches in numerical mathematics and the most classical one. We need to derive error formulae to judge the outcome of a numerical simulation and to measure the difference (distance; see Section 1.3.4) between the numerical solution and the (unknown) exact solution in a certain norm or metric (see again Section 1.3.4)
6. **Efficiency:** In general we can say, the higher the convergence order of an algorithm is, the more efficient the algorithm is. Therefore, we obtain faster the numerical solution to a given problem. But numerical efficiency is not automatically related to resource-effective computing. For instance, developing a parallel code using MPI (message passing interface), hardware-optimization (CPU, GPU), software optimizations (ordering in some optimal way for-loops, arithmetic evaluations, etc.) can further reduce computational costs.
7. **Stability:** Lastly, the robustness of algorithms and implementations with respect to parameter (model, material, numerical) variations, boundary conditions, initial conditions, uncertainties must be studied. Stability relates in the broadest sense to the third condition of Hadamard defined in Section 1.3.1.

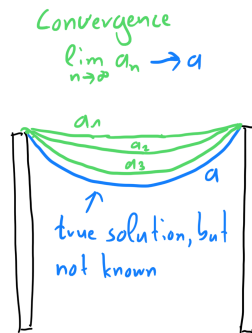
1.2.2 Examples

We illustrate the previous concepts with the help of some examples.

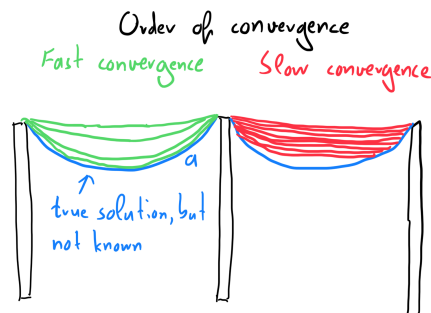
1. **Approximation:** Two approximations of the clothesline¹ problem:



2. **Convergence:** Converging approximations of the clothesline problem:



3. **Order of convergence:** Two different speeds of convergence:



4. **Errors:** We first refer the reader to Section 1.3.3 for a comparison of experimental, theoretical and numerical situations. Second, we sharpen the sense of the influence of different errors. Not

¹Let us explain the clothesline in a very simplified way: wash your clothes in your washing machine. Next, hang your wet clothes on your clothesline, which gets deformed due to gravity. In a great mathematical simplification, we solve the problem $-u''(x) = f$ for $x \in (0, 1)$ with the boundary conditions $u(0) = u(1) = 0$. This model is also known as Poisson problem. To be specific, this is a differential equation, where we search for u according to the given differential equation, the right hand side f (here: gravity) and boundary conditions; see Chapters 16 - 17.

all errors are equally important and sometimes, one might try to ‘optimize’ an error, which has no significant influence. Let’s see this in more detail. Let the errors e_{Model} , $e_{Numerics}$, $e_{Software}$ enter. The total error is defined as

$$e_{Total} = e_{Model} + e_{Numerics} + e_{Software}$$

Let us assume that we have the numbers $e_{Model} = 1000$, $e_{Numerics} = 0.001$, $e_{Software} = 4$, the total error is then given by

$$e_{Total} = 1000 + 0.001 + 4 = 1004.001.$$

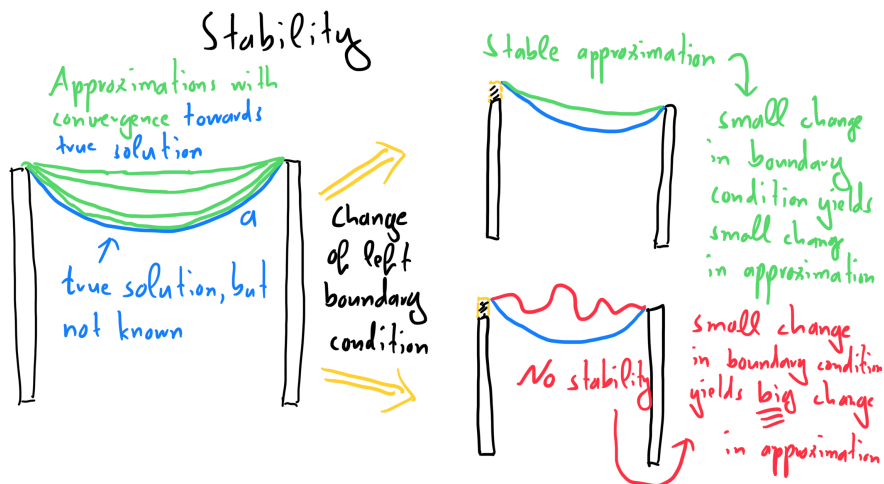
Which error dominates? It is clearly $e_{Model} = 1000$. The relative influence is $e_{Model}/e_{Total} = 0.996$. So, the other two error sources are negligible and would not need further attention in this specific example.

5. **Error estimation:** Error estimation is the process to obtain the concrete numbers 1000, 0.001, 4 in the previous example. Error estimates can be classified into two categories:

- **a priori estimates** include the (unknown) exact solution u , such that $\eta := \eta(u)$, and yield qualitative convergence rates for asymptotic limits. They can be derived before (thus a priori) the approximation is known.
- **a posteriori error estimates** are of the form $\eta := \eta(\tilde{u})$ and explicitly employ the approximation \tilde{u} and therefore yield quantitative information with computable majorants (i.e., bounds) and can be further utilized to design **adaptive schemes**.

6. **Efficiency:** is more or less self-explaining. A first answer is to look at CPU or wall time: how many seconds, minutes, weeks, months does a program need to terminate and yield a result? A second answer is to study ‘iteration numbers’ or arithmetic operations. The latter are often given in terms of the big O notation. Having a linear equation system $Ax = b$ with $A \in \mathbb{R}^{n \times n}$ and $O(n^3)$ complexity means that we need n^3 (cubic in unknowns n) arithmetic operations to calculate the result. For instance, for $n = 100$, we need around 1 000 000 operations. Having another algorithm (yielding the same result of $Ax = b$) with only $O(n)$ operations, means that we only need around 100 operations, which is a great difference. The development of efficient solvers for large linear equations systems is consequently a big branch in numerics and scientific computing.

7. **Stability:** We finally come back to the clothesline problem and change a bit the left boundary condition:



1.3 First related mathematical notions

1.3.1 Well-posedness

The concept of well-posedness, as introduced by Hadamard, is very general and in fact very simple:

- The problem under consideration has a solution;
- This solution is unique;
- The solution depends continuously on the problem data.

The first condition is immediately clear. The second condition is also obvious but often difficult to meet - and in fact many physical processes do not have unique solutions. The last condition says that if a variation of the input data (right hand side, boundary values, initial conditions) is small, then also the (unique) solution should only vary a little.

Remark 1.1. *Problems in which one of the three above conditions is violated are **ill-posed**.*

1.3.2 Measuring errors

Going ahead, we implement an algorithm in a software (for instance in Matlab/Octave, Python, Fortran, C++, Java) using a computer. In the end, several error sources need to be addressed.

1. The set of numbers that can be represented in a computer is finite, therefore a numerical calculation is limited by machine precision, which results in **round-off errors**.
2. The memory of a computer (or cluster) is finite and thus functions and equations can only be represented through approximations. Thus, continuous information has to be represented through discrete information, which results in investigating so-called **discretization errors** or more generally **approximation errors**.
3. **Interpolation errors** appear when complicated functions are interpolated to simpler functions.
4. All further simplifications of a numerical algorithm (in order to solve the discrete problem), with the final goal to reduce the computational time, result into **iteration / truncation errors**. One example is the stopping criterion, which decides after how many steps an iteration is stopped. These can be further divided into **linear** and **nonlinear** iteration errors.
5. **Regularization errors** appear when ill-posed models are modified in order to become well-posed.
6. **Homogenization errors** arise when full-scale models are reduced to larger scale (homogenized) such that they are computationally simpler to solve.
7. **Implementation errors**, better known as **bugs**, appear when mistakes are made implementing algorithms into a software. These errors can simply cause the program to abort. Then it is more or less simple to find the place with the help of a debugger. But there are also subtle errors, which cannot be easily found when the program runs nonetheless, but showing strange results or behaviors.
8. **Model errors**: In order to make a ‘quick guess’ of a possible solution and to start the development of an algorithm aimed at addressing at a later stage a difficult problem, often complicated (nonlinear) equations are reduced to simple (in most cases linear) versions, which results in the so-called **model error**.
9. **Data errors and uncertainties**: the data (e.g., input data, boundary conditions, parameters) are obtained from experimental measurements and may be inaccurate themselves.

10. **Inaccurate experimental setups** that yield wrong reference solutions to which numerical solutions should be compared with.

It is very important to understand that we **never can avoid** all these errors. The important aspect is to **control** these errors and to provide answers if these errors are sufficiently big to influence the interpretation of numerical simulations or if they can be assumed to be small. A big branch of numerical mathematics is to derive error estimates that allow to predict about the size of arising errors.

1.3.3 Illustration of three errors (experimental, theory, numerical) with the clothes-line problem

You need to understand a bit of French for the words, but the figure might be self-explaning.

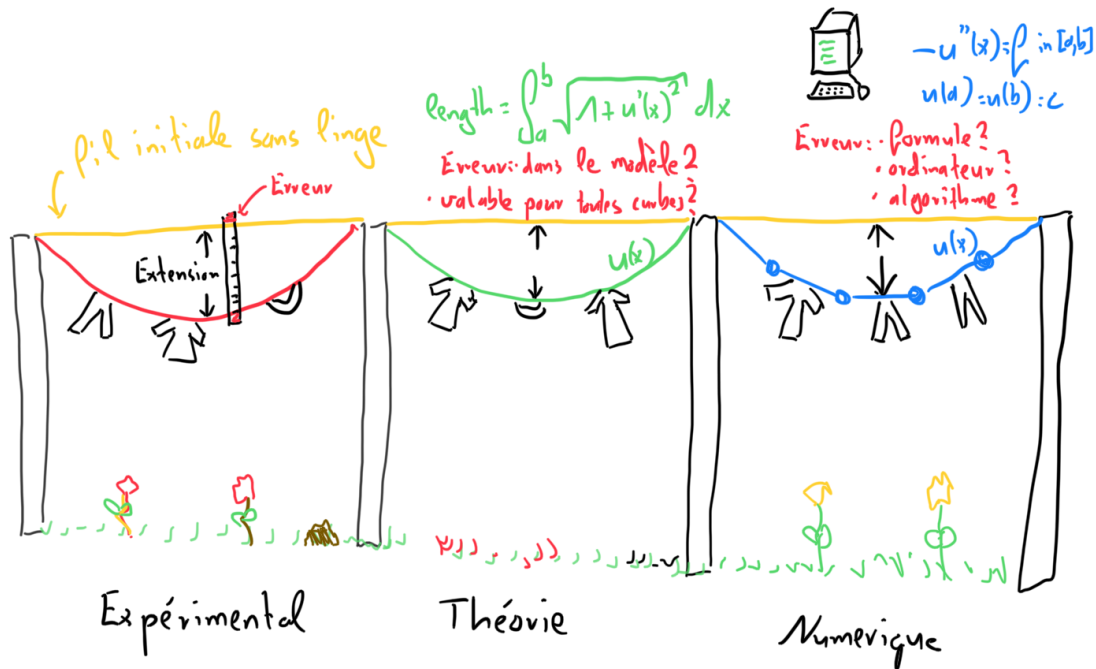


Figure 1.1: Errors in experiments, theory and numerical modeling.

1.3.4 Measuring distances: metrics and norms

To determine the accuracy and quantify the previous errors, we need to frame our problem statements in appropriate spaces that allow us to measure distances. The distance we are interested in is usually between the 'exact' solution u and its numerical approximation \tilde{u} , i.e., it is related to the difference

$$u - \tilde{u}.$$

However, this expression is nearly useless since it can be negative (what does a negative distance mean?) and u and \tilde{u} might be vectors. Consequently, we need something like

$$\|u - \tilde{u}\|.$$

In order to define $\|\cdot\|$, we first introduce metric spaces. A metric space is a set X with a metric on it. This metric associates with any pair of elements (i.e., points), say $a, b \in X$, a distance. The concept of distance is more general than that of norm, in that it does not require the structure of vector space (in particular the difference need not be defined).

Definition 1.2 (Metric space). A metric space is a pair (X, d) where X is a set and d is a metric (or distance) on X . The function d is defined on $X \times X$, where \times denotes the Cartesian product of sets (whose definition is recalled in the next chapter). For all $x, y, z \in X$, it is assumed (axioms defining a distance):

1. d is real-valued, finite and nonnegative;
2. $d(x, y) = 0$ if and only if $x = y$;
3. $d(x, y) = d(y, x)$;
4. $d(x, y) \leq d(x, z) + d(z, y)$.

Example 1.1. On the real line \mathbb{R} the usual metric is defined by

$$d(x, y) = |x - y|.$$

Example 1.2. The Euclidean space \mathbb{R}^3 consists of the set of ordered triples of real numbers $x = (x_1, x_2, x_3)$ and $y = (y_1, y_2, y_3)$ and so forth. The Euclidean metric is defined by

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2}.$$

Example 1.3. On the surface of the earth (which is not a vector space but a manifold) we can define the geodesic distance as the length of the shortest path joining two given points.

In numerical analysis, the framework is usually that of vector spaces and the distances of interest are typically induced by norms.

Definition 1.3 (Normed space). A normed space X is a vector space with a norm. A norm on a real or complex vector space X is a real-valued function on X whose value at $x \in X$ is denoted by $\|x\|$ and which has the following properties:

1. $\|x\| \geq 0$
2. $\|x\| = 0 \Leftrightarrow x = 0$
3. $\|\alpha x\| = |\alpha| \|x\|$
4. $\|x + y\| \leq \|x\| + \|y\|$

where $x, y \in X$ and α is any scalar from the underlying field \mathbb{R} or \mathbb{C} .

Reminders on vector spaces and complements on norms will be given in chapter 6.

Definition 1.4 (Norm vs. metric). A norm on X induces a metric d on X , which is defined by

$$d(x, y) = \|x - y\|.$$

With the help of these definitions it is easy to infer that indeed a norm induces a metric. Therefore, normed spaces are in particular metric spaces. It holds:

Proposition 1.1. The norm is continuous, that is $x \mapsto \|x\|$ is a continuous mapping of $(X, \|\cdot\|)$ into \mathbb{R} .

Distances, induced or not by a norm, can now be employed to study the accuracy of mathematical models and their numerical approximations. Consequently, we are now able to give a meaning to the distance

$$d(u, \tilde{u})$$

between the ‘exact’ solution u and its numerical approximation \tilde{u} .

We finally add an important refinement when errors are evaluated through a norm:

Definition 1.5. Let $\tilde{x} \in X$ be an approximation of $x \in X$. With $\|e\| = \|\tilde{x} - x\|$ we denote the **absolute error**, and with $\|e\|/\|x\|$ we denote the **relative error**.

In most cases, relative errors are of greater importance. For instance, an error of 200m is small if we measure the distance between amphi Painlevé and the moon, but 200m are big if we measure the distance between amphi Painlevé and the CMAP corridor.

1.3.5 Application in numerical mathematics

In numerical mathematics, the key task is often to design appropriate norms for measuring the distance between a numerical approximation u_h and the continuous (exact) solution u . These norms give quantitative information about the error for deciding when an approximation is sufficiently close to our sought solution:

$$\|u - u_h\| \rightarrow 0 \quad (h \rightarrow 0),$$

where h indicates the approximation parameter. In particular, iteration error measurements in certain norms allow to decide when numerical algorithms can be stopped:

$$\|u - u_l\| \rightarrow 0 \quad (l \rightarrow \infty),$$

where l indicates the iteration index. Very often, we have a mixture of situations in which h , representing typically a discretization parameter, and l , the number of iterations, appear simultaneously.

A natural question is whether different norms yield similar answers. Here, one must distinguish between finite-dimensional vector spaces and infinite-dimensional spaces. The latter results into the branch of functional analysis. Therein, norms are in general not equivalent (a precise definition of this notion is given in chapter 6). Consequently, it does matter which norm is adopted. In finite-dimensional spaces, it can be shown, however, that norms are indeed equivalent. We emphasize that we speak here of the vector space in which the exact solution lives, even if the approximate one is sought in a finite-dimensional subspace. Dealing with two equivalent norms implies that convergence for one norm is equivalent to convergence for the other one. However, of course, both errors are not equal, therefore in quantitative error analysis the choice of the norm matters even for intrinsically finite-dimensional problems.

1.3.6 Notation for the complexity and convergence orders: Landau symbols

Definition 1.6 (Landau symbols). (i) Let $g(n)$ be a function with $g \rightarrow \infty$ for $n \rightarrow \infty$. Then $f \in O(g)$ if and only if when

$$\limsup_{n \rightarrow \infty} \left| \frac{f(n)}{g(n)} \right| < \infty$$

and $f \in o(g)$ if and only if

$$\lim_{n \rightarrow \infty} \left| \frac{f(n)}{g(n)} \right| = 0.$$

(ii) Let $g(h)$ be a function with $g(h) \rightarrow 0$ for $h \rightarrow 0$. As before, we define $f \in O(g)$ and $f \in o(g)$:

$$\limsup_{h \rightarrow 0} \left| \frac{f(h)}{g(h)} \right| < \infty \quad \Leftrightarrow \quad f \in O(g),$$

and

$$\lim_{h \rightarrow 0} \left| \frac{f(h)}{g(h)} \right| = 0 \quad \Leftrightarrow \quad f \in o(g).$$

(iii) Specifically:

$$\limsup_{h \rightarrow 0} |f(h)| < \infty \quad \Leftrightarrow \quad f \in O(1),$$

and

$$\lim_{h \rightarrow 0} |f(h)| = 0 \quad \Leftrightarrow \quad f \in o(1).$$

Often, the notation $f = O(g)$ is used rather than $f \in O(g)$ and similarly $f = o(g)$ rather than $f \in o(g)$.

Example 1.4. Seven examples:

1. $\frac{1}{x} = o\left(\frac{1}{x^2}\right) \quad (x \rightarrow 0),$
2. $\frac{1}{x^2} = o\left(\frac{1}{x}\right) \quad (|x| \rightarrow \infty),$

3. $e^{-x} = o(x^{-26})$ ($x \rightarrow \infty$).

4. Let $\varepsilon \rightarrow 0$ and $h \rightarrow 0$. We write

$$h = o(\varepsilon)$$

when

$$\frac{h}{\varepsilon} \rightarrow 0 \quad \text{for } h \rightarrow 0, \quad \varepsilon \rightarrow 0,$$

which means that h tends faster to 0 than ε .

5. Let us assume that we have the error estimate (see for instance Chapter 17)

$$\|y(t_n) - y_n\|_2 = O(k).$$

Here the O notation means nothing else than

$$\frac{\|y(t_n) - y_n\|_2}{k} \rightarrow C \quad \text{for } k \rightarrow 0.$$

6. Let

$$c_1 k^1 + c_2 k^2 + c_3 k^3 + \dots$$

be a development in powers of k . We can shortly write:

$$c_1 k^1 + O(k^2).$$

7. Complexity of algorithms, e.g., LR

$$O(n^3) \quad \text{for } n \rightarrow \infty.$$

Here the fraction converges to a constant C (and not necessarily 0!), which illustrates that $O(\cdot)$ convergence is weaker than $o(\cdot)$ convergence. On the other hand, this should not yield the wrong conclusion that $\|y(t_n) - y_n\|_2$ may not tend to zero. Since $k \rightarrow 0$, also $\|y(t_n) - y_n\|_2 \rightarrow 0$ must hold necessarily.

1.4 Differential equations as an example and guiding questions

As an example of the previously mentioned concepts, we consider differential equations (later in detail in Part IV), which arise in numerous applications in science and engineering and are well-known in continuum mechanics.

1.4.1 Definitions of differential equations

Let us first roughly define the meaning of a differential equation:

Definition 1.7. *A differential equation is a mathematical equation that relates the function with certain of its derivatives.*

Differential equations can be split into two classes:

- **Ordinary differential equations.** An ordinary differential equation (ODE) is an equation (or equation system) involving an unknown function of one independent variable and certain of its derivatives.
- **Partial differential equations.** A partial differential equation (PDE) is an equation (or equation system) involving an unknown function of two or more variables and certain of its partial derivatives.

1.4.2 Some important questions and tasks with regard to numerical concepts

1. What are the underlying physics?
2. What is the mathematical model? Which kind of differential equation are we dealing with? Linear diff. eq.? Coupled system? Nonlinearities? Types of coupling? Order?
3. What discretization scheme is appropriate? Physics-based discretization, i.e., are the physics conserved after discretization in time and/or space?
4. Design of principle algorithms
5. Can we prove that these algorithms really work?
6. Accuracy, efficiency, robustness of algorithms
7. Error-controlled adaptivity or model-order reduction (e.g., with the help of SVD - singular value decomposition; computation of eigenvalues)
8. Linear and/or nonlinear solution schemes? Direct, iterative solution? Multigrid schemes? Computational cost?
9. Analyzing the numerical results: is the numerical solution correct? What does it mean ‘correct’? Is it close enough (appropriate norm!) to some known solution? Comparison with experimental results possible? Comparison with analytical or manufactured solutions u and the numerical solution u_{kh} :

$$\|u - u_{kh}\| \rightarrow 0 \quad \text{for } h, k \rightarrow 0,$$

where k and h are the temporal and spatial discretization parameters, respectively.

10. Advanced postprocessing techniques: quantities of interest rather global norms (see e.g., [20])? Plotting graphical solutions (gnuplot, vtk, ...)? Computational convergence analysis?

The forthcoming sections of these notes aim at refreshing (and perhaps complementing) the mathematical notions needed to address such questions. Applications will mainly be seen through the projects and further courses.

Chapter 2

Basic ingredients: sets and numbers

2.1 Sets

Set theory relies on a series of axioms, which are usually considered as the foundations of mathematics. It includes in particular the definition of natural numbers. This is a very abstract theory, but fortunately it is not needed to know it in details for a practical use of mathematics (rather advanced mathematics existed far before the axiomatic construction). We gather below the main points to be known.

Intuitively, a set is a collection of objects (of any nature) called elements. A set can be finite or not. The empty set, which contains no element, is denoted \emptyset . To write the list of the elements of a set we use the notation

$$A = \{1, 2, 3\}.$$

This means that the set A is made of the numbers 1, 2 and 3. The ordering is irrelevant. To say that 1 belongs to A (or 1 is in A) and 4 does not belong to A we write

$$1 \in A, \quad 4 \notin A.$$

Given two sets A and B , we say that A is included in B , denoted $A \subset B$, if all the elements of A belong to B :

$$A \subset B \iff \forall x \in A, x \in B.$$

The symbol " \forall " means "for all". When $A \subset B$ and $B \subset A$ we say that A and B are equal ($A = B$).

The union of two sets A and B , denoted $A \cup B$, is the set of elements that are either in A or in B :

$$A \cup B = \{x, x \in A \text{ or } x \in B\}.$$

The intersection of two sets A and B , denoted $A \cap B$, is the set of elements that are at the same time in A and in B :

$$A \cap B = \{x, x \in A \text{ and } x \in B\}.$$

The difference of two sets A and B , denoted $A \setminus B$, is the set of elements that are in A but not in B :

$$A \setminus B = \{x, x \in A \text{ and } x \notin B\}.$$

The Cartesian product of two sets A and B , denoted $A \times B$, is the set of pairs (x, y) with $x \in A$ and $y \in B$:

$$A \times B = \{(x, y), x \in A \text{ and } y \in B\}.$$

2.2 Classical sets of numbers

Natural numbers:

$$\begin{aligned}\mathbb{N} &= \{0, 1, 2, \dots\}; \\ \mathbb{N}^* &= \mathbb{N} \setminus \{0\} = \{1, 2, \dots\}.\end{aligned}$$

Integers:

$$\mathbb{Z} = \{\dots - 2, -1, 0, 1, 2, \dots\}.$$

Rational numbers:

$$\mathbb{Q} = \left\{ \frac{p}{q}, p \in \mathbb{Z}, q \in \mathbb{N}^* \right\}.$$

Real numbers: they are the classical numbers we deal with. The set of real numbers is usually denoted by \mathbb{R} . It includes the rational numbers but not only. An example is $\sqrt{2}$. The set \mathbb{R} is constructed by "completion" of \mathbb{Q} . For further definitions and derivations, we refer to classical textbooks on real-valued analysis or calculus.

Complex numbers:

$$\mathbb{C} = \{x + iy, (x, y) \in \mathbb{R} \times \mathbb{R}\}.$$

In the next sections, we recapitulate real and complex numbers. In addition, we also recall some useful things to know when dealing with numbers with the help of computers.

2.3 Real numbers and their usage in computers: floating-point number system

2.3.1 Round-off errors and machine precision

In numerical mathematics it is important to understand that real numbers cannot be represented in their full length by a computer. They are cut (i.e., rounded) according to the so-called machine precision. The **floating-point system** is based on the IEEE-754 standard set in the year 1985. It consists of floating numbers, zero, Inf and NaN.

Example 2.1. *To warm up, please have a try yourself using for instance MATLAB, Octave, Python and type in some real numbers. For instance 1/9 and see what happens.*

Floating-point numbers have usually the form

$$x = (-1)^s \cdot (0.a_1a_2 \dots a_t) \cdot \beta^e = (-1)^s \cdot m \cdot \beta^{e-t}, \quad a_1 \neq 0,$$

which exactly means

$$x = (-1)^s (a_1\beta^{-1} + a_2\beta^{-2} + \dots + a_t\beta^{-t})\beta^e = (-1)^s (a_1\beta^{t-1} + a_2\beta^{t-2} + \dots + a_t\beta^0)\beta^{e-t}.$$

Therein, we have:

1. **sign:** s is either 0 or 1
2. **basis:** $\beta \geq 2$ (being an integer). Computers work with $\beta = 2$ and our usual decimal system has $\beta = 10$
3. **mantissa:** $m = a_1\beta^{t-1} + a_2\beta^{t-2} + \dots + a_t\beta^0 \in \mathbb{N}$. It has length t , which is the maximum number of digits a_i , with $0 \leq a_i \leq \beta - 1$ that can be stored
4. **exponent:** $e \in \{e_{\min}, \dots, e_{\max}\} \subset \mathbb{Z}$

An important quantity is the so-called machine precision ϵ , which is defined as

$$\epsilon = \beta^{1-t}$$

and provides the distance between 1 and the closest floating-point number greater than 1. The machine precision depends on the operating system (i.e., your computer) and not on your software.

Example 2.2. *An example of machine (double) precision is $\epsilon = 2.2204e - 16$. This can be obtained for instance in octave with*

eps

or in python with

```
import numpy
print(numpy.finfo(float).eps)
```

2.3.2 Influence of machine precision in numerical mathematics

The previous explanations need considerations once algorithms are implemented in a software and computer simulations are carried out. Why? Often, we deal with large numbers, e.g., Young's modulus of steel is $E = 190PA = 190 \times 10^9 PA$. Here, we have numbers of order 10^9 . For linear equations system due to condition numbers (see below), we also easily deal with very large numbers. Therefore it is important to know what are the biggest numbers the computer can deal with. As **overflow** we denote the situation when the absolute value of a number is greater than the greatest machine number of the computer.

On the other hand, we are concerned with **underflow**, which denotes the situation when a number unequal to zero is rounded to zero. Many numerical algorithms are requested to terminate once the distance between numerical approximation and sought solution is (very) small, e.g., $\|\tilde{x} - x\| \approx 0$. But zero is too small due to machine precision. What to choose else? We define a tolerance TOL such that

$$\|\tilde{x} - x\| < TOL.$$

In order to avoid difficulties due to machine precision we must choose $TOL \gg \epsilon$. For instance:

- we have for machine epsilon $\epsilon \approx 10^{-16}$ in double precision
- Reasonable tolerances are in the range of $TOL = 10^{-8}, \dots, 10^{-12}$.

Finally, due to round-off some mathematical laws may become violated on a computer. Due to the finite representation of numbers in the computer, the mathematical law of associativity

$$a + (b + c) = (a + b) + c$$

may be violated; for instance in the case of overflow or underflow.

Example 2.3. *In Section 12.5.3, we provide numerical solutions to Newton's method for computing a root-finding problem. Therein, the tolerance TOL must be chosen. We show in the following the interaction of ϵ and TOL. For a classical choice $TOL = 1e - 12$ we obtain*

```
Iter x          f(x)
5    1.414214e+00 4.751755e-14
```

This means that we need 5 iterations to converge to a root value $|f(x)| = 4.75e - 14$. For $TOL = 1e - 12$

```
Iter x          f(x)
4    1.414214e+00 6.156754e-07
```

This means that we need 4 iterations (the scheme is more efficient!) to converge to a root value $|f(x)| = 1.53e-07$ (the final root is less accurate!). Here, we already see the trade-off between efficiency and accuracy. Now let us choose $TOL = 1e-16$, a value of the order of the machine precision. Then:

```

Iter x          f(x)
...
6502 1.414214e+00 -4.440892e-16
6503 1.414214e+00 4.440892e-16
6504 1.414214e+00 -4.440892e-16
6505 1.414214e+00 4.440892e-16
...
16368 1.414214e+00 -4.440892e-16
16369 1.414214e+00 4.440892e-16
16370 1.414214e+00 -4.440892e-16
16371 1.414214e+00 4.440892e-16
...

```

The computation did not stop (endless loop) because of round-off errors due to machine precision such that the tolerance cannot be met. This shows indeed that we must stay away with TOL from ϵ and should choose $TOL > \epsilon$.

Example 2.4. Previously, we wrote that ϵ depends on the hardware. Nonetheless, the actual precision in a program can be changed. The standard machine epsilon is double precision. But we can change to **single precision**. The value for single precision can be obtained in python with

```
print(numpy.finfo(numpy.float32).eps)
```

and yields $\epsilon = 1.1920929e-07$. If we re-do the above calculations, then we even should not obtain the result for $TOL = 1e-12$. Indeed

```

Iter x          f(x)
...
12270 1.414214e+00 -1.192093e-07
12271 1.414214e+00 -1.192093e-07
12272 1.414214e+00 -1.192093e-07
12273 1.414214e+00 -1.192093e-07
...

```

we obtain an endless loop. Convergence is again achieved for $TOL > \epsilon$ as for example $TOL = 1e-5$:

```

Iter x          f(x)
4      1.414214e+00 7.152557e-07

```

Whether such low tolerances are sufficiently accurate depends on the problem and on the other hand on the memory (of course single precision needs less memory to store numbers). For instance in machine learning algorithms, often single precision is used.

2.3.3 Stability and condition number

Related to the machine precision is the stability of algorithms. Due to round-off errors, numerical algorithms can yield wrong results. How severe they are can be estimated with the so-called **condition number**.

Consider a function $f : \mathbb{R} \rightarrow \mathbb{R}$ and let us evaluate its values in the standard way:

$$x \mapsto f(x).$$

Due to round-off (or other, i.e., discretization, iterations, etc.) errors the function $f(x)$ is not evaluated at x , but at $\tilde{x} = x + \Delta x$. The question is how this input error influences the result, i.e., the output error. We denote by

$$\Delta y = f(x + \Delta x) - f(x)$$

the absolute error in the function f . If $f \in C^1(\mathbb{R})$ (one times continuously differentiable), we obtain with the mean value theorem:

$$\Delta y = f(x + \Delta x) - f(x) = f'(\xi)\Delta x$$

where $\xi \in [x, x + \Delta x]$. We observe that $f'(\xi)$ plays the role of an (absolute) amplification factor between Δy and Δx . However, in applications, the relative error is by far more important. Consequently, a representation for the relative amplification factor is:

$$\frac{\Delta y}{y} \approx f'(x) \frac{\Delta x}{f(x)} = \left(f'(x) \frac{x}{f(x)} \right) \frac{\Delta x}{x}.$$

Definition 2.1 (Condition number). *The number $\kappa_{abs} = |f'(x)|$ is the absolute condition number of the problem $x \mapsto f(x)$. For $xf(x) \neq 0$ the relative condition is defined as*

$$\kappa_{rel} = \left| \frac{f'(x)x}{f(x)} \right|.$$

The condition numbers describe the amplification of input errors of a given problem in relation to the output errors. If

$$\kappa_{rel} \gg 1$$

then the problem is ill-conditioned. If $\kappa_{rel} \approx 1$ the problem is well conditioned.

Example 2.5. 1. Addition, subtraction: $f(x, y) = x + y$ (here we evaluate in relation to x):

$$\kappa_{f,x} = \left| \frac{x}{x+y} \right| = \left| \frac{1}{1 + \frac{y}{x}} \right|.$$

If $x \approx -y$, then the (relative) condition number of the addition ($x \approx y$ for subtraction) can become arbitrarily big. An example is:

$$x = -1.019, \quad y = 1.021 \quad \Rightarrow \quad x + y = 0.002.$$

Let $\tilde{x} = -1.020$ be a perturbed argument. The relative error x is very small: $|1.019 - 1.020|/|1.019| \leq 0.1\%$. We obtain the perturbed result

$$\tilde{x} + y = 0.001$$

and an error of 100%. The big amplification of the addition of two numbers with the same absolute value is the so-called cancellation.

2. Multiplication: $f(x, y) = x \cdot y$. This operation is always well-conditioned:

$$\kappa_{f,x} = \left| y \frac{x}{xy} \right| = 1.$$

3. Division: $f(x, y) = x/y$. Same conclusion:

$$\kappa_{f,x} = \left| \frac{1}{y} \frac{x}{\frac{x}{y}} \right| = 1, \quad \kappa_{f,y} = \left| \frac{x}{y^2} \frac{y}{\frac{x}{y}} \right| = 1.$$

4. Square root: $f(x) = \sqrt{x}$:

$$\kappa_{f,x} = \left| \frac{1}{2\sqrt{x}} \frac{x}{\sqrt{x}} \right| = \frac{1}{2}.$$

An input error is even reduced.

2.4 Complex numbers: a brief reminder

The set of complex numbers is denoted by \mathbb{C} . A complex number is written as

$$z = x + iy$$

where $x = \Re(z), y = \Im(z) \in \mathbb{R}$ are the real part and the imaginary part of z , respectively. The imaginary unit is defined as $i^2 = -1$. The main reason to define complex numbers was for solving the quadratic equation

$$x^2 + px + q = 0 \tag{2.1}$$

with coefficients $p, q \in \mathbb{R}$. For values only in \mathbb{R} , this quadratic equation has not always solutions. With complex numbers, we can compute as with real numbers for addition, subtraction, multiplication and division. Let us define

$$\begin{aligned} a &= x + iy \\ b &= u + iv. \end{aligned}$$

Then, the addition gives us:

$$a + b = x + iy + (u + iv) = x + u + i(y + v),$$

with $x, y, u, v \in \mathbb{R}$. The subtraction yields:

$$a - b = x + iy - (u + iv) = x - u + i(y - v).$$

The multiplication yields:

$$a * b = (x + iy) * (u + iv) = xu + ixv + iyu + i^2yv = xu - yv + i(xv + yu).$$

For the division we first introduce the conjugate

$$\bar{a} = x - iy$$

and the squared modulus

$$|a|^2 = a\bar{a} = x^2 + y^2.$$

Then, we can write the reciprocal as

$$\frac{1}{b} = \frac{\bar{b}}{b\bar{b}} = \frac{u - iv}{u^2 + v^2}$$

and the division as

$$\frac{a}{b} = \frac{a\bar{b}}{b\bar{b}} = \frac{1}{u^2 + v^2}((ux + vy) + i(uy - vx)).$$

The modulus (or absolute value) has been defined by

$$r = |a| = \sqrt{x^2 + y^2}.$$

The notation r indicates that this quantity can be identified as the radius of a with respect to the origin in the complex plane. To this end, a can be represented with the help of trigonometric functions:

$$a = r(\cos \phi + i \sin \phi),$$

where ϕ , the argument of a , denotes the angle of 0 to a on the positive real axis. With the help of Euler's formula,

$$e^{i\alpha} = \cos \alpha + i \sin \alpha$$

for $\alpha \in \mathbb{R}$, we obtain the following representation for a complex number a :

$$a = re^{i\phi}.$$

This notation is convenient to handle multiplications, powers and divisions. Indeed trigonometric rules entail

$$e^{i\phi}e^{i\theta} = e^{i(\phi+\theta)}, \quad (e^{i\phi})^n = e^{in\phi}, \quad \frac{e^{i\phi}}{e^{i\theta}} = e^{i(\phi-\theta)}.$$

Going back to (2.1), when the discriminant satisfies $\Delta := p^2 - 4q \leq 0$ we have the complex solutions

$$x = \frac{-p \pm i\sqrt{|\Delta|}}{2}.$$

In general a computer stores a complex number as a pair of floating-point numbers.

Chapter 3

Exercises

Concepts in numerical mathematics

Exercise 3.1. *What is the concept for measuring distances? Recapitulate the properties and provide an example of (function) set.*

Exercise 3.2. *Explain the concepts of approximation, stability and efficiency in your own words.*

Exercise 3.3. *We defined different sources of errors. Which errors are specifically related to computer simulations?*

Exercise 3.4. *Let $C[a, b]$ be the space of continuous real-valued functions defined over the closed interval $[a, b]$. Define a metric on $C[a, b]$ and verify the metric properties.*

Sets

Exercise 3.5. *Let A, B, C be three sets. Show that*

1. $(A \cap B) \cup C = (A \cup C) \cap (B \cup C)$;
2. $(A \cup B) \cap C = (A \cap C) \cup (B \cap C)$;
3. $A \cup B = A \cap B \Rightarrow A = B$.

Real numbers in computers

Exercise 3.6. *In this exercise, we give some ideas for testing how numbers are treated in a computer.*

1. *Test the law of associativity on your computer (e.g., Python, C++, Octave/MATLAB,...) taking $a = 1.0e + 308, b = 1.1e + 308, c = -1.001e + 308$.*
2. *Compute $((1 + x) - 1)/x$ for $x = 1e - 15$ (something close to machine precision) and observe if there is a cancellation of significant digits.*
3. *Find out the machine precision of your personal computer.*

Complex numbers

Exercise 3.7. *Let $a = 8 + 6i$ and $b = 1 + 2i$ be given. Compute (by hand!)*

1. $a + b$
2. $a - b$

3. $a \cdot b$
4. a/b
5. \bar{a}
6. $|a|$
7. a^2
8. *the two square roots of a (complex numbers c such that $c^2 = a$)*
9. *the exponential representation of $z = a + 2b$*
10. *z^2 using the Cartesian (real and imaginary parts) and exponential representations*

Part II

Linear algebra and related numerical notions

Chapter 4

Linear systems and matrices

The numerical solution of large linear systems is one of the two main problems of numerical linear algebra, the other problem being the computation of eigenvalues discussed in chapter 5. Linear systems occur in particular after the approximation (the so-called discretization procedure) of linear differential equations (an introduction to this very wide topic is given in part IV). Even nonlinear problems are usually solved through the construction of a sequence of linear problems, typically by the Newton method (see section 12.4). The precision of the approximation is directly related to the number of unknowns, since it is usually associated with spatial or time resolution. It is common to address linear systems with several thousands (or even millions) of unknowns.

The numerical solution of linear systems is classified in two categories. The first one gathers the so-called direct methods such as Gaussian elimination and LU or Cholesky decompositions. Such methods are exact up to round-off errors made during the calculations. The second category contains the iterative methods as we discuss in Section 7.2. They consist in the construction of a sequence of approximate solutions, converging to an exact solution only when the number of iterations tends to infinity. Some methods are mostly appropriate for systems with specific properties (symmetry, sparsity, block structure, very large systems, systems with indirect knowledge of the coefficients...), or for special computing architectures. The Gaussian elimination method is often considered as the most competitive for general "not too large" systems. The LU decomposition is a kind of reformulation of the Gaussian elimination procedure which is mostly of interest when several systems with the same matrix need to be solved. The Cholesky decomposition is a simplification of the LU decomposition for symmetric (or Hermitian) positive definite matrices. Iterative methods must be used for very large systems (because of the cost complexity and memory requirements), or when one is interested in a quick approximate solution. Moreover, iterative methods serve as so-called smoothers in geometric multigrid methods, which may attain optimal cost complexity in solving linear systems; see e.g. [16].

In the following, unless otherwise specified, the letter \mathbb{K} is used to designate either the set \mathbb{R} of real numbers or the set \mathbb{C} of complex numbers.

4.1 Gaussian elimination for solving linear systems

4.1.1 General definitions

Definition 4.1. *A linear system of m equations and n unknowns is a set of equations of form*

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 & (R_1) \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 & (R_2) \\ \vdots & \vdots \\ \vdots & \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n = b_m & (R_m) \end{cases} . \quad (4.1)$$

The numbers $a_{ij} \in \mathbb{K}$ are the coefficients of the system. The numbers $b_i \in \mathbb{K}$ are the coefficients of the right hand side. The numbers $x_j \in \mathbb{K}$ are the unknowns of the system.

2. Then we replace each row R_p , $p \geq 2$, by $R_p - \frac{a'_{p1}}{a'_{11}}R_1$ (elimination). This leads to an equivalent system of form

$$\left\{ \begin{array}{l} a_{11}^{(1)}x_1 + a_{12}^{(1)}x_2 + \cdots + a_{1n}^{(1)}x_n = b_1^{(1)} \quad (R_1^{(1)}) \\ a_{22}^{(1)}x_2 + \cdots + a_{2n}^{(1)}x_n = b_2^{(1)} \quad (R_2^{(1)}) \\ \vdots \\ \vdots \\ a_{m2}^{(1)}x_2 + \cdots + a_{mn}^{(1)}x_n = b_m^{(1)} \quad (R_m^{(1)}) \end{array} \right. . \quad (4.4)$$

Step 2 consists in applying the above procedure to rows $R_2^{(1)}, \dots, R_m^{(1)}$, obviously choosing the pivot in the second column.

We continue in this way until obtaining a full row echelon system. If this row echelon system is compatible, with last row of form (4.2), then we write the r unknowns x_1, \dots, x_r in terms of the other ones, going upwards from the last row to the first one (backward substitutions).

In the numerical implementation of the Gaussian elimination method, small (in absolute value / modulus compared to other coefficients) pivots are undesirable, because they may lead to significant errors in the floating point arithmetic (an illustration is given as exercise). Therefore, pivoting is often performed so as to select the largest available pivot. To improve accuracy, column swapping is generally also performed. This procedure is called total pivoting.

4.2 Vectors and matrices

4.2.1 Basic definitions

Matrices

Definition 4.7. A matrix with coefficients in \mathbb{K} of type (m, n) is a table with m rows and n columns represented as

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

The numbers $a_{ij} \in \mathbb{K}$ are the coefficients of the matrix. Sometimes we denote $A = (a_{ij})$.

Definition 4.8. We denote by $\mathcal{M}_{mn}(\mathbb{K})$, or $\mathbb{K}^{m \times n}$, the set of matrices with coefficients in \mathbb{K} of type (m, n) .

Definition 4.9. Two matrices $A = (a_{ij}) \in \mathcal{M}_{mn}(\mathbb{K})$ and $B = (b_{ij}) \in \mathcal{M}_{mn}(\mathbb{K})$ are equal if they have the same coefficients, i.e., $a_{ij} = b_{ij} \forall (i, j) \in \{1, \dots, m\} \times \{1, \dots, n\}$.

Definition 4.10. A zero matrix, denoted by 0 , is a matrix with all coefficients equal to 0. Hence, if $A = (a_{ij}) \in \mathcal{M}_{mn}(\mathbb{K})$, then

$$A = 0 \iff a_{ij} = 0 \forall (i, j) \in \{1, \dots, m\} \times \{1, \dots, n\}.$$

Special matrices

- a) Row matrix, also called row vector

$$A = (a_{11} \quad \cdots \quad a_{1n}) \in \mathcal{M}_{1n}(\mathbb{K})$$

b) Column matrix, also called column vector

$$A = \begin{pmatrix} a_{11} \\ \vdots \\ a_{m1} \end{pmatrix} \in \mathcal{M}_{m1}(\mathbb{K})$$

c) Square matrix

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix} \in \mathcal{M}_{nn}(\mathbb{K}) =: \mathcal{M}_n(\mathbb{K})$$

Some square matrices have themselves special features.

i) Lower triangular matrix

$$A = \begin{pmatrix} a_{11} & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & 0 \\ a_{n1} & \cdots & \cdots & a_{nn} \end{pmatrix}$$

ii) Upper triangular matrix

$$A = \begin{pmatrix} a_{11} & \cdots & \cdots & a_{1n} \\ 0 & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & a_{nn} \end{pmatrix}$$

iii) Diagonal matrix

$$A = \begin{pmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a_{nn} \end{pmatrix} = \text{diag}(a_{11}, \dots, a_{nn})$$

iv) Identity matrix

$$I_n = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{pmatrix} = \text{diag}(1, \dots, 1) \in \mathcal{M}_n(\mathbb{K})$$

In contrast to vectors and matrices, the elements of \mathbb{K} are called scalars.

Submatrix

Definition 4.11. A submatrix of a matrix A is a matrix obtained by removing some rows and some columns from A .

4.2.2 Operations on matrices

Addition, multiplication by a scalar

Definition 4.12. Let $A = (a_{ij})$ and $B = (b_{ij})$ be two matrices in $\mathcal{M}_{mn}(\mathbb{K})$. The matrix $A + B$ is defined by

$$A + B = (a_{ij} + b_{ij}) \in \mathcal{M}_{mn}(\mathbb{K}).$$

Definition 4.13. Let $A = (a_{ij}) \in \mathcal{M}_{mn}(\mathbb{K})$ and $\alpha \in \mathbb{K}$. We define the matrix αA by

$$\alpha A = (\alpha a_{ij}) \in \mathcal{M}_{mn}(\mathbb{K}).$$

Proposition 4.2. Let $A, B, C \in \mathcal{M}_{mn}(\mathbb{K})$ and $\alpha \in \mathbb{K}$. We have

$$\begin{aligned} A + B &= B + A \\ (A + B) + C &= A + (B + C) \\ \alpha(A + B) &= \alpha A + \alpha B. \end{aligned}$$

Multiplication of matrices

Definition 4.14. Let $A = (a_{ij}) \in \mathcal{M}_{mn}(\mathbb{K})$ and $B = (b_{ij}) \in \mathcal{M}_{np}(\mathbb{K})$. We define the matrix $C = AB \in \mathcal{M}_{mp}(\mathbb{K})$ by $C = (c_{ij})$ with

$$c_{ij} = \sum_{k=1}^n a_{ik}b_{kj}.$$

Remark 4.2. 1. Be careful with the compatibility of dimensions: the number of columns of the left matrix must be equal to the the number of rows of the second matrix.

2. The product AB may be defined whereas the product BA is not. Even if AB and BA are both defined, in general $AB \neq BA$. If $AB = BA$, then we say that A and B commute.

3. It may happen that $AB = 0$ while $A \neq 0$ and $B \neq 0$!

Proposition 4.3. Let A, B, C be three matrices. We have, under the condition of compatible dimensions, the associativity

$$(AB)C = A(BC).$$

Proposition 4.4. Let $A \in \mathcal{M}_{mn}(\mathbb{K})$. We have $AI_n = I_m A = A$.

Distributivity

Proposition 4.5. Let A, B, C be three matrices and $\alpha \in \mathbb{K}$. We have, under the condition of compatible dimensions,

$$\begin{aligned} A(B + C) &= AB + AC \\ (A + B)C &= AC + BC \\ \alpha(AB) &= (\alpha A)B = A(\alpha B). \end{aligned}$$

Transposition

Definition 4.15. Let $A = (a_{ij}) \in \mathcal{M}_{mn}(\mathbb{K})$. The transpose of A , denoted by A^T , is the matrix

$$A^T = (a_{ji}) \in \mathcal{M}_{nm}(\mathbb{K}).$$

Example 4.1.

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}^T = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$$

Definition 4.16. A matrix $A \in \mathcal{M}_n(\mathbb{R})$ such that $A = A^T$ is called a symmetric matrix. A matrix $A \in \mathcal{M}_n(\mathbb{C})$ such that $A = \bar{A}^T$, where \bar{A} is the complex conjugate matrix of A (obtained by taking the conjugate of each coefficient), is called a Hermitian matrix.

Proposition 4.6. 1. Let $A, B \in \mathcal{M}_{mn}(\mathbb{K})$. We have $(A + B)^T = A^T + B^T$.

2. Let $A \in \mathcal{M}_{mn}(\mathbb{K})$ and $\alpha \in \mathbb{K}$. We have $(\alpha A)^T = \alpha A^T$.

	$k = 0$	$k = 1$	$k = 2$	$k = 3$
$p = 0$	1			
$p = 1$	1	1		
$p = 2$	1	2	1	
$p = 3$	1	3	3	1

Table 4.1: Binomial coefficients C_p^k for $p \leq 3$ (Pascal's triangle)

3. Let $A \in \mathcal{M}_{mn}(\mathbb{K})$ and $B \in \mathcal{M}_{np}(\mathbb{K})$. We have $(AB)^T = B^T A^T$.

4. Let $A \in \mathcal{M}_{mn}(\mathbb{K})$. We have $(A^T)^T = A$.

Definition 4.17. Let $A \in \mathcal{M}_n(\mathbb{K})$ be a Hermitian (symmetric if $\mathbb{K} = \mathbb{R}$) matrix. We say that A is positive semi-definite if $\bar{x}^T A x \geq 0$ for all $x \in \mathbb{K}^n$. We denote $A \geq 0$. We say that A is positive definite if $\bar{x}^T A x > 0$ for all $x \in \mathbb{K}^n \setminus \{0\}$. We denote $A > 0$.

Trace of a square matrix

Definition 4.18. Let $A = (a_{ij}) \in \mathcal{M}_n(\mathbb{K})$. The trace of A , denoted by $\text{tr } A$, is the sum of the diagonal coefficients of A , i.e.

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

Proposition 4.7. 1. Let $A, B \in \mathcal{M}_n(\mathbb{K})$. We have $\text{tr}(A + B) = \text{tr } A + \text{tr } B$.

2. Let $A \in \mathcal{M}_n(\mathbb{K})$ and $\alpha \in \mathbb{K}$. We have $\text{tr}(\alpha A) = \alpha \text{tr } A$.

3. Let $A \in \mathcal{M}_{mn}(\mathbb{K})$, $B \in \mathcal{M}_{nm}(\mathbb{K})$. We have $\text{tr}(AB) = \text{tr}(BA)$.

Power of a square matrix

Definition 4.19. Let $A \in \mathcal{M}_n(\mathbb{K})$ and $p \in \mathbb{N}$. We define

$$A^p = \underbrace{A \times A \times \cdots \times A}_{p \text{ factors}},$$

with the convention $A^0 = I_n$.

Remark 4.3. Be careful not to confuse the matrix power with the power of the coefficients. However, if A is diagonal, then $(\text{diag}(\lambda_1, \dots, \lambda_n))^p = \text{diag}(\lambda_1^p, \dots, \lambda_n^p)$.

For two commuting matrices, we have the following binomial formula.

Proposition 4.8. Let $A, B \in \mathcal{M}_n(\mathbb{K})$ be two matrices such that $AB = BA$ and p be a natural number. The binomial formula

$$(A + B)^p = \sum_{k=0}^p C_p^k A^k B^{p-k}$$

holds, where $C_p^k = \binom{p}{k} = \frac{p!}{(p-k)!k!}$ are the binomial coefficients.

The binomial formula is formally the same as for scalars. Remind that, due to the property $C_{p+1}^{k+1} = C_p^{k+1} + C_p^k$, the binomial coefficients can be obtained from Pascal's triangle, see Table 4.1.

It follows that A is invertible with inverse

$$A^{-1} = \begin{pmatrix} 2 & -1 \\ -5 & 3 \end{pmatrix}.$$

The equivalent matrix formulation is

$$\begin{pmatrix} 3 & 1 \\ 5 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \quad \text{often denoted} \quad \left(\begin{array}{cc|cc} 3 & 1 & 1 & 0 \\ 5 & 2 & 0 & 1 \end{array} \right).$$

After Gaussian elimination this is found equivalent to

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 & -1 \\ -5 & 3 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \quad \text{denoted} \quad \left(\begin{array}{cc|cc} 1 & 0 & 2 & -1 \\ 0 & 1 & -5 & 3 \end{array} \right),$$

from which A^{-1} is readily obtained.

Some algebraic rules

Proposition 4.10. 1. I_n is invertible of inverse I_n .

2. If A is invertible then A^{-1} is invertible and $(A^{-1})^{-1} = A$.

3. If A, B are invertible then AB is invertible and $(AB)^{-1} = B^{-1}A^{-1}$.

4. If A is invertible then A^T is invertible and $(A^T)^{-1} = (A^{-1})^T =: A^{-T}$.

5. If $A = \text{diag}(\lambda_1, \dots, \lambda_n)$, then A is invertible if and only if all coefficients λ_i do not vanish. In this case $A^{-1} = \text{diag}(\lambda_1^{-1}, \dots, \lambda_n^{-1})$.

4.3 Linear subspaces of \mathbb{K}^n

4.3.1 Definitions and first properties

Definition 4.21. A set $F \subset \mathbb{K}^n$ is a linear subspace of \mathbb{K}^n if

1. $0 \in F$,
2. $\forall v, w \in F, v + w \in F$,
3. $\forall v \in F, \forall \lambda \in \mathbb{K}, \lambda v \in F$.

Such sets appear naturally when solving underdetermined homogeneous linear systems, as stated below.

Proposition 4.11. The set of solutions of a linear homogeneous system of the unknown $v = (v_1, \dots, v_n) \in \mathbb{K}^n$ (here v_1, \dots, v_n denote the components of the vector v) is a linear subspace of \mathbb{K}^n .

In the following, elements of \mathbb{K}^n are called vectors, whatever they are considered as row vectors or column vectors.

Hereafter, be careful with the number of vectors indexed by m and the corresponding number of components of each vector that is indicated by n .

Definition 4.22. Let w, v_1, \dots, v_m be vectors from \mathbb{K}^n . Each of these vectors has n components. We say that w is a linear combination of v_1, \dots, v_m if there exist scalars $\lambda_1, \dots, \lambda_m$ such that

$$w = \lambda_1 v_1 + \dots + \lambda_m v_m.$$

Proposition 4.12. Let v_1, \dots, v_m be vectors. The set of linear combinations of v_1, \dots, v_m is a linear subspace of \mathbb{K}^n . It is called the linear subspace spanned by v_1, \dots, v_m , it is denoted by $\text{Span}(v_1, \dots, v_m)$.

4.3.2 Bases

Definition 4.23. A family (v_1, \dots, v_m) of vectors is said to be linearly independent if, for all families $(\lambda_1, \dots, \lambda_m)$ of scalars,

$$\lambda_1 v_1 + \dots + \lambda_m v_m = 0 \implies \lambda_1 = \dots = \lambda_m = 0.$$

Otherwise it is said to be linearly dependent.

Definition 4.24. Let F be a linear subspace of \mathbb{K}^n and (v_1, \dots, v_m) be a family of vectors. If $\text{Span}(v_1, \dots, v_m) = F$ we say that (v_1, \dots, v_m) spans F , or that it is a spanning set of F .

Example. The vectors (e_1, \dots, e_n) defined by $e_i = (e_i^1, \dots, e_i^n)$, $e_i^j = 1$ if $i = j$, $e_i^j = 0$ if $i \neq j$ obviously span \mathbb{K}^n . For instance for $n = 3$:

$$e_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \quad e_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \quad e_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}.$$

Proposition 4.13. Let F be a linear subspace of \mathbb{K}^n . If (v_1, \dots, v_m) spans F then all family of $m + 1$ vectors of F is linearly dependent.

Corollary 4.14. A family of $n + 1$ vectors of \mathbb{K}^n is linearly dependent.

Definition 4.25. Let F be a linear space of \mathbb{K}^n . We say that a family (v_1, \dots, v_m) of vectors is a basis of F if it is linearly independent and a spanning set of F .

Theorem 4.15. All linear subspace of \mathbb{K}^n different from $\{0\}$ admit a basis.

Example. The family (e_1, \dots, e_n) defined above is a basis of \mathbb{K}^n , called *canonical basis*.

Proposition 4.16. Let F be a linear subspace of \mathbb{K}^n and (v_1, \dots, v_m) be a basis of F . Any vector x of F can be decomposed in a unique way as $x = \lambda_1 v_1 + \dots + \lambda_m v_m$ with $\lambda_1, \dots, \lambda_m \in \mathbb{K}$. The scalars $\lambda_1, \dots, \lambda_m$ are called the coordinates of x in the basis (v_1, \dots, v_m) .

4.3.3 Dimension

Here is the "theorem of the dimension".

Theorem 4.17. Let F be a linear subspace of \mathbb{K}^n different from $\{0\}$. All the bases of F contain the same number of vectors.

Definition 4.26. Let F be a linear subspace of \mathbb{K}^n different from $\{0\}$. We call dimension of F , denoted by $\dim F$, the number of vectors of any basis of F . By convention, $\dim\{0\} = 0$.

Proposition 4.18. We have $\dim \mathbb{K}^n = n$ and $\dim F \leq n$ for all linear subspace F of \mathbb{K}^n .

Proposition 4.19. Let F be a linear subspace of \mathbb{K}^n such that $\dim F = m$.

1. If p vectors of F are linearly independent then $p \leq m$.
2. If m vectors of F are linearly independent then they form a basis of F .
3. If p vectors of F span F then $p \geq m$.
4. If m vectors of F span F then they form a basis of F .

Corollary 4.20. Let F be a linear subspace of \mathbb{K}^n such that $\dim F = n$. Then $F = \mathbb{K}^n$.

4.4 Determinant

The mathematical construction of the notion of determinant is a lengthy procedure. Here we only recall the main properties of the determinant of a square matrix.

4.4.1 Expansion with respect to a row or a column

The determinant of a square matrix $A \in \mathcal{M}_n(\mathbb{K})$ is a scalar $\det A \in \mathbb{K}$. We adopt a practical definition by induction on n .

We begin with the trivial case $n = 1$.

Definition 4.27. For $A = (a) \in \mathcal{M}_1(\mathbb{K})$ we set $\det A = a$.

We now assume that the determinant is defined for all matrix of type $(n - 1, n - 1)$.

Definition 4.28. Let $A \in \mathcal{M}_n(\mathbb{K})$ et $i, j \in \{1, \dots, n\}$. We call

- minor (i, j) of A the scalar

$$\Delta_{ij}(A) = \det \tilde{A}$$

where the \tilde{A} is the submatrix of A obtained by removing the row i and the column j ;

- cofactor (i, j) of A the scalar

$$\text{Cof}_{ij}(A) = (-1)^{i+j} \Delta_{ij}(A).$$

Definition 4.29. Let $A = (a_{ij}) \in \mathcal{M}_n(\mathbb{K})$. Given $i_0 \in \{1, \dots, n\}$ we call expansion of the determinant of A with respect to the row i_0 the scalar

$$R_{i_0}(A) = \sum_{j=1}^n a_{i_0,j} \text{Cof}_{i_0,j}(A).$$

Likewise, given $j_0 \in \{1, \dots, n\}$ we call expansion of the determinant of A with respect to the column j_0 the scalar

$$C_{j_0}(A) = \sum_{i=1}^n a_{i,j_0} \text{Cof}_{i,j_0}(A).$$

Definition 4.30. Let $A \in \mathcal{M}_n(\mathbb{K})$. We set

$$\det A = R_1(A).$$

The definitions 4.27 to 4.30 define unambiguously $\det A$ for any matrix $A \in \mathcal{M}_n(\mathbb{K})$. The following theorem permits to sometimes simplify its calculation.

Theorem 4.21. Let $A = (a_{ij}) \in \mathcal{M}_n(\mathbb{K})$. For any $i_0, j_0 \in \{1, \dots, n\}$ we have

$$\det A = R_{i_0}(A) = C_{j_0}(A).$$

Notation. If $A = (a_{ij}) \in \mathcal{M}_n(\mathbb{K})$, we usually write

$$\det A = \begin{vmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \cdots & a_{nn} \end{vmatrix}.$$

Example 4.3.

$$A = \begin{vmatrix} a & b \\ c & d \end{vmatrix}.$$

We have the minors

$$\Delta_{11}(A) = d, \quad \Delta_{12}(A) = c, \quad \Delta_{21}(A) = b, \quad \Delta_{22}(A) = a,$$

the cofactors

$$\text{Cof}_{11}(A) = d, \quad \text{Cof}_{12}(A) = -c, \quad \text{Cof}_{21}(A) = -b, \quad \text{Cof}_{22}(A) = -a,$$

the determinant

$$\det(A) = R_1(A) = R_2(A) = C_1(A) = C_2(A) = ad - bc.$$

4.4.2 Properties

Proposition 4.22. *Let $A \in \mathcal{M}_n(\mathbb{K})$. We denote by A_1, \dots, A_n the columns of A and we write $A = (A_1, \dots, A_n)$.*

1. *Swapping two columns of A multiplies its determinant by -1 .*
2. *If C is an arbitrary column vector then*

$$\begin{aligned} \det(A_1, \dots, A_{i-1}, A_i + C, A_{i+1}, \dots, A_n) \\ = \det(A_1, \dots, A_{i-1}, A_i, A_{i+1}, \dots, A_n) + \det(A_1, \dots, A_{i-1}, C, A_{i+1}, \dots, A_n). \end{aligned}$$

3. *If $\lambda \in \mathbb{K}$ then*

$$\det(A_1, \dots, A_{i-1}, \lambda A_i, A_{i+1}, \dots, A_n) = \lambda \det(A_1, \dots, A_{i-1}, A_i, A_{i+1}, \dots, A_n).$$

Remark 4.5. 1. *If A and B are two matrices then in general $\det(A + B) \neq \det A + \det B$.*

2. *If $\lambda \in \mathbb{K}$ then $\det(\lambda A) = \lambda^n \det A$.*

Corollary 4.23. *Let $A \in \mathcal{M}_n(\mathbb{K})$.*

1. *If a column of A is zero then $\det A = 0$.*
2. *If two columns of A are equal then $\det A = 0$.*
3. *If we add to a column of A a linear combination of the other columns of A then the determinant is unchanged.*
4. *If the columns of A are linearly dependent then $\det A = 0$.*

Proposition 4.24. *Let $A \in \mathcal{M}_n(\mathbb{K})$. We have*

$$\det A = \det A^T.$$

Consequently, the assertions of proposition 4.22 and corollary 4.23 apply also row-wise.

Calculating a determinant via recursive expansions is usually tedious. Only special matrices permit an easy calculation.

Proposition 4.25. *The determinant of a (lower or upper) triangular matrix is equal to the product of its diagonal coefficients.*

In particular, this applies to diagonal matrices. Notably, $\det I_n = 1$.

4.4.3 Determinant of a product, characterization of invertible matrices

The following algebraic property is remarkable.

Theorem 4.26. *Let $A, B \in \mathcal{M}_n(\mathbb{K})$. We have*

$$\det(AB) = \det A \det B.$$

It leads to a very useful characterization of invertible matrices.

Corollary 4.27. *Let $A \in \mathcal{M}_n(\mathbb{K})$. Then A is invertible if and only if $\det A \neq 0$. In this case we have*

$$\det(A^{-1}) = \frac{1}{\det A}.$$

4.5 LU decomposition (complement)

4.5.1 Matrix representation of Gaussian elimination

In this section, we formulate the Gaussian elimination method as a series of matrix operations. These concepts are very useful in computer implementations.

We recall the Gauss elimination principle:

$$\begin{pmatrix} * & * & * & \cdots & * \\ * & * & * & & * \\ * & * & * & & * \\ \vdots & & & \ddots & \vdots \\ * & * & * & \cdots & * \end{pmatrix} \rightarrow \begin{pmatrix} * & * & * & \cdots & * \\ 0 & * & * & & * \\ 0 & * & * & & * \\ \vdots & & & \ddots & \vdots \\ 0 & * & * & \cdots & * \end{pmatrix} \rightarrow \dots \rightarrow \begin{pmatrix} * & * & * & \cdots & * \\ 0 & * & * & & * \\ 0 & 0 & * & & * \\ \vdots & & & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & * \end{pmatrix}$$

Using the upper triangular matrix $U \in \mathcal{M}_n(\mathbb{K})$ ($U = (u_{ij})_{i,j=1}^n$ with $u_{ij} = 0$ for $i > j$), we can solve the reduced system

$$Ux = \tilde{b},$$

using backward substitution. Here, \tilde{b} is the modification of the right hand side b through the elimination process. Let us consider the backward substitution in more detail:

Listing 4.1: Backward substitution

Let $U \in \mathcal{M}_n(\mathbb{K})$ be a regular upper triangular matrix, i.e., $u_{ii} \neq 0$. The solution $x \in \mathbb{K}^n$ of $Ux = b$ is given by:

```
| Set |  $x_n = u_{nn}^{-1}b_n$ .
| For |  $i$  | from |  $n - 1$  | until |  $1$ 
 $x_i = u_{ii}^{-1} \left( b_i - \sum_{j=i+1}^n u_{ij}x_j \right)$ 
```

It holds

Proposition 4.28 (Backward substitution). *Let $U \in \mathcal{M}_n(\mathbb{K})$ be an upper triangular matrix with $u_{ii} \neq 0$. Then, the matrix U is regular and the backward substitution requires*

$$N_U(n) = \frac{n^2}{2} + O(n)$$

elementary operations on numbers.

Proof. It holds $\det(U) = \prod u_{ii} \neq 0$. Also the matrix U is regular.

Each step of the backward substitution consists of additions, multiplications and divisions of the diagonal elements. For $u_{ii} \neq 0$ each step is well-defined.

For the computation of x_i , we need $n - i$ multiplications and additions. Furthermore, we have one division per step. This yields

$$n + \sum_{i=1}^{n-1} (n - i) = n + (n - 1)n - \frac{(n - 1)n}{2} = \frac{n^2}{2} + \frac{n}{2}.$$

□

The transformation of A to a triangular structure is obtained through row-based elimination:

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} \\ a_{31} & a_{32} & a_{33} & \ddots & a_{3n} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn} \end{pmatrix} \rightarrow \begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & \cdots & a_{2n}^{(1)} \\ 0 & a_{32}^{(1)} & a_{33}^{(3)} & \ddots & a_{3n}^{(1)} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & a_{n2}^{(1)} & a_{n3}^{(1)} & \cdots & a_{nn}^{(1)} \end{pmatrix} \rightarrow$$

$$\rightarrow \begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ 0 & a_{22}^{(1)} & a_{23}^{(1)} & \cdots & a_{2n}^{(1)} \\ 0 & 0 & a_{33}^{(2)} & \ddots & a_{3n}^{(2)} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & a_{n3}^{(2)} & \cdots & a_{nn}^{(2)} \end{pmatrix} \rightarrow \cdots \rightarrow A^{(n-1)} =: U.$$

Starting from $A^{(0)} := A$, we obtain successively the matrices $A^{(i)}$ with $A^{(n-1)} =: U$. Therein, in step i , we eliminate the i th column of $A^{(i-1)}$ below the diagonal. This is obtained through subtraction of the $g_k^{(i)}$ -fold of the i th row from the k th row. Here, it holds for $k = i + 1, \dots, n$

$$g_k^{(i)} := \frac{a_{ki}^{(i-1)}}{a_{ii}^{(i-1)}},$$

and the new row is computed as

$$a_{kj}^{(i)} = a_{kj}^{(i-1)} - g_k^{(i)} a_{ij}^{(i-1)}, \quad k = i + 1, \dots, n, \quad j = i, \dots, n.$$

These instructions are well-defined if $a_{ii}^{(i-1)} \neq 0$. This follows not necessarily from the regularity of A and all $A^{(i)}$. First of all, we assume that $a_{ii}^{(i-1)} \neq 0$ holds true, and discuss later the general case. In the i th elimination step, the first $i - 1$ rows and columns remain unchanged. The i th elimination step is written in compact form with the help of a matrix-matrix multiplication

$$A^{(i)} = F^{(i)} A^{(i-1)}.$$

Here, the elimination matrix is

$$F^{(i)} := \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & -g_{i+1}^{(i)} & \ddots & \\ & & \vdots & \ddots & \\ & & -g_n^{(i)} & & 1 \end{pmatrix}, \quad g_k^{(i)} := \frac{a_{ki}^{(i-1)}}{a_{ii}^{(i-1)}}.$$

Therein, all non-specified entries are zero. Multiple applications yield

$$U = A^{(n-1)} = F^{(n-1)} A^{(n-2)} = F^{(n-1)} F^{(n-2)} A^{(n-3)} = \underbrace{F^{(n-1)} \cdots F^{(1)}}_{=: F} A^{(0)} = FA. \quad (4.6)$$

Matrices $F^{(i)}$ with this structure are called Frobenius matrices. It holds

Proposition 4.30 (LU decomposition). *Let $A \in \mathcal{M}_n(\mathbb{K})$ be a regular matrix. We assume that all diagonal entries $a_{ii}^{(i-1)}$ arising through the elimination process are nonzero. Then, there exists a unique decomposition $A = LU$ with an upper triangular matrix $U \in \mathcal{M}_n(\mathbb{K})$ and a lower triangular matrix $L \in \mathcal{M}_n(\mathbb{K})$ with diagonal entries 1. The computational cost to construct the LU decomposition requires*

$$\frac{1}{3}n^3 + O(n^2)$$

elementary operations.

Proof. (i) *Uniqueness.* Let us assume there exist two LU decompositions

$$A = L_1U_1 = L_2U_2 \quad \leftrightarrow \quad L_2^{-1}L_1 = U_2U_1^{-1}.$$

The product of triangular matrices is again a triangular matrix. The product $L_2^{-1}L_1$ has only values of 1 on the diagonal. It follows

$$L_2^{-1}L_1 = U_2U_1^{-1} = I,$$

and therefore $L_1 = L_2$ and $U_1 = U_2$.

(ii) *Feasibility.* Each step of the elimination is well-defined as long as we do not need to divide by $a_{ii}^{(i-1)} = 0$. The matrix F is per construction regular and therefore the matrix L exists.

(iii) *Computational cost.* In the i th elimination step

$$A^{(i)} = F^{(i)}A^{(i-1)}$$

we have first of all $n - i$ arithmetic operations for the computation of $g_j^{(i)}$ for $j = i + 1, \dots, n$. The matrix-matrix multiplication applies to all elements a_{kl} with $k > i$ and $l > i$. It holds

$$a_{kl}^{(i)} = a_{kl}^{(i-1)} - g_k^{(i)}a_{il}^{(i-1)}, \quad k, l = i + 1, \dots, n.$$

Here, we need $(n - i)^2$ arithmetic operations. In total, the computational cost in the $n - 1$ steps sums-up to

$$N_{LU}(n) = \sum_{i=1}^{n-1} \{n - i + (n - i)^2\} = \sum_{i=1}^{n-1} \{i + i^2\},$$

and with the well-known summation rules, it follows

$$N_{LU}(n) = \frac{n^3}{3} + \frac{n}{3}.$$

□

The LU decomposition can be used to solve linear equation system:

Listing 4.2: Solving $Ax = b$ with the LU decomposition

Let $A \in \mathcal{M}_n(\mathbb{R})$ be a regular matrix which admits an LU decomposition.

Construct the LU decomposition $A = LU$

Forward substitution: $Ly = b$

Backward substitution: $Ux = y$

The forward substitution is similar to Algorithm 4.1 and requires $O(n^2)$ operations according to Proposition 4.28. This is an interesting observation. The solution process is much cheaper than constructing the LU decomposition itself. For this reason, if possible, the LU decomposition is only constructed when needed, but in many applications it can be used several times (to solve several systems with the same matrix). This should be considered in implementations and reduces significantly the computational cost for large problems.

Remark 4.6 (Practical aspects). *The matrix L is a lower triangular matrix with values 1 on the diagonal. These values are not required to be stored. Similarly, the zero-elements of $A^{(i)}$ below the diagonal are not required to be stored either. Consequently, both matrices L and U can be stored in a common square matrix. In step i , it holds*

$$\tilde{A}^{(i)} = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & \cdots & \cdots & a_{1n} \\ \mathbf{l_{21}} & a_{22}^{(1)} & a_{23}^{(1)} & & & & \vdots \\ \mathbf{l_{31}} & \mathbf{l_{32}} & a_{33}^{(2)} & & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & & & \vdots \\ \vdots & & & \mathbf{l_{i+1,i}} & a_{i+1,i+1}^{(i)} & \cdots & a_{i+1,n}^{(i)} \\ \vdots & & & \vdots & & \ddots & \vdots \\ \mathbf{l_{n1}} & \mathbf{l_{n2}} & \cdots & \mathbf{l_{n,i}} & a_{n,i+1}^{(i)} & \cdots & a_{nn}^{(i)} \end{pmatrix}.$$

The bold entries are the entries of L . The values above the line do not change anymore during the algorithm and form already the respective entries of L and U .

The algorithm reads:

Listing 4.3: LU decomposition of a matrix A without additional memory

Let $A = (a_{ij})_{ij} \in \mathcal{M}_n(\mathbb{K})$ be a regular matrix for which the LU decomposition exists

```
| For | i | from | 1 | until | n
  | For | k | from | i + 1 | until | n
    aki = aki/aii
  | For | j | from | i + 1 | until | n
    akj = akj - aki · aij
```

The element $a_{ii}^{(i-1)}$ is the so-called *Pivot element*. Until now, this element was assumed to be non-zero. For regular matrices this is however in general not the case. Let us consider the following example:

$$A := \begin{pmatrix} 1 & 4 & 2 \\ 2 & 8 & 1 \\ 1 & 2 & 1 \end{pmatrix}.$$

In the first step of the construction, we have $a_{11}^{(0)} = 1$ and it holds

$$A^{(1)} = F^{(1)}A = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 4 & 2 \\ 2 & 8 & 1 \\ 1 & 2 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 4 & 2 \\ 0 & 0 & -3 \\ 0 & -2 & -1 \end{pmatrix}$$

At this step, the algorithm aborts, because we have $a_{22}^{(1)} = 0$. We could have continued the algorithm with the choice $a_{32}^{(1)} = -2$ as a new Pivot element. This approach can be applied in a systematic manner and is called pivoting. In the i th step of the scheme, we search first for an appropriate Pivot element a_{ki} in the i th column with $k \geq i$. The k th and i th rows are switched and the LU decomposition can

We illustrate this process with the help of a simple example:

$$\begin{aligned}
 U &= F^{(3)} P^{(3)} F^{(2)} P^{(2)} F^{(1)} P^{(1)} A \\
 &= F^{(3)} P^{(3)} F^{(2)} \underbrace{P^{(3)} P^{(3)} P^{(2)} F^{(1)}}_{=I} \underbrace{P^{(2)} P^{(3)} P^{(3)} P^{(2)}}_{=I} P^{(1)} A \\
 &= \underbrace{F^{(3)}}_{=\tilde{F}^{(3)}} \underbrace{P^{(3)} F^{(2)} P^{(3)}}_{=\tilde{F}^{(2)}} \underbrace{P^{(3)} P^{(2)} F^{(1)} P^{(2)} P^{(3)}}_{=\tilde{F}^{(1)}} \underbrace{P^{(3)} P^{(2)} P^{(1)}}_{=P} A
 \end{aligned}$$

With $\tilde{L}^{(i)} = [\tilde{F}^{(i)}]^{-1}$ it holds

$$\underbrace{\tilde{L}^{(1)} \dots \tilde{L}^{(n-1)}}_{=: \tilde{L}} U = PA.$$

Since $\tilde{L}^{(i)}$ are again Frobenius matrices, it holds further with Proposition 4.29:

$$\begin{aligned}
 \tilde{L} &= \tilde{L}^{(1)} \dots \tilde{L}^{(n-1)} \\
 &= \sum_{i=1}^{n-1} \tilde{L}^{(i)} - (n-2)I \\
 &= P^{(n-1)} \left(L^{(n-1)} + P^{(n-1)} \left(L^{(n-2)} + \dots \right. \right. \\
 &\quad \left. \left. + \dots + P^{(2)} F^{(1)} P^{(2)} \right) \dots P^{(n-2)} \right) P^{(n-1)} - (n-2)I.
 \end{aligned}$$

For the construction of the LU decomposition not only do we need to permute the $A^{(i)}$, but we also need to change the $L^{(i)}$ that have been computed so far. We summarize:

Proposition 4.32 (LU decomposition with pivoting). *Let $A \in \mathcal{M}_n(\mathbb{K})$ be a regular matrix. Then, there exists an LU decomposition*

$$PA = LU,$$

where P is a product of Pivot matrices, L is a lower triangular matrix with values 1 on the diagonal and U an upper triangular matrix. The LU decomposition without pivoting $P = I$ is unique in case it exists.

Proof. The proof is left for the reader. □

We recall that pivoting yields a well-defined construction of the LU decomposition. But there is a second advantage. With the help of pivoting, we can increase the numerical stability of Gaussian elimination. Through the choice of a Pivot element a_{ki} with maximal relative value (with respect to the row), we can reduce the risk of cancellation.

Example 4.4 (LU decomposition without pivoting). *Let:*

$$A = \begin{pmatrix} 2.3 & 1.8 & 1.0 \\ 1.4 & 1.1 & -0.7 \\ 0.8 & 4.3 & 2.1 \end{pmatrix}, \quad b = \begin{pmatrix} 1.2 \\ -2.1 \\ 0.6 \end{pmatrix},$$

and the solution of $Ax = b$ is given by (five-digit accuracy):

$$x \approx \begin{pmatrix} +0.34995 \\ -0.98023 \\ +2.1595 \end{pmatrix}.$$

For the matrix A it holds $\text{cond}_\infty(A) = \|A\|_\infty \|A^{-1}\|_\infty \approx 7.2 \cdot 1.2 \approx 8.7$. This problem is well-conditioned. The amplification factor is 8.7 and we expect an amplification of the error of about one

digit. First, we construct the LU decomposition (three-digit accuracy). The entries of L are in bold letters:

$$F^{(1)} = \begin{pmatrix} 1 & 0 & 0 \\ -\frac{1.4}{2.3} & 1 & 0 \\ -\frac{0.8}{2.3} & 0 & 1 \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 \\ -0.609 & 1 & 0 \\ -0.348 & 0 & 1 \end{pmatrix},$$

$$[L^{(1)}, A^{(1)}] \approx \begin{pmatrix} 2.3 & 1.8 & 1.0 \\ 0.609 & 0.0038 & -1.31 \\ 0.348 & 3.67 & 1.75 \end{pmatrix}.$$

In the second step, we obtain

$$F^{(2)} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{3.67}{0.0038} & 1 \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -966 & 1 \end{pmatrix},$$

$$[L^{(2)}L^{(1)}, A^{(2)}] \approx \begin{pmatrix} 2.3 & 1.8 & 1.0 \\ 0.609 & 0.0038 & -1.31 \\ 0.348 & 966 & 1270 \end{pmatrix}.$$

The LU decomposition is then

$$L = \begin{pmatrix} 1 & 0 & 0 \\ 0.609 & 1 & 0 \\ 0.348 & 966 & 1 \end{pmatrix}, \quad U := \begin{pmatrix} 2.3 & 1.8 & 1.0 \\ 0 & 0.0038 & -1.31 \\ 0 & 0 & 1270 \end{pmatrix}.$$

We solve the linear equation system through forward and backward substitution:

$$A\tilde{x} = L \underbrace{U\tilde{x}}_{=y} = b$$

First, we have

$$y_1 = 1.2, \quad y_2 = -2.1 - 0.609 \cdot 1.2 \approx -2.83,$$

$$y_3 = 0.6 - 0.348 \cdot 1.2 + 966 \cdot 2.83 \approx 2730.$$

And finally, we obtain

$$\tilde{x}_3 = \frac{2730}{1270} \approx 2.15,$$

$$\tilde{x}_2 = \frac{-2.83 + 1.31 \cdot 2.15}{0.0038} \approx -3.55,$$

$$\tilde{x}_1 = \frac{1.2 + 1.8 \cdot 3.55 - 1 \cdot 2.15}{2.3} \approx 2.37.$$

For the solution \tilde{x} it holds

$$\tilde{x} = \begin{pmatrix} 2.37 \\ -3.55 \\ 2.15 \end{pmatrix}, \quad \frac{\|\tilde{x} - x\|_2}{\|x\|_2} \approx 1.4,$$

which is a relative error of 140%. Here, we only considered round-off errors and not yet perturbed entries!

The previous example shows the significance of pivoting. In the second step, we chose as Pivot element 0.0038, which is close to 0. For this reason, we obtained values of different orders in the matrices L and U . This yields instabilities in the computation.

Example 4.5 (LU decomposition with pivoting). We continue the previous example in step 2 and search now first the Pivot element:

$$[L^{(1)}, A^{(1)}] = \begin{pmatrix} 2.3 & 1.8 & 1.0 \\ 0.609 & 0.0038 & -1.31 \\ 0.348 & \boxed{3.67} & 1.75 \end{pmatrix}, \quad P^{(2)} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$

Then:

$$[\tilde{L}^{(1)}, \tilde{A}^{(1)}] = \begin{pmatrix} 2.3 & 1.8 & 1.0 \\ 0.348 & 3.67 & 1.75 \\ 0.609 & 0.0038 & -1.31 \end{pmatrix}.$$

It follows

$$F^{(2)} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{0.0038}{3.67} & 1 \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -0.00104 & 1 \end{pmatrix},$$

$$[L^{(2)}\tilde{L}^{(1)}, \tilde{A}^{(2)}] \approx \begin{pmatrix} 2.3 & 1.8 & 1.0 \\ 0.348 & 3.67 & 1.75 \\ 0.609 & 0.00104 & -1.31 \end{pmatrix}.$$

We obtain the decomposition $\tilde{L}U = PA$ as

$$\tilde{L} := \begin{pmatrix} 1 & 0 & 0 \\ 0.348 & 1 & 0 \\ 0.609 & 0.00104 & 1 \end{pmatrix}, \quad U := \begin{pmatrix} 2.3 & 1.8 & 1.0 \\ 0 & 3.67 & 1.75 \\ 0 & 0 & -1.31 \end{pmatrix}, \quad P := \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$

The linear equation system, we solve in the form:

$$PAx = \tilde{L} \underbrace{Ux}_{=y} = Pb.$$

For the right hand side, it holds $\tilde{b} = Pb = (1.2, 0.6, -2.1)^T$ and forward substitution in $\tilde{L}y = \tilde{b}$ yields

$$\begin{aligned} y_1 &= 1.2, \\ y_2 &= 0.6 - 0.348 \cdot 1.2 \approx 0.182, \\ y_3 &= -2.1 - 0.609 \cdot 1.2 - 0.00104 \cdot 0.182 \approx -2.83. \end{aligned}$$

As approximation \tilde{x} , we obtain

$$\tilde{x} = \begin{pmatrix} 0.350 \\ -0.980 \\ 2.160 \end{pmatrix}$$

with a relative error

$$\frac{\|\tilde{x} - x\|_2}{\|x\|_2} \approx 0.0002,$$

which is now only 0.02% rather than 140% as in the previous example.

These examples show that the constructed LU decomposition is not really a truly decomposition, but only an approximation of $A \approx LU$ due to round-off errors. Computing LU easily yields the test and we can compute the error $A - LU$.

As we emphasized at the beginning of this section, the LU decomposition is one of the most important schemes for solving linear equation systems. Due to third order computational cost, however, computing time exceeds easily the capacities of modern computers - despite faster and bigger machines! For instance, we obtain as illustration:

Within the discretization of differential equations we usually obtain very large linear equation systems with $n \sim 10^6 - 10^9$. The good news are that these systems have a special structure, such as symmetry or a sparse pattern with just a few non-zero entries. This is for instance the case for the finite element discretization of Poisson's problem (Chapter 17). Here, in contrast to the previous table, we may have $n = 1\,000\,000$, but thanks to these structural properties, we can solve such systems within one minute.

n	Operations ($\approx \frac{1}{3}n^3$)	Time LU decomposition
100	$300 \cdot 10^3$	30 μ s
1 000	$300 \cdot 10^6$	30 ms
10 000	$300 \cdot 10^9$	30 s
100 000	$300 \cdot 10^{12}$	10 h
1 000 000	$300 \cdot 10^{15}$	1 year

Table 4.2: Computational time for constructing the LU decomposition of a non-sparse matrix $A \in \mathcal{M}_n(\mathbb{K})$ on a computer with 10 GigaFLOPS with optimal load-balancing.

4.5.2 LU decomposition for diagonal-dominant matrices

Proposition 4.32 showed us that the LU decomposition for arbitrary regular matrices with pivoting exists. On the other hand, there are many matrices for which the LU decomposition is stable without pivoting. Examples are positive definite or diagonal dominant matrices:

Definition 4.31 (Diagonal dominant matrices). *A matrix $A \in \mathcal{M}_n(\mathbb{K})$ is called diagonal dominant, if*

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|, \quad i = 1, \dots, n.$$

A diagonal dominant matrix has the largest (absolute) value on the diagonal. For regular matrices the diagonal is moreover always nonzero.

Proposition 4.33 (LU decomposition of diagonal dominant matrices). *Let $A \in \mathcal{M}_n(\mathbb{K})$ be a regular and diagonal dominant matrix. Then, the LU decomposition can be constructed without pivoting and all Pivot elements $a_{ii}^{(i-1)}$ are non-zero.*

4.5.3 Case of Hermitian positive definite matrices: Cholesky decomposition

In the case of a Hermitian positive definite matrix A , a variant of the LU decomposition can be obtained, namely the Cholesky factorization. It has the form

$$A = \bar{C}^T C,$$

where C is an upper triangular matrix with real positive diagonal entries.

4.6 Canonical inner product, Euclidean norm, matrix norm

4.6.1 Canonical inner product and Euclidean norm

The general definition of a norm has been given in Definition 1.3, and the general definition of an inner product will be given in chapter 6. In this chapter dedicated to linear algebra in \mathbb{K}^n we restrict ourselves to the canonical inner product and the associated Euclidean norm which are the most standard ones.

Definition 4.32. *The canonical inner (or scalar) product of two vectors $(x_1, \dots, x_n) \in \mathbb{K}^n$ and $(y_1, \dots, y_n) \in \mathbb{K}^n$ is defined by*

$$\langle x, y \rangle := \overline{x_1}y_1 + \dots + \overline{x_n}y_n.$$

The Euclidean norm of $(x_1, \dots, x_n) \in \mathbb{K}^n$ is defined by

$$\|x\|_2 := \sqrt{\langle x, x \rangle} = \sqrt{|x_1|^2 + \dots + |x_n|^2}.$$

Obviously, for real vectors, conjugacy and moduli can be dropped.

Next, we have the classical Cauchy-Schwarz inequality for vectors of \mathbb{K}^n :

Proposition 4.34. *For all $x, y \in \mathbb{K}^n$ we have*

$$|\langle x, y \rangle| \leq \|x\|_2 \|y\|_2.$$

4.6.2 Matrix norm

Definition 4.33. *A matrix norm on $\mathcal{M}_n(\mathbb{K})$ is a norm $\|\cdot\|$ on the vector space $\mathcal{M}_n(\mathbb{K})$ that satisfies*

$$\|AB\| \leq \|A\| \|B\| \quad \forall A, B \in \mathcal{M}_n(\mathbb{K}).$$

As example we first mention the Frobenius norm, which is associated with an inner product (see chapter 6).

Definition 4.34. *The inner product of $A = (a_{ij}) \in \mathcal{M}_n(\mathbb{K})$ and $B = (b_{ij}) \in \mathcal{M}_n(\mathbb{K})$ is defined by*

$$\langle A, B \rangle := \text{tr}(\bar{A}^T B) = \sum_{i=1}^n \sum_{j=1}^n \bar{a}_{ij} b_{ij}.$$

The Frobenius norm of $A = (a_{ij}) \in \mathcal{M}_n(\mathbb{K})$ is defined by

$$\|A\|_F := \sqrt{\langle A, A \rangle} = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2}.$$

There are also the so-called induced norms, which are more appropriate for some applications. These norms are constructed from a vector norm.

Definition 4.35. *Let $\|\cdot\|$ be a vector norm on \mathbb{K}^n . The induced norm on $\mathcal{M}_n(\mathbb{K})$ is defined by*

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}. \quad (4.8)$$

It is an easy exercise to check that such norms are indeed matrix norms. However they are not associated with an inner product on matrices, even if the Euclidean vector norm is considered.

Solving the maximization problem (4.8) is in principle not an easy task. However, when the Euclidean vector norm is used, it can be reformulated as an eigenvalue problem.

Proposition 4.35. *The matrix norm induced by the Euclidean norm satisfies*

$$\begin{aligned} \|A\|_2 &= \sqrt{\varrho(\bar{A}^T A)} \\ &= \varrho(A) \text{ if } A \text{ is Hermitian (or real symmetric),} \end{aligned}$$

where ϱ stands for the spectral radius (Definition 5.4 in the next chapter). It is called the spectral norm.

By definition an induced norm satisfies

$$\|Ax\| \leq \|A\| \|x\| \quad \forall x \in \mathbb{K}^n.$$

When a matrix norm and a vector norm satisfy the above relation they are said to be compatible.

Chapter 5

Diagonalization of square matrices and applications

Diagonalization of matrices has many applications. The major purpose is to simplify some given matrix (in case it is diagonalizable) to a diagonal matrix by conserving important properties such as the rank, characteristic polynomial, determinant, trace, and eigenvalues and their algebraic multiplicities. Mathematically, diagonalization is the process to construct a similar matrix. Similar matrices represent the same linear mapping under two different bases.

5.1 Eigenvalues, eigenvectors, diagonalization

5.1.1 Eigenvalues and eigenvectors

Let $A \in \mathcal{M}_n(\mathbb{K})$ with in general $\mathbb{K} = \mathbb{C}$.

Definition 5.1. Let $\lambda \in \mathbb{K}$. We say that λ is an eigenvalue of A if there exists a column vector $x \in \mathbb{K}^n$, $x \neq 0$, such that

$$Ax = \lambda x.$$

Such an x is called an eigenvector associated with the eigenvalue λ .

Proposition 5.1. The function

$$\begin{aligned} \mathbb{K} &\rightarrow \mathbb{K} \\ \lambda &\mapsto P_A(\lambda) := \det(A - \lambda I_n) \end{aligned}$$

is a polynomial of degree n . It is called the characteristic polynomial of A .

Proposition 5.2. Let $\lambda \in \mathbb{K}$. Then λ is an eigenvalue of A if and only if $P_A(\lambda) = 0$.

Determining the eigenvalues of a matrix amounts to finding the roots of its characteristic polynomial. However this is not an easy task, since characteristic polynomials have no special property¹. Explicit formulas exist only for $n \leq 4$. For the numerical approximation of eigenvalues iterative methods are used. They are usually not based on the characteristic polynomial. The simplest method, the power method, is described in section 5.3.

Definition 5.2. Let λ be an eigenvalue of A . We call an eigenspace of A associated with the eigenvalue λ the set

$$E_A(\lambda) = \{x \in \mathbb{K}^n \text{ s.t. } Ax = \lambda x\}.$$

It is a linear subspace of \mathbb{K}^n .

In other words, $E_A(\lambda)$ is the set of eigenvectors associated with λ and the zero vector.

¹Given an arbitrary polynomial P one can always construct a matrix, called companion matrix, whose characteristic polynomial is P , up to a multiplicative constant.

5.1.2 Diagonalization

Definition 5.3. A matrix $A \in \mathcal{M}_n(\mathbb{K})$ is diagonalizable if there exists an invertible matrix $P \in \mathcal{M}_n(\mathbb{K})$ such that $D := P^{-1}AP$ is diagonal.

Remark 5.1. The matrix $D := P^{-1}AP$ corresponds to a change of coordinate system (or change of basis) in the map $x \mapsto Ax$. Indeed, if $x = P\tilde{x}$, $y = P\tilde{y}$ and $y = Ax$ then $\tilde{y} = P^{-1}AP\tilde{x}$. Therefore P is called a change of basis matrix (or transition matrix). Its columns contain the coordinates of the new basis vectors in the former (here canonical) basis. The relation $D = P^{-1}AP$ is straightforwardly inverted into $A = PDP^{-1}$.

Proposition 5.3. A matrix $A \in \mathcal{M}_n(\mathbb{K})$ is diagonalizable if and only if there exists a basis $\mathcal{B} = (v_1, \dots, v_n)$ such that each v_i is an eigenvector associated with an eigenvalue λ_i of A . In this case, if P is the matrix made columnwise of the vectors (v_1, \dots, v_n) , then

$$P^{-1}AP = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} =: \text{diag}(\lambda_1, \dots, \lambda_n).$$

Therefore it is to be remembered that the columns of a transition matrix associated with a diagonalization form a basis of eigenvectors.

Remark 5.2. The diagonal coefficients λ_i need not be distinct.

Proposition 5.4. Eigenvectors associated with distinct eigenvalues are linearly independent.

Corollary 5.5. If $A \in \mathcal{M}_n(\mathbb{K})$ admits n distinct eigenvalues then A is diagonalizable.

The fact that a matrix be diagonalizable is related to the number of its eigenvalues and their order of multiplicity². We do not enter into these developments, which are extremely classical and well-documented. We only stress that a real matrix may be diagonalizable on \mathbb{C} while it is not on \mathbb{R} , due to a possibly larger number of eigenvalues. The case of symmetric and Hermitian matrices is of special interest.

5.1.3 Case of Hermitian matrices

Theorem 5.6. If A is Hermitian then A is diagonalizable with the help of real eigenvalues. Moreover the transition matrix P can be chosen such that $P^{-1} = \bar{P}^T$ (this is called a unitary matrix).

This theorem has a direct consequence on the characterization of Hermitian positive (semi-) definite matrices.

Corollary 5.7. If $A \in \mathcal{M}_n(\mathbb{K})$ be a Hermitian matrix, then A is positive semi-definite (resp. positive definite) if and only if all its eigenvalues are ≥ 0 (resp. > 0).

5.2 Applications

The diagonalization has many applications. One of them, the expression of matrix power, turns out to be very important in the analysis of iterative methods.

Proposition 5.8. If $A = PDP^{-1}$ with $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ then, for all $k \in \mathbb{N}$, we have $A^k = PD^kP^{-1}$ with $D^k = \text{diag}(\lambda_1^k, \dots, \lambda_n^k)$.

²The notion of order of multiplicity is related, but not equivalent, to the fact that an eigenvalue may appear in several places in the diagonal form.

When k goes to infinity the behavior of A^k is governed by the λ_i^k 's. If $|\lambda_i| < 1$ for all i then clearly $A^k \rightarrow 0$ when $k \rightarrow +\infty$. If $|\lambda_{i_0}| > 1$ for some index i_0 then, calling v_{i_0} an eigenvector associated with λ_{i_0} , we have $A^k v_{i_0} = \lambda_{i_0}^k v_{i_0}$, whose norm goes to infinity when $k \rightarrow +\infty$.

Definition 5.4. *The largest modulus of the eigenvalues (on \mathbb{C}) of a (real or complex) matrix A is called spectral radius and is denoted by $\rho(A)$ with $\rho = \max_{\lambda \in \sigma(A)} |\lambda|$.*

A fairly similar reasoning can be made concerning the analysis of first order linear differential systems with constant coefficients, where diagonalization uncouples the system, leading to explicit solutions:

$$y' = Ay \Leftrightarrow (P^{-1}y)' = D(P^{-1}y) \Leftrightarrow \begin{cases} y = Pz, & z = (z_1, \dots, z_n)^T \\ z_i' = \lambda_i z_i & \forall i = 1, \dots, n \end{cases}.$$

Here the signs of the real parts of the eigenvalues govern the asymptotic behavior of the solutions.

5.3 The power method (complement)

There exist many more or less sophisticated methods for the numerical approximation of all or some eigenvalues of a given matrix. The extremal eigenvalues (in modulus or in real part) often carry the most important information, as mentioned in section 5.2. The power method is an elementary method aiming at computing the largest modulus of the eigenvalues.

Let $A \in \mathcal{M}_n(\mathbb{R})$. The power method consists in defining two sequences (x_n) and (y_n) by

- $x_0 \in \mathbb{R}^n$ given,
- for all $k \in \mathbb{N}^*$

$$y_k = Ax_{k-1}, \quad x_k = \frac{y_k}{\|y_k\|}.$$

Here, $\|\cdot\|$ is an arbitrary norm on \mathbb{R}^n , typically the Euclidean norm is used.

Theorem 5.9. *We assume that A is diagonalizable (possibly on \mathbb{C}), with basis of eigenvectors (e_1, \dots, e_n) , chosen to be of unit norm, associated with the eigenvalues $(\lambda_1, \dots, \lambda_n)$. We assume further that*

$$|\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_{n-1}| < |\lambda_n|.$$

If $x_0 = \sum_{i=1}^n \beta_i e_i$ with $\beta_n \neq 0$, then

$$\lim_{k \rightarrow +\infty} \|y_k\| = |\lambda_n|.$$

Moreover, if λ_n is real and positive, $e_n \in \mathbb{R}^n$, then $\lim_{k \rightarrow +\infty} x_k = \pm e_n$.

Remark 5.3. 1. *There is a priori no way to ensure that x_0 admits a nonzero coordinate along e_n , since by definition e_n is unknown. However, because of round-off errors in a numerical implementation, such a component will automatically appear in the vector x_k after some iterations.*

2. *A variant, the inverse power method, aims at computing the smallest modulus of the eigenvalues.*

Chapter 6

Vector spaces and elements of topology

Here again, \mathbb{K} stands either for \mathbb{R} or for \mathbb{C} .

6.1 Vector spaces

6.1.1 General definition

We have defined linear subspaces of \mathbb{K}^n in section 4.3, however there is a more general and axiomatic definition of the concept of vector space.

Definition 6.1. *A vector space on \mathbb{K} is a set X , whose elements are called vectors, equipped with two operations*

$$(x, y) \in X^2 \mapsto x + y \in X,$$

$$(\lambda, x) \in \mathbb{K} \times X \mapsto \lambda x \in X,$$

satisfying the following properties:

$$\forall x, y, z \in X, \quad (x + y) + z = x + (y + z),$$

$$\forall x, y \in X, \quad x + y = y + x,$$

there exists a vector denoted by 0 such that $0 + x = x$ for all $x \in X$,

for all $x \in X$ there exists a vector denoted by $-x$ such that $x + (-x) = 0$,

$$\forall (\lambda, \mu, x) \in \mathbb{K} \times \mathbb{K} \times X, \quad (\lambda\mu)x = \lambda(\mu x),$$

$$\forall x \in X, \quad 1x = x,$$

$$\forall (\lambda, x, y) \in \mathbb{K} \times X \times X, \quad \lambda(x + y) = \lambda x + \lambda y,$$

$$\forall (\lambda, \mu, x) \in \mathbb{K} \times \mathbb{K} \times X, \quad (\lambda + \mu)x = \lambda x + \mu x.$$

Clearly, \mathbb{K}^n equipped with its standard operations is a vector space on \mathbb{K} . This also holds true for all its subspaces. We speak of real vector space when $\mathbb{K} = \mathbb{R}$, and of complex vector space when $\mathbb{K} = \mathbb{C}$.

6.1.2 Vector spaces of finite dimension

Definition 6.2. Let X be a vector space and $(u_i)_{i \in I}$ be a family (not necessarily finite) of elements of X . A vector $x \in X$ is said to be a linear combination of $(u_i)_{i \in I}$ if there exist a finite subfamily $(u_i)_{i \in J}$ and a family $(\alpha_i)_{i \in J}$ of elements of \mathbb{K} such that

$$x = \sum_{i \in J} \alpha_i u_i.$$

Definition 6.3. A vector space X is said to be of finite dimension if there exists a finite family (u_1, \dots, u_n) such that every vector of X is a linear combination of (u_1, \dots, u_n) .

Of course, \mathbb{K}^n is of finite dimension. But **function spaces are usually not of finite dimension.** For instance $\mathcal{C}([a, b], \mathbb{R})$, the set of continuous real-valued functions on $[a, b]$ is a vector space but it is not of finite dimension. There are exceptions, however. For instance the space of polynomial functions of degree $\leq n$ is of finite dimension: it is spanned by the monomials $(x^k)_{0 \leq k \leq n}$.

6.2 Basic topology

6.2.1 Norms and balls

Definition 6.4. Let X be a vector space. A norm on X is a function

$$x \in X \mapsto \|x\| \in \mathbb{R}_+$$

that satisfies:

$$\begin{aligned} \forall x \in X, \quad \|x\| = 0 &\Leftrightarrow x = 0, \\ \forall (x, \alpha) \in X \times \mathbb{K}, \quad \|\alpha x\| &= |\alpha| \|x\|, \\ \forall (x, y) \in X^2, \quad \|x + y\| &\leq \|x\| + \|y\|. \end{aligned} \tag{6.1}$$

A vector space endowed with a norm is called a normed vector space.

On $X = \mathbb{K}$, the canonical norm is the modulus (or absolute value) $|\cdot|$.

The triangle inequality (6.1) implies other inequality

$$\forall (x, y) \in X^2, \quad \left| \|x\| - \|y\| \right| \leq \|x - y\|. \tag{6.2}$$

Due to this property, we say that the norm is a 1-Lipschitz (or non-expansive) function.

When there is no confusion the notation $\|x\|$ will be systematically used for the norm of x on its underlying space.

Example 6.1. On \mathbb{K}^n the most classical norms are defined, for $x = (x_1, \dots, x_n)$, by

$$\|x\|_1 = \sum_{i=1}^n |x_i|, \quad \|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}, \quad \|x\|_\infty = \max_{i=1, \dots, n} |x_i|.$$

For instance, for $x = (-1, 2)$ we obtain $\|x\|_1 = 3$, $\|x\|_2 = \sqrt{5}$, $\|x\|_\infty = 2$.

The norm $\|\cdot\|_2$ has already been mentioned, and it is the most classical norm because it has a noticeable property: it is associated with the canonical inner product of \mathbb{K}^n . Actually, to check that the norms from example 6.1 are indeed norms according to Definition 6.4, the only nontrivial task is to check the triangle inequality (6.1) for $\|\cdot\|_2$. For this we use the expression of the norm in terms of the canonical inner product and the Cauchy-Schwarz inequality from Proposition 4.34 to obtain

$$\begin{aligned} \|x + y\|_2^2 &= \langle x + y, x + y \rangle = \|x\|_2^2 + \|y\|_2^2 + 2\Re\langle x, y \rangle \leq \|x\|_2^2 + \|y\|_2^2 + 2|\langle x, y \rangle| \\ &\leq \|x\|_2^2 + \|y\|_2^2 + 2\|x\|_2\|y\|_2 = (\|x\|_2 + \|y\|_2)^2. \end{aligned}$$

Definition 6.5. Let X be a normed vector space, $z \in X$ and $\rho > 0$. The open ball of center z and radius ρ is the set

$$B_\rho(z) = \{x \in X \text{ s.t. } \|x - z\| < \rho\}.$$

The closed ball of center z and radius ρ is the set

$$\bar{B}_\rho(z) = \{x \in X \text{ s.t. } \|x - z\| \leq \rho\}.$$

When the radius ρ is chosen equal to 1 and the center is fixed at the origin, the balls $B_1(0)$ and $\bar{B}_1(0)$ are called the (open or closed) unit balls. See Fig. 6.1.

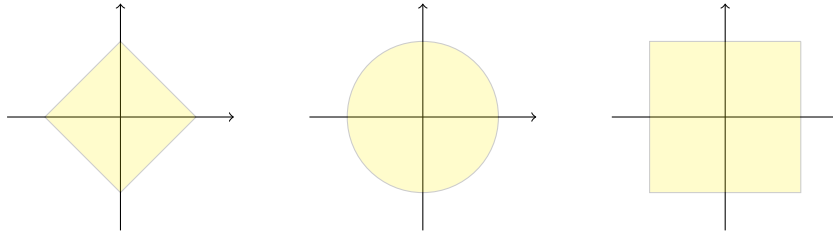


Figure 6.1: Unit balls of \mathbb{R}^2 for the norms (from left to right): $\|\cdot\|_1$, $\|\cdot\|_2$ and $\|\cdot\|_\infty$.

Definition 6.6. Let X be a normed vector space and $z \in X$. A subset U of X is said to be a neighborhood of z if there exists $\rho > 0$ such that $B_\rho(z) \subset U$.

6.2.2 Open sets, closed sets (complement)

Definition 6.7. Let X be a normed vector space. A subset A of X is said to be open if

$$\forall z \in A \exists \rho > 0 \text{ s.t. } B_\rho(z) \subset A,$$

i.e., A is a neighborhood of all its points.

More generally, a subset $A \subset E \subset X$ is said to be open relatively to the set E if

$$\forall z \in A \exists \rho > 0 \text{ s.t. } B_\rho(z) \cap E \subset A,$$

i.e., A is a relative neighborhood of all its points.

Definition 6.8. Let X be a normed vector space. A subset A of X is said to be closed if its complementary set $X \setminus A$ is open.

A subset $A \subset E \subset X$ is said to be closed relatively to E if $E \setminus A$ is open relatively to E .

It is a classical exercise to show that open balls are indeed open sets, and that closed balls are closed sets.

On \mathbb{R} (equipped with the absolute value), all intervals of the form $]a, b[$, $] - \infty, b[$, $]a, +\infty[$ are open, and all intervals of the form $[a, b]$, $] - \infty, b]$, $[a, +\infty[$ are closed. Intervals of the form $[a, b[$ or $]a, b]$ are neither open nor closed. The interval $] - \infty, +\infty[= \mathbb{R}$ is open and closed.

The following proposition is straightforward and also applies to relative open / closed sets.

Proposition 6.1. An arbitrary union of open sets is an open set. A finite intersection of open sets is an open set.

An arbitrary intersection of closed sets is a closed set. A finite union of closed sets is a closed set.

6.2.3 Closure, interior (complement)

Definition 6.9. Let X be a normed vector space and A be a subset of X . The closure of A , denoted by \bar{A} , is the subset of X defined by

$$\forall x \in X, \quad x \in \bar{A} \Leftrightarrow \forall \rho > 0, B_\rho(x) \cap A \neq \emptyset.$$

The closure of a set A is always closed, in fact it is the smallest (in the sense of inclusion) closed set containing A . If A is closed then $\bar{A} = A$.

Definition 6.10. Let X be a normed vector space and A be a subset of X . The interior of A , denoted by $\text{int } A$ is the subset of X defined by

$$\forall x \in X, \quad x \in \text{int } A \Leftrightarrow \exists \rho > 0, B_\rho(x) \subset A.$$

The interior of a set A is always open, in fact it is the largest open set contained in A . If A is open then $\text{int } A = A$.

It is classical to show that $X \setminus \bar{A} = \text{int}(X \setminus A)$ and $X \setminus \text{int } A = \overline{X \setminus A}$. In addition, the closure of an open ball is the closed ball of the same center and radius. The interior of a closed ball is the open ball of the same center and radius.

6.2.4 Converging sequences (complement)

The notion of convergence in normed vector spaces extends that of real-valued sequences, only replacing the absolute value by the norm.

Definition 6.11. Let X be a normed vector space and (x_n) be a sequence of elements of X . We say that (x_n) is converging to $x \in X$ if

$$\forall \varepsilon > 0 \exists N \in \mathbb{N} \text{ s.t. } n \geq N \Rightarrow \|x_n - x\| < \varepsilon.$$

This can be phrased as: x_n becomes arbitrarily close to x provided that n be large enough.

6.2.5 Equivalent norms

Definition 6.12. Two norms $\|\cdot\|_a$ and $\|\cdot\|_b$ on X are said to be equivalent if there exists $\alpha, \beta > 0$ such that

$$\alpha\|x\|_a \leq \|x\|_b \leq \beta\|x\|_a \quad \forall x \in X.$$

Proposition 6.2. If two norms $\|\cdot\|_a$ and $\|\cdot\|_b$ on X are equivalent then a subset of X is open for the norm $\|\cdot\|_a$ if and only if it is open for the norm $\|\cdot\|_b$. Of course, the same holds for closed sets. Likewise, if the norms $\|\cdot\|_a$ and $\|\cdot\|_b$ are equivalent, a sequence of X is converging for the norm $\|\cdot\|_a$ if and only if it converges for the norm $\|\cdot\|_b$.

Theorem 6.3. If X is a vector space of finite dimension then all norms on X are equivalent.

This property is convenient, however it only holds in finite dimensional spaces. Be careful!

6.3 Special normed spaces (complement)

6.3.1 Inner product spaces

Definition 6.13. Let X be a vector space. A function

$$(x, y) \in X^2 \mapsto \langle x, y \rangle \in \mathbb{K}$$

is said to be an inner product (or scalar product) if

$$\forall (x, y) \in X^2, \quad \langle x, y \rangle = \overline{\langle y, x \rangle},$$

$$\begin{aligned} \forall(x, x', y) \in X^3, \forall \alpha \in \mathbb{K}, \quad \langle \alpha x + x', y \rangle &= \bar{\alpha} \langle x, y \rangle + \langle x', y \rangle, \\ \forall x \in X \setminus \{0\}, \quad \langle x, x \rangle &> 0. \end{aligned}$$

A vector space endowed with an inner product is called inner product space.

Observe that two above properties imply

$$\langle x, \alpha y + y' \rangle = \alpha \langle x, y \rangle + \langle x, y' \rangle.$$

Therefore one gets the standard linearity of the inner product with respect to the second argument, while one sometimes speaks of antilinearity with respect to the first argument because of the conjugacy. Prescribing the antilinearity with respect to the first argument, hence the linearity with respect to the second one, is a matter of (non-universal) convention. Of course, when $\mathbb{K} = \mathbb{R}$ this discussion is irrelevant.

Typical inner products are those defined in section 4.6 (in the case of matrices, considering $\mathcal{M}_n(\mathbb{K})$ as a vector space).

Theorem 6.4. *If X is an inner product space then it is a normed space for the norm defined by*

$$\|x\| = \sqrt{\langle x, x \rangle}.$$

In addition we have the Cauchy-Schwarz inequality

$$\forall(x, y) \in X^2, \quad |\langle x, y \rangle| \leq \|x\| \|y\|.$$

The norm defined above will be considered as canonical in any inner product space.

If x, y belong to an inner product space then simple algebra yields the expansion formula

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2 + 2\langle x, y \rangle,$$

the parallelogram equality

$$\|x + y\|^2 + \|x - y\|^2 = 2(\|x\|^2 + \|y\|^2),$$

and the polarization identity

$$\langle x, y \rangle = \frac{1}{4} (\|x + y\|^2 - \|x - y\|^2).$$

If A is an arbitrary subset of X then we define the orthogonal of A by

$$A^\perp = \{x \in X \text{ s.t. } \langle x, y \rangle = 0 \forall y \in A\}.$$

It is always a closed linear subspace of X . When A is a singleton we often write a^\perp for $\{a\}^\perp$.

6.3.2 Banach and Hilbert spaces

Let us first extend the notion of real-valued Cauchy sequence.

Definition 6.14. *Let X be a normed vector space and (x_n) be a sequence of elements of X . We say that (x_n) is a Cauchy sequence if*

$$\forall \varepsilon > 0 \exists N \in \mathbb{N} \text{ s.t. } n, m \geq N \Rightarrow \|x_n - x_m\| < \varepsilon.$$

Definition 6.15. *A normed vector space such that every Cauchy sequence is converging is said to be a Banach space. When it is at the same time an inner product space, it is called a Hilbert space.*

Note that a Banach space remains a Banach space when changing the norm to an equivalent one. But in general, the choice of the norm is crucial when defining a Banach space.

Theorem 6.5. *A normed vector space of finite dimension is a Banach space.*

6.3.3 Complement on matrix norms

In order to establish the convergence of sequences of vectors of \mathbb{K}^n the choice of the norm is irrelevant, as explained in section 6.2.5. When such sequences arise from matrix operations it comes the question of the choice of a compatible matrix norm (see section 4.6.2). Optimal convergence results take advantage of this freedom thanks to the following convenient theorem. It highlights the prominent role played by the spectral radius (see Definition 5.4).

Theorem 6.6. *Let A be a square matrix with coefficients in \mathbb{K} .*

1. *It holds for every matrix norm $\|\cdot\|$*

$$\varrho(A) \leq \|A\|.$$

2. *For all $\varepsilon > 0$ there exists an induced matrix norm $\|\cdot\|$ such that*

$$\|A\| \leq \varrho(A) + \varepsilon.$$

As example of application we immediately see that the sequence (A^n) goes to 0 as soon as $\rho(A) < 1$, since $\|A^n\| \leq \|A\|^n \leq (\varrho(A) + \varepsilon)^n$ and ε can be taken arbitrarily small. The if and only if statement holds also true by definition of an eigenvalue. We underline that it is not requested that A be diagonalizable.

Chapter 7

Linear equation systems and iterative solvers (complement)

7.1 Stability analysis of linear equation systems

7.1.1 Motivation

Previously, we have shown how to solve linear equation systems. Tacitly we assumed that we deal with ‘nice’ numbers. However, in practice the computer will solve these systems for us and in most cases we deal with floating point numbers from \mathbb{R} . As we discussed in the introduction (Chapters 1 and 2), we now deal with round-off errors.

The question comes, how such round-off errors may influence the solution of $Ax = b$. As an example, let us consider

$$\begin{pmatrix} 0.988 & 0.960 \\ 0.992 & 0.963 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0.084 \\ 0.087 \end{pmatrix}$$

with the solution $(x, y)^T = (3, -3)^T$. We compute the solution with a precision of three correct digits using Gauss elimination:

$$\left(\begin{array}{cc|c} 0.988 & 0.960 & 0.084 \\ 0.992 & 0.963 & 0.087 \end{array} \right) \quad \times 0.988/0.992$$

$$\left(\begin{array}{cc|c} 0.988 & 0.960 & 0.084 \\ 0.988 & 0.959 & 0.0866 \end{array} \right) \quad \downarrow -$$

$$\left(\begin{array}{cc|c} 0.988 & 0.959 & 0.084 \\ 0 & 0.001 & -0.0026 \end{array} \right) \quad .$$

We now have the typical triangular system of A . The right hand side b was accordingly modified. Through backward substitution we obtain

$$\begin{aligned} 0.001y &= -0.0026 \Rightarrow y = -2.6 \\ 0.988x &= 0.087 - 0.959 \cdot (-2.6) \approx 2.58 \Rightarrow x = 2.61. \end{aligned}$$

Here, we have $(x, y) = (2.61, -2.60)$. The relative error of the numerical approximation is more than 10%. It seems that the solution of linear equation systems is either ill-conditioned (see the general concept of condition numbers in the chapter 2) or the Gaussian elimination is numerically unstable. The conditioning and stability will be addressed of the following subsection. In addition, in many applications, we deal with very large linear equation systems, where $A \in \mathbb{R}^{n \times n}$ with $n \gg 10^6$ and larger. For this reason, the efficiency and memory requirements play a crucial role and therefore, direct methods (such as Gaussian elimination) are unfortunately not always feasible, and rather iterative methods come into play.

7.1.2 Stability analysis

We are given $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$ and consider

$$Ax = b.$$

Through numerical errors, round-off errors, false input data or experimental errors, A and b may be perturbed:

$$\tilde{A}\tilde{x} = \tilde{b}.$$

Here, $\tilde{A} = A + \delta A$ and $\tilde{b} = b + \delta b$ with the perturbations δA and δb , respectively.

We derive now results for the conditioning and error amplification for the solution of linear equation systems. Errors can be in the matrix A as well as in the right hand side b .

We first consider errors in the right hand side vector.

Proposition 7.1 (Perturbation of the right hand side b). *Let $A \in \mathbb{R}^{n \times n}$ be a regular matrix and $b \in \mathbb{R}^n$. With $x \in \mathbb{R}^n$ we denote the solution of $Ax = b$. Let δb be a perturbation of the right hand side $\tilde{b} = b + \delta b$ and \tilde{x} the solution of the perturbed system $A\tilde{x} = \tilde{b}$. Furthermore, we denote with $\|\cdot\|$ a matrix norm that is compatible with the vector norm $\|\cdot\|$. Then, it holds*

$$\frac{\|\delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\delta b\|}{\|b\|},$$

with the condition number of the matrix, i.e.,

$$\text{cond}(A) = \|A\| \cdot \|A^{-1}\|.$$

Proof. Let $\|\cdot\|$ be an arbitrary matrix norm with compatible vector norm $\|\cdot\|$. For the solution $x \in \mathbb{R}^n$ and the perturbed solution $\tilde{x} \in \mathbb{R}^n$, we have

$$\tilde{x} - x = A^{-1}(A\tilde{x} - Ax) = A^{-1}(\tilde{b} - b) = A^{-1}\delta b$$

Thus:

$$\frac{\|\delta x\|}{\|x\|} \leq \|A^{-1}\| \frac{\|\delta b\|}{\|x\|} \cdot \frac{\|b\|}{\|b\|} = \|A^{-1}\| \frac{\|\delta b\|}{\|b\|} \cdot \frac{\|Ax\|}{\|x\|} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{=: \text{cond}(A)} \frac{\|\delta b\|}{\|b\|}.$$

□

Remark 7.1. *The condition number of a matrix is very important in numerical linear algebra. Let us also consider the matrix-vector multiplication $y = Ax$. With perturbed input $\tilde{x} = x + \delta x$, we obtain, since obviously $\text{cond}(A) = \text{cond}(A^{-1})$:*

$$\frac{\|\delta y\|}{\|y\|} \leq \text{cond}(A) \frac{\|\delta x\|}{\|x\|}.$$

Note that by definition the condition number of a matrix depends on the chosen norm. are also equivalent. For symmetric matrices it can be inferred that $\text{cond}_2(A)$:

$$\text{cond}_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2 = \frac{\max\{|\lambda|, \lambda \text{ eigenvalue of } A\}}{\min\{|\lambda|, \lambda \text{ eigenvalue of } A\}}.$$

In the second step, we consider now a perturbation δA of the matrix A . First, we need to ensure that $\tilde{A} = A + \delta A$ is still regular.

Lemma 7.2. *Let $\|\cdot\|$ be the matrix norm induced by the vector norm. Let $B \in \mathbb{R}^{n \times n}$ be a matrix with $\|B\| < 1$. Then, the matrix $I + B$ is regular and it holds the estimate*

$$\|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|}.$$

Proof. It holds

$$\|(I + B)x\| \geq \|x\| - \|Bx\| \geq (1 - \|B\|)\|x\|.$$

Since $1 - \|B\| > 0$, the mapping $I + B$ is one-to-one. Thus, $I + B$ is regular. Furthermore, we have

$$\begin{aligned} 1 = \|I\| &= \|(I + B)(I + B)^{-1}\| = \|(I + B)^{-1} + B(I + B)^{-1}\| \\ &\geq \|(I + B)^{-1}\| - \|B\| \|(I + B)^{-1}\| = \|(I + B)^{-1}\|(1 - \|B\|) > 0. \end{aligned}$$

□

With this result, we address now the perturbation of a matrix A :

Proposition 7.3 (Perturbation of the matrix A). *Let $A \in \mathbb{R}^{n \times n}$ be a regular matrix, and $b \in \mathbb{R}^n$. Furthermore, let $x \in \mathbb{R}^n$ be the solution of $Ax = b$ and let $\tilde{A} = A + \delta A$ a perturbed matrix with $\|\delta A\| \leq \|A^{-1}\|^{-1}$. For the perturbed solution $\tilde{x} = x + \delta x$ of $\tilde{A}\tilde{x} = b$, it holds*

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A)\|\delta A\|/\|A\|} \frac{\|\delta A\|}{\|A\|}.$$

Proof. It holds for x and the perturbed solution \tilde{x} and the error $\delta x := \tilde{x} - x$:

$$\begin{aligned} (A + \delta A)\tilde{x} &= b \\ (A + \delta A)x &= b + \delta Ax \end{aligned} \quad \Rightarrow \quad \delta x = -[A + \delta A]^{-1}\delta Ax.$$

According to our assumption $\|A^{-1}\delta A\| \leq \|A^{-1}\| \|\delta A\| < 1$, it follows with the Lemma 7.2:

$$\begin{aligned} \|\delta x\| &\leq \|[I + A^{-1}\delta A]^{-1}A^{-1}\delta A\| \|x\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\delta A\|} \|\delta A\| \|x\| \\ &\leq \frac{\text{cond}(A)}{1 - \|A^{-1}\| \|\delta A\|} \frac{\|\delta A\|}{\|A\|} \|x\|. \end{aligned}$$

We establish the assertion with the expansion $\|A\|/\|A\|$. □

The combination of the two previous perturbation propositions yields the main result:

Theorem 7.4 (Perturbation theorem for linear equation systems). *Let $A \in \mathbb{R}^{n \times n}$ be a regular matrix and $b \in \mathbb{R}^n$ a right hand side vector. Furthermore, let $x \in \mathbb{R}^n$ be the solution of $Ax = b$ with the regular matrix $A \in \mathbb{R}^{n \times n}$. For the solution $\tilde{x} \in \mathbb{R}^n$ of the perturbed system $\tilde{A}\tilde{x} = \tilde{b}$ with $\delta b = \tilde{b} - b$ and $\delta A = \tilde{A} - A$ and the assumption*

$$\|\delta A\| < \frac{1}{\|A^{-1}\|}$$

it holds the estimate

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A)\|\delta A\|/\|A\|} \left(\frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right)$$

with the condition number

$$\text{cond}(A) = \|A\| \|A^{-1}\|.$$

Proof. We combine the results from Proposition 7.1 and 7.3. To this end, let x be the solution of $Ax = b$, \tilde{x} the perturbed solution $\tilde{A}\tilde{x} = \tilde{b}$ and \hat{x} the solution to the perturbed right hand side $A\hat{x} = \tilde{b}$. Then, we obtain

$$\|x - \tilde{x}\| \leq \|x - \hat{x}\| + \|\hat{x} - \tilde{x}\| \leq \text{cond}(A) \frac{\|\delta b\|}{\|b\|} \|x\| + \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\delta A\|}{\|A\|}} \frac{\|\delta A\|}{\|A\|} \|x\|.$$

Taking into account $\|\delta A\| < \|A^{-1}\|^{-1}$, we have

$$0 \leq \text{cond}(A) \frac{\|\delta A\|}{\|A\|} < \|A\| \|A^{-1}\| \frac{1}{\|A\| \|A^{-1}\|} = 1.$$

Therefore,

$$\text{cond}(A) \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\delta A\|}{\|A\|}},$$

which shows the assertion. \square

With this theorem, we re-consider our motivating example from the beginning of this section. We recall

$$A = \begin{pmatrix} 0.988 & 0.959 \\ 0.992 & 0.963 \end{pmatrix} \\ \Rightarrow A^{-1} \approx \begin{pmatrix} 8302 & -8267 \\ -8552 & 8517 \end{pmatrix}.$$

Employing the maximum row sum $\|\cdot\|_\infty$ it holds

$$\|A\|_\infty = 1.955, \quad \|A^{-1}\|_\infty \approx 17069, \quad \text{cond}_\infty(A) \approx 33370.$$

Solving linear equation systems is indeed often ill-conditioned. As a consequence, round-off errors may be significantly amplified. This also shows that the difficulty is an intrinsic part of the problem statement, but not necessarily a stability issue with the Gaussian elimination (provided appropriate pivoting is performed!). These are partially good news and the reason why Gaussian elimination is often used for moderate sizes of the matrix A .

7.2 Basic solvers

For a moderate number of unknowns a so-called direct solver (LU decomposition, Cholesky) is a good choice to solve $AU = B$. Such methods are always available in most software packages such as Matlab, octave, python, etc:

```
U = sp.sparse.linalg.spsolve(A,B) // in Python
U = A\B // in octave / Matlab
```

For big systems, **iterative solvers** are the methods of choice because they require less memory and less computational cost than direct solvers. In addition, they serve as smoothers in so-called multigrid methods, which may yield optimal complexities $O(n)$, where n is the number of unknowns. We provide a brief introduction to iterative solvers in the following section.

7.2.1 Fixed-point solvers: Richardson, Jacobi, Gauss-Seidel

A large class of schemes is based on so-called **fixed point** methods for solving:

$$f(x) = x.$$

We provide in the following a brief introduction that is based on [15]. First, we have

Definition 7.1. Let $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$ and $C \in \mathbb{R}^{n \times n}$. For an initial guess $x^0 \in \mathbb{R}^n$ we iterate for $k = 1, 2, \dots$:

$$x^k = x^{k-1} + C(b - Ax^{k-1}).$$

Please be careful that k does not denote the power, but the current iteration index. Furthermore, we introduce:

$$B := I - CA \quad \text{and} \quad c := Cb.$$

Then:

$$x^k = Bx^{k-1} + c.$$

Thanks to the construction of

$$g(x) = Bx + c = x + C(b - Ax)$$

it is trivial to see that in the limit $k \rightarrow \infty$, it holds

$$g(x) = x$$

with the solution

$$Ax = b.$$

Remark 7.2. Thanks to Banach's fixed point theorem (see section 12.3), we can investigate under which conditions the above scheme will converge. Actually it should hold

$$\|B\| < 1,$$

for a matrix norm such that

$$\|g(x) - g(y)\| \leq \|B\| \|x - y\|.$$

A similar condition appears in the stability analysis of the schemes for solving ODEs, as reported in the numerical tests of Section 17.5 (blow-up, zig-zag solution and converged solution). A big problem (which is also true for the ODE cases) is that different norms may predict different results. For instance it may happen that

$$\|B\|_2 < 1 \quad \text{but} \quad \|B\|_\infty > 1.$$

Discussions can be found for instance in [15]. In particular, if the spectral radius of B is $\rho(B) < 1$ then $\|B\| < 1$ for some induced matrix norm and $\rho(B) > 1$ implies $\|B\| > 1$ for all matrix norm.

We concentrate now on the algorithmic aspects. The two fundamental requirements for the matrix C (defined above) are:

- It should hold $C \approx A^{-1}$ and therefore $\|I - CA\| \ll 1$;
- It should be simple to construct C .

Of course, we easily see that these two requirements are conflicting statements. As always in numerics we need to find a trade-off that is satisfying for the developer and the computer.

Definition 7.2 (Richardson iteration). The simplest choice of C is the identity matrix, i.e.,

$$C = I$$

Then, we obtain the Richardson iteration

$$x^k = x^{k-1} + \omega(b - Ax^{k-1})$$

with a relaxation parameter $\omega > 0$.

Further schemes require more work and we need to decompose the matrix A first:

$$A = L + D + U.$$

Here, L is a lower-triangular matrix, D a diagonal matrix, and U an upper-triangular matrix. In more detail:

$$A = \underbrace{\begin{pmatrix} 0 & & \dots & 0 \\ a_{21} & \ddots & & \\ \vdots & \ddots & \ddots & \\ a_{n1} & \dots & a_{n,n-1} & 0 \end{pmatrix}}_{=:L} + \underbrace{\begin{pmatrix} a_{11} & & \dots & 0 \\ & \ddots & & \\ & & \ddots & \\ 0 & \dots & & a_{nn} \end{pmatrix}}_{=:D} + \underbrace{\begin{pmatrix} 0 & a_{12} & \dots & a_{1n} \\ & \ddots & \ddots & \vdots \\ & & \ddots & a_{n-1,n} \\ 0 & \dots & & 0 \end{pmatrix}}_{=:U}.$$

With this, we can now define two very classical schemes:

Definition 7.3 (Jacobi method). To solve $Ax = b$ with $A = L + D + U$ let $x^0 \in \mathbb{R}^n$ be an initial guess. We iterate for $k = 1, 2, \dots$

$$x^k = x^{k-1} + D^{-1}(b - Ax^{k-1})$$

or in other words $J := -D^{-1}(L + U)$:

$$x^k = Jx^{k-1} + D^{-1}b.$$

Definition 7.4 (Gauss-Seidel method). To solve $Ax = b$ with $A = L + D + U$ let $x^0 \in \mathbb{R}^n$ be an initial guess. We iterate for $k = 1, 2, \dots$

$$x^k = x^{k-1} + (D + L)^{-1}(b - Ax^{k-1})$$

or in other words $H := -(D + L)^{-1}U$:

$$x^k = Hx^{k-1} + (D + L)^{-1}b.$$

To implement these two schemes, we provide the presentation in index-notation:

Theorem 7.5 (Index-notation of the Jacobi and Gauss-Seidel methods). One step of the Jacobi method and Gauss-Seidel method, respectively, can be carried out in $n^2 + O(n)$ operations. For each step, in index-notation for each entry it holds:

$$x_i^k = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^{k-1} \right), \quad i = 1, \dots, n,$$

i.e., (for the Gauss-Seidel method):

$$x_i^k = \frac{1}{a_{ii}} \left(b_i - \sum_{j < i} a_{ij} x_j^k - \sum_{j > i} a_{ij} x_j^{k-1} \right), \quad i = 1, \dots, n.$$

7.2.2 Gradient descent

An alternative class of methods is based on so-called **descent** or **gradient** directions, which further improve the previously introduced methods. So far, we have:

$$x^{k+1} = x^k + d^k, \quad k = 1, 2, 3, \dots$$

where d^k denotes the **direction** in which we go at each step. For instance:

$$d^k = D^{-1}(b - Ax^k), \quad d^k = (D + L)^{-1}(b - Ax^k)$$

for the Jacobi and Gauss-Seidel methods, respectively. To improve these kind of iterations, we have two possibilities:

- introduce a relaxation (or so-called damping) parameter $\omega^k > 0$ (possibly adapted at each step) such that

$$x^{k+1} = x^k + \omega^k d^k,$$

- and/or to improve the search direction d^k such that we reduce the error as best as possible.

We restrict our attention to positive definite symmetric matrices as they appear naturally many problems. A key point is another view on the problem by regarding it as a minimization problem for which $Ax = b$ is the first-order necessary condition and consequently the sought solution. Imagine for

simplicity that we want to minimize $f(x) = \frac{1}{2}ax^2 - bx$. The first-order necessary condition is nothing else than the derivative $f'(x) = ax - b$. We find a possible minimum via $f'(x) = 0$, namely

$$ax - b = 0 \quad \Rightarrow \quad x = a^{-1}b, \quad \text{if } a \neq 0.$$

That is exactly the same how we would solve a linear matrix system $Ax = b$. By regarding it as a minimum problem we understand better the purpose of our derivations: How does minimizing a function $f(x)$ work in terms of an iteration? Well, we try to minimize f at each step k :

$$f(x^0) > f(x^1) > \dots > f(x^k).$$

This means that the direction d^k (to determine $x^{k+1} = x^k + \omega^k d^k$) should be a descent direction. This idea can be applied for solving linear equation systems. We first define the quadratic form

$$Q(y) = \frac{1}{2}(Ay, y) - (b, y),$$

where (\cdot, \cdot) is the Euclidian scalar product. Then, we can define

Algorithm 7.6 (Descent method - basic idea). *Let $A \in \mathbb{R}^{n \times n}$ be positive definite and $x^0, b \in \mathbb{R}^n$. Then for $k = 0, 1, 2, \dots$*

- Compute d^k ;
- Determine ω^k as minimum of $\omega^k = \operatorname{argmin} Q(x^k + \omega^k d^k)$;
- Update $x^{k+1} = x^k + \omega^k d^k$.

For instance d^k can be determined via the Jacobi or Gauss-Seidel methods.

Another possibility is the gradient method in which we use the gradient to obtain search directions d^k . This brings us to the gradient method:

Algorithm 7.7 (Gradient descent). *Let $A \in \mathbb{R}^{n \times n}$ positive definite and the right hand side $b \in \mathbb{R}^n$. Let the initial guess be $x^0 \in \mathbb{R}^n$ and the initial search direction $d^0 = b - Ax^0$. Then $k = 0, 1, 2, \dots$*

- Compute the vector $r^k = Ad^k$;
- Compute the relaxation

$$\omega^k = \frac{\|d^k\|_2^2}{(r^k, d^k)};$$
- Update the solution vector $x^{k+1} = x^k + \omega^k d^k$;
- Update the search direction vector $d^{k+1} = d^k - \omega^k r^k$.

One can show that the gradient method converges to the solution of the linear equation system $Ax = b$ (see for instance [15]).

7.3 Conjugate gradient method

The previous gradient descent method may converge slowly in most cases because search directions are in general $d^k \not\perp d^{k+2}$. An improvement is the conjugate gradient method in which all directions are conjugate in a certain sense. During this procedure another subspace is constructed for which this method belongs to so-called Krylov subspace methods. The standard textbook with all information is from Saad [16]. In the following, we closely follow to the second author's co-authored book [15][Section 7.8]

7.3.1 Formulation of the CG scheme.

In order to enhance the performance of gradient descent, the conjugate gradient (CG) scheme was developed. Here, the search directions $\{d^0, \dots, d^{k-1}\}$ are pairwise orthogonal. The measure of orthogonality is achieved by using the A scalar product:

$$(Ad^r, d^s) = 0 \quad \forall r \neq s.$$

At step k , we seek the approximation $x^k = x^0 + \sum_{i=0}^{k-1} \alpha_i d^i$ as the minimum of all $\alpha = (\alpha_0, \dots, \alpha_{k-1})$ with respect to $Q(x^k)$:

$$\min_{\alpha \in \mathbb{R}^k} Q \left(x^0 + \sum_{i=0}^{k-1} \alpha_i d^i \right) = \min_{\alpha \in \mathbb{R}^k} \left\{ \frac{1}{2} \left(Ax^0 + \sum_{i=0}^{k-1} \alpha_i Ad^i, x^0 + \sum_{i=0}^{k-1} \alpha_i d^i \right) - \left(b, x^0 + \sum_{i=0}^{k-1} \alpha_i d^i \right) \right\}.$$

The stationary point is given by

$$0 \stackrel{!}{=} \frac{\partial}{\partial \alpha_j} Q(x^k) = \left(Ax^0 + \sum_{i=0}^{k-1} \alpha_i Ad^i, d^j \right) - (b, d^j) = - (b - Ax^k, d^j), \quad j = 0, \dots, k-1.$$

Therefore, the new residual $b - Ax^k$ is perpendicular to all search directions d^j for $j = 0, \dots, k-1$. The resulting linear equation system

$$(b - Ax^k, d^j) = 0 \quad \forall j = 0, \dots, k-1 \tag{7.1}$$

has the feature of "Galerkin orthogonality", similar to FEM schemes.

While constructing the CG method, new search directions should be linearly independent of the current d^j . Otherwise, the space would not become larger and consequently, the approximation cannot be improved.

Definition 7.5 (Krylov space). *We choose an initial approximation $x^0 \in \mathbb{R}^n$ and set $d^0 := b - Ax^0$. The Krylov space $K_k(d^0, A)$ is*

$$K_k(d^0, A) := \text{span}\{d^0, Ad^0, \dots, A^{k-1}d^0\}.$$

Here, A^k means the k -th power of A .

It holds:

Lemma 7.8. *Assume $A^k d^0 \in K_k(d^0, A)$. Then, the solution $x \in \mathbb{R}^n$ of $Ax = b$ is an element of $K_k(d_0, A)$.*

Proof. Let $x^k \in x_0 + K_k(d^0, A)$ be the best approximation, which fulfills the Galerkin equation (7.1). Let $r^k := b - Ax^k$. Since

$$r^k = b - Ax^k = \underbrace{b - Ax^0}_{=d^0} + \underbrace{A(x^0 - x^k)}_{\in K_k(d^0, A)} \in d^0 + AK_k(d^0, A)$$

it holds $r^k \in K_{k+1}(d^0, A)$. Supposing that $K_{k+1}(d^0, A) \subset K_k(d^0, A)$, we obtain $r^k \in K_k(d^0, A)$. The Galerkin equation yields $r^k \perp K_k(d^0, A)$, from which we obtain $r^k = 0$ and $Ax^k = b$. \square

If the CG scheme aborts since it cannot find new search directions, the solution is found. Let us assume that the A -orthogonal search directions $\{d^0, d^1, \dots, d^{k-1}\}$ have been found, then we can compute the next CG approximation using the basis representation $x^k = x^0 + \sum \alpha_i d^i$ and employing the Galerkin equation:

$$\left(b - Ax^0 - \sum_{i=0}^{k-1} \alpha_i Ad^i, d_j \right) = 0 \quad \Rightarrow \quad (b - Ax^0, d^j) = \alpha_j (Ad^j, d^j) \quad \Rightarrow \quad \alpha_j = \frac{(d^0, d^j)}{(Ad^j, d^j)}.$$

The A -orthogonal basis $\{d^0, \dots, d^{k-1}\}$ of the Krylov space $K_k(d^0, A)$ can be computed with the Gram-Schmidt procedure. However, this procedure has a high computational cost; see e.g., [15]. A better procedure is a two-step recursion formula, which is efficient and stable:

Lemma 7.9 (Two-step recursion formula). *Let $A \in \mathbb{R}^{n \times n}$ symmetric positive definite and $x^0 \in \mathbb{R}^n$ and $d^0 := b - Ax^0$. Then, for $k = 1, 2, \dots$, the iteration*

$$r^k := b - Ax^k, \quad \beta_{k-1} := -\frac{(r^k, Ad^{k-1})}{(d^{k-1}, Ad^{k-1})}, \quad d^k := r^k - \beta_{k-1}d^{k-1}$$

constructs an A -orthogonal basis with $(Ad^r, d^s) = 0$ for $r \neq s$. Here x^k in step k defines the new Galerkin solution $(b - Ax^k, d^j) = 0$ for $j = 0, \dots, k-1$.

Proof. See [15]. □

We collect now all ingredients to construct the CG scheme. Let x^0 be an initial guess and $d^0 := b - Ax^0$ the resulting defect. Suppose that $K_k := \text{span}\{d^0, \dots, d^{k-1}\}$ and $x^k \in x^0 + K_k$ and $r^k = b - Ax^k$ have been computed. Then, we can compute the next iterate d^k according to Lemma 7.9:

$$\beta_{k-1} = -\frac{(r^k, Ad^{k-1})}{(d^{k-1}, Ad^{k-1})}, \quad d^k = r^k - \beta_{k-1}d^{k-1}. \quad (7.2)$$

For the new coefficient α_k in $x^{k+1} = x^0 + \sum_{i=0}^k \alpha_i d^i$ holds with testing in the Galerkin equation (7.1) with d^k :

$$\left(\underbrace{b - Ax^0}_{=d^0} - \sum_{i=0}^k \alpha_i Ad^i, d^k \right) = (b - Ax^0, d^k) - \alpha_k (Ad^k, d^k) = (b - Ax^0 + \underbrace{A(x^0 - x^k)}_{\in K_k}, d^k) - \alpha_k (Ad^k, d^k).$$

That is

$$\alpha_k = \frac{(r^k, d^k)}{(Ad^k, d^k)}, \quad x^{k+1} = x^k + \alpha_k d^k. \quad (7.3)$$

This allows to compute the new residual r^{k+1} :

$$r^{k+1} = b - Ax^{k+1} = b - Ax^k - \alpha_k Ad^k = r^k - \alpha_k Ad^k. \quad (7.4)$$

We summarize (7.2 – 7.4) and formulate the classical CG scheme:

Algorithm 7.10. *Let $A \in \mathbb{R}^{n \times n}$ symmetric positive definite and $x^0 \in \mathbb{R}^n$ and $r^0 = d^0 = b - Ax^0$ be given. Iterate for $k = 0, 1, \dots$:*

1. $\alpha_k = \frac{(r^k, d^k)}{(Ad^k, d^k)}$
2. $x^{k+1} = x^k + \alpha_k d^k$
3. $r^{k+1} = r^k - \alpha_k Ad^k$
4. $\beta_k = \frac{(r^{k+1}, Ad^k)}{(d^k, Ad^k)}$
5. $d^{k+1} = r^{k+1} - \beta_k d^k$

Without round-off errors, the CG scheme yields after (at most) n steps the solution of a n -dimensional problem and is in this sense a direct method rather than an iterative scheme. However, in practice for huge n , the CG scheme is usually stopped earlier, yielding an approximate solution.

Proposition 7.11 (CG as a direct method). *Let $x^0 \in \mathbb{R}^n$ be any initial guess. Assuming no round-off errors, the CG scheme terminates after (at most) n steps with $x^n = x$. At each step, we have:*

$$Q(x^k) = \min_{\alpha \in \mathbb{R}} Q(x^{k-1} + \alpha d^{k-1}) = \min_{y \in x^0 + K_k} Q(y)$$

i.e.,

$$\|b - Ax^k\|_{A^{-1}} = \min_{y \in x^0 + K_k} \|b - Ay\|_{A^{-1}}$$

with the norm

$$\|x\|_{A^{-1}} = (A^{-1}x, x)^{\frac{1}{2}}.$$

Proof. That the CG scheme is a direct scheme follows from Lemma 7.8.

The iterate is given by

$$Q(x^k) = \min_{y \in x^0 + K_k} Q(y),$$

with is equivalent to (7.1). The ansatz

$$x^k = x^0 + \sum_{k=0}^{k-1} \alpha_k d^{k-1} = x^0 + \underbrace{y^{k-1}}_{\in K_{t-1}} + \alpha_{t-1} d^{k-1}$$

yields

$$(b - Ax^k, d^j) = (b - Ay^{k-1}, d_j) - \alpha_{t-1} (Ad^{k-1}, d^j) = 0 \quad \forall j = 0, \dots, t-1$$

that is

$$(b - Ay^{k-1}, d_j) = 0 \quad \forall j = 0, \dots, t-2,$$

and therefore $y^{k-1} = x^{k-1}$ and

$$Q(x^k) = \min_{\alpha \in \mathbb{R}} Q(x^{k-1} + \alpha d^{k-1}).$$

Finally, employing symmetry $A = A^T$, we obtain:

$$\begin{aligned} \|b - Ay\|_{A^{-1}}^2 &= (A^{-1}[b - Ay], b - Ay) = (Ay, y) - (A^{-1}b, Ay) - (y, b) \\ &= (Ay, y) - 2(b, y), \end{aligned}$$

i.e., the relationship $\|b - Ay\|_{A^{-1}}^2 = 2Q(y)$. □

Remark 7.3 (The CG scheme as iterative scheme). *As previously mentioned, in practice the CG scheme is (always) used as iterative method rather than a direct method. Due to round-off errors the search directions are never 100% orthogonal.*

7.3.2 Convergence analysis of the CG scheme

We now turn our attention to the convergence analysis, which is a nontrivial task. The key is the following characterization of one iteration $x^k = x^0 + K_k$ by

$$x^k = x^0 + p_{k-1}(A)d^0,$$

where $p_{k-1} \in P_{k-1}$ is a polynomial in A :

$$p_{k-1}(A) = \sum_{i=0}^{k-1} \alpha_i A^i$$

The characterization as minimization of Proposition 7.11 can be written as:

$$\|b - Ax^k\|_{A^{-1}} = \min_{y \in x^0 + K_k} \|b - Ay\|_{A^{-1}} = \min_{q \in P_{k-1}} \|b - Ax^0 - Aq(A)d^0\|_{A^{-1}}.$$

When we employ the $\|\cdot\|_A$ norm, we obtain with $d^0 = b - Ax^0 = A(x - x^0)$

$$\|b - Ax^k\|_{A^{-1}} = \|x - x^k\|_A = \min_{q \in P_{k-1}} \|(x - x^0) - q(A)A(x - x^0)\|_A,$$

that is

$$\|x - x^k\|_A = \min_{q \in P_{k-1}} \|[I - q(A)A](x - x^0)\|_A.$$

In the sense of the best approximation property, we can formulate this task as:

$$p \in P_{k-1} : \|[I - p(A)A](x - x^0)\|_A = \min_{q \in P_{k-1}} \|[I + q(A)A](x - x^0)\|_A. \quad (7.5)$$

The characterization as best approximation is key in the convergence analysis of the CG scheme. Let $q(A)A \in P_k(A)$. We seek a polynomial $q \in P_k$ with $q(0) = 1$, such that

$$\|x^k - x\|_A \leq \min_{q \in P_k, q(0)=1} \|q(A)\|_A \|x - x^0\|_A. \quad (7.6)$$

The convergence of the CG method is related to the fact whether we can construct a polynomial $q \in P_k$ with $p(0) = 1$ such that the A norm is as small as possible. First, we have:

Lemma 7.12 (Bounds for matrix polynomials). *Let $A \in \mathbb{R}^{n \times n}$ symmetric positive definite with the eigenvalues $0 < \lambda_1 \leq \dots \leq \lambda_n$, and $p \in P_k$ a polynomials with $p(0) = 1$:*

$$\|p(A)\|_A \leq M, \quad M := \min_{p \in P_k, p(0)=1} \sup_{\lambda \in [\lambda_1, \lambda_n]} |p(\lambda)|.$$

Proof. See [15]. □

Employing the previous result and the error estimate (7.6), we can now derive a convergence result for the CG scheme.

Proposition 7.13 (Convergence of the CG scheme). *Let $A \in \mathbb{R}^{n \times n}$ be symmetric positive definite. Let $b \in \mathbb{R}^n$ a right hand side vector and let $x^0 \in \mathbb{R}^n$ be an initial guess. Then:*

$$\|x^k - x\|_A \leq 2 \left(\frac{1 - 1/\sqrt{\kappa}}{1 + 1/\sqrt{\kappa}} \right)^k \|x^0 - x\|_A, \quad k \geq 0,$$

with the spectral condition number $\kappa = \text{cond}_2(A)$ of the matrix A .

Remark 7.4. *We see immediately that a large condition number $\kappa \gg 1$ yields*

$$\frac{1 - 1/\sqrt{\kappa}}{1 + 1/\sqrt{\kappa}} \rightarrow 1$$

and deteriorates significantly the convergence rate of the CG scheme. This is the key reason why preconditioners of the form $P^{-1} \approx A^{-1}$ are introduced that re-scale the system; see Section 7.4:

$$\underbrace{P^{-1}A}_{\approx I} x = P^{-1}b.$$

Proof of Prop. 7.13. From the previous lemma and the estimate (7.6) it follows

$$\|x^k - x\|_A \leq M \|x^0 - x\|_A$$

with

$$M = \min_{q \in P_k, q(0)=1} \max_{\lambda \in [\lambda_1, \lambda_n]} |q(\lambda)|.$$

We have to find a sharp estimate of the size of M . That is to say, we seek a polynomial $q \in P_k$ which takes at the origin the value 1, i.e., $q(0) = 1$ and which simultaneously has values near 0 in the maximum norm in the interval $[\lambda_1, \lambda_n]$. To this end, we work with the Tschebyscheff approximation (see e.g., [15] and references therein for the original references). We seek a best approximation $p \in P_k$ of the zero function on $[\lambda_1, \lambda_n]$. Such a polynomial should have the property $p(0) = 1$. For this reason, the trivial solution $p = 0$ is not valid. A Tschebyscheff polynomial reads:

$$T_k = \cos(k \arccos(x))$$

and has the property:

$$2^{-k-1} \max_{[-1,1]} |T_k(x)| = \min_{\alpha_0, \dots, \alpha_{k-1}} \max_{[-1,1]} |x^k + \sum_{i=0}^{k-1} \alpha_i x^i|.$$

We choose now the transformation:

$$x \mapsto \frac{\lambda_n + \lambda_1 - 2t}{\lambda_n - \lambda_1}$$

and obtain with

$$p(t) = T_k \left(\frac{\lambda_n + \lambda_1 - 2t}{\lambda_n - \lambda_1} \right) T_k \left(\frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} \right)^{-1}$$

a polynomial of degree k , which is minimal on $[\lambda_1, \lambda_n]$ and can be normalized by

$$p(0) = 1.$$

It holds:

$$\sup_{t \in [\lambda_1, \lambda_n]} |p(t)| = T_k \left(\frac{\lambda_n + \lambda_1}{\lambda_n - \lambda_1} \right)^{-1} = T_k \left(\frac{\kappa + 1}{\kappa - 1} \right)^{-1} \quad (7.7)$$

with the spectral condition:

$$\kappa := \frac{\lambda_n}{\lambda_1}.$$

We now employ the Tschebyscheff polynomials outside of $[-1, 1]$:

$$T_n(x) = \frac{1}{2} \left[(x + \sqrt{x^2 - 1})^n + (x - \sqrt{x^2 - 1})^n \right].$$

For $x = \frac{\kappa+1}{\kappa-1}$, it holds:

$$\frac{\kappa + 1}{\kappa - 1} + \sqrt{\left(\frac{\kappa + 1}{\kappa - 1} \right)^2 - 1} = \frac{\kappa + 2\sqrt{\kappa} + 1}{\kappa - 1} = \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1}$$

and therefore

$$\frac{\kappa + 1}{\kappa - 1} - \sqrt{\left(\frac{\kappa + 1}{\kappa - 1} \right)^2 - 1} = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}.$$

Using this relationship, we can estimate (7.7):

$$T_k \left(\frac{\kappa + 1}{\kappa - 1} \right) = \frac{1}{2} \left[\left(\frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^k + \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \right] \geq \frac{1}{2} \left(\frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^k.$$

It follows that

$$\sup_{t \in [\lambda_1, \lambda_n]} T_k \left(\frac{\kappa + 1}{\kappa - 1} \right)^{-1} \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k = 2 \left(\frac{1 - \frac{1}{\sqrt{\kappa}}}{1 + \frac{1}{\sqrt{\kappa}}} \right)^k.$$

This finishes the proof. \square

7.4 Preconditioning

Very often in practice the matrix A has a bad condition number for which iterative solvers fail to converge or converge very slowly. Therefore, one can precondition the system with a matrix P^{-1} that has something to do with the inverse of A . In practice left- and right-preconditioning are known, respectively:

$$\begin{aligned} P^{-1}Ax &= P^{-1}b, \\ AP^{-1}u &= b, \quad u = Px. \end{aligned}$$

The key question is how to design P^{-1} . The two extreme cases are

$$\begin{aligned} P^{-1} &= I \\ P^{-1} &= A^{-1} \end{aligned}$$

Therefore, the ‘best’ lies somewhere in between, since $P^{-1} = I$ will not do the job and $P^{-1} = A^{-1}$ means that we actually would construct the inverse of A , which means that we can ‘solve’ the system by Gaussian elimination or LU and do not need any iterative solver.

In the following, we concentrate first on the CG scheme and in the section later, we will also see preconditioned GMRES solvers. For instance, we recall that for second-order operators (such as Laplace) we have a dependence on the mesh size $O(h^{-2}) = O(N)$ (in 2D). For the CG scheme it holds:

$$\rho_{CG} = \frac{1 - \frac{1}{\sqrt{\kappa}}}{1 + \frac{1}{\sqrt{\kappa}}} = 1 - \frac{2}{\sqrt{\kappa}} + O\left(\frac{1}{\kappa}\right).$$

Preconditioning reformulates the original system with the goal of obtaining a moderate condition number for the modified system. Let $P \in \mathbb{P}^{n \times n}$ be a matrix with

$$P = KK^T.$$

Then:

$$Ax = b \quad \Leftrightarrow \quad \underbrace{K^{-1}A(K^T)^{-1}}_{=: \tilde{A}} \underbrace{K^T x}_{=: \tilde{x}} = \underbrace{K^{-1}b}_{=: \tilde{b}},$$

which is

$$\tilde{A}\tilde{x} = \tilde{b}.$$

In the case of

$$\text{cond}_2(\tilde{A}) \ll \text{cond}_2(A)$$

and if the application of K^{-1} is cheap, then the consideration of a preconditioned system $\tilde{A}\tilde{x} = \tilde{b}$ yields a much faster solution of the iterative scheme. The condition $P = KK^T$ is necessary such that the matrix \tilde{A} keeps its symmetry.

The preconditioned CG scheme (PCG) can be formulated as:

Algorithm 7.14. *Let $A \in \mathbb{R}^{n \times n}$ symmetric positive definite and $P = KK^T$ a symmetric preconditioner. Choosing an initial guess $x^0 \in \mathbb{R}^n$ yields:*

1. $r^0 = b - Ax^0$
2. $Pp^0 = r^0$
3. $d^0 = p^0$
4. For $k = 0, 1, \dots$
 - (a) $\alpha_k = \frac{(r^k, d^k)}{(Ad^k, d^k)}$
 - (b) $x^{k+1} = x^k + \alpha_k d^k$
 - (c) $r^{k+1} = r^k - \alpha_k Ad^k$
 - (d) $Pp^{k+1} = r^{k+1}$
 - (e) $\beta_k = \frac{(r^{k+1}, p^{k+1})}{(r^k, g^k)}$
 - (f) $d^{k+1} = p^{k+1} + \beta_k d^k$

At each step, we have as additional cost the application of the preconditioner P . We recall that P allows the decomposition into K and K^T even if they are not explicitly used.

We seek P such that

$$P \approx A^{-1}.$$

On the other hand

$$P \approx I,$$

such that the construction of P is not too costly. Obviously, these are two conflicting requirements. Typical preconditioners are:

- *Jacobi preconditioning*

We choose $P \approx D^{-1}$, where D is the diagonal part of A . It holds

$$D = D^{\frac{1}{2}}(D^{\frac{1}{2}})^T,$$

which means that for $D_{ii} > 0$, this preconditioner is admissible. For the preconditioned matrix, it holds

$$\tilde{A} = D^{-\frac{1}{2}}AD^{-\frac{1}{2}} \Rightarrow \tilde{a}_{ii} = 1$$

- *SSOR preconditioning*

The SSOR scheme is a symmetric variant of the SOR method (successive over-relaxation) and is based on the decomposition:

$$P = (D + \omega L)D^{-1}(D + \omega U) = \underbrace{(D^{\frac{1}{2}} + \omega LD^{-\frac{1}{2}})}_K \underbrace{(D^{\frac{1}{2}} + \omega D^{-\frac{1}{2}}U)}_{=K^T}.$$

For instance, for the Poisson problem, we can find an optimal ω (which is a non-trivial task) such that

$$\text{cond}_2(\tilde{A}) = \sqrt{\text{cond}_2(A)}$$

can be shown. Here, the convergence improves significantly. The number of necessary steps to achieve a given error reduction by a factor of ϵ improves to

$$t_{CG}(\epsilon) = \frac{\log(\epsilon)}{\log(1 - \kappa^{-\frac{1}{2}})} \approx -\frac{\log(\epsilon)}{\sqrt{\kappa}}, \quad \tilde{t}_{CG}(\epsilon) = \frac{\log(\epsilon)}{\log(1 - \kappa^{-\frac{1}{4}})} \approx \frac{\log(\epsilon)}{\sqrt[4]{\kappa}}.$$

Rather than having 100 steps, we only need 10 steps for instance, in case an optimal ω can be found.

7.5 GMRES - generalized minimal residual method

The iterative GMRES algorithm [17, 16] is suited for solving nonsymmetric linear systems. Let $A \in \mathbb{R}^{n \times n}$ be a regular matrix, but not necessarily symmetric. We demonstrate first an option, which turns out to be not feasible when n is large. A symmetric version of the problem

$$Ax = b$$

can be achieved by multiplication with A^T :

$$A^T Ax = A^T b.$$

The matrix $B = A^T A$ is positive definite since

$$(Bx, x)_2 = (A^T Ax, x)_2 = (Ax, Ax)_2 = \|Ax\|_2.$$

In principle, we could now apply the CG scheme to $A^T A$. Instead of one matrix-vector multiplication, we would need two such multiplications per step. However, using $A^T A$, the convergence rate will deteriorate since

$$\kappa(B) = \text{cond}_2(A^T A) = \text{cond}_2(A)^2.$$

For this reason, the CG scheme is not really an option.

Pure GMRES

Let us now prepare the ingredients for the GMRES scheme. The main part is Arnoldi's method [1] in which the (modified) Gram-Schmidt method (see Numerik 1, e.g. [15] for the modified Gram-Schmidt see e.g., [16]) is used for orthonormalization of the basis $\{v_1, \dots, v_k\}$ of the Krylov space $K_k = \text{span}\{v_1, Av_1, \dots, A^{k-1}v_1\}$.

Algorithm 7.15 (Arnoldi). *Choose some initial vector v_1 with $\|v_1\| = 1$. Iterate for $j = 1, 2, \dots$, do*

1. $h_{i,j} = (Av_j, v_i)$, $i = 1, 2, \dots, j$
2. $\hat{v}_{j+1} = Av_j - \sum_{i=1}^j h_{i,j}v_i$
3. $h_{j+1,j} = \|\hat{v}_{j+1}\|$
4. $v_{j+1} = \hat{v}_{j+1}/h_{j+1,j}$

We then define

$$V_k = \{v_1, \dots, v_k\} \in \mathbb{R}^{n \times k}$$

and $H_k = V_k^T AV_k \in \mathbb{R}^{k \times k}$ is the upper Hessenberg matrix, with entries h_{ij} just computed by Arnoldi's algorithm.

Recall the task: Solve

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}$$

Using a Galerkin method with an l_2 -orthonormal basis (produced by Arnoldi's algorithm), we seek x_k of the form

$$x_k = x_0 + z_k$$

where x_0 is some initial guess and z_k comes from the Krylov space $K_k = \text{span}\{r_0, Ar_0, \dots, A^{k-1}r_0\}$ with the usual initial residual $r_0 = b - Ax_0$.

In order to prepare for GMRES, we assume that after k steps of Arnoldi's algorithm, we have

- an orthonormal system $V_{k+1} = \{v_1, \dots, v_{k+1}\}$
- a matrix $\bar{H}_k \in \mathbb{R}^{(k+1) \times k}$ with the non-zero elements h_{ij} .

We have

$$AV_k = V_{k+1}\bar{H}_k. \tag{7.8}$$

With this, we are interested in solving the following least squares problem:

$$\min_{z \in K_k} \|b - A[x_0 + z]\| = \min_{z \in K_k} \|r_0 - Az\|. \tag{7.9}$$

Introducing the variable $z = V_k y$ we obtain the following minimization problem:

$$\min J(y) = \min \|\beta v_1 - AV_k y\|, \quad \beta = \|r_0\|.$$

We notice that (using (7.8))

$$b - Ax = b - A(x_0 + V_k y) = r_0 - AV_k y = \beta v_1 - V_{k+1}\bar{H}_k y = V_{k+1}(\beta e_1 - \bar{H}_k y)$$

This brings us to

$$J(y) = \|V_{k+1}[\beta e_1 - \bar{H}_k y]\|$$

with the unit vector $e_1 \in \mathbb{R}^{(k+1) \times (k+1)}$. Since V_{k+1} is orthonormal we have

$$J(y) = \|\beta e_1 - \bar{H}_k y\|, \quad y \in \mathbb{R}^k.$$

The solution to (7.9) is given by

$$x_k = x_0 + V_k y_k$$

where y_k minimizes $J(y)$. This minimizer y_k is inexpensive since it requires to solve a $(k+1) \times k$ least-squares problem, where k is typically small.

Algorithm 7.16. *GMRES:*

1. Choose the initial guess x_0 . Compute $r_0 = b - Ax_0$ and $v_1 = r_0/\|r_0\|$.
2. Iterate: For $j = 1, 2, \dots, k, \dots$ until some stopping criterion:
 - (a) $h_{i,j} = (Av_j, v_i)_2, \quad i = 1, 2, \dots, j$
 - (b) $\hat{v}_{j+1} = Av_j - \sum_{i=1}^j h_{i,j}v_i$
 - (c) $h_{j+1,j} = \|\hat{v}_{j+1}\|$. If $h_{j+1,j} = 0$ set $k := j$ and go to Step 3 (Hessenberg matrix)
 - (d) $v_{j+1} = \hat{v}_{j+1}/h_{j+1,j}$
3. Define the $(k+1) \times k$ Hessenberg matrix $\bar{H}_k = \{h_{ij}\}_{1 \leq i \leq k+1, 1 \leq j \leq k}$;
4. Compute minimizer y_k of $J(y)$;
5. Set as approximate solution $x_k = x_0 + V_k y_k$.

Unfortunately, for increasing k the storage of vectors becomes an issue and increases like k . Moreover, the arithmetic cost increases as $1/2k^2n$. Consequently, there exists a restarted version, say after m steps, denoted by GMRES(m):

Algorithm 7.17. *GMRES(m):*

1. Choose the initial guess x_0 . Compute $r_0 = b - Ax_0$ and $v_1 = r_0/\|r_0\|$.
2. Iterate: For $j = 1, 2, \dots, m$ do:
 - (a) $h_{i,j} = (Av_j, v_i)_2, \quad i = 1, 2, \dots, j$
 - (b) $\hat{v}_{j+1} = Av_j - \sum_{i=1}^j h_{i,j}v_i$
 - (c) $h_{j+1,j} = \|\hat{v}_{j+1}\|$
 - (d) $v_{j+1} = \hat{v}_{j+1}/h_{j+1,j}$
3. Set as approximate solution $x_m = x_0 + V_m y_k$ where y_m minimizes $\|\beta e_1 - \bar{H}_m y\|$ for $y \in \mathbb{R}^m$.
4. Restart:
 - Compute $r_m = b - Ax_m$; if satisfied, stop
 - else compute $x_0 := x_m, r_1 := r_m/\|r_m\|$ and go to step 2.

Remark 7.5. Of course in both versions, still the minimization problem $J(y)$ must be solved, which can be done by the QR method. For more details, we refer the reader to [16].

Left-preconditioned GMRES

The left-preconditioned scheme reads

$$P^{-1}Ax = P^{-1}b.$$

Then, we have

Algorithm 7.18. *Left-preconditioned GMRES:*

1. Compute $r_0 = P^{-1}(b - Ax_0)$, $\beta = \|r_0\|_2$, and $v_1 = r_0/\beta$
2. For $j = 1, \dots, k$ do
 - (a) Compute $w := P^{-1}Av_j$
 - (b) For $i = 1, \dots, j$ do

- i. $h_{i,j} := (w, v_i)$
- ii. $w := w - h_{i,j}v_i$
- (c) Compute $h_{j+1,j} = \|w\|_2$ and $v_{j+1} = w/h_{j+1,j}$
- 3. Define $V_k := \{v_1, \dots, v_k\}$ and $\bar{H}_k = \{h_{i,j}\}_{1 \leq i \leq j+1, 1 \leq j \leq k}$
- 4. Compute $y_k = \operatorname{argmin}_y \|\beta e_1 - \bar{H}_k y\|_2$
- 5. Set $x_k = x_0 + V_k y_k$
- 6. If satisfied stop (solution found), else set $x_0 := x_k$ and go to Step 1.

Right-preconditioned GMRES

The right-preconditioned scheme reads

$$AP^{-1}u = b, \quad u = Px.$$

Then, we have

Algorithm 7.19. *Right-preconditioned GMRES:*

- 1. Compute $r_0 = b - Ax_0$, $\beta = \|r_0\|_2$, and $v_1 = r_0/\beta$
- 2. For $j = 1, \dots, k$ do
 - (a) Compute $w := AP^{-1}v_j$
 - (b) For $i = 1, \dots, j$ do
 - i. $h_{i,j} := (w, v_i)$
 - ii. $w := w - h_{i,j}v_i$
 - (c) Compute $h_{j+1,j} = \|w\|_2$ and $v_{j+1} = w/h_{j+1,j}$
- 3. Define $V_k := \{v_1, \dots, v_k\}$ and $\bar{H}_k = \{h_{i,j}\}_{1 \leq i \leq j+1, 1 \leq j \leq k}$
- 4. Compute $y_k = \operatorname{argmin}_y \|\beta e_1 - \bar{H}_k y\|_2$
- 5. Set $x_k = x_0 + P^{-1}V_k y_k$
- 6. If satisfied stop (solution found), else set $x_0 := x_k$ and go to Step 1.

The difference to the left-preconditioned version is that the residual norm is taken with respect to the initial system $b - Ax_k$ (without P^{-1}). The reason is that the algorithm obtains the residual implicitly $b - Ax_k = b - AP^{-1}u_k$. Moreover, the right-preconditioned version can be extended to FGMRES (where ‘F’ stands for ‘flexible’) in which the preconditioner P^{-1} can be changed at each iteration step [16].

Chapter 8

Convex sets and systems of linear inequalities

Throughout this chapter, all vector spaces will be considered as real.

8.1 Convex sets: definition and first properties

Definition 8.1. Let X be a (real) vector space. A set $A \subset X$ is said to be convex if

$$\forall (x, y) \in A^2, \forall \theta \in [0, 1], \quad \theta x + (1 - \theta)y \in A.$$

The interpretation is the following. The set $\{\theta x + (1 - \theta)y, \theta \in [0, 1]\}$ is the line segment joining x to y . The definition means that this line segment is contained in A as soon as x and y belong to A , see fig. 8.1.

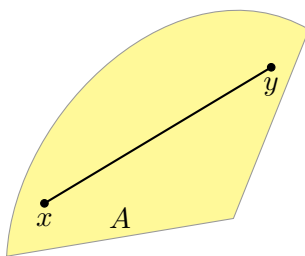


Figure 8.1: Convex set

The following fact is a straightforward consequence of the definition.

Proposition 8.1. An (arbitrary) intersection of convex sets of X is convex.

On the contrary, a union of convex sets is not convex in general, see fig. 8.2.

Due to the triangle inequality:

Proposition 8.2. If X is a normed vector space then all balls (open or closed) of X are convex.

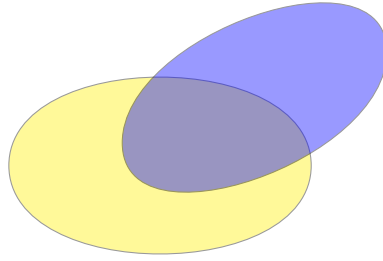


Figure 8.2: Intersection vs union of convex sets

8.2 Convex combinations and convex hull (complement)

8.2.1 Convex combinations

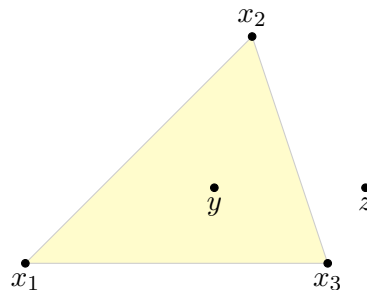
Definition 8.2. Let X be a vector space. A convex combination of m elements $x_1, \dots, x_m \in X$ is an element $y \in X$ that can be decomposed as

$$y = \sum_{i=1}^m \theta_i x_i \quad \text{with } (\theta_1, \dots, \theta_m) \in (\mathbb{R}_+)^m, \sum_{i=1}^m \theta_i = 1.$$

In particular a convex combination is a linear combination, but we ask more. In fig. 8.3, the convex combinations of $\{x_1, x_2, x_3\}$ form the triangle.

Proposition 8.3. If A is convex then any convex combination of elements of A belongs to A .

Proposition 8.3 is illustrated by fig. 8.3: take A as the closed triangle and consider the points x_1, x_2, x_3 .

Figure 8.3: y is a convex combination of $\{x_1, x_2, x_3\}$, but not z .

Theorem 8.4 (Carathéodory). In a vector space of dimension n , all convex combinations of m elements, $m > n + 1$, can be written as a convex combination of at most $n + 1$ of these elements.

Carathéodory's theorem in dimension $n = 2$ is illustrated in fig. 8.4.

8.2.2 Convex hull

Since an intersection of convex sets is a convex set, the following definition makes sense.

Definition 8.3. Let X be a vector space and A be a subset of X . The convex hull of A , denoted by $\text{conv } A$, is the smallest (in the sense of inclusion) convex subset of X containing A . It is the intersection of all convex sets containing A .

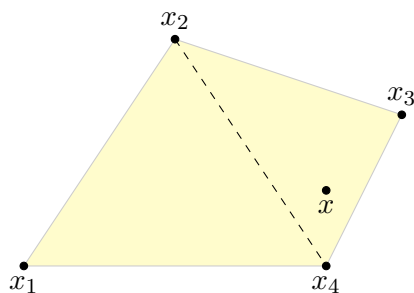


Figure 8.4: Carathéodory's theorem: x is a convex combination of $\{x_1, x_2, x_3, x_4\}$, but it is also a convex combination of $\{x_2, x_3, x_4\}$, or $\{x_1, x_3, x_4\}$.

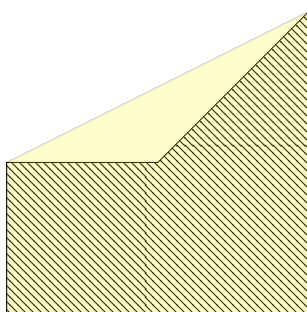


Figure 8.5: A set (hatched) and its convex hull (yellow)

The above construction (see fig. 8.5) can be qualified as *external*. It admits an *internal* counterpart:

Proposition 8.5. *The convex hull of A is the set of all convex combinations of elements of A .*

Proposition 8.5 can be combined with Carathéodory's theorem: if A is a subset of a vector space of dimension n , then the convex hull of A is the set of all convex combinations of at most $n + 1$ elements of A , see fig. 8.6.

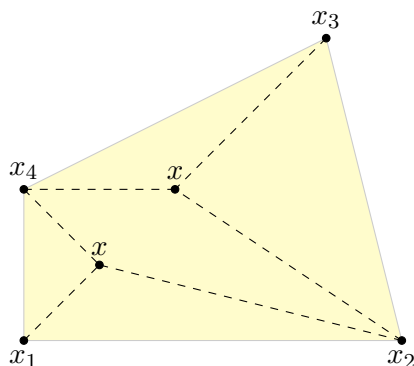


Figure 8.6: The convex hull of $\{x_1, x_2, x_3, x_4\}$ in dimension 2 is the set of all convex combinations of at most 3 vectors taken among $\{x_1, x_2, x_3, x_4\}$.

Proposition 8.6. *Let X, Y be two vector spaces, $A \subset X$, $B \subset Y$. We have*

$$\text{conv}(A \times B) = \text{conv } A \times \text{conv } B.$$

Proposition 8.6 is illustrated in fig. 8.7. A straightforward application is that

$$\text{conv}(\{0, 1\}^n) = (\text{conv}\{0, 1\})^n = [0, 1]^n.$$

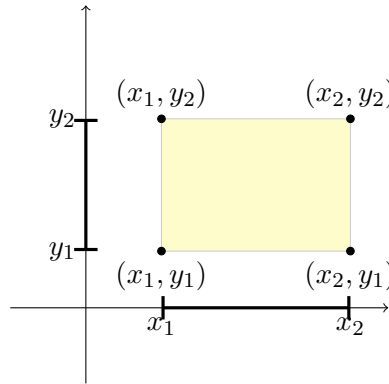


Figure 8.7: The convex hull of $\{x_1, x_2\} \times \{y_1, y_2\} = \{(x_1, y_1), (x_2, y_1), (x_1, y_2), (x_2, y_2)\}$ is $[x_1, x_2] \times [y_1, y_2]$.

8.3 Projection

Theorem 8.7. *Let X be a Hilbert space (typically \mathbb{R}^n equipped with the canonical inner product and the associated Euclidean norm) and A be a nonempty closed convex subset of X . For all $x \in X$ there exists a unique $y \in A$ such that*

$$\|x - y\| = \min_{z \in A} \|x - z\|. \quad (8.1)$$

This element is characterized by the variational inequality (see fig. 8.8)

$$\langle x - y, z - y \rangle \leq 0 \quad \forall z \in A. \quad (8.2)$$

It is called the projection of x onto A , denoted by $y = P_A(x)$.

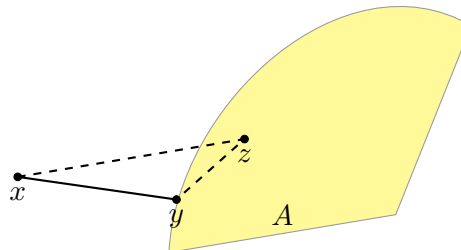


Figure 8.8: Projection onto a convex set

Computing a projection is often difficult (one has to solve an optimization problem), however we discuss below two easy and very useful cases, where the variational inequality (8.2) is easily solved. The first one deals with linear subspaces and corresponds to the standard orthogonal projection, see fig. 8.9..

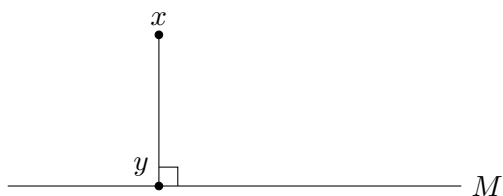


Figure 8.9: Orthogonal projection onto a linear subspace

Proposition 8.8. *Let X be a Hilbert space and M be a closed linear subspace of X (typically $X = \mathbb{R}^n$ and M is a linear subspace of \mathbb{R}^n). The projection $y = P_M(x) \in M$ is characterized by*

$$\langle x - y, v \rangle = 0 \quad \forall v \in M. \quad (8.3)$$

The second case is related to the positive orthant

$$\mathbb{R}_+^n = \{x = (x_1, \dots, x_n) \in \mathbb{R}^n : x_i \geq 0 \forall i\},$$

and corresponds to the positive part.

Proposition 8.9. *It holds for all $x = (x_1, \dots, x_n) \in \mathbb{R}^n$*

$$P_{\mathbb{R}_+^n}(x) = x^+ = (x_1^+, \dots, x_n^+),$$

with $x_i^+ = \max(0, x_i)$.

8.4 Cones

8.4.1 Definitions

Definition 8.4. *A subset C of a vector space X is said to be a cone if*

$$\forall(\alpha, x) \in \mathbb{R}_+ \times C, \alpha x \in C.$$

Cones can be convex or not, see fig. 8.10.

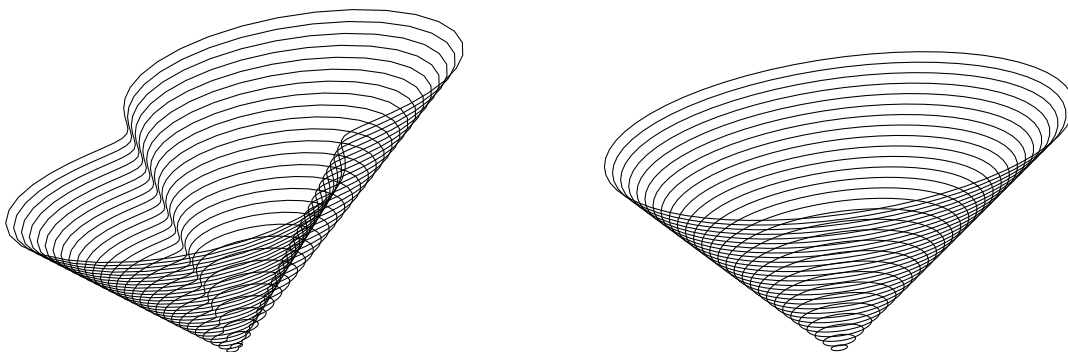


Figure 8.10: Non-convex cone (left) and convex cone (right)

The positive orthant \mathbb{R}_+^n is an example of convex cone.

8.4.2 Normal cone (complement)

Definition 8.5. Let A be a convex set of a Hilbert space X and let $a \in A$. The normal cone of A at point a is the set

$$N_A(a) = \{x \in X \text{ s.t. } \langle x, y - a \rangle \leq 0 \forall y \in A\}.$$

It generalizes the notion of outward normal, see Fig. 8.11. Clearly, if a is in the interior of A then $N_A(a)$ is reduced to $\{0\}$.

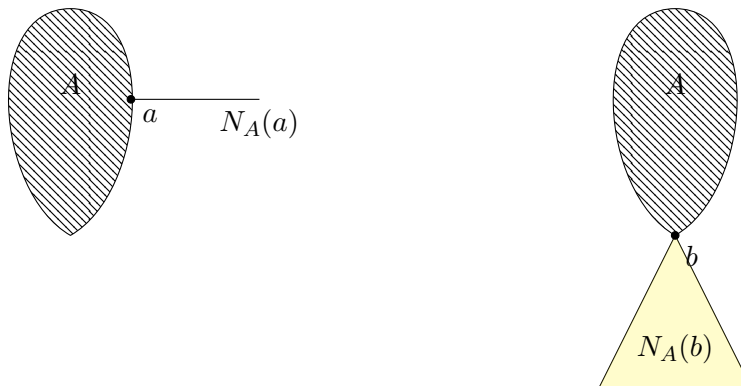


Figure 8.11: Normal cones

8.4.3 Polar cones (complement)

Definition 8.6. Let C be a cone of a Hilbert space X , with inner product $\langle \cdot, \cdot \rangle$. The positive and negative polar (or dual) cones of C are respectively defined by (see fig. 8.12)

$$C^+ = \{x \in X \text{ s.t. } \langle x, y \rangle \geq 0 \forall y \in C\}, \quad C^- = \{x \in X \text{ s.t. } \langle x, y \rangle \leq 0 \forall y \in C\}.$$

Of course those two cones have similar properties. The negative polar cone is somehow more standard, we will simply call it polar cone. It is clear from the definition that a polar cone is always closed and convex.

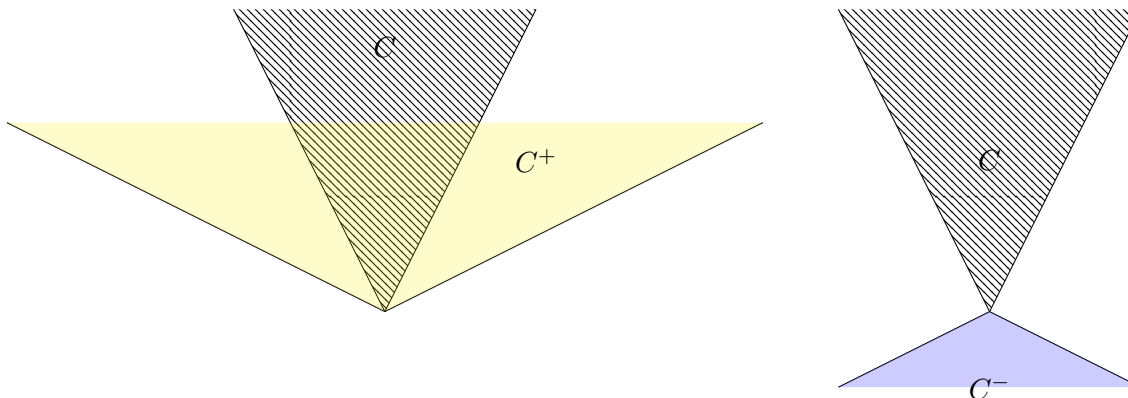


Figure 8.12: Positive and negative polar cones

Example 8.1. For the positive orthant it is easily checked that

$$(\mathbb{R}_+^n)^+ = \mathbb{R}_+^n, \quad (\mathbb{R}_+^n)^- = \mathbb{R}_-^n.$$

Let us now consider reiterated polarity. The definition shows that $C \subset C^{--}$ for every cone C , and even that $\overline{C} \subset C^{--}$ since C^{--} is closed. The following theorem deals with the other inclusion.

Theorem 8.10. *If C is a convex cone of a Hilbert space X then*

$$C^{--} = \overline{C}.$$

8.5 Systems of linear inequalities, introduction to linear programming

8.5.1 Linear inequality systems: Fourier-Motzkin elimination

Here we discuss a practical way of solving linear inequality systems of form $Ax \leq b$, with unknown $x \in \mathbb{R}^n$ and matrix $A \in \mathbb{R}^{m \times n}$. The inequality is meant componentwise. This kind of problem is often known as a feasibility problem: the inequalities are interpreted as constraints and the questions are

1. to know whether it is possible to satisfy all these constraints simultaneously;
2. to describe the feasible set, when it is not empty.

Let

$$C = \{x \in \mathbb{R}^n \mid Ax \leq b\}.$$

As an intersection of m affine half-spaces (corresponding to each inequality), the set C is convex. Specifically it is called a convex polyhedron. When $b = 0$, C is clearly a cone and we call this a polyhedral cone.

One possibility to compute feasible solutions (in case they exist) is the Fourier-Motzkin algorithm, which is a basic method.

The Fourier-Motzkin elimination works like Gaussian elimination: eliminate from n -variables to $n - 1$ -variables in an equivalent fashion. This procedure is then recursively applied until 1-variable. Here, we can compute all numbers between the given bounds. By backward substitution we then compute the values of the other variables.

Example 8.2. *Consider the system*

$$\begin{cases} x_1 + x_2 \leq 1 \\ 2x_1 - x_2 \leq 4 \\ x_1 + 2x_2 \leq 8. \end{cases} \quad (8.4)$$

In order to eliminate x_2 we write the equivalent system

$$\begin{cases} x_2 \leq 1 - x_1 \\ x_2 \geq 2x_1 - 4 \\ x_2 \leq 4 - \frac{1}{2}x_1 \end{cases} \iff 2x_1 - 4 \leq x_2 \leq \min(1 - x_1, 4 - \frac{1}{2}x_1).$$

This is possible if and only if

$$\begin{cases} 2x_1 - 4 \leq 1 - x_1 \\ 2x_1 - 4 \leq 4 - \frac{1}{2}x_1 \end{cases} \iff x_1 \leq \frac{5}{3}.$$

We obtain the (nonempty) set of solutions

$$\left\{ (x_1, x_2) \mid x_1 \leq \frac{5}{3}, 2x_1 - 4 \leq x_2 \leq \min(1 - x_1, 4 - \frac{1}{2}x_1) \right\}.$$

See Fig. 8.13.

Remark 8.1. *If we now consider the 'positively homogeneous' system*

$$\begin{cases} x_1 + x_2 \leq 0 \\ 2x_1 - x_2 \leq 0, \end{cases}$$

then we get the set of solutions

$$\{(x_1, x_2) \mid x_1 \leq 0, 2x_1 \leq x_2 \leq -x_1\}.$$

This is clearly a cone of \mathbb{R}^2 , as expected from the structure of the system: a solution remains a solution when it is multiplied by a nonnegative number. See Fig. 8.14

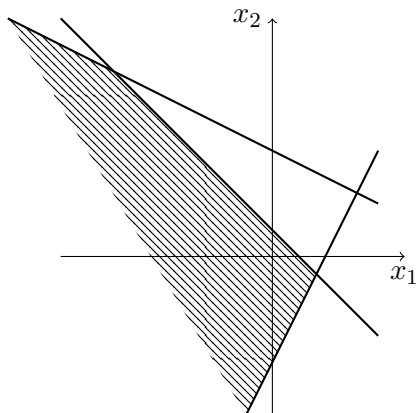


Figure 8.13: Set of solutions in example 8.2 (it is not bounded on the left side).

8.5.2 Conical hull and Farkas lemma (complement)

Here are a few useful results.

Lemma 8.11. *Let $a_1, \dots, a_m \in \mathbb{R}^n$ be given vectors. Then, the set*

$$\text{cone}(a_1, \dots, a_m) := \{x_1 a_1 + \dots + x_m a_m \mid x_i \geq 0 \forall i = 1, \dots, m\}$$

is a convex cone in \mathbb{R}^n . It is the conical hull of a_1, \dots, a_m . The vectors a_1, \dots, a_m are generators of the cone.

In the example from Remark 8.1, the vectors $(-1, 1)$ and $(-1, -2)$ are generators of the set of solutions.

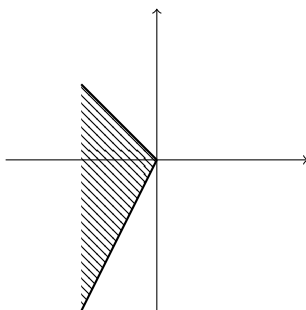


Figure 8.14: Cone of solutions in Remark 8.1

Lemma 8.12. *Let $A \in \mathbb{R}^{m \times n}$. Then, the set*

$$K := \{y \in \mathbb{R}^n \mid y = A^T x, x \geq 0\}$$

(the inequality is meant component-wise) is a non-empty, closed and convex cone.

The following is a consequence of Theorem 8.10.

Lemma 8.13 (Farkas). *Let $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ be given. Then, the following two statements are equivalent:*

1. *The system $A^T x = b, x \geq 0$ has a solution.*
2. *The inequality $b^T d \geq 0$ holds true for all $d \in \mathbb{R}^n$ with $Ad \geq 0$.*

8.5.3 Introduction to linear programming (complement)

Linear programming deals with optimization problems with linear objective function and linear (or affine) constraints. We provide a brief introduction following [11].

Definition 8.7 (Standard form). *A linear program in standard form is:*

$$\text{minimize } c^T x, \quad \text{s.t. } Ax = b, x \geq 0,$$

where $c \in \mathbb{R}^n, x \in \mathbb{R}^n, b \in \mathbb{R}^m$ and $A \in \mathbb{R}^{m \times n}$.

In the above problem the function $x \mapsto c^T x$ is the objective function, the expression $Ax = b$ is the equality constraint and the expression $x \geq 0$ is the inequality constraint.

Other problem statements can be transformed into this form. For instance

$$\text{minimize } c^T x, \quad \text{s.t. } Ax \leq b,$$

is an example in which the constraint $Ax \leq b$ represents a linear inequality system. By using the so-called *slack variable* z inequality constraints can be converted into equalities. Then, we obtain

$$\text{minimize } c^T x, \quad \text{s.t. } Ax + z = b, z \geq 0.$$

However, we do not yet deal with the standard form, because x has no sign. To this end, we must split the vector x into nonnegative and nonpositive parts:

$$x = x^+ - x^-$$

with $x^+ = \max(x, 0)$ and $x^- = \max(-x, 0)$. With these notations, we finally arrive at the following standard form:

$$\text{minimize } \begin{bmatrix} c \\ -c \\ 0 \end{bmatrix}^T \begin{bmatrix} x^+ \\ x^- \\ z \end{bmatrix} \quad \text{s.t. } [A \quad -A \quad I] \begin{bmatrix} x^+ \\ x^- \\ z \end{bmatrix} = b, \quad \begin{bmatrix} x^+ \\ x^- \\ z \end{bmatrix} \geq 0.$$

Remark 8.2. *Further inequality constraints of the form $x \leq u$ or $Ax \geq 0$ can be converted by similar derivations while introducing again slack variables.*

Chapter 9

Exercises

Linear systems

Exercise 9.1. Solve the following linear systems using the Gauss elimination method.

$$(a) \begin{cases} 8x - 2y + z = 3 \\ 12x + y - 2z = 1 \\ 96x - 16y + 5z = 29, \end{cases}$$

$$(b) \begin{cases} 2x - 2y + z - t + u = 1 \\ x + 2y - z + t - 2u = 1 \\ 4x - 10y + 5z - 5t + 7u = 1 \\ 2x - 14y + 7z - 7t + 11u = -1, \end{cases}$$

$$(c) \begin{cases} 2x_1 + 3x_2 + 5x_3 = 3 \\ x_1 - 3x_2 - 6x_3 = 2 \\ 3x_1 - x_3 = 4, \end{cases}$$

$$(d) \begin{cases} 2x_1 + 3x_2 - x_3 = -1 \\ x_1 - 5x_2 + x_3 = 6 \\ 3x_1 + 3x_2 + 7x_3 = 0 \\ x_1 + 2x_2 - 4x_3 = -1. \end{cases}$$

Exercise 9.2. We consider the system

$$\begin{cases} 10^{-4}x + y = 1 \\ -x + y = 2. \end{cases}$$

Solve this system in floating point arithmetic with 3 digits (for the mantissa, in basis 10)

1. without pivoting,
2. with row pivoting.

Conclusion?

Matrix algebra

Exercise 9.3. Compute the matrix multiplication AB with

$$A = \begin{pmatrix} 3 & 1 & 2 \\ 4 & 1 & 1 \\ 5 & 2 & 2 \end{pmatrix} \quad \text{and} \quad B = \begin{pmatrix} 7 & 2 & 3 \\ 2 & 2 & 2 \\ 4 & 5 & 6 \end{pmatrix}.$$

Exercise 9.4. Compute the matrix multiplication AB with

$$A = \begin{pmatrix} 1 & -1 & 2 \\ 3 & -2 & 4 \end{pmatrix} \quad \text{and} \quad B = \begin{pmatrix} 1 & 2 & 11 & 4 \\ -2 & 3 & 0 & 2 \\ 3 & 1 & 4 & 0 \end{pmatrix}.$$

Exercise 9.5. Consider the matrix

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

- 1) Calculate $B \in M_3(\mathbb{R})$ such that $A = 2I + B$.
- 2) Calculate B^2 and B^3 .
- 3) Deduce an expression of A^n for all $n \in \mathbb{N}$.

Exercise 9.6. Let m be a real number. When it is possible calculate the inverse of the matrix

$$A = \begin{pmatrix} 1 & m & -2 \\ 1 & m+1 & m-2 \\ 2 & 2m+1 & 2m-4 \end{pmatrix}.$$

Exercise 9.7. Consider the matrices

$$P = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, \quad A = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}.$$

- 1) Calculate P^{-1} .
- 2) Calculate $B = P^{-1}AP$.
- 3) Show that $A^n = PB^nP^{-1}$ for all $n \in \mathbb{N}$. Deduce an expression of A^n .

Linear spaces

Exercise 9.8. Let (e_1, e_2, e_3) be the canonical basis of \mathbb{R}^3 . Let a_1, a_2, a_3 be three vectors of \mathbb{R}^3 defined by $a_1 = (0, -2, 3)$, $a_2 = (1, 2, 1)$, $a_3 = (3, 0, -4)$.

- 1) Show that (a_1, a_2, a_3) is a basis of \mathbb{R}^3 .
- 2) Let u be the vector of coordinates $(1, 1, 1)$ in the basis (a_1, a_2, a_3) . What are the coordinates of u in the basis (e_1, e_2, e_3) ?
- 3) Let v be the vector of coordinates $(1, 1, 1)$ in the basis (e_1, e_2, e_3) . What are the coordinates of v in the basis (a_1, a_2, a_3) ?

Determinants

Exercise 9.9. Calculate the following determinants and give a necessary and sufficient condition on the real numbers a, b, c for d_2 being equal to 0:

$$d_1 = \begin{vmatrix} 1 & 2 & 3 \\ 2 & 3 & 1 \\ 3 & 1 & 2 \end{vmatrix}, \quad d_2 = \begin{vmatrix} 1 & 1 & 1 \\ a & b & c \\ a^2 & b^2 & c^2 \end{vmatrix}.$$

Exercise 9.10. Calculate the determinant of order n :

$$\Delta_n = \begin{vmatrix} 1 & 1 & \dots & \dots & 1 \\ 1 & 1 & 0 & \dots & 0 \\ \vdots & 0 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & 1 & 0 \\ 1 & 0 & \dots & 0 & 1 \end{vmatrix}.$$

Exercise 9.11. Let $a \neq b$ be two real numbers. Calculate the determinant of order n

$$\Delta_n = \begin{vmatrix} a+b & ab & 0 & \dots & 0 \\ 1 & a+b & ab & \ddots & \vdots \\ 0 & 1 & a+b & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & ab \\ 0 & \dots & 0 & 1 & a+b \end{vmatrix}.$$

Indication: show by induction that $\Delta_n = \frac{a^{n+1} - b^{n+1}}{a-b}$.

LU decomposition (complement)

Exercise 9.12. Let the matrix A and the right hand side vector b given as follows:

$$A = \begin{pmatrix} 2 & 1 & 0 \\ 4 & 3 & 2 \\ 0 & 1 & 1 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 5 \\ 2 \end{pmatrix}.$$

1. Construct the LU decomposition of the matrix A .
2. Solve $Ax = b$ with the LU decomposition.
3. Use the Jacobi and Gauss-Seidel iterative methods and evaluate by hand the first two steps. Take as initial guess vector $x^0 = (0, 0, 0)$.

Diagonalization

Exercise 9.13. Let

$$A = \begin{pmatrix} 1 & 2 \\ -1 & 4 \end{pmatrix}.$$

- 1) Determine the eigenvalues of A and show that A is diagonalizable on \mathbb{R} .
- 2) Diagonalize A (give the diagonal form and the transition matrix).
- 3) Calculate A^n for all $n \in \mathbb{N}$.

Exercise 9.14. 1) Calculate the complex eigenvalues of the matrices

$$A = \begin{pmatrix} 0 & 1 & -6 \\ 1 & 0 & 2 \\ 1 & 0 & 2 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{pmatrix}, \quad C = \begin{pmatrix} 2 & 1 & 0 \\ 0 & 1 & -1 \\ 0 & 2 & 4 \end{pmatrix}.$$

- 2) Show that A is diagonalizable on \mathbb{C} and B is diagonalizable on \mathbb{R} .
- 3) What are the spectral radii of A , B and C ?

Normed spaces

Exercise 9.15 (complement). We define on $\mathcal{M}_n(\mathbb{R})$ the map

$$A = (a_{ij}) \mapsto N(A) = n \max_{ij} |a_{ij}|.$$

Show that this defines a matrix norm.

Exercise 9.16 (complement). Let $A \in \mathcal{M}_n(\mathbb{R})$ be symmetric positive definite. We define the map

$$(x, y) \in \mathbb{R}^n \times \mathbb{R}^n \mapsto \langle x, y \rangle_A = \langle Ax, y \rangle,$$

where $\langle \cdot, \cdot \rangle$ is the canonical inner product of \mathbb{R}^n .

1. Show that this defines an inner product on \mathbb{R}^n .

2. We denote by $\|\cdot\|_A$ the associated norm. Give tight constants α, β , in terms of the eigenvalues of A , such that

$$\alpha\|x\|_2 \leq \|x\|_A \leq \beta\|x\|_2 \quad \forall x \in \mathbb{R}^n.$$

Exercise 9.17. We consider the space $\mathcal{C}([0, 1], \mathbb{R})$ of continuous functions from $[0, 1]$ into \mathbb{R} . Show that the quantities

$$\|f\|_1 = \int_0^1 |f(t)| dt, \quad \|f\|_2 = \sqrt{\int_0^1 |f(t)|^2 dt}, \quad \|f\|_\infty = \sup_{t \in [0, 1]} |f(t)|$$

define norms such that

$$\|f\|_1 \leq \|f\|_2 \leq \|f\|_\infty.$$

Remarks. 1. One can show that these norms are nevertheless not equivalent.

2. Knowing how to prove the triangle inequality for $\|\cdot\|_2$ is not mandatory.

Stability of linear systems and iterative methods (complement)

Exercise 9.18. We perform a stability analysis of the following linear system. Let

$$A = \begin{pmatrix} 6 & -2 \\ 11.5 & -3.85 \end{pmatrix}, \quad b = \begin{pmatrix} 10 \\ 17 \end{pmatrix} \quad \text{and} \quad \tilde{A} = \begin{pmatrix} 6 & -2 \\ 11.5 & -3.84 \end{pmatrix}.$$

In this exercise you may **perform the calculations using a computer** and a software such as Python, Matlab, or Octave.

1. Compute the solutions x and \tilde{x} of the systems

$$Ax = b, \quad \tilde{A}\tilde{x} = b.$$

2. What do you observe by comparing x and \tilde{x} and comparing the matrices A and \tilde{A} ?

3. We know the stability theorem:

$$\frac{\|x - \tilde{x}\|}{\|\tilde{x}\|} \leq \text{cond}(A) \frac{\|A - \tilde{A}\|}{\|A\|}. \quad (9.1)$$

Compute the factors

$$\text{cond}_2(A) = \|A^{-1}\|_2 \|A\|_2$$

and

$$\frac{\|A - \tilde{A}\|_2}{\|A\|_2}.$$

4. Compute the right hand side of the estimate (9.1) and interpret the result in view of the relative error of the solutions x and \tilde{x} , i.e.,

$$\frac{\|x - \tilde{x}\|_2}{\|\tilde{x}\|_2}.$$

5. Is the lastly obtained solution plausible?

6. **(optional alternative by hand)** You may replace the spectral norm by the maximal row sum norm, which is defined as $\|A\|_\infty = \max_{i=1,\dots,n} \sum_{j=1}^n |a_{ij}|$ for $A \in \mathbb{R}^{n \times n}$. Then, you can re-do steps No. 1-5 again by hand. You may finally also compare the results using the $\|\cdot\|_2$ and $\|\cdot\|_\infty$.

Exercise 9.19. Implement (in a software) the descent method, Jacobi and Gauss-Seidel approaches for solving

$$A = \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{pmatrix}, \quad b = \begin{pmatrix} 2 \\ -3 \\ 4 \end{pmatrix}$$

As initial guess take $x^0 = 0$ and compute for each method 10 steps. As information, the exact solution is $x = (1, 0, 2)^T$.

Convex sets

Exercise 9.20. Show that the simplex

$$S_n = \left\{ x = (x_1, \dots, x_n) \in \mathbb{R}_+^n \text{ s.t. } \sum_{i=1}^n x_i = 1 \right\},$$

(see fig. 9.1) is a convex set.

Complement: show that $S_n = \text{conv}\{e_1, \dots, e_n\}$, where (e_1, \dots, e_n) is the canonical basis of \mathbb{R}^n .

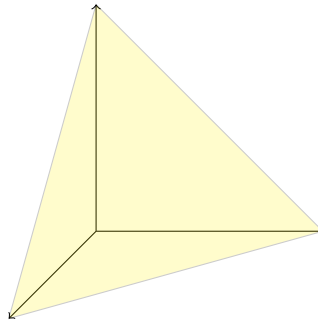


Figure 9.1: Simplex of \mathbb{R}^3

Exercise 9.21. Let v_1, \dots, v_n be arbitrary vectors of a vector space. Show that

$$K := \left\{ \sum_{i=1}^n \alpha_i v_i, \alpha_i \geq 0 \forall i \right\}.$$

is a convex cone.

Complement: Show that

$$K = \text{conv} \left(\bigcup_{i=1}^n h_i \right)$$

with the half-lines

$$h_i = \mathbb{R}_+ v_i = \{t v_i, t \geq 0\}.$$

Exercise 9.22. Let

$$M = \left\{ x = (x_1, \dots, x_n) \in \mathbb{R}^n : \sum_{i=1}^n x_i = 0 \right\}.$$

Describe the projection operator onto M .

Linear inequalities

Exercise 9.23. *Using the Fourier-Motzkin elimination, find out whether these systems admit solutions:*

$$(a) \begin{cases} x_1 + 2x_2 \leq 4 \\ -2x_1 - x_2 \leq 2 \\ x_1 - x_2 \leq 1, \end{cases} \quad (b) \begin{cases} x_1 + 2x_2 \leq 2 \\ -2x_1 - x_2 \leq 4 \\ 4x_1 + x_2 \leq -16. \end{cases}$$

Part III

Functions and related numerical notions

Chapter 10

Functions of one or several variables

10.1 Basic concepts

A real-valued or complex-valued function on a set U is a mapping f which associates to each $x \in U$ a unique number $f(x) \in \mathbb{K}$, $\mathbb{K} = \mathbb{R}$ or $\mathbb{K} = \mathbb{C}$. We write

$$\begin{aligned} f : U &\rightarrow \mathbb{K} \\ x &\mapsto f(x). \end{aligned}$$

The set U is the definition set and $f(U) := \{f(x) \in \mathbb{K} : x \in U\}$ is the image set. The graph of f is the set

$$G(f) := \{(x, f(x)) : x \in U\} \subset U \times \mathbb{K}.$$

We will mainly work in the situation where U is a subset of \mathbb{R}^d , hence the variable is represented by a point in \mathbb{R}^d , denoted by

$$x = (x_1, \dots, x_d),$$

actually there are d independent space variables. For instance they can be space and time variables. The canonical inner product in \mathbb{R}^d is denoted by $\langle x, y \rangle = x \cdot y = \sum_{i=1}^d x_i y_i$.

We can similarly consider vector-valued functions, namely functions taking values in an arbitrary vector space. Actually, studying

$$\left\{ \begin{array}{l} f : U \rightarrow \mathbb{R}^m \text{ (or } \mathbb{C}^m) \\ x \mapsto f(x) = (f_1(x), \dots, f_m(x)) \end{array} \right.$$

amounts to studying the m functions f_1, \dots, f_m . This is why the number of variables deserves much more mathematical care than the number of values, at least when this number is finite. We also insist on the fact that we will only consider real variables. Dealing with complex variables is a completely different story!

10.2 Differentiation

10.2.1 Partial derivatives

Let f be a (scalar, vector or matrix-valued) function of a variable $x = (x_1, \dots, x_d) \in \mathbb{R}^d$. We recall that the partial derivative of f with respect to (w.r.t.) one of these variables is obtained after differentiating f (i.e. each component or entry of f) with respect to this variable, the other variables being fixed. We frequently use the notations:

$$\frac{\partial f}{\partial x_i} = \partial_{x_i} f \quad (\text{partial derivative w.r.t. } x_i),$$

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \partial_{x_i x_j} f \quad (\text{second partial derivative w.r.t. } x_j \text{ then } x_i),$$

$$\frac{\partial^2 f}{\partial x_i \partial x_i} = \partial_{x_i}^2 f \quad (\text{second partial derivative w.r.t. } x_i).$$

10.2.2 Fréchet derivative, Jacobian matrix

There are several notions of differentiability, here we will merely use the most classical notion of Fréchet differentiability. This concept can be defined for functions from an arbitrary real vector space to an arbitrary real or complex vector space, however for simplicity we restrict ourselves here to functions from \mathbb{R}^n to \mathbb{R}^m .

Definition 10.1. Let $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $x \in U$. We say that f is Fréchet differentiable at point x if there exists $L \in \mathcal{M}_{mn}(\mathbb{R})$ such that

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

In this case the matrix L is unique, it is called the Jacobian matrix of f at x , denoted by $Df(x)$ or $df(x)$ or $Jf(x)$. The linear map $y \in \mathbb{R}^n \mapsto Ly \in \mathbb{R}^m$ is the Fréchet derivative of f at point x .

If the map $x \in U \mapsto Df(x)$ is continuous we say that f is continuously Fréchet differentiable.

Note that due to the equivalence of norms (Theorem 6.3), any norm on \mathbb{R}^m and \mathbb{R}^n can be indifferently chosen.

The Jacobian matrix can be inferred from the partial derivatives.

Proposition 10.1. 1. If $f : x = (x_1, \dots, x_n) \in U \subset \mathbb{R}^n \mapsto f(x) = (f_1(x), \dots, f_m(x)) \in \mathbb{R}^m$ is Fréchet differentiable at $x \in U$ then

$$Df(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(x) & \dots & \frac{\partial f_1}{\partial x_n}(x) \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1}(x) & \dots & \frac{\partial f_m}{\partial x_n}(x) \end{pmatrix}.$$

2. If all partial derivatives of f exist and are continuous functions on U , then f is continuously Fréchet differentiable.

10.2.3 Chain rule

Next, we introduce the chain rule for differentiating compositions of functions. Essentially, it says that the Jacobian matrix of a composition is the product of the Jacobian matrices.

Theorem 10.2. Let $f : U \subset \mathbb{R}^n \rightarrow V \subset \mathbb{R}^m$ and $g : V \rightarrow \mathbb{R}^p$ be Fréchet differentiable in a and $b = f(a)$, respectively. Set $h = g \circ f$, i.e. $h(x) = g(f(x))$. Then h is Fréchet differentiable in a with

$$Dh(a) = Dg(b)Df(a).$$

Here is a special case.

Corollary 10.3. Let $f : I \subset \mathbb{R} \rightarrow V \subset \mathbb{R}^m$ be differentiable in a and $g : V \rightarrow \mathbb{R}$ be differentiable in $b = f(a)$. Then, $h := g \circ f$ is differentiable in a with the derivative

$$h'(a) = Dg(b)f'(a) = \sum_{i=1}^m \partial_i g(b) f'_i(a).$$

For instance in physics, if a quantity depends on the temperature and the pressure, say $g = g(T, P)$, while T, P both depend on the space position x (say in 1d for simplicity), then we write $h(x) = g(T(x), P(x))$ and

$$\frac{dh}{dx} = \frac{\partial g}{\partial T} \frac{dT}{dx} + \frac{\partial g}{\partial P} \frac{dP}{dx}.$$

10.2.4 Schwarz' theorem

We now provide a result for second-order partial derivatives and when they can be interchanged.

Theorem 10.4 (Schwarz). *Let the function f have in a neighborhood of $a \in \mathbb{R}^n$ the partial derivatives $\partial_i f$, $\partial_j f$ and $\partial_{ji} f$. Let $\partial_{ji} f$ be continuous in a . Then, there exists $\partial_{ij} f$ and it holds*

$$\partial_{ij} f(a) = \partial_{ji} f(a).$$

Example 10.1. *Given $f(x, y) = x^3 y + e^{xy^2}$. Compute $\partial_x f$, $\partial_y f$, $\partial_{xx} f$, $\partial_{xy} f$, $\partial_{yx} f$, $\partial_{yy} f$. For instance, for the mixed second-order derivatives it holds*

$$f_{xy} = f_{yx} = 3x^2 + 2ye^{xy^2}(1 + xy^2).$$

This shows that Schwarz' theorem holds true in this example.

10.2.5 Multiindex notation (complement)

For a general description of ODEs and PDEs the multiindex notation is commonly used.

- A multiindex is a vector $\alpha = (\alpha_1, \dots, \alpha_n)$, where each component $\alpha_i \in \mathbb{N}_0$. The order is

$$|\alpha| = \alpha_1 + \dots + \alpha_n.$$

- For a given multiindex we define the partial derivative:

$$D^\alpha u(x) := \partial_{x_1}^{\alpha_1} \dots \partial_{x_n}^{\alpha_n} u.$$

If $k \in \mathbb{N}_0$, we define the set of all partial derivatives of order k :

$$D^k u(x) := \{D^\alpha u(x) : |\alpha| = k\}.$$

Example 10.2. *Let the problem dimension $n = 3$. Then, $\alpha = (\alpha_1, \alpha_2, \alpha_3)$. For instance, let $\alpha = (2, 0, 1)$. Then $|\alpha| = 3$ and $D^\alpha u(x) = \partial_x^2 \partial_z^1 u(x)$.*

10.3 Classical differential operators

Well-known in physics, it is convenient to work with the **nabla-operator** to define derivative expressions:

$$\nabla = \begin{pmatrix} \partial_{x_1} \\ \vdots \\ \partial_{x_d} \end{pmatrix}.$$

Very often in applications, when space and time variables are involved, the classical differential operators apply to the space variables only.

The gradient of a single-valued function is the vector of partial derivatives

$$\nabla u = \begin{pmatrix} \partial_{x_1} u \\ \vdots \\ \partial_{x_d} u \end{pmatrix}.$$

The gradient of a vector-valued function $v : \mathbb{R}^d \rightarrow \mathbb{R}^n$ reads:

$$\nabla v = \begin{pmatrix} \partial_{x_1} v_1 & \dots & \partial_{x_d} v_1 \\ \vdots & & \vdots \\ \partial_{x_1} v_n & \dots & \partial_{x_d} v_n \end{pmatrix}.$$

It corresponds exactly with the Jacobian matrix! Note the (universal but a little confusing) fact that for single-valued functions the gradient is the transpose of the Jacobian matrix (a row matrix).

The divergence is defined for vector-valued functions $v : \mathbb{R}^d \rightarrow \mathbb{R}^d$ by

$$\operatorname{div} v := \nabla \cdot v := \nabla \cdot \begin{pmatrix} v_1 \\ \vdots \\ v_d \end{pmatrix} = \sum_{k=1}^d \partial_{x_k} v_k.$$

The divergence for a matrix-valued function $\sigma : \mathbb{R}^d \rightarrow \mathcal{M}_d(\mathbb{R})$ is defined as:

$$\nabla \cdot \sigma = \left(\sum_{j=1}^d \frac{\partial \sigma_{ij}}{\partial x_j} \right)_{1 \leq i \leq d}.$$

The Laplace operator of a scalar-valued function $u : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as

$$\Delta u = \sum_{k=1}^d \partial_{x_k x_k} u.$$

For a vector-valued function $u : \mathbb{R}^d \rightarrow \mathbb{R}^n$, we define the Laplace operator component-wise as

$$\Delta u = \Delta \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} = \begin{pmatrix} \sum_{k=1}^d \partial_{x_k x_k} u_1 \\ \vdots \\ \sum_{k=1}^d \partial_{x_k x_k} u_n \end{pmatrix}.$$

Let us recall the **cross product** of two vectors $u, v \in \mathbb{R}^3$:

$$\begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} \times \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = \begin{pmatrix} u_2 v_3 - u_3 v_2 \\ u_3 v_1 - u_1 v_3 \\ u_1 v_2 - u_2 v_1 \end{pmatrix}.$$

With the help of the cross product, we can define the **rotational** of $v : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ by

$$\operatorname{curl} v = \nabla \times v = \begin{pmatrix} \partial_{x_1} \\ \partial_{x_2} \\ \partial_{x_3} \end{pmatrix} \times \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = \begin{pmatrix} \partial_{x_2} v_3 - \partial_{x_3} v_2 \\ \partial_{x_3} v_1 - \partial_{x_1} v_3 \\ \partial_{x_1} v_2 - \partial_{x_2} v_1 \end{pmatrix}.$$

10.4 Taylor expansions, Hessian matrix

If $f : \mathbb{R} \rightarrow \mathbb{R}$ is of class \mathcal{C}^n around some point a then we have the classical Taylor-Young formula:

$$f(x) = \sum_{j=0}^n \frac{f^{(j)}(a)}{j!} (x-a)^j + o(|x-a|^n).$$

The remainder $o(|x-a|^n)$ (Landau notation) is a function such that

$$\lim_{x \rightarrow a} \frac{o(|x-a|^n)}{|x-a|^n} = 0.$$

There is also the Taylor-Lagrange formula, if f is of class \mathcal{C}^{n+1} :

$$f(x) = \sum_{j=0}^n \frac{f^{(j)}(a)}{j!} (x-a)^j + \frac{f^{(n+1)}(y)}{(n+1)!} (x-a)^{n+1},$$

for some intermediate point y between a and x .

For functions of several variables we limit ourselves to the Taylor-Young formula of order 2.

Proposition 10.5. Let $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}$ be a C^2 function. Let $a, x \in U$ points such that the line segment joining a and x is in U . Then

$$f(x) = f(a) + Df(a)(x - a) + \frac{1}{2}(x - a)^T D^2 f(a)(x - a) + o(\|x - a\|^2),$$

with $Df(a)$ the Jacobian (row) matrix of f at a and $D^2 f(a)$ (or $\nabla^2 f(a)$, or $f''(a)$) the Hessian (square) matrix of f at a given by

$$D^2 f(a) = \begin{pmatrix} \partial_{x_1 x_1} f(a) & \dots & \partial_{x_1 x_n} f(a) \\ \vdots & & \vdots \\ \partial_{x_n x_1} f(a) & \dots & \partial_{x_n x_n} f(a) \end{pmatrix}.$$

Remark 10.1. 1. Under the assumptions made the Hessian matrix is symmetric (Schwarz theorem 10.4).

2. The first order term can be rewritten as

$$Df(a)(x - a) = \langle \nabla f(a), x - a \rangle.$$

10.5 Convex functions

10.5.1 Definitions

A geometrical way to define the convexity of a function is through the notion of epigraph. This point of view is very useful in convex analysis, a branch of mathematics related to optimization. Again we place ourselves in the framework of real vector spaces, mostly \mathbb{R}^n .

Definition 10.2. Let X be a vector space and $f : X \rightarrow \mathbb{R}$. The epigraph of f is the set (see fig. 10.1)

$$\text{epi } f = \{(x, \xi) \in X \times \mathbb{R} \text{ s.t. } \xi \geq f(x)\}.$$

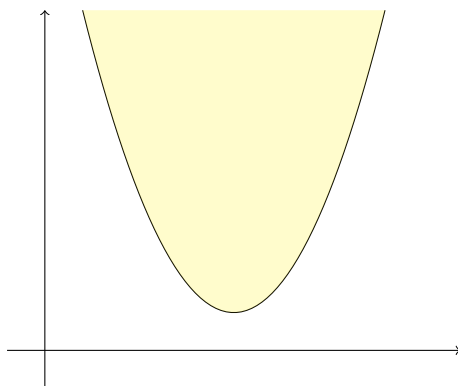


Figure 10.1: Epigraph

Definition 10.3. Let X be a vector space and $f : X \rightarrow \mathbb{R}$. The function f is said to be convex if its epigraph is convex.

The following characterization is extremely standard, see fig. 10.2.

Proposition 10.6. A function $f : X \rightarrow \mathbb{R}$ is convex if and only if

$$\forall x, y \in X, \forall \theta \in (0, 1), \quad f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y). \quad (10.1)$$

A little refinement is the notion of strict convexity.

Definition 10.4. A function $f : X \rightarrow \mathbb{R}$ is strictly convex if and only if

$$\forall x, y \in X, x \neq y, \forall \theta \in (0, 1), \quad f(\theta x + (1 - \theta)y) < \theta f(x) + (1 - \theta)f(y).$$

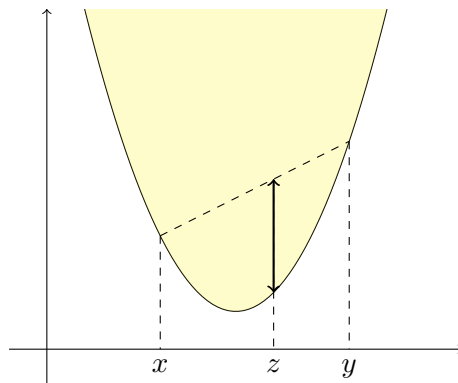


Figure 10.2: Algebraic characterization of convex functions: if $z = \theta x + (1 - \theta)y$ for $0 < \theta < 1$ then $f(z) \leq \theta f(x) + (1 - \theta)f(y)$.

10.5.2 Characterization

Proposition 10.7. *Let $f : X = \mathbb{R}^n \rightarrow \mathbb{R}$ be a Fréchet differentiable function. The following assertions are equivalent:*

- (i) f is convex;
- (ii) $f(y) - f(x) \geq \langle \nabla f(x), y - x \rangle \quad \forall (x, y) \in X^2$;
- (iii) $\langle \nabla f(y) - \nabla f(x), y - x \rangle \geq 0 \quad \forall (x, y) \in X^2$;
- (iv) $D^2 f(x) \geq 0$ (positive semi-definite) $\quad \forall x \in X$, when f is \mathcal{C}^2 .

10.6 Integration by parts

10.6.1 Substitution rule

One of the most important formulas in calculus, numerical mathematics, calculus of variations, and continuum mechanics is the substitution rule that allows to transform integrals from one domain to another.

In 1D

Proposition 10.8. *Let $I = [a, b]$ be given. To transform this interval to a new interval, we use a mapping $T : [a, b] \rightarrow [\alpha, \beta]$ with $T(a) = \alpha$ and $T(b) = \beta$. If T is of class \mathcal{C}^1 we have the transformation rule:*

$$\int_{\alpha}^{\beta} f(y) dy = \int_{T(a)}^{T(b)} f(y) dy = \int_a^b f(T(x)) T'(x) dx.$$

Remark 10.2. *In case that $T(a) = \beta$ and $T(b) = \alpha$ the previous proposition still holds true, but with a negative sign:*

$$\int_{\alpha}^{\beta} f(y) dy = \int_{T(b)}^{T(a)} f(y) dy = - \int_{T(a)}^{T(b)} f(y) dy = \int_a^b f(T(x)) (-T'(x)) dx.$$

Both cases can be summarized as:

Theorem 10.9. *Let $I = [a, b]$ be given. To transform this interval to a new interval $[\alpha, \beta]$, we employ a monotone mapping T such that $T([a, b]) = [\alpha, \beta]$. If T is of class \mathcal{C}^1 , it holds:*

$$\int_{T(I)} f(y) dy := \int_{\alpha}^{\beta} f(y) dy = \int_a^b f(T(x)) |T'(x)| dx = \int_I f(T(x)) |T'(x)| dx.$$

Remark 10.3. We observe the formal relation between the integration increments:

$$dy = |T'(x)| dx.$$

In higher dimensions

We have the following generalization of the substitution rule in arbitrary dimension (also known as **change of variables** under the integral):

Theorem 10.10. Let $\Omega \subset \mathbb{R}^n$ be an open, measurable, domain. Let the function $T : \Omega \rightarrow \mathbb{R}$ be of class C^1 , one-to-one (injective) and Lipschitz continuous. Then:

- The domain $\hat{\Omega} := T(\Omega)$ is measurable.
- The function $f(T(\cdot))|\det T'(\cdot)| : \Omega \rightarrow \mathbb{R}$ is (Riemann)-integrable.
- For all measurable subdomains $M \subset \Omega$ it holds the substitution rule:

$$\int_{T(M)} f(y) dy = \int_M f(T(x))|\det T'(x)| dx,$$

and in particular as well for $M = \Omega$.

Remark 10.4. The determinant of the Jacobian matrix is sometimes called *Jacobian determinant*. The formal relation between increments then reads $dy = |\det T'(x)|dx$. This formula is usually stated and proved in the framework of measure theory, thus it involves particular notions of regularity. Here and in the next sections we do not specify the regularity assumptions, which are anyway very mild.

10.6.2 Integration by parts and Green's formulae

Let $\Omega \subset \mathbb{R}^d$ be a bounded domain with 'smooth' boundary $\partial\Omega$. We denote by n its outward unit normal vector. One of the most important formula in applied mathematics, physics and continuum mechanics is **integration by parts**, which for $u, v : \mathbb{R}^d \rightarrow \mathbb{R}$ reads:

$$\int_{\Omega} \nabla u v dx = - \int_{\Omega} u \nabla v dx + \int_{\partial\Omega} u v n ds.$$

We can obtain further results, which are very useful, directly based on the integration by parts. For instance the following one is at the basis of the finite element method:

$$\int_{\Omega} \nabla u \cdot \nabla v dx = - \int_{\Omega} \Delta u v dx + \int_{\partial\Omega} v \partial_n u ds.$$

Here, $\partial_n u := \nabla u \cdot n$ is called normal derivative. Very much related and well-known in physics is the divergence formula, valid for $u : \mathbb{R}^d \rightarrow \mathbb{R}^d$:

$$\int_{\Omega} \operatorname{div} u dx = \int_{\partial\Omega} u \cdot n ds.$$

Chapter 11

Introduction to nonlinear optimization

11.1 General concepts

11.1.1 Brief classification of optimization problems

Optimization aims at finding points that minimize or maximize a given function $f : U \rightarrow V$. We distinguish between single objective optimization when $V = \mathbb{R}$ and multi-objective optimization when V is a more general vector space equipped with an order relation. We will place ourselves in the first class, and up to a change of sign we will focus on minimization problems. As to the admissible (or feasible) set U , it is usually either a "smooth enough" subset of a vector space (*continuous optimization*), or a countable set (*discrete optimization*). We focus on the first category. Further, we speak of *unconstrained optimization* if U is the full vector space, and of *constrained optimization* otherwise. We mainly develop the first one. In the case where f is a convex function and U is a convex set we are in the field of convex optimization.

11.1.2 Problem setting

As said above, we consider an unconstrained continuous optimization problem of form

$$\underset{x \in X}{\text{minimize}} \quad f(x) \tag{11.1}$$

where $f : X \rightarrow \mathbb{R}$ is an arbitrary function and X is a normed vector space. This function is usually called *cost function*, or *objective function*, or *criterion*. Solving (11.1) not only means finding the value of the minimum (or infimum), but also finding minimizers, if there are some, namely points where the minimum is attained. In some cases, which will be specified, we will assume that $X = \mathbb{R}^n$.

11.1.3 Global and local minimizers

Definition 11.1. We say that $a \in X$ is a global minimizer of $f : X \rightarrow \mathbb{R}$ if

$$f(a) \leq f(x) \quad \forall x \in X.$$

Definition 11.2. We say that $a \in X$ is a local minimizer of $f : X \rightarrow \mathbb{R}$ if there exists $\varepsilon > 0$ such that

$$f(a) \leq f(x) \quad \forall x \in B_\varepsilon(a),$$

where $B_\varepsilon(a)$ is the (open) ball of center a and radius ε . In words, a minimizes f over a neighborhood of a . An illustration is given in fig. 11.1.

Obviously the problem stated in (11.1) deals with global minima. However, the conditions of optimality we are going to write as well as the output of most algorithms do not allow to distinguish between global and local minima. Finding global minima is often a very difficult task, requiring

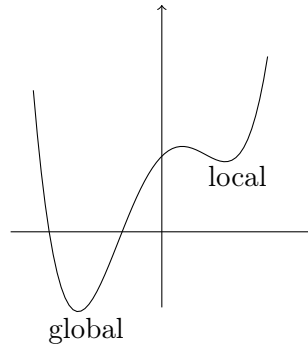


Figure 11.1: Global vs local minimizer

extremely costly numerical methods. We will see later that these bothering facts disappear under convexity assumptions.

More finely, when in definitions 11.1 and 11.2 the inequalities can be made strict for $x \neq a$ we say that a is a strict global / local minimizer. In Figure 11.2 we plot a function which has all these different minima:

$$f(x) = \begin{cases} \sin(4x) + \sin(x) + \frac{1}{4}x^2 & x < x_* \\ \sin(4x_*) + \sin(x_*) + \frac{1}{4}x_*^2 & x_* \leq x \leq x_* + \frac{1}{2} \\ \sin(4(x - \frac{1}{4})) + \sin(x - \frac{1}{4}) + \frac{1}{4}(x - \frac{1}{4})^2 & x > x_* \end{cases}$$

with $x_* \approx 2.72086$.

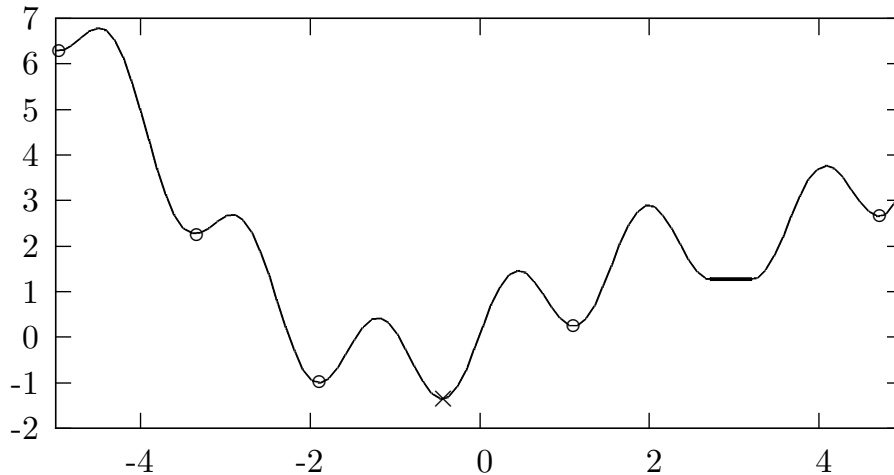


Figure 11.2: Function with different minima. With \circ , we denote strict local minima, with \times the global (strict) minimum and the fat region characterizes a non-strict local minimum.

11.1.4 Optimality conditions

For the construction and the justification of numerical algorithms, we introduce the optimality conditions. They are essentially consequences of the Taylor-Young formula (Proposition 10.5). Here it is assumed that $X = \mathbb{R}^n$.

Theorem 11.1 (Necessary optimality conditions of first and second order).

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuously differentiable and $x \in \mathbb{R}^n$ be local minimizer of f . Then it holds the first-order necessary condition

$$\nabla f(x) = 0.$$

We say that x is a critical or stationary point.

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a two times continuously differentiable function. Let $x \in \mathbb{R}^n$ be a local minimizer of f . Then it holds the necessary second order optimality condition:

$$D^2f(x) \geq 0$$

(the Hessian is positive semi-definite).

We next state sufficient optimality conditions:

Theorem 11.2 (Sufficient optimality condition of second order). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a two times continuously differentiable function and $x \in \mathbb{R}^n$. If $\nabla f(x) = 0$ (critical point) and $D^2f(x) > 0$ (the Hessian is positive definite) then x is a strict local minimizer.*

Definition 11.3. *A stationary point which is neither a local minimizer nor a local maximizer is called local saddle point.*

Recall that, given a symmetric matrix $A \in \mathbb{R}^{n \times n}$, it holds

- $A > 0 \Leftrightarrow$ all eigenvalues are > 0
- $A < 0 \Leftrightarrow$ all eigenvalues are < 0
- $A \geq 0 \Leftrightarrow$ all eigenvalues are ≥ 0
- $A \leq 0 \Leftrightarrow$ all eigenvalues are ≤ 0 .

In dimension is $n = 2$ we have the following criterion. Let

$$A = \begin{pmatrix} a & b \\ b & c \end{pmatrix}.$$

The two eigenvalues λ_1, λ_2 satisfy $\lambda_1 \lambda_2 = \det A$ and $\lambda_1 + \lambda_2 = \operatorname{tr} A$, hence:

- $A > 0 \Leftrightarrow \det(A) > 0$ and $\operatorname{tr} A > 0$
- $A < 0 \Leftrightarrow \det(A) > 0$ and $\operatorname{tr} A < 0$
- $A \geq 0 \Leftrightarrow \det(A) \geq 0$ and $\operatorname{tr} A \geq 0$
- $A \leq 0 \Leftrightarrow \det(A) \geq 0$ and $\operatorname{tr} A \leq 0$.

11.2 Convex case (complement)

11.2.1 Specific aspects of convex optimization

A major feature of convexity is that it makes local minimizers automatically global ones.

Proposition 11.3. *Let $f : X \rightarrow \mathbb{R}$ be convex. Every local minimizer of f is a global minimizer.*

Proposition 11.4. *Let $f : X \rightarrow \mathbb{R}$ be strictly convex. If f admits a global minimizer, then it is unique.*

Another outstanding property of convex functions is that the stationarity condition becomes sufficient. We will state this next in a more general form, with possible convex constraints.

11.2.2 Towards constrained optimization

Proposition 11.5. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be Fréchet differentiable and convex and \mathcal{U} be a closed convex subset of \mathbb{R}^n . Then a point $a \in \mathcal{U}$ is a global minimized of f over \mathcal{U} if and only if*

$$-\nabla f(a) \in N_{\mathcal{U}}(a),$$

where $N_{\mathcal{U}}(a)$ is the normal cone of \mathcal{U} in a .

Of course, if $\mathcal{U} = \mathbb{R}^n$, then $N_{\mathcal{U}}(a) = \{0\}$ and the optimality condition boils down to the stationarity $\nabla f(a) = 0$.

But in general, describing the normal cone is likely to be a highly non-trivial task. This is typically achieved through Lagrange multipliers and needs to be studied in specialized courses.

11.3 Steepest descent method

11.3.1 Descent methods

A general algorithm for the minimization of a continuous function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is:

Listing 11.1: **Descent for minimization of continuous functions**

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ a continuous function.

Choose an initial guess $x^0 \in \mathbb{R}^n$

| For | k | from | $1, 2, \dots$

 Determine descent direction $d^k \in \mathbb{R}^n$

 Determine a step length $s_k \in \mathbb{R}$

 Compute new approximation

$$x^k = x^{k-1} + s_k d^k$$

 with

$$f(x^k) < f(x^{k-1})$$

This is a general abstract scheme, which must be specified for the search direction $d^k \in \mathbb{R}^n$ and the step length $s_k \in \mathbb{R}$. The iteration number in x^k and d^k is here written as a superscript to avoid confusion with the vectors components. A descent direction is a vector d^k such that

$$f(x^{k-1} + s d^k) < f(x^{k-1})$$

as soon as $s > 0$ is small enough. If f is continuously differentiable, it holds the Taylor expansion

$$f(x^{k-1} + s d^k) = f(x^{k-1}) + s \nabla f(x^{k-1}) \cdot d^k + o(s).$$

We infer the (sufficient) condition for d^k to be a descent direction:

$$\nabla f(x^{k-1}) \cdot d^k < 0. \tag{11.2}$$

11.3.2 Steepest descent

A steepest descent direction at a point x is a direction d of unit norm that minimizes the expression $\nabla f(x) \cdot d$. Whereas it heavily depends on the chosen norm, for the Euclidean norm the following holds.

Theorem 11.6 (Optimal descent direction). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function. Let $x \in \mathbb{R}^n$ with $\nabla f(x) \neq 0$. The direction of steepest descent in x is uniquely determined as*

$$d = -\frac{\nabla f(x)}{\|\nabla f(x)\|}.$$

Proof. It holds for arbitrary $d \in \mathbb{R}^n$ with $\|d\| = 1$

$$\nabla f(x) \cdot d \geq -\|\nabla f(x)\| \|d\| = -\|\nabla f(x)\|.$$

In the case $d = -\nabla f(x)/\|\nabla f(x)\|$ it holds an equality

$$\nabla f(x) \cdot d = -\|\nabla f(x)\|.$$

Therefore, $-\nabla f(x)/\|\nabla f(x)\|$ characterizes a unique normalized direction of steepest descent. \square

With this, we obtain in a natural way the *steepest descent or gradient descent* scheme, which is a descent algorithm with optimal descent direction:

Listing 11.2: Minimization with gradient descent Let $f : X \rightarrow \mathbb{R}$ continuously differentiable.

Choose $x^0 \in X$

| For | k | from | $1, 2, \dots$

| Abort | , if $\|\nabla f(x^{k-1})\| = 0$

Determine descent direction $d^k = -\nabla f(x^{k-1})/\|\nabla f(x^{k-1})\|$

Choose a step length $s_k \in \mathbb{R}$

Determine new approximation

$$x^k = x^{k-1} + s_k d^k$$

with

$$f(x^k) < f(x^{k-1})$$

11.3.3 Line search

We have now a more specific scheme, but the step length is still an open question. Let us recall descent schemes for solving linear equation systems. Given a symmetric positive semidefinite matrix A , solving

$$Ax = b$$

is equivalent to the minimization of the quadratic functional (just double-check yourself that $\nabla Q(x) = Ax - b$)

$$Q(x) = \frac{1}{2}(Ax, x) - (b, x).$$

For such a quadratic function, the optimal step length $s_k \in \mathbb{R}$ with given direction d^k is the minimum of the scalar-valued function

$$h(s) := Q(x^{k-1} + sd^k),$$

which is just obtained through a second order algebraic equation. For a general function, this procedure is not as simple. The step length is rather determined by another numerical procedure known as line search, also used in Newton-type methods (e.g., damped Newton method). The general idea is to make trials in an organized way.

A classical line search technique is the Armijo rule (see e.g., [11]).

Listing 11.3: Armijo step length rule

Let $\beta, \gamma \in (0, 1)$ be given. We denote the previous step by $x^{k-1} \in \mathbb{R}^n$ and by $d^k \in \mathbb{R}^n$ the chosen descent direction.

Initialize $s_k^{(0)}$

| For | l | from | $0, 1, 2, \dots$

| Abort | , if $f(x^{k-1} + s_k^{(l)} d^k) < f(x^{k-1}) + \gamma s_k^{(l)} \nabla f(x^{k-1}) \cdot d^k$

$$s_k^{(l+1)} = \beta s_k^{(l)}$$

k	x^k	y^k	$\ \nabla f(x^k)\ $	$f(x^k)$	k	x^k	y^k	$\ \nabla f(x^k)\ $	$f(x^k)$
0	-1.	1.5	356					...	
1	-1.212	1.3675	4.9958	9.5597	100	1.0050	1.0100	2.55e-05	0.012493
2	-1.0898	1.3940	4.7929	6.3382	101	1.0048	1.0101	2.52e-05	0.012492
3	-1.1373	1.3533	4.6039	1.9585	102	1.0049	1.0099	2.48e-05	0.012491
4	-1.0384	1.2768	4.5494	5.7600	103	1.0047	1.0099	2.44e-05	0.012491
5	-1.0836	1.2337	4.3770	1.9861	104	1.0049	1.0097	2.41e-05	0.012490
		...			105	1.0046	1.0098	2.37e-05	0.012489

Table 11.1: Results for the approximation of the Rosenbrock function. The table is taken from [15].

The Armijo rule ensures that we really approach the minimum. One typically chooses $\beta = \gamma = 1/2$. It can be shown that this algorithm terminates after a finite number of steps and always finds a suitable step length:

Theorem 11.7 (Armijo step length rule). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuously differentiable. Let $\gamma \in (0, 1)$ and $d \in \mathbb{R}^n$ a descent direction at a point $x \in X$, just satisfying*

$$\nabla f(x) \cdot d < 0.$$

For any $s_0 > 0$ there exists $m \in \mathbb{N}$ such that

$$f(x + \beta^m s_0 d) - f(x) \leq \gamma \beta^m s_0 \nabla f(x) \cdot d.$$

Proof. It holds

$$\frac{f(x + sd) - f(x)}{s} - \gamma \nabla f(x) \cdot d \rightarrow (1 - \gamma) \nabla f(x) \cdot d < 0 \quad (s \rightarrow 0).$$

It follows from the continuity of ∇f that there exists an $\epsilon > 0$ such that for all $s < \epsilon$:

$$\frac{f(x + sd) - f(x)}{s} - \gamma \nabla f(x) \cdot d \leq 0 \quad \forall |s| < \epsilon.$$

Given $\beta \in (0, 1)$ we find a minimal $m \in \mathbb{N}$ such that $s = \beta^m s_0 < \epsilon$ holds true. \square

The gradient descent scheme converges slowly, and is therefore not very efficient (see our numerical concepts at the beginning of these lecture notes). However it is very robust and easy to implement, hence it is often used for complicated or high dimensional problems where more sophisticated methods are hardly applicable.

In Figure 11.3, we show the convergence behavior with the help for the so-called Rosenbrock function. This example is a typical test case in numerical optimization:

$$\min_{x, y \in \mathbb{R}} \{f(x, y) = 10(y - x^2)^2 + (1 - x)^2\}.$$

We choose as initial guesses $x^0 = (-1, 1.5)$ and $x^0 = (2, -2)$. For the first test case, we summarize the convergence behavior in Table 11.1.

11.4 Gradient descent with projection (complement)

Here we discuss the problem:

$$\underset{x \in \mathcal{U}}{\text{minimize}} f(x) \tag{11.3}$$

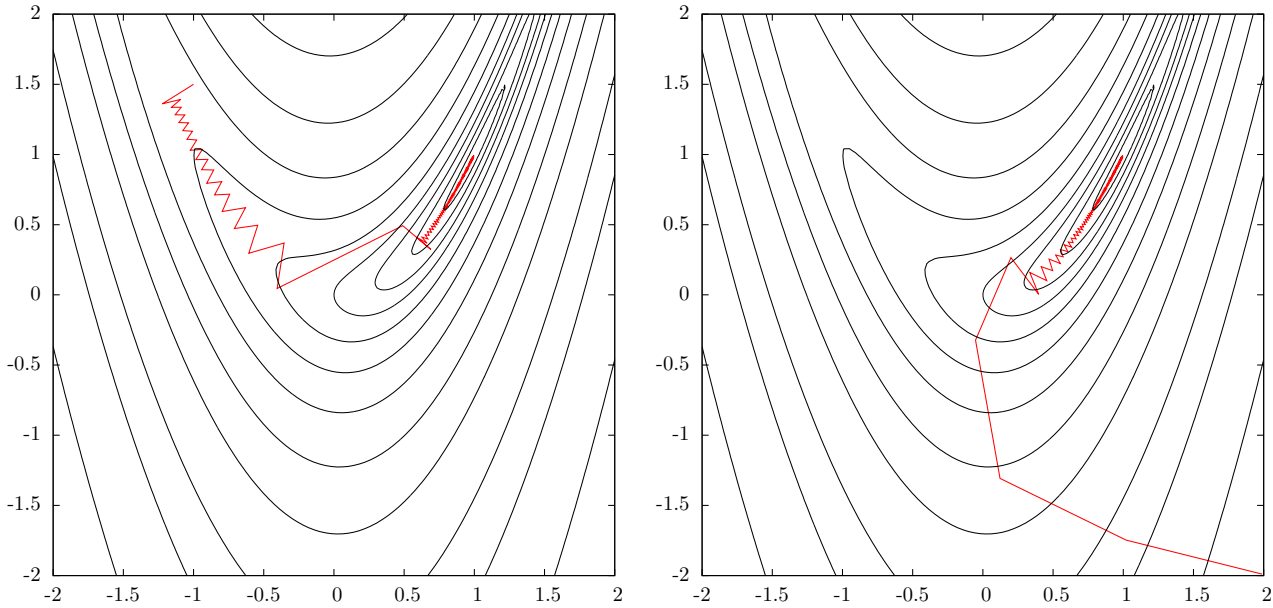


Figure 11.3: Convergence behavior of the gradient scheme for the approximation of the minimum of the Rosenbrock function $f(x, y) = 10(y - x^2)^2 + (1 - x)^2$ with different initial values. Figure taken from [15].

where \mathcal{U} is a (closed) convex subset of \mathbb{R}^n , and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a two-times continuously differentiable function. For clarity and simplicity we also assume that f is convex, but this assumption is not crucial. One way to cope with the constraint $x \in \mathcal{U}$ is to replace the iteration $x^k = x^{k-1} + s_k d^k$ by

$$x^k = P_{\mathcal{U}}(x^{k-1} + s_k d^k),$$

where $P_{\mathcal{U}}$ is the projection operator onto \mathcal{U} , see section 8.3. This method is appropriate if we have at disposal an efficient means to compute the projection. In this case we speak of simple constraints, examples of which have been discussed in section 8.3.

In order to guaranty the descent property

$$f(P_{\mathcal{U}}(x^{k-1} + s d^k)) < f(x^{k-1}) \text{ for small } s > 0,$$

the condition (11.2) turns out not to be sufficient. Luckily:

Proposition 11.8. *The direction $d^k = -\nabla f(x^{k-1})/\|\nabla f(x^{k-1})\|$ is a descent direction, unless x^{k-1} is a minimizer.*

Proof. We have for any $s > 0$

$$f(P_{\mathcal{U}}(x^{k-1} + s d^k)) - f(x^{k-1}) = \nabla f(x^{k-1}) \cdot (P_{\mathcal{U}}(x^{k-1} + s d^k) - x^{k-1}) + O(\|P_{\mathcal{U}}(x^{k-1} + s d^k) - x^{k-1}\|^2).$$

Setting $y_s = x^{k-1} + s d^k$ the first order term rewrites as

$$\begin{aligned} \nabla f(x^{k-1}) \cdot (P_{\mathcal{U}}(y_s) - x^{k-1}) &= -\frac{\|\nabla f(x^{k-1})\|}{s} (y_s - x^{k-1}) \cdot (P_{\mathcal{U}}(y_s) - x^{k-1}) \\ &= -\frac{\|\nabla f(x^{k-1})\|}{s} \left((y_s - P_{\mathcal{U}}(y_s)) \cdot (P_{\mathcal{U}}(y_s) - x^{k-1}) + \|P_{\mathcal{U}}(y_s) - x^{k-1}\|^2 \right) \\ &\leq -\frac{\|\nabla f(x^{k-1})\|}{s} \|P_{\mathcal{U}}(y_s) - x^{k-1}\|^2, \end{aligned}$$

using the variational inequality (8.2). It follows that

$$f(P_{\mathcal{U}}(x^{k-1} + sd^k)) - f(x^{k-1}) \leq \|P_{\mathcal{U}}(y_s) - x^{k-1}\|^2 \left(-\frac{\|\nabla f(x^{k-1})\|}{s} + O(1) \right).$$

Moreover, using again (8.2) we obtain

$$P_{\mathcal{U}}(y_s) = x^{k-1} \Leftrightarrow (y_s - x^{k-1}) \cdot (z - x^{k-1}) \leq 0 \quad \forall z \in \mathcal{U} \Leftrightarrow -\nabla f(x^{k-1}) \in N_{\mathcal{U}}(x^{k-1}),$$

which is the optimality condition stated in Proposition 11.3. We conclude the desired strict inequality $f(P_{\mathcal{U}}(x^{k-1} + sd^k)) - f(x^{k-1}) < 0$ for small s , unless x^{k-1} is a minimizer. \square

11.5 Application of Newton-type methods (complement)

We go back to the unconstrained problem (11.1) with $X = \mathbb{R}^n$. There are many more methods for increasing the efficiency of the gradient descent scheme. We simply refer to the standard literature [11] and discuss the Newton method and a few variants.

11.5.1 Newton's method for optimization

If the governing function $f(x)$ is two-times continuously differentiable, then the Newton scheme, see section 12.4 for a general presentation, can be adopted for the numerical solution of the minimization problem. As previously discussed, the necessary optimality condition is given by

$$F(x) := \nabla f(x) \stackrel{!}{=} 0.$$

The Newton iteration for solving the above system from an initial guess $x^0 \in X$, reads for $k = 1, 2, 3, \dots$

$$DF(x^k)(x^{k+1} - x^k) = -F(x^k).$$

Here it holds

$$DF(x^k) = D^2 f(x^k),$$

where $D^2 f(x^k)$ denotes the Hessian of f . We arrive at

$$x^{k+1} = x^k - D^2 f(x^k)^{-1} \nabla f(x^k).$$

In order to apply Newton schemes, the cost functional $f(x)$ must be twice continuously differentiable. In order to apply the classical Newton-Kantorovich theorem (see later) the second order derivative must be even Lipschitz continuous.

Remark 11.1. *Often in practice the strict Newton assumptions cannot be made or it is even clear that these do not hold true. Nonetheless, Newton's method can be used. Indeed, it often still yields satisfactory results. This shows that between mathematical rigorous results and practical applications might be often a gap.*

In Figure 11.4 we plot the convergence behavior of Newton's method for the Rosenbrock function. Our findings are summarized in Table 11.2.

In both cases, the minima could be approximated in only a very few steps. In the values of the table, we also observe that Newton's method is in general not a descent method. The iterates $f(x^k)$ do not monotonically converge to the minimum.

In numerical optimization (but also mildly and highly nonlinear PDE problems) the biggest problem is often the small convergence radius of Newton's method; i.e., the initial guess must be quite close to the solution to expect convergence. Here, we need so-called globalization strategies such as line search method or trust region concepts [11] or see also [4] for PDE problems. Another alternative is a mixture of gradient descent methods and Newton type methods: we first approach a good initial guess by using the robust gradient descent scheme. Once we are in the convergence region, we switch to Newton's method.

k	x^k	y^k	$\ \nabla f(x^k)\ $	$f(x^k)$	k	x^k	y^k	$\ \nabla f(x^k)\ $	$f(x^k)$
0	-1	1.5	18.86	6.5	0	2	-2	496.71	361
1	-1.222	1.444	6.92	4.963	1	1.992	3.966	1.988	0.983
2	-0.104	-1.239	26.08	16.84	2	1.001	0.021	43.91	9.620
3	-0.063	0.002	2.128	1.127	3	1.001	1.002	2.57e-03	1.65e-06
4	0.963	-0.123	45.51	11.03	4	1.000	1.000	7.41e-05	2.74e-11
5	0.965	0.931	7.05e-02	0.00124					
6	0.999	0.999	5.56e-02	1.55e-05					
7	1.000	1.000	9.67e-08	2.34e-15					

Table 11.2: Convergence behavior of Newton's method for the Rosenbrock function. The results are taken from [15].

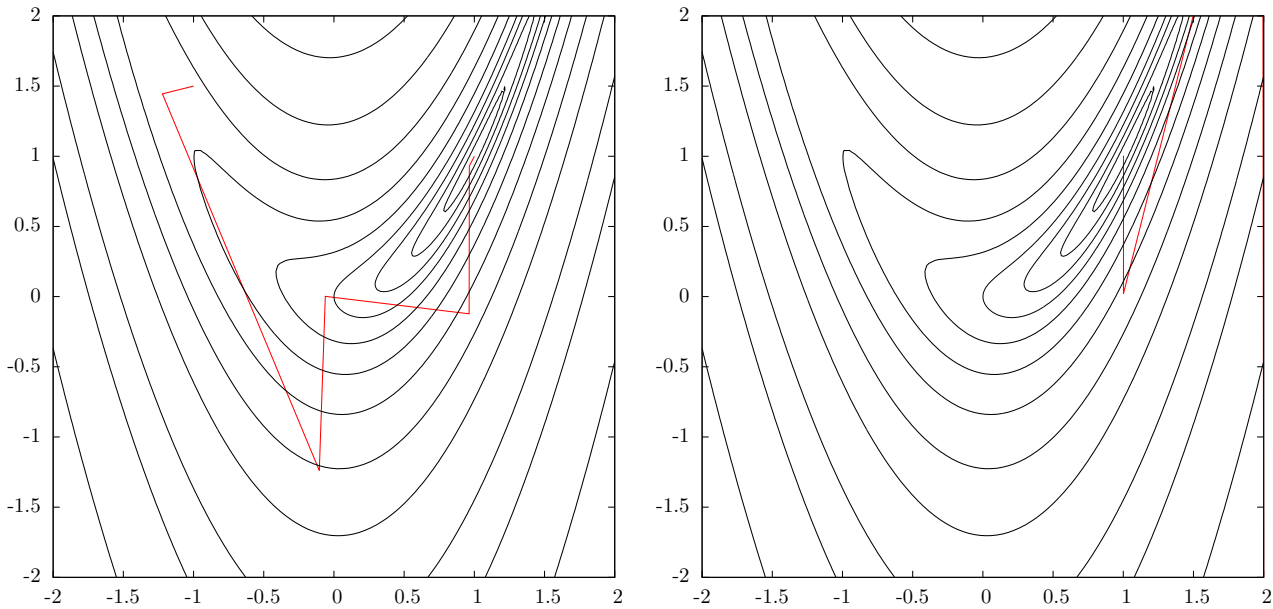


Figure 11.4: Convergence behavior of Newton's method for the approximation of the minimum of the Rosenbrock function $f(x, y) = 10(y - x^2)^2 + (1 - x)^2$. These figures are taken from [15].

11.5.2 Quasi-Newton methods

For problems of large dimension, or when the Hessian matrix is out of reach, we can use quasi-Newton methods which mimic Newton's methods without computing the exact Hessian. They extend the secant method to higher dimension. They have the general expression

$$\begin{cases} d^k = -H^{k-1}\nabla f(x^{k-1}), \\ s_k \text{ minimizes } \alpha \mapsto f(x^{k-1} + \alpha d^k), \\ x^k = x^{k-1} + s_k d^k, \end{cases}$$

where H^{k-1} is some symmetric positive definite approximation of $\nabla^2 f(x^{k-1})^{-1}$. The DFP method starts with some H^0 then uses the construction:

$$\begin{cases} \delta^k = x^k - x^{k-1} = s_k d^k, \\ \gamma^k = \nabla f(x^k) - \nabla f(x^{k-1}), \\ H^k = H^{k-1} + \frac{\delta^k (\delta^k)^T}{(\delta^k)^T \gamma^k} - \frac{H^{k-1} \gamma^k (\gamma^k)^T H^{k-1}}{(\gamma^k)^T H^{k-1} \gamma^k}. \end{cases}$$

For the BFGS variant, which is very popular, we modify the above by:

$$H^k = H^{k-1} + \left(1 + \frac{(\gamma^k)^T H^{k-1} \gamma^k}{(\delta^k)^T \gamma^k}\right) \frac{\delta^k (\delta^k)^T}{(\delta^k)^T \gamma^k} - \frac{\delta^k (\gamma^k)^T H^{k-1} + H^{k-1} \gamma^k (\delta^k)^T}{(\delta^k)^T \gamma^k}.$$

Of course, the line search is done approximately, using e.g. Armojo's rule or simply requiring a decrease of the cost function.

Chapter 12

Solving nonlinear equations

Solving nonlinear equations is one of the most-found tasks in numerical mathematics, but also in the mathematical analysis via fixed-point theorems. Examples are

- Root-finding of nonlinear functions, e.g., $f(x) = ax^3 + bx^2 + cx + d$
- Solution of (nonlinear) differential equations
- Solution of nonlinear optimization problems, e.g., classical tasks or more recently in machine learning.

Analytical results are again rarely achievable and for this reason, numerical algorithms are used. These are iterations in which an initial guess is chosen and from this guess, further approximations are calculated.

12.1 Introduction to iterative methods

Consider a function $f : I \subset \mathbb{R} \rightarrow \mathbb{R}$. We search, if it exists, some point $x \in I$ such that $f(x) = 0$. In general, such an equation can only be solved approximately. Iterative methods consist in the construction of a sequence (x_n) such that $x_n \rightarrow x$ when $n \rightarrow +\infty$ with $f(x) = 0$. Moreover, the speed of convergence of the sequence is of interest. It can be represented by the notion of order of convergence.

Let (x_n) be a sequence of real numbers that converges to some x , and let $r \geq 1$ be a real number.

Definition 12.1. (i) *The convergence of the sequence (x_n) is said to be of order at least r if there exist $C \geq 0$ and $N \in \mathbb{N}$ such that*

$$|x_{n+1} - x| \leq C|x_n - x|^r \quad \forall n \geq N.$$

(ii) *The convergence is said to be of order r exactly if, in addition, there exist $c > 0$ and $N' \in \mathbb{N}$ such that*

$$|x_{n+1} - x| \geq c|x_n - x|^r \quad \forall n \geq N'.$$

It is immediately seen that, if there exists $\lambda > 0$ such that

$$\frac{|x_{n+1} - x|}{|x_n - x|^r} \rightarrow \lambda,$$

then the convergence is of order r exactly.

Obviously, the larger the order of convergence is, the better it is. In the case of $r = 1$, however, the constant $0 < C < 1$ must be bounded. When the convergence is of (exact) order 1, we also speak of linear convergence. For $C := C_n \rightarrow 0$ for $n \rightarrow \infty$, we have so-called superlinear convergence. When it is of order 2 we speak of quadratic convergence.

12.2 The bisection method

Let $f : [a, b] \rightarrow \mathbb{R}$ be a continuous function such that $f(a)f(b) < 0$. The bisection method consists in the construction of three sequences (a_n) , (b_n) and (x_n) according to the following procedure

1. Initialization: set $a_0 = a$, $b_0 = b$, $x_0 = (a + b)/2$.

2. Iteration k : set

$$\begin{aligned} a_{k+1} &= a_k, \quad b_{k+1} = x_k && \text{if } f(x_k)f(a_k) < 0, \\ a_{k+1} &= x_k, \quad b_{k+1} = b_k && \text{if } f(x_k)f(b_k) < 0, \end{aligned}$$

$$x_{k+1} = \frac{a_{k+1} + b_{k+1}}{2}.$$

If $f(x_{k+1}) = 0$ then an exact solution is obtained and the iterations stop. Otherwise the iterations are continued.

It is not difficult to establish the following result.

Proposition 12.1 (Bisection). *Let $f \in C[a, b]$ with $f(a)f(b) < 0$. Let a_n, b_n, x_n the sequence constructed through the above algorithm. It holds*

$$x_n \rightarrow z, \quad f(z) = 0,$$

and the error estimate.

$$|x_n - z| \leq \frac{b - a}{2^{n+1}}.$$

Proof. (i) For $a < b$, it holds

$$a \leq a_1 \leq \cdots \leq a_n \leq a_{n+1} \leq x_n \leq b_{n+1} \leq b_n \leq \cdots \leq b_1 \leq b.$$

Thus, the scheme can be constructed. The sequence x_n is bounded through a_n and b_n . Therefore (classical result in Analysis), the sequence x_n converges. The algorithm terminates when $f(x_n) = 0$, i.e., $x_n = z$ is the sought root.

(ii) It holds

$$b_{n+1} - a_{n+1} = \begin{cases} x_n - a_n = \frac{b_n - a_n}{2} & f(a_n)f(x_n) < 0 \\ b_n - x_n = \frac{b_n - a_n}{2} & f(a_n)f(x_n) > 0 \end{cases}.$$

In both cases, it follows

$$|b_n - a_n| = \frac{|b_0 - a_0|}{2^n}. \tag{12.1}$$

Due to

$$a_n \leq x_n \leq b_n$$

it follows the convergence of $x_n \rightarrow z$ towards $z \in [a_n, b_n] \subset [a, b]$.

(iii) For the limit z it holds thanks to continuity $f(\cdot)$

$$0 \leq f(z)^2 = \lim_{n \rightarrow \infty} f(a_n)f(b_n) \leq 0 \quad \Rightarrow \quad f(z) = 0.$$

Therefore, the limit is the sought root. Finally, we can derive from (12.1) the error estimate:

$$|x_n - z| \leq \frac{|b_n - a_n|}{2} \leq \frac{|b - a|}{2^{n+1}}.$$

□

Remark 12.1 (A priori error estimate). *In Proposition 12.1 we have the a priori error estimate*

$$|x_n - z| \leq \frac{b - a}{2^{n+1}}.$$

Such an estimate can be evaluated before the computation. For instance, we obtain a guess for the required number of iterations n .

Exercise 12.1 (Bisection). *We consider $I = [0, 1]$. The a priori error estimate yields for different $n \in \mathbf{N}$ the bounds*

$$|x_9 - z| < 10^{-3}, \quad |x_{19} - z| < 10^{-6}, \quad |x_{29} - z| < 10^{-9}.$$

This means that in iteration 29 we can expect an accuracy of $TOL = 10^{-9}$. Compare this result to Newton's method (no matter its construction) in which a similar accuracy is obtained in approx. 4 steps. This shows that Newton's method can be more efficient than bisection. More specifically, bisection has a linear order of convergence and Newton's method it quadratically convergent.

Finally, we discuss an a posteriori result for which f must be continuously differentiable.

Proposition 12.2 (A posteriori error estimate). *Let $f \in C^1[a, b]$ with $f(a)f(b) < 0$ and $f(z) = 0$ for $z \in [a, b]$. The derivative is assumed to be bounded:*

$$0 < m \leq |f'(x)| \leq M < \infty$$

For the root z of f and the sequence $(x_n)_{n \in \mathbf{N}}$, it holds

$$\frac{|f(x_n)|}{M} \leq |x_n - z| \leq \frac{|f(x_n)|}{m}.$$

Proof. With Taylor expansion, it holds

$$f(x_n) = f(z) + f'(\xi)(x_n - z) = f'(\xi)(x_n - z) \quad (12.2)$$

with an intermediate value $\xi \in [a_n, b_n]$. It follows

$$\frac{|f(x_n)|}{\max |f'(x)|} \leq |x_n - z| \leq \frac{|f(x_n)|}{\min |f'(x)|}.$$

□

Exercise 12.2. *Given $f(x) = (x - 3)^2 - 1$. Take $a = 0$ and $b = 3.5$ and compute the root via bisection. Make first a sketch of the function curve. Check whether bisection can be applied. Compute three steps per hand. Optional: implement the code. Complement: is there only one root? If not, how can be approximate the other(s)?*

Exercise 12.3. *Given $f(x) = (x - 3)^2$. Take again $a = 0$ and $b = 3.5$. Can bisection be applied?*

12.3 Fixed points

12.3.1 Reminder: the (Banach) fixed point theorem

Definition 12.2. *A function $F : U \subset \mathbb{R} \rightarrow \mathbb{R}$ is said to be Lipschitz continuous with Lipschitz constant $k \geq 0$ if*

$$|F(x) - F(y)| \leq k|x - y| \quad \forall x, y \in U.$$

If F is Lipschitz continuous with Lipschitz constant $k < 1$ then F is said to be contracting.

Definition 12.3. *We say that x is a fixed point of F if $F(x) = x$.*

Note that the equation $f(x) = 0$ can always be equivalently rewritten as $F(x) = x$, for instance with $F(x) = f(x) + x$.

Theorem 12.3 (fixed point). *Let $F : U \rightarrow U$ be contracting, with $U \subset \mathbb{R}$ closed (for instance a closed interval $[a, b]$). Then F admits a unique fixed point x^* in U . Moreover, the sequence (x_n) defined by $x_0 \in U$ and $x_{n+1} = F(x_n)$ converges to x^* , and we have*

$$|x_n - x^*| \leq \frac{k^n}{1-k} |x_1 - x_0| \quad \forall n \in \mathbb{N}.$$

12.3.2 Attractive and repulsive fixed points

Let I be an open interval of \mathbb{R} , $F : I \rightarrow \mathbb{R}$ be a function of class \mathcal{C}^1 , and $x \in I$ be a fixed point of F .

With the mean value theorem in analysis, we have

$$\left| \frac{x_{n+1} - x}{x_n - x} \right| = \left| \frac{F(x_{n+1}) - F(x)}{x_n - x} \right| \rightarrow |F'(x)| \quad \text{for } n \rightarrow \infty$$

From this it follows that the asymptotic convergence rate is linear and the limit is $|F'(x)|$. Specifically, in the case of $|F'(x)| = 0$, we have at least superlinear convergence. From these considerations we obtain

Definition 12.4. *We say that*

- x is an attractive fixed point if $|F'(x)| < 1$,
- x is a repulsive fixed point if $|F'(x)| > 1$.

This notion gives an information on whether F is contracting or not in a neighborhood of x . It is therefore not surprising to have the following two results.

Proposition 12.4. *Let x be an attractive fixed point of F . There exists $h > 0$ such that $F([x-h, x+h]) \subset [x-h, x+h] \subset I$ and that all sequence defined by $x_0 \in [x-h, x+h]$ and $x_{n+1} = F(x_n)$ converges to x . Moreover, for such a sequence, the convergence is of order at least 1.*

Proposition 12.5. *Let x be a repulsive fixed point of F and (x_n) be a sequence such that $x_{n+1} = F(x_n) \forall n \in \mathbb{N}$. Then*

- either there exists an integer N such that $x_n = x \forall n \geq N$,
- or (x_n) does not converge to x .

12.3.3 Calculating fixed points

The algorithm to compute fixed-points is simple:

$$x_{n+1} = F(x_n), \quad n = 0, 1, 2, 3, \dots$$

where x_0 is the initial guess, therein F is the so-called iteration function.

Exercise 12.4. *Take for instance $F(x) = \cos(x)$ as an example for an iteration function.*

Exercise 12.5 (Root-finding problem). *Fixed-point problems are naturally root-finding problems since $x = F(x)$ is equivalent to $x - F(x) = 0$.*

Exercise 12.6. *Given again $f(x) = (x-3)^2 - 1$. Take $x_0 = 0$ as initial guess and design first the fixed point iteration function $F(x)$. Then introduce the iteration scheme and a relaxation parameter $\omega > 0$. The latter one is required such that we deal with an attractive fixed point with $|F'(x)| < 1$. Compute for $\omega = 1$ and $\omega = 0.1$ the first three iterations x_1, x_2, x_3 . Justify theoretically for which ω the condition $|F'(x)| < 1$ holds true.*

12.4 The Newton method

There exist many variants and simplifications, globalizations (Newton's method is a priori only a local method) as well as different names depending on we work with scalar equations, in \mathbb{R}^n , or even in infinite-dimensional function spaces.

Let us come back to the root-finding problem

$$f(x) = 0.$$

We assume that f is differentiable. We construct a sequence of iterates $(x_n)_{n \in \mathbb{N}}$ and hopefully reach at some point

$$|f(x_n)| < TOL, \quad \text{where } TOL \text{ is again small, e.g., } TOL = 10^{-10}.$$

One has to start with a Taylor expansion. In our lecture we do this as follows. Let us assume that we are at x_n and can evaluate $f(x_n)$. Now we want to compute this next iterate x_{n+1} with the unknown value $f(x_{n+1})$. Taylor gives us:

$$f(x_{n+1}) = f(x_n) + f'(x_n)(x_{n+1} - x_n) + o(x_{n+1} - x_n)^2$$

We assume that $f(x_{n+1}) = 0$ (or very close to zero $f(x_{n+1}) \approx 0$). Then, x_{n+1} is the sought root and neglecting the higher-order terms we obtain:

$$0 = f(x_n) + f'(x_n)(x_{n+1} - x_n),$$

i.e.,

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}.$$

If (x_n) converges, the Taylor-Young formula yields $f(x_{n+1}) = o(x_{n+1} - x_n)$, hence a fast convergence of $f(x_n)$ to 0 is expected.

Exercise 12.7. Given again $f(x) = (x-3)^2 - 1$. Take $x_0 = 0$ as initial guess and use Newton's method to compute the first three iterates x_1, x_2, x_3 . Compare with your results the speed of convergence for bisection, fixed-point and Newton.

Exercise 12.8 (Newton's method as fixed-point iteration). The previous algorithm can be stated in terms of a fixed-point iteration with the iteration function

$$F(x) = x - \frac{f(x)}{f'(x)}.$$

Theorem 12.6. Let I be an open interval and $f : I \rightarrow \mathbb{R}$ be a function of class \mathcal{C}^2 . We assume that there exist positive real numbers m, k such that

$$|f'(y)| \geq m \quad \forall y \in I,$$

$$|f''(y)| \leq k \quad \forall y \in I.$$

Let $x \in I$ be such that $f(x) = 0$. If $0 < h < 2m/k$, $x_0 \in]x - h, x + h[\subset I$ and

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad \forall n \in \mathbb{N},$$

then

1. $x_n \in]x - h, x + h[\quad \forall n \in \mathbb{N}$,
2. $x_n \rightarrow x$,
3. the convergence is of order at least 2.

In the practical use of the Newton method, the constants m and k are usually not known (it is even not always clear that such constants exist). Therefore it is very difficult to ensure that x_0 is taken sufficiently close to x (recall also that x is by definition also unknown) in order to have convergence. This is the major drawback of the Newton method. However, when it converges, the convergence is fast. In contrast, the bisection method guarantees the convergence under the condition $f(a)f(b) < 0$, which is easy to check, but the convergence is slow. Efficient methods often combine several techniques, with Newton-type iterations performed as last stage.

12.5 Generalization to higher dimensions (complement)

The fixed point and Newton methods admit rather straightforward extensions to higher dimensions, i.e., for solving $F(x) = 0$ or $f(x) = x$ with $F, f : \mathbb{R}^n \rightarrow \mathbb{R}^n$. In Newton's method, the derivative is replaced by the Jacobian matrix $DF(x)$, hence the iteration becomes

$$x_{n+1} = x_n - DF(x_n)^{-1}F(x_n).$$

It is actually not needed to compute the full inverse of $DF(x_n)$, but only the solution of the system $DF(x_n)h = F(x_n)$ has to be found to obtain $DF(x_n)^{-1}F(x_n)$. Still, for large problems, this may yield a significant computational cost, more or less balanced by the fast convergence of the sequence.

12.5.1 Newton's method: going from \mathbb{R} to higher dimensions

Overview:

- Newton-Raphson (1D), find $x \in \mathbb{R}$ via iterating $k = 0, 1, 2, \dots$ such that $x_k \approx x$ via:

$$\begin{aligned} \text{Find } \delta x \in \mathbb{R} : \quad & f'(x_k)\delta x = -f(x_k), \\ \text{Update:} \quad & x_{k+1} = x_k + \delta x. \end{aligned}$$

- Newton in \mathbb{R}^n , find $x \in \mathbb{R}^n$ via iterating $k = 0, 1, 2, \dots$ such that $x_k \approx x$ via:

$$\begin{aligned} \text{Find } \delta x \in \mathbb{R}^n : \quad & F'(x_k)\delta x = -F(x_k), \\ \text{Update:} \quad & x_{k+1} = x_k + \delta x. \end{aligned}$$

Here we need to solve a linear equation system to compute the update $\delta x \in \mathbb{R}^n$.

Often, Newton's method is formulated in terms of a defect-correction scheme. To illustrate this more clearly we consider the problem $f(x) = y$ instead of $f(x) = 0$.

Definition 12.5 (Defect). *Let $\tilde{x} \in \mathbb{R}$ an approximation of the solution $f(x) = y$. The defect (or similarly the residual) is defined as*

$$d(\tilde{x}) = y - f(\tilde{x}).$$

Definition 12.6 (Newton's method as defect-correction scheme).

$$\begin{aligned} f'(x_k)\delta x &= d_k, \quad d_k := y - f(x_k), \\ x_{k+1} &= x_k + \delta x, \quad k = 0, 1, 2, \dots \end{aligned}$$

The iteration is finished with the same stopping criterion as for the classical scheme. To compute the update δx we need to invert $f'(x_k)$:

$$\delta x = (f'(x_k))^{-1}d_k.$$

This step seems trivial but is the most critical one if we deal with problems in \mathbb{R}^n with $n > 1$ or in function spaces. Because here, the derivative becomes a matrix. Therefore, the problem results in solving a linear equation system of the type $A\delta x = b$.

Remark 12.2. *This previous form of Newton's method is already very close to the schemes that are used in research. In higher dimensions it applies to deal with problems such as nonlinear PDEs or optimization. The 'only' aspects that are however big research topics are the choice of*

- *good initial Newton guesses;*
- *globalization techniques.*

Two very good books on these topics, including further materials as well, are [4, 11].

12.5.2 A basic algorithm for a residual-based Newton method

In this type of methods, the main criterion is a decrease of the residual in each step. We denote $F(x) = -d(x) = f(x) - y$.

Algorithm 12.7 (Residual-based Newton's method). *Given an initial guess x_0 . Iterate for $k = 0, 1, 2, \dots$ such that*

$$\begin{aligned} \text{Find } \delta x \in \mathbb{R}^n : \quad & F'(x_k)\delta x_k = -F(x_k), \\ \text{Update:} \quad & x_{k+1} = x_k + \lambda_k \delta x_k, \end{aligned}$$

with $\lambda_k \in (0, 1]$ (see the next sections how λ_k can be determined). A full Newton step corresponds to $\lambda_k = 1$. The criterion for convergence is the contraction of the residuals measured in terms of a discrete vector norm:

$$\|F(x_{k+1})\| < \|F(x_k)\|,$$

here $\|\cdot\|$ is a norm associated to the function space, here \mathbb{R}^n , for instance $\|\cdot\|_{l_2}$. In order to save some computational cost, close to the solution x^ , intermediate simplified Newton steps can be used. In the case of $\lambda_k = 1$ we observe*

$$\theta_k = \frac{\|F(x_{k+1})\|}{\|F(x_k)\|} < 1.$$

If $\theta_k < \theta_{max}$, e.g., $\theta_{max} = 0.1$, then the old Jacobian $F'(x_k)$ is kept and used again in the next step $k + 1$. Otherwise, if $\theta_k > \theta_{max}$, the Jacobian will be assembled. Finally the stopping criterion is one of the following (relative preferred!):

$$\begin{aligned} \|F(x_{k+1})\| &\leq \text{TOL}_N \quad (\text{absolute}) \\ \|F(x_{k+1})\| &\leq \text{TOL}_N \|F(x_0)\| \quad (\text{relative}) \end{aligned}$$

If fulfilled, set $x^ := x_{k+1}$ and the (approximate) root x^* of the problem $F(x) = 0$ is found.*

12.5.3 Example of the basic Newton method

We illustrate the basic Newton schemes with the help a numerical example in which we want to solve

$$x^2 = 2 \quad \text{in } \mathbb{R}$$

We reformulate this equation as root-finding problem:

$$f(x) = 0,$$

where

$$f(x) = x^2 - 2.$$

For Newton's method we need to compute the first derivative, which is trivial here,

$$f'(x) = 2x,$$

but can become a real challenge for complicated problems. To start Newton's method (or any iterative scheme), we need an initial guess (similar to ODE-IVP). Let us choose $x_0 = 3$. Then we obtain as results:

Iter	x	f(x)
0	3.000000e+00	7.000000e+00
1	1.833333e+00	1.361111e+00
2	1.462121e+00	1.377984e-01
3	1.414998e+00	2.220557e-03
4	1.414214e+00	6.156754e-07
5	1.414214e+00	4.751755e-14

We see that Newton's method converges very fast, i.e., quadratically, which means that the correct number of digits roughly doubles at each iteration step.

12.5.4 Example using a Newton defect-correction scheme including line search

In this second example, we present Newton's method as defect correction scheme as introduced in Definition 12.6. Additionally we introduce a line search parameter $\omega \in [0, 1]$. If the initial Newton guess is not good enough we can enlarge the convergence radius of Newton's method by choosing $\omega < 1$. Of course for the given problem here, Newton's method will always converge, because the underlying function is convex. Despite the fact that ω is not really necessary in this example, we can highlight another feature very well. In the optimal case (as seen above in the results), Newton's method will converge quadratically. Any adaptation will deteriorate the performance. Thus, let us choose $\omega = 0.9$. Then we obtain as result only linear convergence:

Iter	x	f(x)
0	3.000000e+00	7.000000e+00
1	1.950000e+00	1.802500e+00
2	1.534038e+00	3.532740e-01
3	1.430408e+00	4.606670e-02
4	1.415915e+00	4.816699e-03
5	1.414385e+00	4.840133e-04
6	1.414231e+00	4.842504e-05
7	1.414215e+00	4.842742e-06
8	1.414214e+00	4.842766e-07
9	1.414214e+00	4.842768e-08
10	1.414214e+00	4.842768e-09
11	1.414214e+00	4.842771e-10
12	1.414214e+00	4.842748e-11

And if we choose $\omega = 0.5$ we only obtain (perfect) linear convergence:

Iter	x	f(x)
0	3.000000e+00	7.000000e+00
1	2.416667e+00	3.840278e+00
2	2.019397e+00	2.077962e+00
3	1.762146e+00	1.105159e+00
4	1.605355e+00	5.771631e-01
5	1.515474e+00	2.966601e-01
...		
32	1.414214e+00	2.286786e-09
33	1.414214e+00	1.143393e-09


```

34  1.414214e+00 5.716965e-10
35  1.414214e+00 2.858482e-10
36  1.414214e+00 1.429239e-10
37  1.414214e+00 7.146195e-11
=====

```

In summary, using a full Newton method we only need 5 iterations until a tolerance of $TOL_N = 10^{-10}$ is reached. Having linear convergence with a good convergence factor still 12 iterations are necessary. And a linear scheme does need 37 iterations. This shows nicely the big advantage of Newton's method (i.e., quadratic convergence) in comparison to linear or super-linear convergence.

12.5.5 Newton's method in higher dimensions and the Newton-Kantorovich theorem

In this section, we now study Newton's method in \mathbb{R}^n from the theoretical side including proofs and further extensions such as globalization techniques.

Let $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$. We seek a solution to the root-finding $f(x) = (f_1(x), \dots, f_n(x)) = 0$. We take an initial guess $x^{(0)} \in D$ and iterate $x^{(k)} \rightarrow x^{(k+1)}$. The derivation is based on the one-dimensional situation. Around the iteration $x^{(k)}$ we linearize the function $f(\cdot)$ with the help of the Taylor up to the first order. In an arbitrary direction, it holds $w^{(k)} \in \mathbb{R}^n$

$$f(x^{(k)} + w^{(k)}) = f(x^{(k)}) + \sum_{j=1}^n \frac{\partial f}{\partial x_j}(x^{(k)})w_j^{(k)} + O(|w^{(k)}|^2).$$

The new iterate $x^{(k+1)} = x^{(k)} + w^{(k)}$ is defined as root of the linearization. Therefore, we seek a solution $w^{(k)} \in \mathbb{R}^n$ of the linear equation system

$$\sum_{j=1}^n \frac{\partial f_i}{\partial x_j}(x^{(k)})w_j^{(k)} = -f_i(x^{(k)}), \quad i = 1, \dots, n.$$

The Jacobian $Df : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$ reads

$$Df = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}$$

and consequently, we obtain in compact notation as linear equation system

$$Df(x^{(k)})w^{(k)} = -f(x^{(k)}).$$

Consequently, the Newton iteration with initial guess $x^{(0)} \in D$ for $k = 0, 1, 2, \dots$ reads

$$x^{(0)} \in D, \quad Df(x^{(k)})w^{(k)} = -f(x^{(k)}), \quad x^{(k+1)} = x^{(k)} + w^{(k)}. \quad (12.3)$$

The biggest difference to the one-dimensional case is that we now deal with a matrix for the derivative. Specifically, we need to compute n^2 derivatives. Rather than just dividing, we need to solve a linear equation system with the coefficient matrix $Df(x_k) \in \mathbb{R}^{n \times n}$.

The central convergence results of this Newton method is the Newton-Kantorovich theorem. Let $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a differentiable mapping. With $\|\cdot\|$ we denote the Euclidian vector norm and the induced matrix norm; here the spectral norm. We search a root $z \in D$ such that $f(z) = 0$.

Theorem 12.8 (Newton-Kantorovich). *Let $D \subset \mathbb{R}^n$ be an open and convex set. Moreover, let $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ be continuously differentiable.*

(1) *Let the Jacobian matrix Df be uniformly Lipschitz-continuous for all $x, y \in D$:*

$$\|Df(x) - Df(y)\| \leq L\|x - y\|, \quad x, y \in D, \quad (12.4)$$

with $L < \infty$.

(2) *Let the Jacobian on D be a uniformly bounded inverse*

$$\|Df(x)^{-1}\| \leq \beta, \quad x \in D, \quad (12.5)$$

with $\beta < \infty$.

(3) *For the initial guess $x^{(0)} \in D$, assume*

$$q := \alpha\beta L < \frac{1}{2}, \quad \alpha := \|Df(x^{(0)})^{-1}f(x^{(0)})\|. \quad (12.6)$$

(4) *For $r := 2\alpha$ suppose that the closed ball*

$$B_r(x_0) := \{x \in \mathbb{R}^n : \|x - x_0\| \leq r\}$$

is contained in the set D .

Then, the function f has a unique root $z \in B_r(x_0)$ and the Newton iteration

$$Df(x^{(k)})\delta x = -f(x^{(k)}), \quad x^{(k+1)} = x^{(k)} + \delta x, \quad k = 0, 1, 2, \dots,$$

converges quadratically to the root z . Furthermore, we have the a priori error estimate

$$\|x^{(k)} - z\| \leq 2\alpha q^{2^k - 1}, \quad k = 0, 1, \dots$$

Proof. We translate the proof from [15], which was based on [14], which is itself close to Kantorovich [8]. Alternative proofs can be found in [4] for instance.

The proof is complicated and we first make an outline:

- (i) Derivation of some auxiliary lemma.
 - (ii) We then show that all iterates $x^{(k)}$ are contained in the ball $B_r(x_0)$. Furthermore, we then obtain the a priori error estimate.
 - (iii) We show that the iterates $(x^{(k)})_{k \in \mathbb{N}}$ form a Cauchy sequence and have consequently a limit z .
 - (iv) Next, we establish that this limit z is a root of f .
 - (v) Finally, we show that we have one and only one root in $B_r(x_0)$.
- (i) Let $x, y, z \in D$. Since D is convex, it holds for all $x, y \in D$

$$f(x) - f(y) = \int_0^1 Df(y + s(x - y))(x - y) ds.$$

We subtract on both sides $Df(z)(x - y)$:

$$f(x) - f(y) - Df(z)(x - y) = \int_0^1 ((Df(y + s(x - y)) - Df(z))(x - y) ds$$

With the Lipschitz continuity of the Jacobian Df , it follows

$$\begin{aligned} \|f(y) - f(x) - Df(z)(y - x)\| &\leq L\|y - x\| \int_0^1 \|s(x - z) + (1 - s)(y - z)\| ds \\ &\leq \frac{L}{2}\|y - x\|(\|x - z\| + \|y - z\|). \end{aligned}$$

For the choice $z = x$, we obtain

$$\|f(y) - f(x) - Df(x)(y - x)\| \leq \frac{L}{2} \|y - x\|^2, \quad \forall x, y \in D. \quad (12.7)$$

And for $z = x_0$, it yields

$$\|f(y) - f(x) - D(x_0)(y - x)\| \leq rL \|y - x\|, \quad \forall x, y \in B_r(x_0). \quad (12.8)$$

(ii) Next, we show that all iterates are contained in the ball $B_r(x^{(0)})$. These derivations also yield as byproduct the a priori error estimate. We establish the proof using induction and show that

$$\|x^{(k+1)} - x^{(0)}\| \leq r, \quad \|x^{(k+1)} - x^{(k)}\| \leq \alpha q^{2^{(k)}-1}, \quad k = 1, 2, \dots \quad (12.9)$$

Let us begin with $k = 0$. It holds for the Newton iteration $x^{(1)} - x^{(0)} = -Df(x^{(0)})^{-1}f(x^{(0)})$ with the condition (12.6)

$$\|x^{(1)} - x^{(0)}\| = \|Df(x^{(0)})^{-1}f(x^{(0)})\| = \alpha = \frac{r}{2} < r,$$

d.h. $x^{(1)} \in B_r(x^{(0)})$, and it holds the estimate

$$\|x^{(1)} - x^{(0)}\| \leq \alpha = \alpha q^{2^0-1}.$$

Induction step from $k \rightarrow k + 1$. According to the induction assumption, both equations (12.9) are true for $k \geq 0$. Consequently, it holds $x^{(k)} \in B_r(x^{(0)})$ such that the Newton iterate $x^{(k+1)}$ is well-defined. Then, we infer with the condition (12.5), with the estimate (12.7) and the induction assumption (12.9) with the definition of q the following estimates:

$$\begin{aligned} \|x^{(k+1)} - x^{(k)}\| &= \|Df(x^{(k)})^{-1}f(x^{(k)})\| \\ &\leq \beta \|f(x^{(k)})\| \\ &= \beta \| \underbrace{f(x^{(k)}) - f(x^{(k-1)}) - Df(x^{(k-1)})(x^{(k)} - x^{(k-1)})}_{=0} \| \\ &\leq \frac{\beta L}{2} \underbrace{\|x^{(k)} - x^{(k-1)}\|^2}_{\text{Induction}} \leq \frac{\beta L}{2} (\alpha q^{2^{(k-1)}-1})^2 \\ &= \frac{\alpha}{2} q^{2^{(k)}-1} < \alpha q^{2^{(k)}-1} \end{aligned}$$

Hence, we have

$$\begin{aligned} \|x^{(k+1)} - x^{(0)}\| &\leq \|x^{(k+1)} - x^{(k)}\| + \dots + \|x^{(1)} - x^{(0)}\| \\ &\leq \alpha(1 + q + q^3 + q^7 + \dots + q^{2^{(k)}-1}) \\ &\leq \frac{\alpha}{1 - q} \leq 2\alpha = r. \end{aligned}$$

Thus, $x^{(k+1)} \in B_r(x^{(0)})$. This shows that induction step from $k \rightarrow k + 1$, i.e., the two inequalities (12.9) hold true for $k + 1$.

(iii) We now show that the iterates $x^{(k)} \in B_r(x^{(0)})$ form a Cauchy sequence. Let $m > 0$. Since $q < \frac{1}{2}$, it holds

$$\begin{aligned} \|x^{(k)} - x^{k+m}\| &\leq \|x^{(k)} - x^{(k+1)}\| + \dots + \|x^{k+m-1} - x^{k+m}\| \\ &\leq \alpha(q^{2^{(k)}-1} + q^{2^{(k+1)}-1} + \dots + q^{2^{m+k-1}-1}) \\ &= \alpha q^{2^{(k)}-1} (1 + q^{2^{(k)}} + \dots + (q^{2^{(k)}})^{2^{m-1}-1}) \\ &\leq 2\alpha q^{2^{(k)}-1}. \end{aligned} \quad (12.10)$$

This shows that indeed $(x^{(k)}) \subset D$ form a Cauchy sequence because $q < \frac{1}{2}$. In the Banach space \mathbb{R}^n , the limit

$$z = \lim_{k \rightarrow \infty} x^{(k)} \in \mathbb{R}^n$$

exists. Passing to the limit $k \rightarrow \infty$, we obtain with (12.9)

$$\|z - x^{(0)}\| \leq r,$$

such that $z \in B_r(x_0)$. For the limit $m \rightarrow \infty$ in (12.10) we establish the error estimate

$$\|x^{(k)} - z\| \leq 2\alpha q^{2^{(k)}-1}, \quad k = 0, 1, \dots$$

(iv) It remains to show that $z \in B_r(x^{(0)})$ is a root of f . The Newton iteration rule as well as the condition (12.4) yield

$$\begin{aligned} \|f(x^{(k)})\| &= \|Df(x^{(k)})(x^{(k)} - x^{(k-1)})\| \\ &\leq \|Df(x^{(k)}) - Df(x^{(0)}) + Df(x^{(0)})\| \|x^{(k+1)} - x^{(k)}\| \\ &\leq \left(L\|x^{(k)} - x^{(0)}\| + \|Df(x^{(0)})\| \right) \|x^{(k+1)} - x^{(k)}\| \rightarrow 0 \quad (k \rightarrow \infty). \end{aligned}$$

Hence, it holds

$$f(x^{(k)}) \rightarrow 0, \quad k \rightarrow \infty.$$

The continuity of f implies $f(z) = 0$.

(v) Finally, we show that the root $z \in B_r(x_0)$ is unique. This is obtained with the help of the contraction property and the formulation of the Newton scheme as fixed-point iteration. Each root of $f(\cdot)$ is a fixed-point of the simplified Newton iteration

$$g(x) := x - Df(x^{(0)})^{-1}f(x)$$

The fixed-point iteration function $g(\cdot)$ is Lipschitz continuous. With (12.8), it holds

$$\begin{aligned} g(x) - g(y) &= x - y - Df(x^{(0)})^{-1}f(x) + Df(x^{(0)})^{-1}f(y) \\ &= Df(x^{(0)})^{-1}(f(y) - f(x) - Df(x^{(0)})(y - x)) \\ &\leq \underbrace{\|Df(x^{(0)})^{-1}\|}_{\leq \beta} rL \|y - x\| \\ &\leq \beta Lr \|y - x\|. \end{aligned}$$

With $r = 2\alpha$, it follows $\beta Lr \leq 2\alpha\beta L \leq 2q < 1$. This yields that g is a contraction. Banach's fixed-point theorem yields that we have only one fixed-point and consequently only one root. \square

Remark 12.3. *The Newton-Kantorovich theorem differs in several points from well-posedness results for one-dimensional Newton schemes. For instance, the root is not an assumption, but follows within the proof. Moreover, the Newton-Kantorovich theorem only requires a Lipschitz continuous first derivative, but not a twice differentiable function.*

From Theorem 12.8, the following convergence result can be inferred:

Corollary 12.9. *Let $D \subset \mathbb{R}^n$ be open and $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ be twice continuously differentiable. We assume that $z \in D$ is a root with regular Jacobian $Df(z)$. Then, Newton's method is locally convergent, i.e., there exists a neighborhood B around z such that the Newton scheme converges for all starting values $x^{(0)} \in B$.*

Example 12.1 (Newton's method in \mathbb{R}^n). We seek the root of the function

$$f(x_1, x_2) = \begin{pmatrix} 1 - x_1^2 - x_2^2 \\ (x_1 - 2x_2)/(1/2 + x_2) \end{pmatrix}$$

with the Jacobian

$$Df(x) = \begin{pmatrix} -2x_1 & -2x_2 \\ \frac{2}{1+2x_2} & -\frac{4+4x_1}{(1+2x_2)^2} \end{pmatrix}.$$

The roots of f are given by

$$x \approx \pm(0.894427, 0.447214).$$

We start the iteration with $x^{(0)} = (1, 1)^T$ and obtain the iterates

$$\begin{aligned} x^1 &\approx \begin{pmatrix} 1.14286 \\ 0.357143 \end{pmatrix}, & x^2 &\approx \begin{pmatrix} 0.92659 \\ 0.442063 \end{pmatrix}, \\ x^3 &\approx \begin{pmatrix} 0.894935 \\ 0.447349 \end{pmatrix}, & x^4 &\approx \begin{pmatrix} 0.894427 \\ 0.447214 \end{pmatrix}. \end{aligned}$$

After only four iterations, the first six digits are exact.

12.5.6 Globalization of Newton's method

Each step of Newton's method consists of the parts

1. Computation of the Jacobian;
2. Solving the linear equation system;
3. Update and monitoring the convergence.

In contrast to the one-dimensional case, the computation of the derivative in \mathbb{R}^n may be cumbersome, since n^2 partial derivatives must be calculated:

$$(Df)_{ij} = \frac{\partial f_i}{\partial x_j}, \quad i, j = 1, \dots, n.$$

This can be a non-negligible computational cost. In the second step of Newton's method, we must solve a linear equation system. Using for instance LU decomposition (as we discussed before in detail), it requires $O(n^3)$ operations. The solution requires then $O(n^2)$ operations.

As well we can work with simplified Newton methods. If the derivative $f'(x^{(k)})$ is not computed according to the current right hand side, but to some other value $f'(c)$, we still have convergence, but only with a linear order. In the higher-dimensional case, we obtain

Listing 12.1: (Simplified Newton's method)

Let $f \in C^1(D)$ be given, an initial guess $x^{(0)} \in D$ and $c \in D$.

Compute the Jacobian $Df(c)$

Construct LU decomposition $Df(c) = L(c)R(c)$

Iterate for $k = 1, 2, \dots$

$$L(c)R(c)w^{(k)} = -f(x^{(k-1)})$$

$$x^{(k)} = x^{(k-1)} + w^{(k)}$$

The two expensive steps of Newton's method are only required once. The iteration can be performed efficiently with forward and backward substitution. This can be further improved when the linearization point $c \in \mathbb{R}^n$ is selected in a dynamic way, for instance at $c = x^{(3)}$, $c = x^{(6)}$, \dots or when the convergence behavior deteriorates. As monitor, we can observe the residual of two subsequent iterates

$$\rho_k = \frac{|f(x^{(k)})|}{|f(x^{(k-1)})|}.$$

If $\rho_k > \rho_0$ (bad convergence), then with $c = x^{(k)}$, we perform an update of the Jacobian.

The second challenge is the choice of a good initial guess. Enlarging the convergence radius, i.e., globalization, can be obtained with damping. Here, step 5) in Algorithm 12.1 is replaced by

$$x^{(k)} = x^{(k-1)} + \omega_k w^{(k)}$$

with a damping parameter $\omega_k \in (0, 1]$. For the Newton iteration, this reads

$$x^{(k+1)} = x^{(k)} - \omega_k Df(x^{(k)})^{-1} f(x^{(k)}).$$

Proposition 12.10 (Damped Newton's method). *Let the assumptions from Theorem 12.8 hold true except for the condition on the initial value (12.6). The damped Newton iteration*

$$Df(x^{(k)})\delta x = -f(x^{(k)}), \quad x^{(k+1)} = x^{(k)} + \omega_k \delta x$$

with

$$\omega_k := \min\left\{1, \frac{1}{\alpha_k \beta L}\right\}, \quad \alpha_k := \|Df(x^{(k)})^{-1} f(x^{(k)})\|$$

(ω_k is the so-called line search parameter) generates a sequence $(x^{(k)})_{k \in \mathbb{N}}$, for which after k_* steps

$$q_* := \alpha_{k_*} \beta L < \frac{1}{2}$$

is fulfilled. For $k > k_*$ the sequence $x^{(k)}$ converges quadratically and we have the a priori error estimate

$$\|x^{(k)} - z\| \leq \frac{\alpha}{1 - q_*} q_*^{2^{(k)} - 1}, \quad k \geq k_*.$$

Proof. We show the monotone convergence in the sense of $\|f(x^{(k+1)})\| \leq \|f(x^{(k)})\|$ of the damped Newton iteration. For $x^{(0)}$, the set

$$D_0 := \{x \in D : \|f(x)\| \leq \|f(x^{(0)})\|\}$$

is closed and non-empty. We define the damped iteration

$$x_\omega = g_\omega(x) := x - \omega Df(x)^{-1} f(x).$$

For $x \in D_0$, let $\omega_{max} \leq 1$ be the maximal damping parameter such that

$$f(x_\omega) \leq f(x^{(0)}) \quad 0 \leq \omega \leq \omega_{max} \leq 1.$$

We define

$$h(\omega) := f(x_\omega)$$

with $h(0) = f(x)$. For $h(\cdot)$ it holds further

$$h'(\omega) = \frac{d}{d\omega} f(x - \omega Df(x)^{-1} f(x)) = -Df(x_\omega) Df(x)^{-1} f(x).$$

With $h'(0) = -f(x) = -h(0)$, it follows

$$\begin{aligned} f(x_\omega) - f(x) &= h(\omega) - h(0) = \int_0^\omega h'(s) ds = - \int_0^\omega Df(x_s) Df(x)^{-1} f(x) ds \\ &= - \int_0^\omega \left\{ (Df(x_s) - Df(x)) Df(x)^{-1} f(x) + f(x) \right\} ds. \end{aligned}$$

Hence,

$$f(x_\omega) = \left(- \int_0^\omega (Df(x_s) - Df(x)) Df(x)^{-1} ds \right) f(x) + (1 - \omega) f(x)$$

and we estimate

$$\begin{aligned} \|f(x_\omega)\| &\leq (1 - \omega)\|f(x)\| + \beta L \int_0^\omega \|x_s - x\| ds \|f(x)\| \\ &= (1 - \omega)\|f(x)\| + \beta L \int_0^\omega s \|Df(x)^{-1}f(x)\| ds \|f(x)\| \\ &\leq \left(1 - \omega + \alpha_x \beta L \frac{\omega^2}{2}\right) \|f(x)\|, \quad \alpha_x = \|Df(x)^{-1}f(x)\|. \end{aligned}$$

Now, we determine $0 \leq \omega \leq \omega_{max} \leq 1$ such that the expression

$$\rho(\omega) := \left|1 - \omega + \alpha_x \beta L \frac{\omega^2}{2}\right|$$

is minimal. The minimum is obtained for

$$\omega_{min} = \min \left\{1, \frac{1}{\alpha_x \beta L}\right\}.$$

Then, it holds

$$1 - \omega_{min} + \frac{\omega_{min}^2}{2} \alpha_x \beta L \leq 1 - \frac{1}{2\alpha_x \beta L} < 1.$$

For this choice of the damping parameter, we have monotone convergence:

$$\|f(x^{(k+1)})\| \leq \left(1 - \frac{1}{2\alpha_k \beta L}\right) \|f(x^{(k)})\|$$

After a finite number of steps, it holds

$$\frac{\alpha_k \beta L}{2} \leq \frac{\beta^2 L}{2} \|f(x^{(k)})\| < 1$$

and we arrive at the quadratic convergence radius. \square

Remark 12.4 (Globalization). *The damped Newton scheme is a globalization strategy. The condition for the initial guess*

$$x^{(0)} \in D : \quad \underbrace{\|Df(x^{(0)})^{-1}f(x^{(0)})\|}_{=\alpha_0} \beta L < \frac{1}{2}$$

is not required to be necessarily fulfilled. If this product is too big, we can iterate with the damped iteration until $\alpha_k \beta L$ with $\alpha_k = \|Df(x^{(k)})^{-1}f(x^{(k)})\|$ fulfills the condition. As said, globalization strategies yield convergence in a larger region, but therein, however, the convergence is initially usually slow with the rate

$$\rho \leq 1 - \frac{1}{2\alpha_k \beta L}.$$

For large α_k , we can choose ρ close to 1.

It is important that the switch to fast convergence be in an automated fashion. Here, this is given since

$$\alpha_k \beta L < \frac{1}{2} \quad \Rightarrow \quad \omega_k = \min \left\{1, \frac{2}{\alpha_k \beta L}\right\} = 1.$$

A difficulty remains however as to the choice of the damping parameter ω_k . Indeed the values

$$\alpha_k = \|Df(x^{(k)})^{-1}f(x^{(k)})\|, \quad \beta = \max_{x \in D} \|Df(x)^{-1}\|,$$

and the Lipschitz constant of $Df(x)$ are in general not computable. A common strategy (see also in numerical optimization [11]) is a line-search technique.

Listing 12.2: **(Line-Search)**

We are given $f \in C^1(D)$, and $x^{(0)} \in D$ and a parameter $\sigma \in (0, 1)$ and a maximal number of line search steps $L_{max} \in \mathbb{N}$.

- 1 For $k = 0, 1, 2, \dots$, we iterate
- 2 $Df(x^{(k)})w^{(k)} = -f(x^{(k)})$
- 3 Set $\omega_0^k = 1$
- 4 For $l = 0, 1, \dots, L_{max}$ iterate
- 5 $x^{(k+1)} = x^{(k)} + \omega_l w^{(k)}$
- 6 Stop, if $\|f(x^{(k+1)})\| < \|f(x^{(k)})\|$
- 7 $\omega_{l+1}^{(k)} = \sigma \omega_l^{(k)}$

Using line search we enforce a monotonic contraction of the residuals. We first try a full Newton step with $\omega = 1$. As long as $\|f(x^{(k+1)})\| > \|f(x^{(k)})\|$, we reduce ω . Common values for σ are $\sigma = \frac{1}{2}$ or $\sigma = \frac{1}{4}$. Further modifications are described in numerical optimization [11] and also so-called error-oriented Newton methods going back to [4] and recently applied in coupled variational inequality systems in the monograph [19].

12.6 Newton's method for a coupled, nonlinear system of ODEs (complement)

We discuss an example motivated from the later Chapter 17. Therein, ordinary differential equations (ODEs) shall be solved at some time point t_n . Often, such systems result into nonlinear problems, which can be solved by methods introduced in the current chapter. We neglect all details concerning ODEs and refer the reader to Chapter 17, but concentrate on all relevant steps for the nonlinear solution.

We look for two unknowns $p(t) : [0, T] \rightarrow \mathbb{R}$ and $s(t) : [0, T] \rightarrow \mathbb{R}$, which are determined by a coupled nonlinear system of two ODEs:

$$\begin{aligned} p'(t) &= r_p p(t) - \alpha s(t)p(t) \\ s'(t) &= r_s s(t) + \mu \alpha p(t)s(t) \end{aligned}$$

The derivatives are defined with respect to time, i.e., $p'(t) := \frac{d}{dt}p(t)$ and $s'(t) := \frac{d}{dt}s(t)$. Furthermore, we deal with the parameters $r_p > 0, \alpha > 0, r_s > 0, \mu > 0$, which are independent of the time t .

When we discretize the time-continuous derivatives by an implicit procedure (taking the unknown as well on the right hand side), we obtain

$$\begin{aligned} \frac{p^n - p^{n-1}}{\Delta t} &= r_p p^n - \alpha s^n p^n \\ \frac{s^n - s^{n-1}}{\Delta t} &= r_s s^n + \mu \alpha p^n s^n, \end{aligned}$$

where $\Delta = t_n - t_{n-1}$ is the time step size (difference between two time points). Since we look simultaneously for the time-discretized p^n and s^n , we easily see that the system is nonlinear due to the couplings $s^n p^n$ (1st equation) and $p^n s^n$ (2nd equation).

One possibility to solve this nonlinear system is to formulate a root-finding problem of the form

$$F(p^n, s^n) = 0$$

with

$$F(p^n, s^n) = \begin{pmatrix} p^n - p^{n-1} - \Delta t[r_p p^n - \alpha s^n p^n] \\ s^n - s^{n-1} - \Delta t[r_s s^n + \mu \alpha p^n s^n] \end{pmatrix}.$$

For Newton's method, we introduce the iteration index j . Formally, we have: set $p^{n,0} := p^{n-1}$ and $s^{n,0} := s^{n-1}$ as initial guess. Then, a defect-correction variant reads:

$$\begin{aligned} F'(p^{n,j}, s^{n,j})(w, v)^T &= -F(p^{n,j}, s^{n,j}) \\ (p^{n,j+1}, s^{n,j+1})^T &= (p^{n,j}, s^{n,j})^T + (w, v)^T \end{aligned}$$

until some stopping criterion, e.g.,

$$\max(\|p^{n,j+1} - p^{n,j}\|_2, \|s^{n,j+1} - s^{n,j}\|_2) < TOL$$

with for instance $TOL = 1e - 10$. The derivative F' is the Jacobian matrix; here $F' \in \mathbb{R}^{2 \times 2}$. The crucial point is how the Jacobian reads in detail. Formally, we have

$$F'(p^n, s^n) = \begin{pmatrix} \partial_p F_1, \partial_s F_1 \\ \partial_p F_2, \partial_s F_2 \end{pmatrix}$$

Specifically, for the above system, we obtain

$$\begin{aligned} \partial_p F_1 &= 1 - 0 - \Delta t[r_p - \alpha s^n] \\ \partial_s F_1 &= 0 - 0 - \Delta t[-\alpha p^n] \\ \partial_p F_2 &= 0 - 0 - \Delta t[\mu \alpha s^n] \\ \partial_s F_2 &= 1 - 0 - \Delta t[r_s + \mu \alpha p^n] \end{aligned}$$

The full Newton system at a time point t^n for the iteration j then reads:

$$\begin{pmatrix} 1 - \Delta t[r_p - \alpha s^{n,j}] & -\Delta t[-\alpha p^{n,j}] \\ -\Delta t[\mu \alpha s^{n,j}] & 1 - \Delta t[r_s + \mu \alpha p^{n,j}] \end{pmatrix} \begin{pmatrix} w \\ v \end{pmatrix} = - \begin{pmatrix} p^{n,j} - p^{n-1} - \Delta t[r_p p^{n,j} - \alpha s^{n,j} p^{n,j}] \\ s^{n,j} - s^{n-1} - \Delta t[r_s s^{n,j} + \mu \alpha p^{n,j} s^{n,j}] \end{pmatrix}$$

$$\begin{pmatrix} p^{n,j+1} \\ s^{n,j+1} \end{pmatrix} = \begin{pmatrix} p^{n,j} \\ s^{n,j} \end{pmatrix} + \begin{pmatrix} w \\ v \end{pmatrix}.$$

Until

$$\left\| \begin{pmatrix} p^{n,j+1} \\ s^{n,j+1} \end{pmatrix} - \begin{pmatrix} p^{n,j} \\ s^{n,j} \end{pmatrix} \right\|_2 < TOL$$

With this procedure, we then have obtained the solution for the current time point and set

$$p^n := p^{n,j+1}, \quad s^n := s^{n,j+1}.$$

Afterward, we can proceed to the next time point t_{n+1} ; for details we refer to Chapter 17.

We finally notice as well that a simplified (approximate) Newton scheme is obtained by neglecting the off-diagonal terms. Then, the defect step reads

$$\begin{pmatrix} 1 - \Delta t[r_p - \alpha s^{n,j}] & 0 \\ 0 & 1 - \Delta t[r_s + \mu \alpha p^{n,j}] \end{pmatrix} \begin{pmatrix} w \\ v \end{pmatrix} = - \begin{pmatrix} p^{n,j} - p^{n-1} - \Delta t[r_p p^{n,j} - \alpha s^{n,j} p^{n,j}] \\ s^{n,j} - s^{n-1} - \Delta t[r_s s^{n,j} + \mu \alpha p^{n,j} s^{n,j}] \end{pmatrix}$$

12.7 Iteration schemes in nonlinear optimization (complement)

12.7.1 Linear regression

A typical situation for approximation problems is linear regression. Here, often given data in form of (x_n, p_n) for $n = 1, \dots, N$ is available and a function of the form

$$p(x) = b + mx$$

is sought that best approximates in some mathematical sense (choice of the norm!) these data. The specific goal is to determine the parameters $b \in \mathbb{R}$ and $m \in \mathbb{R}$.

Here, a least-squares approximation is one possibility:

$$S := S(m, b) = \frac{1}{N} \sum_{n=1}^N (p_n - (b - mx_n))^2.$$

The function S is the so-called cost (or objective) functional. The goal is solution of the minimization problem:

$$\min_{m, b} S(m, b).$$

We also refer the reader back to chapter 11 for an introduction to nonlinear optimization. For the numerical solution, we present two options.

The first one is the steepest descent method. We just need to calculate the two partial derivatives with respect to m and b :

$$\frac{\partial S}{\partial m} = \frac{2}{N} \sum_{k=1}^N -x_k (p_k - (mx_k + b)) \quad (12.11)$$

$$\frac{\partial S}{\partial b} = \frac{2}{N} \sum_{k=1}^N -(p_k - (mx_k + b)). \quad (12.12)$$

The final algorithm is then given by: for $l = 1, 2, 3, \dots, N_{max}$, iterate such that

$$\begin{pmatrix} m_{l+1} \\ b_{l+1} \end{pmatrix} = \begin{pmatrix} m_l \\ b_l \end{pmatrix} - \rho \begin{pmatrix} \frac{\partial S(m_l, b_l)}{\partial m} \\ \frac{\partial S(m_l, b_l)}{\partial b} \end{pmatrix}.$$

The second option consists in exploiting the fact the cost function is quadratic, hence the first order optimality condition

$$\nabla S(m, b) = 0$$

is a simple linear system given by (12.11)-(12.12). It can be solved directly.

12.7.2 Neural networks

In the solution of machine learning problems with the area of neural networks, the task is mathematically exactly the same as in the previous subsection, except that linearity is lost. For a concise introduction from a mathematical viewpoint, we refer to [7]. We need to determine weights w and biases b . Determining those two quantities is called training. Let us put both of them into a function $a^{[L]} := a^{[L]}(w, b)$, which is the global activation function at the output layer L . Let N pieces of data (training points) and corresponding target outputs be denoted as $(x^n, y^n), n = 1, \dots, N$. Then, a quadratic cost functional (as before) can be defined:

$$C = \frac{1}{N} \sum_{n=1}^N \frac{1}{2} \|y^n - a^{[L]}(x^n)\|_2^2.$$

This problem can be solved as in the previous section by gradient descent. However, since usually N is very big (i.e., a big data problem) the computational cost may be prohibitively high. For this reason, a so-called stochastic gradient method is often employed in which the individual gradients running over all training points N is replaced by a single randomly chosen, training point. There is a vast literature on machine learning and neural networks because of its enormous applications in science, technical applications, and our society. The important message in the frame of this class is that the concepts, and numerical methods presented here, are already the key ingredients of algorithms used for machine learning.

Let us now have a look at the activation functions that are classically used. To fix ideas, consider a single layer ($L = 1$). Typically, writing as subscripts the components of vectors, one considers

$$(a^{[L]}(x))_i = S \left(\sum_{j=1}^M w_{ij} x_j - b_i \right)$$

with the sigmoid function

$$S(t) = \frac{1}{1 + e^{-t}},$$

see Fig. 12.1. This function approximates the thresholding phenomenon that occurs in natural neurons, while being differentiable (very helpful for optimization!).

In the above notation M is the number of inputs, and i covers the number of outputs.

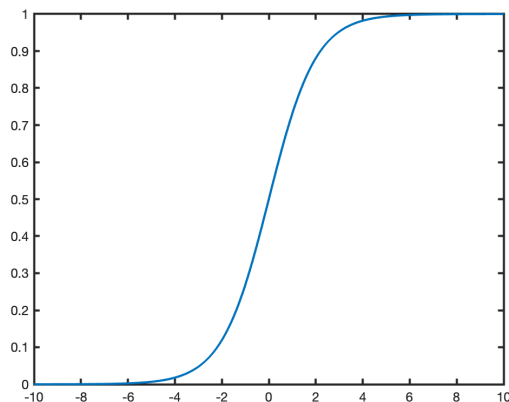


Figure 12.1: Sigmoid function.

We can incorporate intermediate or hidden layers just by composing such operations. The gradient is obtained by the chain rule, and performing Jacobian matrix multiplications in a clever ordering (actually from output to input) is called error back-propagation. Machine learning and neural networks are very active research fields for which we refer to other teaching materials for an introduction and further references [9].

Chapter 13

Interpolation and approximation

In this chapter, we discuss interpolation and approximation. The former consists in constructing a simple function (often a polynomial) through certain given points. With the help of such simple functions further operations such as numerical differentiation and numerical integration can be constructed. Approximation is more general and does not necessarily need to fit in certain given points. One example is linear regression.

13.1 Polynomial interpolation

13.1.1 Introduction

Let (x_0, \dots, x_n) and (y_0, \dots, y_n) be two families of real numbers, with the so-called nodes x_i 's pairwise distinct. We denote by \mathcal{P}_n the set of polynomial functions of degree $\leq n$. We investigate the following issues:

1. Does there exist a polynomial $P \in \mathcal{P}_n$ such that $P(x_i) = y_i \forall i = 0, \dots, n$? What about uniqueness?
2. How to calculate it?
3. How does it behave in between the points x_i ?

13.1.2 Existence, uniqueness, expression

Theorem 13.1. *There exists a unique $P \in \mathcal{P}_n$ such that $P(x_i) = y_i \forall i = 0, \dots, n$. It is called Lagrange interpolating polynomial.*

Let us write the interpolating polynomial as $P(x) = \sum_{k=0}^n a_k x^k$. The coefficients must solve the linear system

$$\begin{cases} a_0 + a_1 x_0 + \dots + a_n x_0^n = y_0 \\ a_0 + a_1 x_1 + \dots + a_n x_1^n = y_1 \\ \vdots \\ a_0 + a_1 x_n + \dots + a_n x_n^n = y_n, \end{cases}$$

which admits the matrix representation

$$\begin{pmatrix} 1 & x_0 & \dots & x_0^n \\ 1 & x_1 & \dots & x_1^n \\ \vdots & \vdots & & \vdots \\ 1 & x_n & \dots & x_n^n \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{pmatrix}.$$

The determinant of the matrix of this system, called Vandermonde determinant (a particular case was given as exercise), is equal to

$$D_n = \prod_{0 \leq i < j \leq n} (x_j - x_i) \neq 0.$$

Hence this system admits a unique solution. In fact, there is a direct expression of the interpolating polynomial, given by

$$P(x) = \sum_{j=0}^n y_j \varphi_j(x), \quad \varphi_j(x) = \prod_{k \neq j} \frac{x - x_k}{x_j - x_k}.$$

Indeed, the functions $\varphi_0, \dots, \varphi_n$ belong to \mathcal{P}_n and satisfy $\varphi_i(x_j) = \delta_{ij}$, with $\delta_{ij} = 1$ if $i = j$, $\delta_{ij} = 0$ if $i \neq j$ (Kronecker symbol).

An example is shown in Fig. 13.1.

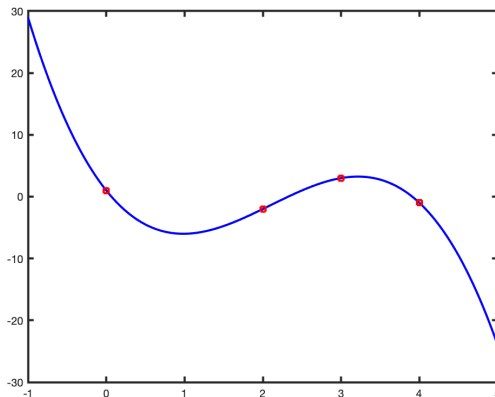


Figure 13.1: Interpolating polynomial of degree 3. The data points are the red squares.

13.1.3 Interpolation error (complement)

For a given function f we now assume that the data points are such that $y_j = f(x_j) \forall j = 0, \dots, n$. A natural question is then: how does the Lagrange interpolating polynomial of the families (x_0, \dots, x_n) and $(f(x_0), \dots, f(x_n))$, called Lagrange interpolating polynomial of f , approximate f ? An answer is provided by the following theorem.

Theorem 13.2. *Let f be a function of class \mathcal{C}^{n+1} on $[a, b]$ and P be its Lagrange interpolating polynomial at the points $x_0, x_1, \dots, x_n \in [a, b]$. Then, for all $x \in [a, b]$, there exists $\xi \in [a, b]$ such that*

$$f(x) - P(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \pi_{n+1}(x), \quad \pi_{n+1} = \prod_{j=0}^n (x - x_j).$$

We obviously always have $|\pi_{n+1}(x)| \leq (b-a)^{n+1}$. If the points are equidistant, i.e. $x_j = a + jh$, $h = (b-a)/n$, then we have the finer estimate $|\pi_{n+1}(x)| \leq n!h^{n+1}/4$. This yields

$$|f(x) - P(x)| \leq \max_{\xi \in [a, b]} |f^{(n+1)}(\xi)| \frac{h^{n+1}}{4(n+1)}.$$

We could expect that $f(x) - P(x) \rightarrow 0$ when $n \rightarrow +\infty$. However this is not always the case since the derivatives $|f^{(n+1)}(\xi)|$ may increase quickly when n grows. It can even happen that $|f(x) - P(x)| \rightarrow +\infty$: this is called the Runge phenomenon, due to oscillations of P between the points x_i .

A classical example is the function $f : t \mapsto \frac{1}{1+25t^2}$ on $[-1, 1]$, see figure 13.2.

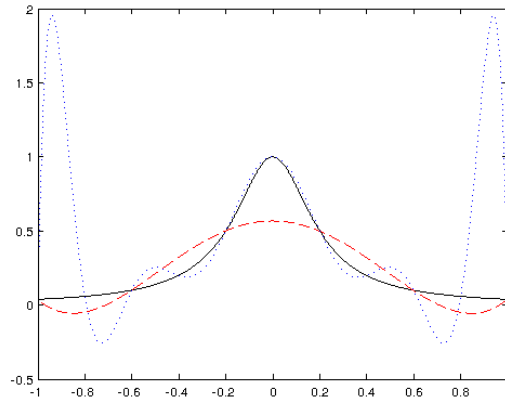


Figure 13.2: Runge phenomenon. The function f is depicted in solid line. The Lagrange interpolating polynomials of f of degree 5 and 10 with equidistant points are depicted in dashed and dotted lines, respectively.

13.2 Numerical differentiation

13.2.1 Approximation of the first derivative

Consider a "sufficiently smooth" function $f : \mathbb{R} \rightarrow \mathbb{R}$ and two prescribed points $x_0, x_1 \in \mathbb{R}$ where f takes known values. Hence we can interpolate f linearly at the nodes x_0 and x_1 by

$$p_1(x) = \frac{x - x_1}{x_0 - x_1} f(x_0) + \frac{x - x_0}{x_1 - x_0} f(x_1).$$

Now we can approximate the derivative $f'(x)$ through the corresponding derivative of the polynomial

$$p_1'(x) = \frac{f(x_1) - f(x_0)}{x_1 - x_0}.$$

Proposition 13.3 (Difference quotient of the first derivative). *Let $f \in C^2[a, b]$. For two points $x_0, x_1 \in [a, b]$ with $h := x_1 - x_0 > 0$ we have the left-sided and right-sided difference quotients*

$$\begin{aligned} \frac{f(x_1) - f(x_0)}{h} &= f'(x_0) + \frac{1}{2} h f''(x_0) + O(h^2), \\ \frac{f(x_1) - f(x_0)}{h} &= f'(x_1) - \frac{1}{2} h f''(x_1) + O(h^2). \end{aligned}$$

For $f \in C^4[a, b]$, it holds $x_0, x_0 - h, x_0 + h \in [a, b]$ for the central difference quotient

$$\frac{f(x_0 + h) - f(x_0 - h)}{2h} = f'(x_0) + \frac{1}{4} h^2 f^{(iv)}(x_0) + O(h^3).$$

Proof. (i) With Taylor expansion, it holds

$$f(x_0 \pm h) = f(x_0) \pm h f'(x_0) + \frac{1}{2} h^2 f''(x_0) \pm \frac{1}{6} h^3 f'''(x_0) + \frac{1}{24} h^4 f^{(iv)}(\xi)$$

with some $\xi \in [a, b]$. It follows for the right-sided difference quotient

$$\frac{f(x_0 + h) - f(x_0)}{h} = \frac{h f'(x_0) + \frac{1}{2} h^2 f''(x_0) + O(h^3)}{h} = f'(x_0) + \frac{1}{2} h f''(x_0) + O(h^2).$$

The proof for the left-sided difference quotient is similar.

(ii) For the central difference quotient, the even powers of h cancel due to symmetry:

$$f(x_0 + h) - f(x_0 - h) = 2f'(x_0) + \frac{1}{3}h^3 f'''(x_0) + O(h^4).$$

This yields the final result. \square

We see in the previous proofs that the choice of the nodal values plays a crucial role. If they are chosen in some clever way, we obtain better convergence than expected. This phenomenon is known as super-convergence, which is the case for the central difference quotient. The key point is symmetry. This can be seen as follows: let the derivative in x_0 be constructed through $x_0 - h$ and $x_0 + 2h$ (no symmetry). Then, Taylor around x_0 yields

$$\begin{aligned} \frac{f(x_0 + 2h) - f(x_0 - h)}{3h} &= f'(x_0) + \left(2hf''(x_0) + \frac{8}{6}h^2 f'''(x_0)\right) - \left(\frac{1}{2}hf''(x_0) - \frac{1}{6}h^2 f'''(x_0)\right) + O(h^3) \\ &= f'(x_0) + \frac{3}{2}hf''(x_0) + O(h^2). \end{aligned}$$

In this case, higher order powers of h do not cancel and we only obtain linear convergence, but not quadratic.

13.2.2 Approximation of the second derivative (complement)

We interpolate f with the help of three points $x_0 - h$, x_0 , $x_0 + h$ using a quadratic polynomial:

$$\begin{aligned} p_2(x) &= f(x_0 - h) \frac{(x_0 - x)(x_0 + h - x)}{2h^2} \\ &\quad + f(x_0) \frac{(x - x_0 + h)(x_0 + h - x)}{h^2} + f(x_0 + h) \frac{(x - x_0 + h)(x - x_0)}{2h^2}. \end{aligned}$$

In $x = x_0$ it holds for the first and second derivative:

$$p_2'(x_0) = \frac{f(x_0 + h) - f(x_0 - h)}{2h}, \quad p_2''(x_0) = \frac{f(x_0 + h) - 2f(x_0) + f(x_0 - h)}{h^2}.$$

For the first derivative, we obtain the same result as in the previous subsection.

With Taylor, we obtain for the second derivative

$$p_2''(x_0) = \frac{-2f(x_0) + f(x_0 - h) + f(x_0 + h)}{h^2} = f''(x_0) + \frac{1}{12}h^2 f^{(iv)}(x_0) + O(h^4).$$

Using p_2 we can also derive a one-sided difference quotient for the second derivative. Thus, we approximate $f''(x_0 - h) \approx p_2''(x_0 - h)$. With Taylor, we obtain

$$p_2''(x_0 - h) = f''(x_0 - h) + hf'''(x_0 - h) + O(h^2),$$

which is only an approximation of first order. Again, the choice of good support points is crucial.

13.3 Numerical integration (complement)

13.3.1 Goal

Let $f : [a, b] \rightarrow \mathbb{R}$ be a continuous function. Our goal is to calculate

$$\int_a^b f(x) dx.$$

In general, we do not know any explicit function F such that $F' = f$. Then we have to content ourselves with an approximation of the integral.

13.3.2 General principle

We first introduce a subdivision $\sigma = (x_0, x_1, \dots, x_N)$ of $[a, b]$ (i.e., $x_0 = a < x_1 < \dots < x_N = b$). The stepsize of the subdivision σ is denoted by :

$$h = \max_{0 \leq k \leq N-1} |x_{k+1} - x_k|.$$

Obviously the integral can be decomposed as

$$\int_a^b f(x) dx = \sum_{k=0}^{N-1} \int_{x_k}^{x_{k+1}} f(x) dx. \quad (13.1)$$

Each integral $\int_{x_k}^{x_{k+1}} f(x) dx$ will be approximated by a so-called elementary quadrature formula.

Note that the change of variable

$$x = x_k + \frac{x_{k+1} - x_k}{2}(t + 1), \quad t \in [-1, 1],$$

yields

$$\int_{x_k}^{x_{k+1}} f(x) dx = \frac{x_{k+1} - x_k}{2} \int_{-1}^1 g_k(t) dt \quad (13.2)$$

with $g_k(t) = f(x_k + \frac{x_{k+1} - x_k}{2}(t + 1))$. Therefore it is sufficient to define elementary quadrature formulas for the standard interval $[-1, 1]$.

13.3.3 Elementary quadrature formula

Let $g : [-1, 1] \rightarrow \mathbb{R}$ be a continuous function. Our goal is to approximate $\int_{-1}^1 g(x) dx$ by a number $J(g)$.

Definition 13.1. *An elementary quadrature formula is defined by*

- points t_0, \dots, t_L such that $-1 \leq t_0 < t_1 < \dots < t_L \leq 1$ called nodes,
- coefficients $w_0, w_1, \dots, w_L \in \mathbb{R}$ called weights.

It reads

$$J(g) = \sum_{j=0}^L w_j g(t_j).$$

13.3.4 Simplest rules: box, mid-point, trapezoidal, Simpson

Let us work on the interval $[a, b]$. The simplest quadrature rules are the left-sided box rule

$$I^0(f) = (b - a)f(a),$$

the right-sided box rule

$$I^0(f) = (b - a)f(b),$$

the mid-point rule

$$I^0(f) = (b - a)f\left(\frac{a + b}{2}\right),$$

the trapezoidal rule

$$I^1(f) = \frac{b - a}{2}(f(a) + f(b)),$$

and the Simpson rule

$$I^2(f) = \frac{b - a}{6}(f(a) + 4f\left(\frac{a + b}{2}\right) + f(b)).$$

13.3.5 Composite quadrature formula

In view of (13.1) and (13.2) it is natural to set the following definition.

Definition 13.2. *The composite quadrature formula associated with the subdivision σ and the elementary quadrature formula $J(g) = \sum_{j=0}^L w_j g(t_j)$ is :*

$$I_{\sigma,J}(f) = \sum_{k=0}^{N-1} \frac{x_{k+1} - x_k}{2} \sum_{j=0}^L w_j f(\xi_{k,j}) \quad \text{with } \xi_{k,j} = x_k + \frac{x_{k+1} - x_k}{2} (t_j + 1).$$

Here are classical examples.

	elementary formula $J(g)$	corresponding composite formula $I_{\sigma,J}(f)$
leftpoint formula	$2g(-1)$	$\sum_{k=0}^{N-1} (x_{k+1} - x_k) f(x_k)$
rightpoint formula	$2g(1)$	$\sum_{k=0}^{N-1} (x_{k+1} - x_k) f(x_{k+1})$
midpoint formula	$2g(0)$	$\sum_{k=0}^{N-1} (x_{k+1} - x_k) f\left(\frac{x_k + x_{k+1}}{2}\right)$
trapezoidal formula	$g(-1) + g(1)$	$\sum_{k=0}^{N-1} \frac{x_{k+1} - x_k}{2} (f(x_k) + f(x_{k+1}))$

13.3.6 Order of a quadrature formula

Definition 13.3. *An elementary quadrature formula J is said to be of order at least n if it is exact for polynomials in \mathcal{P}_n , i.e.,*

$$J(g) = \int_{-1}^1 g(t) dt \quad \forall g \in \mathcal{P}_n.$$

It is of order exactly n if, in addition, there exists $g \in \mathcal{P}_{n+1}$ such that

$$J(g) \neq \int_{-1}^1 g(t) dt.$$

For the previous examples it holds:

	exact order
left/right-point formula	0
midpoint formula	1
trapezoidal formula	1

It might seem surprising that the midpoint and trapezoidal formulas have the same order (super-convergence). This is due to the fact that the midpoint is at a special location that enables to gain orders for a fixed number of nodes. Such a node is called a Gauss point. We will not enter into this (beautiful and classical) theory.

13.3.7 Interpolatory quadrature

All examples of quadrature formulas we have mentioned so far belong to the class of interpolatory quadrature formulas, which is the standard way of constructing quadrature formulas.

Definition 13.4. *An interpolatory quadrature formula is an elementary quadrature formula such that*

$$J(g) = \int_{-1}^1 P_g(t) dt,$$

where P_g is the Lagrange interpolating polynomial of g at the points t_0, \dots, t_N .

Then, from

$$P_g(t) = \sum_{j=0}^L g(t_j) \varphi_j(t), \quad \varphi_j(t) = \prod_{k \neq j} \frac{t - t_k}{t_j - t_k},$$

we obtain the expression

$$J(g) = \sum_{j=0}^L w_j g(t_j), \quad w_j = \int_{-1}^1 \varphi_j(t) dt.$$

By construction, the formula is of order at least L . It is easy to prove the reciprocal statement.

Theorem 13.4. *The elementary quadrature formula $J(g) = \sum_{j=0}^L w_j g(t_j)$ is of order at least L if and only if it is an interpolatory quadrature formula.*

13.3.8 Gauss quadrature

We finally discuss a class of quadrature rules that are often used in practice. So far, the order (of convergence) of a quadrature rule was determined through the number of support points and the polynomial function that was constructed using these points. The question is whether the quadrature order can be further enhanced or whether our developments are optimal.

A careful revision entails that the mid-point rule has a better convergence order than expected despite the fact that the number of support points is the same as for the box rule. The reason is symmetry (see also the section on difference quotients; super-convergence).

This raises the question whether this is a general principle that an optimal distribution of the support points enhances the convergence order. The answer is positive and leads to the so-called Gauss quadrature.

So far a quadrature rule for $n+1$ support points $x_0, \dots, x_n \in [a, b]$ had the order $n+1$. This means that polynomials with degree n are exactly integrated. For instance, the box rule integrates constant polynomials exactly (degree $n=0$). The trapezoidal rule integrates linear polynomials exactly (degree $n=1$).

As intermediate summary, we note: so far the order of a quadrature rule for $n+1$ support points is $n+1$.

An optimal distribution (Gauss quadrature) yields however the much better result:

Proposition 13.5. *An interpolatory quadrature rule with $n+1$ support points has at most the order $2n+2$.*

In practice the key question is how such quadrature rules can be constructed. This can be achieved via orthogonalization. More information about the construction can be found in standard numerical analysis books.

Chapter 14

Trigonometric interpolation, Fourier series and Fourier transform

14.1 Trigonometric interpolation: discrete Fourier transform

We address here a similar problem as in section 13.1, but instead of classical algebraic polynomials we consider trigonometric polynomials. Such polynomials are better suited to the approximation of periodic functions, and more generally oscillatory phenomena.

Let a function $f : [0, 2\pi] \rightarrow \mathbb{C}$ such that $f(0) = f(2\pi)$ be given. A trigonometric polynomial is a function of form

$$P(x) = \sum_{k=-M}^M c_k e^{ikx},$$

for some $M \in \mathbb{N}$ and coefficients $c_{-M}, \dots, c_M \in \mathbb{C}$. Note that P is by construction continuous and 2π -periodic on \mathbb{R} , this is why we have assumed $f(0) = f(2\pi)$. Obviously, the period 2π is chosen for convenience and can be changed by rescaling.

We now define the equidistant nodes $x_0, \dots, x_n \in [0, 2\pi[$ by

$$x_j = \frac{2j\pi}{n+1}.$$

Our goal is to determine P in order to interpolate f at these nodes, i.e., to satisfy

$$P(x_j) = f(x_j) \quad \forall j = 0, \dots, n. \quad (14.1)$$

This yields $n+1$ equations, and since there are $2M+1$ coefficients to fix, it is natural to assume that

$$n = 2M.$$

If n is odd, it is possible to choose $M = [n/2] + 1$ ($[x]$ stands for the greatest integer $\leq x$) and impose an additional condition on the coefficients, see e.g. [13]. For simplicity, we assume that n is even and fix $M = n/2$.

We have the following result.

Proposition 14.1. *The relations (14.1) are satisfied if and only if*

$$c_k = \frac{1}{n+1} \sum_{j=0}^n f(x_j) e^{-ikx_j}, \quad k = -M, \dots, M. \quad (14.2)$$

Proof. Suppose that (14.1) holds true. Multiplying (14.1) by e^{-imx_j} and summing over j yields

$$\sum_{j=0}^n \sum_{k=-M}^M c_k e^{i(k-m)x_j} = \sum_{j=0}^n e^{-mx_j} f(x_j).$$

We can interchange the summations to obtain

$$\sum_{k=-M}^M c_k \sum_{j=0}^n e^{i(k-m)x_j} = \sum_{j=0}^n e^{-mx_j} f(x_j). \quad (14.3)$$

Now, we remark that

$$\sum_{j=0}^n e^{i(k-m)x_j} = \sum_{j=0}^n e^{i(k-m)j \frac{2\pi}{n+1}} = \sum_{j=0}^n \left(e^{i(k-m) \frac{2\pi}{n+1}} \right)^j.$$

This is a geometric sum, equal to

$$\frac{n+1}{1 - e^{i(k-m) \frac{2\pi}{n+1}}} = \frac{1 - e^{i(k-m)2\pi}}{1 - e^{i(k-m) \frac{2\pi}{n+1}}} = 0 \quad \text{if } k \neq m.$$

We infer from (14.3) the claimed expression

$$(n+1)c_m = \sum_{j=0}^n e^{-mx_j} f(x_j).$$

Suppose now (14.2) holds true. We compute

$$P(x_l) = \sum_{k=-M}^M c_k e^{ikx_l} = \sum_{k=-M}^M \frac{1}{n+1} \sum_{j=0}^n f(x_j) e^{-ikx_j} e^{ikx_l} = \frac{1}{n+1} \sum_{j=0}^n f(x_j) \sum_{k=-M}^M e^{ik(x_l - x_j)}.$$

By the same arguments as previously the inner sum is equal to $2M+1 = n+1$ if $j = l$, 0 else. This entails (14.1). \square

The family of coefficients (c_{-M}, \dots, c_M) is called the discrete Fourier transform of f relatively to the nodes x_0, \dots, x_n , and the resulting trigonometric polynomial P is the interpolating trigonometric polynomial of f .

There exist fast algorithms to compute these coefficients. They are termed fast Fourier transform (in short FFT) algorithms. One of them (the Cooley-Tukey algorithm) can be studied as project (see Chapter 19).

Remark 14.1. If f is real-valued, we infer from (14.2) that $c_{-k} = \overline{c_k}$. The Euler formula (Chapter 2) entails

$$P(x) = \frac{a_0}{2} + \sum_{k=1}^M [a_k \cos(kx) + b_k \sin(kx)],$$

with the real coefficients

$$a_k = c_k + c_{-k} = 2\Re(c_k) = \frac{2}{n+1} \sum_{j=0}^n f(x_j) \cos(kx_j),$$

$$b_k = i(c_k - c_{-k}) = -2\Im(c_k) = \frac{2}{n+1} \sum_{j=0}^n f(x_j) \sin(kx_j).$$

An example of trigonometric interpolation with $n = 4$, $M = 2$, is shown in Fig. 14.1.

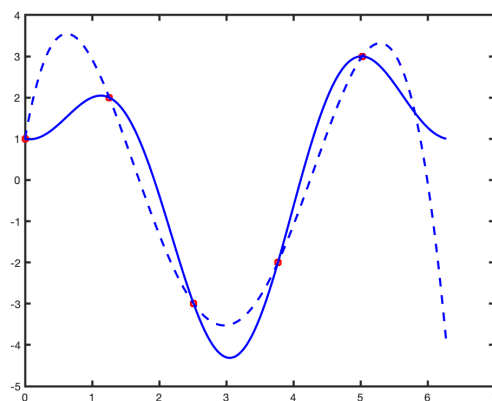


Figure 14.1: Trigonometric interpolation (solid line) vs polynomial interpolation (dashed line). The data points are the red squares. Observe that the trigonometric interpolation extends to a periodic smooth (C^∞) function, unlike the polynomial interpolation.

14.2 Fourier series

In the above considerations, assume that M (thus also n) goes to infinity. Formally it appears that

$$P(x) \rightarrow \sum_{k=-\infty}^{+\infty} c_k e^{ikx}$$

and, recognizing the leftpoint quadrature formula with stepsize $x_{j+1} - x_j = 2\pi/(n+1)$,

$$c_k \rightarrow \frac{1}{2\pi} \int_0^{2\pi} f(x) e^{-ikx} dx.$$

These ideas admit a strong mathematical background: the Fourier series theory.

14.2.1 Series

Definition 14.1. A series is an expression of form

$$S = \sum_{k=0}^{+\infty} a_k, \quad \text{meaning exactly} \quad S = \lim_{n \rightarrow +\infty} \sum_{k=0}^n a_k.$$

By extension,

$$S = \sum_{k=-\infty}^{+\infty} a_k, \quad \text{means} \quad S = \lim_{n \rightarrow +\infty} \sum_{k=-n}^n a_k.$$

For series of functions, several types of convergence can be considered. Here we will only speak of pointwise convergence, i.e., convergence of the series of numbers when the function is evaluated at a fixed (but arbitrary) point.

14.2.2 Convergence of Fourier series

For the pointwise convergence of Fourier series one has the following result, known as Dirichlet theorem.

Theorem 14.2. Let $f : \mathbb{R} \rightarrow \mathbb{C}$ be a continuous and 2π -periodic function. Let $x \in \mathbb{R}$ such that f admits a left and right derivative at x . We have

$$f(x) = \sum_{k=-\infty}^{+\infty} c_k e^{ikx} \tag{14.4}$$

with

$$c_k = \frac{1}{2\pi} \int_0^{2\pi} f(x) e^{-ikx} dx. \quad (14.5)$$

The expression (14.4) is called Fourier series of f and the c_k 's from (14.5) are called Fourier coefficients of f .

Remark 14.2. Similarly to remark 14.1, if f is real-valued, we have $c_{-k} = \overline{c_k}$, whereby

$$f(x) = \frac{a_0}{2} + \sum_{k=1}^{+\infty} [a_k \cos(kx) + b_k \sin(kx)],$$

with

$$a_k = 2\Re(c_k) = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos(kx) dx, \quad b_k = -2\Im(c_k) = \frac{1}{\pi} \int_0^{2\pi} f(x) \sin(kx) dx.$$

Remark 14.3. If f is only piecewise continuous then we have at a discontinuity point x_0

$$\sum_{k=-\infty}^{+\infty} c_k e^{ikx_0} = \frac{1}{2} \left(\lim_{x \rightarrow x_0^-} f(x) + \lim_{x \rightarrow x_0^+} f(x) \right).$$

14.2.3 Parseval's equality

Parseval's equality states that the Fourier coefficients carry the 'energy' of the function.

Theorem 14.3. Let $f : \mathbb{R} \rightarrow \mathbb{C}$ be a continuous and 2π -periodic function (square integrable on $[0, 2\pi]$ is actually sufficient) and call c_k its Fourier coefficients. We have

$$\frac{1}{2\pi} \int_0^{2\pi} |f(x)|^2 dx = \sum_{k=-\infty}^{+\infty} |c_k|^2.$$

14.2.4 Applications

Fourier series are a very important tool in mathematical analysis and signal processing. The main reason of their usefulness is that the basis functions

$$\varphi_k(x) = e^{ikx}$$

are 'eigenfunctions' for the differentiation operator: indeed it holds $\varphi_k' = ik\varphi_k$. Therefore they appear as special solutions of many differential equations. They are also orthogonal, in the sense that

$$\int_0^{2\pi} \varphi_k(x) \varphi_l(x) dx = 0 \quad \forall k \neq l.$$

It is thus natural to decompose a given function on this 'basis' (to be rigorous the notion of Hilbert basis must be considered, allowing for infinite linear combinations).

14.3 Fourier transform

14.3.1 General concept

Fourier series admit a generalization to non-periodic functions defined over the whole real line: the Fourier transform. In this case the frequencies are no longer discrete, and the series becomes an integral. When the integrals below are well defined ¹ one has

$$\hat{f}(\omega) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} f(x) e^{-i\omega x} dx, \quad f(x) = \int_{-\infty}^{+\infty} \hat{f}(\omega) e^{i\omega x} d\omega.$$

¹Basically one needs that $\int_{-\infty}^{+\infty} |f(x)| dx < +\infty$ and $\int_{-\infty}^{+\infty} |\hat{f}(\omega)| d\omega < +\infty$, however, using generalized integrals, $\int_{-\infty}^{+\infty} |f(x)|^2 dx < +\infty$ works also.

We call \hat{f} the Fourier transform of f , and the latter equality is known as the Fourier inversion formula. In some cases the above expressions do not have a classical meaning and can be understood in the framework of distribution theory. The Fourier inversion formula holds true for what we call 'tempered distributions'. This is a fairly involved theory that we will not develop. For practical use it is enough to handle the prototype tempered distribution: the so-called Dirac distribution. It is described in the last paragraph.

Let us go back to classical functions and state two properties.

- If f is square integrable over \mathbb{R} , i.e. $\int_{-\infty}^{+\infty} |f(x)|^2 dx < +\infty$, then Parseval's equality holds:

$$\frac{1}{2\pi} \int_{-\infty}^{+\infty} |f(x)|^2 dx = \int_{-\infty}^{+\infty} |\hat{f}(\omega)|^2 d\omega.$$

- The Fourier transform of a derivative is formally obtained through integration by parts: if f is differentiable and f, f' admit Fourier transforms, then

$$\widehat{f'}(\omega) = i\omega \hat{f}(\omega).$$

This property is also recovered by differentiating the inversion formula. Due to this, as for Fourier series, the Fourier transform is very useful to construct special solutions to differential equations.

Remark 14.4. *There is a noticeable symmetry between the definition of the Fourier transform and the inversion formula. Actually, passing to the conjugate in this latter yields*

$$\overline{f(x)} = \int_{-\infty}^{+\infty} \overline{\hat{f}(\omega)} e^{-i\omega x} d\omega = 2\pi \widehat{\widehat{f}}(x). \quad (14.6)$$

In particular, if f is real-valued as well as its Fourier transform, then $f = 2\pi \widehat{\widehat{f}}$. Note that placing the normalization coefficient $1/2\pi$ in the definition of the Fourier transform rather than in the inversion formula is a matter of (non-universal) convention. Sometimes, $1/\sqrt{2\pi}$ is placed in front of both.

14.3.2 A glimpse on the Dirac distribution

As mentioned above, it is sometimes convenient to apply the Fourier transform to distributions rather than functions. We will not give a rigorous definition of the concept of distribution, we will content ourselves with the formal description of the most classical one: the Dirac distribution. Intuitively, the Dirac distribution δ is a generalization of a function that is equal to 0 everywhere except at point 0, and with 'integral' equal to 1. More precisely, it satisfies

$$\int_{-\infty}^{+\infty} \delta(x)\varphi(x)dx = \varphi(0)$$

for all smooth function φ . Indeed, choosing this (so-called test function) φ identically equal to 1 yields $\int_{-\infty}^{+\infty} \delta(x)dx = 1$, while choosing φ arbitrary but vanishing at 0 yields $\delta(x) = 0$ for all $x \neq 0$. It is to be stressed that the integral notation used above is abusive in full rigor, but it suggests to a large extent correct ideas and calculations. An intuitive understanding of this concept is to see the Dirac distribution as the limit of a sequence of nonnegative functions of integral 1 that concentrate at 0, see Fig. 14.2.

We obtain from the definition and the inversion formula in the form (14.6)

$$\hat{\delta} = \frac{1}{2\pi}, \quad \widehat{\widehat{1}} = \delta.$$

One may remark that the function 1 is neither integrable nor square integrable over \mathbb{R} (in fact, it is another example of tempered distribution). Its Fourier transform has been determined so as to satisfy the inversion formula, which is the main goal of the construction.

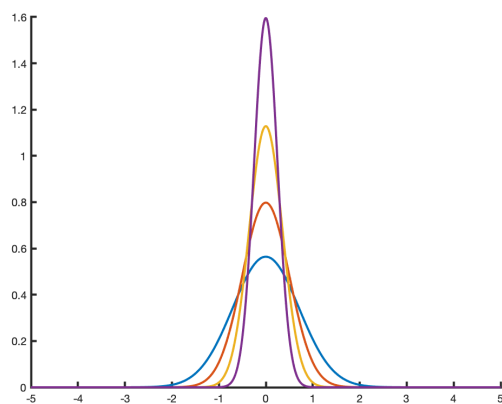


Figure 14.2: Illustration of the Dirac distribution as the limit of a sequence of functions that become thinner and thinner, and higher and higher. Here a Gaussian function with standard deviation going to 0 is considered.

Chapter 15

Exercises

Functions of one or several variables

Exercise 15.1. Calculate the gradient (Jacobian matrix), the divergence, the rotational and the Laplacian of the following vector fields:

$$f(x, y, z) = \begin{pmatrix} 3x^2 + y + z \\ x^2 - y^2 - z \\ x + z^2 \end{pmatrix}, \quad g(x, y, z) = \begin{pmatrix} y^2 + z^2 \\ x^2 + z^2 \\ x^2 + y^2 \end{pmatrix}.$$

Exercise 15.2. For the following scalar fields calculate the gradient, the Hessian matrix, the critical points and their nature (local minimum/maximum...):

$$f(x, y) = 3x^2 + 2y^2 + 2xy + 1,$$

$$g(x, y) = x^2 - 2y^2 + xy - 9x.$$

Exercise 15.3 (Integration). Evaluate

$$\int_0^1 \exp(x) \cdot x \, dx.$$

Evaluate

$$\int_0^2 2x \exp(x^2) \, dx.$$

Hint: Use integration by parts and substitution rules.

Exercise 15.4 (Taylor). Write down the Taylor polynomials of degree 4 in $x_0 = 0$ for

1. $f(x) = \sin(x)$,

and for $x_0 = 1$ for

2. $f(x) = \exp(2x)$.

Exercise 15.5. Let us define

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

Is this function continuous or discontinuous in $(0, 0)$. Justify your answer.

Nonlinear equations

Exercise 15.6. Find the fixed points of the following functions and give their nature (attractive, repulsive). Give a graphical illustration of the convergence (or not) of a sequence of type $x_{n+1} = f(x_n)$.

1. $f : x \in \mathbb{R}_+^* \mapsto f(x) = \sqrt{x}$,

2. $f : x \in \mathbb{R} \mapsto f(x) = 1 - x^2$.

Exercise 15.7 (complement). Describe the sequence generated by Newton's method for solving $\log x = 0$, with initial point $x_0 > 0$. Show that it converges if and only if $x_0 \in]0, e[$.

Exercise 15.8. Let the function $f(x, y) = y^2(x - 1) + x^2(x + 1)$ be given. Compute the stationary points and check which kind of extremum point we deal with.

Exercise 15.9. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ with $f(x) = \sqrt{x^2 + 1}$.

1. Compute the minimum of $f(x)$.
2. That is to solve $f'(x) = 0$.
3. Write down Newton's method.
4. Perform some steps (per hand) for different initial guesses $x_0 = 0.5$ and $x_0 = 2$ respectively.
5. (optional) Implement Newton's method and compute the solution \hat{x} to a tolerance $TOL = 1e - 10$.

Exercise 15.10. Let $L : \mathbb{R}^3 \rightarrow \mathbb{R}$ with

$$L(x, y, \lambda) = f(x) - \lambda c(x) = -x - 0.5y^2 - \lambda(1 - x^2 - y^2)$$

be given.

1. Find the stationary point of L . What needs to be done?
2. Formulate Newton's method (formally; without implementation) for the previous minimization problem.

Exercise 15.11 (This exercise is half of a project because implementation is required). Develop a Newton scheme in \mathbb{R}^2 to find the root of the problem:

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \quad f(x, y) = \left(2xay^2, 2(x^2 + \kappa)ay \right)^T,$$

where $\kappa = 0.01$ and $a = 5$.

1. Justify first that integration of f yields $F(x, y) = (x^2 + \kappa)ay^2$. What is the relation between f and F ?
2. Compute the root of f by hand. Derive the derivative f' and study its properties.
3. Finally, design the requested Newton algorithm. As initial guess, take $(x_0, y_0) = (4, -5)$.
4. What do you observe with respect to the number of Newton iterations?
5. How could we reduce the number of Newton steps?

Polynomial interpolation

Exercise 15.12. Consider the points $x_n = n$ for all $n \in \mathbb{N}$. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous function such that $f(x_n) = \frac{1}{2}(1 + (-1)^{n+1})$ for all $n \in \mathbb{N}$. Calculate the Lagrange interpolating polynomial $P \in \mathcal{P}_4$ of f at the points x_0, x_1, x_2, x_3, x_4 .

Exercise 15.13. Let the following data set be given:

i	0	1	2
x_i	-1	0	2
y_i	-1	0	1

Construct the piece-wise linear and quadratic interpolation polynomials, respectively.

Exercise 15.14 (complement). Let $T > 0$. We consider the function defined by $f(x) = \sin(\frac{2\pi}{T}x)$ for all $x \in [0, 1]$.

Let $n \in \mathbb{N}^*$. We denote by $P_n \in \mathcal{P}_n$ the Lagrange interpolating polynomial of f at the $n + 1$ points $0, 1/n, 2/n, \dots, 1$.

1. Determine an upper bound of $|f^{(n+1)}(x)|$. Deduce an upper bound of the error $|f(x) - P_n(x)|$ for all $x \in [0, 1]$.
2. Let $\epsilon > 0$ and $T > 0$ be given. Provide a condition on n ensuring that the error is less than ϵ for all $x \in [0, 1]$. Make explicit this condition for $T = 2\pi$ and $\epsilon = 0.1$.

Numerical integration (complement)

Exercise 15.15. We place ourselves on the reference interval $[-1, 1]$ and we consider the elementary quadrature formula

$$J(g) = w_0g(-1) + w_1g(0) + w_2g(1).$$

Fix the weights w_0, w_1, w_2 so that the formula is of maximal order (this is the Simpson formula). What is this order?

Exercise 15.16. Evaluate the integral

$$\int_0^1 \frac{1}{1+x} dx$$

with the mid-point rule, trapezoidal rule and the Simpson rule, respectively.

Trigonometric interpolation, Fourier series, Fourier transform

Exercise 15.17. Let the function $f : [0, 2\pi] \rightarrow \mathbb{R}$ be defined by $f(x) = \sin(2x)$. For $M \in \mathbb{N}$, $n = 2M$, and

$$x_j = \frac{2j\pi}{n+1} \quad j = 0, \dots, n,$$

let P_M be the interpolating trigonometric polynomial of f associated with these nodes.

1. Calculate P_0 and P_1 . Without calculation, determine P_M for any $M \geq 2$.
2. For $f(x) = \sin(mx)$, indicate a relation between the pulsation m and the number of nodes needed to guaranty an exact representation of f by its trigonometric interpolant. This is a simplified version of Shannon's theorem.

Exercise 15.18. Let the function $f : [0, 2\pi] \rightarrow \mathbb{R}$ be defined by

$$\begin{aligned} f(0) &= f(2\pi) = 1, \\ f(x) &= 0 \text{ otherwise.} \end{aligned}$$

Let $M \in \mathbb{N}$, $n = 2M$, and

$$x_j = \frac{2j\pi}{n+1} \quad j = 0, \dots, n.$$

1. Calculate the discrete Fourier transform of f associated with the above nodes.
2. Show that the interpolating trigonometric polynomial of f admits the expression

$$P_M(x) = \frac{1}{2M+1} \frac{\cos(Mx) - \cos((M+1)x)}{1 - \cos x}.$$

Plot this function for $M = 8$ and $M = 16$. What do you observe concerning the minimal value of P_M ?

3. (optional) Calculate $P_M(\frac{3\pi}{2M+1})$ and its limit when $M \rightarrow +\infty$. What do you conclude?

Exercise 15.19 (complement). Let the function $f : [0, 2\pi] \rightarrow \mathbb{R}$ be defined by

$$f(x) = x(2\pi - x).$$

1. Calculate the Fourier coefficients c_k , a_k , b_k of f .
2. Deduce the value of

$$\sum_{k=1}^{+\infty} \frac{1}{k^2}.$$

3. (optional) Using Parseval's equality calculate

$$\sum_{k=1}^{+\infty} \frac{1}{k^4}.$$

Exercise 15.20. Consider the function $r : \mathbb{R} \rightarrow \mathbb{R}$ defined by

$$r(x) = \begin{cases} 1 & \text{if } -1 \leq x \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

1. Calculate the Fourier transform \hat{r} of r .
2. (optional) Formulate Parseval's equality.
3. (optional) Show that

$$\int_{-\infty}^{+\infty} |\omega \hat{r}(\omega)|^2 = +\infty. \quad (15.1)$$

This shows in a quantitative way that the spectrum of r decays slowly at infinity. How do you qualitatively interpret the high frequencies? Interpret (15.1) using Parseval's equality.

4. Provide the expression of the spatial Fourier transform $\hat{u}(t, \omega)$ of the function $u(t, x)$ solution of the heat equation (see also next chapter for a review on differential equations)

$$\begin{cases} \frac{\partial u}{\partial t}(t, x) - \frac{\partial^2 u}{\partial x^2}(t, x) = 0 & \forall (t, x) \in \mathbb{R}_+ \times \mathbb{R} \\ u(0, x) = r(x) & \forall x \in \mathbb{R}. \end{cases}$$

What about the decay of the spectrum? Give a physical interpretation.

5. Same question when r is replaced by the Dirac distribution δ . Optional: give an expression of the function u itself.

Optimization

Exercise 15.21. We consider the minimization problem:

$$\text{minimize}_{u \in \mathbb{R}^n} J(u),$$

where $J : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuously differentiable function.

We will specifically work with the 3 following functions defined on \mathbb{R}^2 :

- the quadratic function

$$J_1(x, y) = (x - 1)^2 + p(y - 1)^2,$$

- the Rosenbrock function

$$J_2(x, y) = (x - 1)^2 + p(x^2 - y)^2,$$

- the radial function

$$J_3(r) = e^{-r} + r, \quad r = \sqrt{x^2 + y^2}.$$

In the first 2 cases the positive parameter p can be chosen equal to 10. Of course those are academic examples for which the solution is obvious.

1. Define the above functions, for instance in Python:

```
[ def J(x,y):
  return (x-1)**2+p*(y-1)**2
```

2. Plot isovalues, for instance:

```
[ import numpy as np
  import matplotlib.pyplot as plt
  h=0.025
  xx=np.arange(-1,1.5,h)
  yy=np.arange(-0.5,2,h)
  X,Y=np.meshgrid(xx,yy)
  Z=J(X,Y)
  fig, ax=plt.subplots()
  CS=ax.contour(X,Y,Z,20)
  ax.clabel(CS,inline=1,fontsize=10)
```

3. Implement the steepest descent method with the following special case of Armijo's line search:

```
| Set  $d^k = -\nabla J(u^{k-1})$ ,  $\alpha_0 = 1$ 
| While  $J(u^{k-1} + \alpha_i d^k) > J(u^{k-1}) + \frac{1}{2} \alpha_i \nabla J(u^{k-1}) \cdot d^k$ 
|    $\alpha_{i+1} = \alpha_i / 2$ 
| end
| Set  $s_k = \alpha_i$ 
```

The stopping criterion of the main loop will be the satisfaction of one of the conditions:

- an optimality criterion: $\|\nabla J(u^k)\| \leq \varepsilon$ where $\varepsilon = 10^{-7} \|\nabla J(u^0)\|$,
- a maximal number of iterations $\text{maxiter} = 5000$, in order to stop the algorithm if it does not converge.

Display the obtained minimizers and plot the trajectories.

4. Play with the initial guess and the algorithmic parameters (such as the stopping criterion and parameters within Armijo's rule). What are your conclusions?

Part IV

Differential equations

Chapter 16

Introduction to ODEs (ordinary differential equations)

This chapter is on purpose short. We introduce a few examples of ODEs and some analytical properties. For more detailed introductions, we refer the reader to [2] or [5].

16.1 An introductory example

Let us compute the growth of a species (for example human beings) with a very simple (and finally not that realistic) model. But this shows that reality can be represented to some extent at least by simple models, and that continuous comparisons with other data is also necessary. Furthermore, this also shows that mathematical modeling often starts with a simple equation and is then continuously enriched with further terms and coefficients. In the end we arrive at complicated formulae.

To get started, let us assume that the population number is $y = y(t)$ at time t . Furthermore, we assume constant growth g and mortalities rates m , respectively. In a short time frame dt we have a relative increase of

$$\frac{dy}{y} = (g - m)dt$$

of the population. Re-arranging and taking the limit $dt \rightarrow 0$ yields:

$$\lim_{dt \rightarrow 0} \frac{dy}{dt} = (g - m)y$$

and thus

$$y' = (g - m)y.$$

This is the so-called **Malthusian law of growth**. For this ordinary differential equation (ODE), we can even explicitly compute the solution:

$$y(t) = c \exp((g - m)(t - t_0)).$$

This can be formally obtained through:

$$y' = \frac{dy}{dt} = (g - m)y \tag{16.1}$$

$$\Rightarrow \int \frac{dy}{y} = \int_{t_0}^t (g - m)dt \tag{16.2}$$

$$\Rightarrow \log |y| + C = (g - m)(t - t_0) \tag{16.3}$$

$$\Rightarrow y = \exp[C] \cdot \exp[(g - m)(t - t_0)]. \tag{16.4}$$

In order to work with this ODE and to compute the future development of the species we need an initial value at some starting point t_0 :

$$y(t_0) = y_0.$$

With this value, we can further work to determine the constant $c = \exp(C)$:

$$y(t_0) = c \exp[(g - m)(t_0 - t_0)] = c = y_0. \quad (16.5)$$

Let us say in the year $t_0 = 2021$ there have been two members of this species: $y(2021) = 2$. Supposing a growth rate of 25 per cent per year yields $g = 0.25$. Let us say $m = 0$ - nobody will die. In the following we compute two estimates of the future evolution of this species: for the year $t = 2024$ and $t = 2032$. We first obtain:

$$y(2024) = 2 \exp(0.25 * (2024 - 2021)) = 4.117 \approx 4.$$

Thus, four members of this species exist after three years. Secondly, we want to give a 'long term' estimate for the year $t = 2022$ and calculate:

$$y(2022) = 2 \exp(0.25 * (2032 - 2021)) = 31.285 \approx 31.$$

In fact, this species has an increase of 29 members within 11 years. If you translate this to human beings, we observe that the formula works quite well for a short time range but becomes somewhat unrealistic for long-term estimates though.

Remark 16.1. *An improvement of the basic population model is the **logistic differential equation** that was proposed by Verhulst; see e.g., [2].*

16.2 The model ODE

In the previous example (setting $a = g - m$), we introduced one of the most important ODEs:

$$y' = ay, \quad y(t_0) = y_0. \quad (16.6)$$

This ODE serves often as model problem and important concepts such as existence, uniqueness, stability are usually introduced in terms of this ODE. It is linear with a possibly variable coefficient a . The unique solution can be expressed as

$$y(t) = y_0 \exp\left(\int_{t_0}^t a(s) ds\right). \quad (16.7)$$

For a non-homogeneous ODE of form $y' = ay + b$, the general expression of the solution becomes

$$y(t) = \bar{y}(t) + \alpha \exp\left(\int_{t_0}^t a(s) ds\right),$$

where $\bar{y}(t)$ is any particular solution and α is a constant. The classical way to find a particular solution is the method of 'variation of the constant', which means searching a solution in the form

$$\bar{y}(t) = \alpha(t) \exp\left(\int_{t_0}^t a(s) ds\right).$$

Then one finds $\bar{y}'(t) = a(t)\bar{y}(t) + \alpha'(t) \exp(\int_{t_0}^t a(s) ds)$, whereby α can be obtained through $\alpha'(t) = b(t) \exp(-\int_{t_0}^t a(s) ds)$.

For nonlinear ODEs of form $y' = f(t, y)$ the existence of solutions is much less trivial and often relies on the Cauchy-Lipschitz theorem. We will not enter into these notions. There are also higher order differential equations, which involve second or higher order derivatives of the unknown. Introducing some of these derivatives as additional unknowns, they can be reduced to a system of first order differential equations. Therefore we will subsequently mostly focus our attention to first order ODEs. We only give below a short reminder on exact solutions for second order linear ODE's with constant coefficients.

16.3 Well-posedness

In ODE lectures, we usually learn first the theorem of Peano, which ensures existence (but not uniqueness) of an ODE solution. In these notes, we consider directly the Cauchy-Lipschitz (sometimes called Picard-Lindelöf) theorem, which establishes existence and uniqueness.

Consider the n -dimensional case:

$$y'(t) = f(t, y(t))$$

where $y(t) = (y_1, \dots, y_n)^T$ and $f(t, x) = (f_1, \dots, f_n)^T$. Let an initial point $(t_0, y_0) \in \mathbb{R} \times \mathbb{R}^n$ be given. Furthermore, define

$$D = I \times \Omega \subset \mathbb{R}^1 \times \mathbb{R}^n$$

and let $f(t, y)$ be continuous on D . We seek a solution $y(t)$ on a time interval $I = [t_0, t_0 + T]$ that satisfies $y(t_0) = y_0$.

A crucial aspect for uniqueness is a Lipschitz condition on the right hand side $f(t, y)$: The function $f(t, y)$ on D is said to be (uniformly) Lipschitz continuous if for $L(t) > 0$ it holds

$$\|f(t, x_1) - f(t, x_2)\| \leq L(t)\|x_1 - x_2\|, \quad (t, x_1), (t, x_2) \in D.$$

The function is said to be (locally) Lipschitz continuous if the previous statement holds on every bounded subset of D .

Theorem 16.1 (Cauchy-Lipschitz). *Let $f : D \rightarrow \mathbb{R}^n$ be Lipschitz-continuous. Then there exists for each $(t_0, y_0) \in D$ an $\varepsilon > 0$ and a solution $y : I := [t_0 - \varepsilon, t_0 + \varepsilon] \rightarrow \mathbb{R}^n$ of the IVP such that*

$$y'(t) = f(t, y(t)), \quad t \in I, \quad y(t_0) = y_0.$$

16.4 Second order linear ODE's with constant coefficients

We consider the differential equation

$$y'' + ay' + by = 0, \tag{16.8}$$

where a, b are fixed real numbers. We call \mathcal{S} the set of solutions.

Definition 16.1. *The algebraic equation*

$$r^2 + ar + b = 0 \tag{16.9}$$

is called characteristic equation of (16.8).

Theorem 16.2. *1. If the characteristic equation (16.9) admits two distinct real solutions α and β , then*

$$\mathcal{S} = \left\{ t \mapsto \lambda e^{\alpha t} + \mu e^{\beta t}, \lambda, \mu \in \mathbb{R} \right\}.$$

2. If the characteristic equation (16.9) admits a double root α then

$$\mathcal{S} = \left\{ t \mapsto (\lambda t + \mu) e^{\alpha t}, \lambda, \mu \in \mathbb{R} \right\}.$$

3. If the characteristic equation (16.9) admits two distinct conjugate complex roots $\rho + i\omega$ and $\rho - i\omega$ then

$$\mathcal{S} = \left\{ t \mapsto e^{\rho t} [\lambda \cos(\omega t) + \mu \sin(\omega t)], \lambda, \mu \in \mathbb{R} \right\}.$$

Examples.

If $a = 0$ and $b = \omega^2 > 0$ then $\mathcal{S} = \{t \mapsto \lambda \cos(\omega t) + \mu \sin(\omega t), \lambda, \mu \in \mathbb{R}\}$.

If $a = 0$ and $b = -\omega^2 < 0$ then $\mathcal{S} = \{t \mapsto \lambda e^{\omega t} + \mu e^{-\omega t}, \lambda, \mu \in \mathbb{R}\}$.

Chapter 17

Finite differences for ODEs

Finite differences have been very popular in the past and are still used to solve initial-value problems (IVP) and boundary value problems (BVP). Often finite differences are used in combination with other discretization methods as for example finite elements. For example in a time-dependent PDE, we need to discretize in time and space. Using the Rothe method (horizontal method of lines) temporal discretization is often based on finite differences and spatial discretization on finite elements. In these notes we restrict ourselves to single-step methods.

17.1 Problem statement of an IVP (initial value problem)

We consider ODE initial value problems of the form: find a differentiable function $y(t)$ for $0 \leq t < T < \infty$ such that

$$\begin{aligned}y'(t) &= f(t, y(t)), \\ y(0) &= y_0.\end{aligned}$$

The second condition is the so-called initial condition. Furthermore, $y'(t) = \frac{d}{dt}y$ is the time derivative with respect to time.

The model problem of an ODE is

$$y' = \lambda y, \quad y(0) = y_0, \quad y_0, \lambda \in \mathbb{R}. \quad (17.1)$$

This ODE has the solution:

$$y(t) = \exp(\lambda t)y_0.$$

Furthermore, Problem (17.1) is an example of an autonomous and linear ODE. It is autonomous because the right hand side does not explicitly depend on the variable t , i.e., $f(t, y) = f(y) = \lambda y$. And this ODE is linear because the coefficient λ is independent of the solution y . One further classification can be made in case the right hand side $f = 0$. Then we say that the ODE is homogeneous. Lastly, the coefficient λ is here a constant: we say that it is a linear ODE with constant coefficient. It is a particular case of the ODE with variable coefficient briefly discussed in (16.6), (16.7).

17.2 Stiff problems

An essential difficulty in developing stable numerical schemes, is associated with **stiffness**, which we shall define in the following. Stiffness is very important in both ODE and (time-dependent) PDE applications.

Stiff problems are characterized as ODE solutions that contain various components with (significant) different evolutions over time. A very illustrative example can be found in [18] in which for instance a solution to an ODE problem may look like:

$$y(t) = \exp(-t) + \exp(-99t).$$

Here, the first term defines the trajectory of the solution whereas the second term requires very small time steps in order to resolve the very rapid decrease of that function.

Such behaviors specifically appear in ODE systems in which we seek $y(t) \in \mathbb{R}^n$ rather than $y(t) \in \mathbb{R}$. Then the following definition can be adopted. An IVP is called **stiff** (along a solution $y(t)$) if the eigenvalues $\lambda(t)$ of the Jacobian $f'(t, y(t))$ of f w.r.t. y yield the stiffness ratio:

$$\kappa(t) := \frac{\max_{\Re \lambda(t) < 0} |\Re \lambda(t)|}{\min_{\Re \lambda(t) < 0} |\Re \lambda(t)|} \gg 1.$$

17.3 One-step schemes

17.3.1 The Euler method

The Euler method is the simplest scheme. It is an explicit scheme and also known as forward Euler method.

We consider again the IVP from above with right hand side that satisfies again a Lipschitz condition. According to Theorem 16.1, there exists a unique solution on some neighborhood of t_0 . Here we suppose in addition the existence of a global solution on $[t_0, t_0 + T]$.

For a numerical approximation of the IVP, we first select a sequence of discrete (time) points:

$$t_0 < t_1 < \dots < t_N = t_0 + T.$$

Furthermore we set

$$I_n = [t_{n-1}, t_n], \quad k_n = t_n - t_{n-1}, \quad k := \max_{1 \leq n \leq N} k_n.$$

The derivation of the Euler method is as follows: approximate the derivative with a forward difference quotient (we sit at the time point t_{n-1} and look forward in time):

$$y'(t_{n-1}) \approx \frac{y(t_n) - y(t_{n-1})}{k_n}.$$

Yet $y'(t_{n-1}) = f(t_{n-1}, y_{n-1}(t_{n-1}))$. Then the ODE can be approximated as:

$$\frac{y(t_n) - y(t_{n-1})}{k_n} \approx f(t_{n-1}, y_{n-1}(t_{n-1})).$$

Setting the approximation $y_n \approx y(t_n)$ we obtain the scheme:

Definition 17.1 (Euler method). *For a given starting point $y_0 := y(0) \in \mathbb{R}^n$, the Euler method generates a sequence $\{y_n\}_{n \in \mathbb{N}}$ through*

$$y_n = y_{n-1} + k_n f(t_{n-1}, y_{n-1}), \quad n = 1, \dots, N.$$

17.3.2 Implicit schemes

With the same notation as introduced in Section 17.3.1, we define two further schemes. Beside the Euler method, low-order simple schemes are the implicit Euler and the trapezoidal rule. The main difference is that in general a nonlinear equation system needs to be solved in order to compute the solution. On the other hand we have better numerical stability properties in particular for stiff problems (an analysis will be undertaken in Section 17.4).

The derivation of the implicit (or backward) Euler method is derived as follows: approximate the derivative with a backward difference quotient (we sit at t_n and look back to t_{n-1}):

$$y'(t_n) \approx \frac{y(t_n) - y(t_{n-1})}{k_n}.$$

Consequently, we take the right hand side f at the current time step $y'(t_n) = f(t_n, y_n(t_n))$ and obtain as approximation

$$\frac{y(t_n) - y(t_{n-1})}{k_n} = f(t_n, y_n(t_n)).$$

Consequently, we obtain a scheme in which the right hand side is unknown itself, which leads to a formulation of a nonlinear system:

Definition 17.2 (Implicit Euler). *The implicit Euler scheme is defined as:*

$$\begin{aligned} y_0 &:= y(0), \\ y_n - k_n f(t_n, y_n) &= y_{n-1}, \quad n = 1, \dots, N. \end{aligned}$$

In contrast to the Euler method, the ‘right hand side’ function f now depends on the unknown solution y_n . Thus the computational cost is (much) higher than for the (forward) Euler method. But on the contrary, the method does not require a time step restriction as we shall see in Section 17.4.1.

To derive the trapezoidal rule (or Crank-Nicolson scheme), we take again the difference quotient on the left hand side but approximate the right hand side through its mean value:

$$\frac{y_n - y_{n-1}}{k_n} = \frac{1}{2} \left(f(t_n, y_n(t_n)) + f(t_{n-1}, y_{n-1}(t_{n-1})) \right),$$

which yields:

Definition 17.3 (Trapezoidal rule). *The trapezoidal rule reads:*

$$\begin{aligned} y_0 &:= y(0), \\ y_n &= y_{n-1} + \frac{1}{2} k_n \left(f(t_n, y_n) + f(t_{n-1}, y_{n-1}) \right), \quad n = 1, \dots, N. \end{aligned}$$

It can be shown that the trapezoidal rule is of second order, which means that halving the step size k_n leads to an error that is four times smaller. This is illustrated with the help of a numerical example in Section 17.5.

17.4 Numerical analysis

In the previous section, we have constructed algorithms that yield a sequence of discrete solution $\{y_n\}_{n \in \mathbb{N}}$. In the numerical analysis our goal is to derive a convergence result of the form

$$\|y_n - y(t_n)\| \leq Ck^\alpha$$

where α is the order of the scheme. This result will tell us that the discrete solution y^n really approximates the exact solution y and if we come closer to the exact solution at which rate we come closer.

The model problem

To derive error estimates we work with the model scalar problem:

$$y' = \lambda y, \quad y(0) = y_0, \quad y_0, \lambda \in \mathbb{R}. \tag{17.2}$$

Splitting into stability and consistency

For linear numerical schemes the convergence is composed by **stability** and **consistency**. First of all we have from the previous section that y_n is obtained for the forward Euler method as:

$$\begin{aligned} y_n &= (1 + k\lambda)y_{n-1} \\ &= B_E y_{n-1}, \quad B_E := 1 + k\lambda. \end{aligned}$$

Let us write the error at each time point t_n as:

$$e_n := y_n - y(t_n) \quad \text{for } 1 \leq n \leq N.$$

It holds:

$$\begin{aligned} e_n &= y_n - y(t_n) \\ &= B_E y_{n-1} - y(t_n) \\ &= B_E (e_{n-1} + y(t_{n-1})) - y(t_n) \\ &= B_E e_{n-1} + B_E y(t_{n-1}) - y(t_n) \\ &= B_E e_{n-1} + \frac{k(B_E y(t_{n-1}) - y(t_n))}{k} \\ &= B_E e_{n-1} - k \underbrace{\frac{y(t_n) - B_E y(t_{n-1})}{k}}_{=: \eta_{n-1}}. \end{aligned}$$

Therefore, the error can be split into two parts:

Definition 17.4 (Error splitting of the model problem). *The error at step n can be decomposed as*

$$e_n := \underbrace{B_E e_{n-1}}_{\text{Stability}} - \underbrace{k \eta_{n-1}}_{\text{Consistency}}. \quad (17.3)$$

The first term, namely the stability, provides an idea how the previous error e_{n-1} is propagated from t_{n-1} to t_n . The second term η_{n-1} is the so-called truncation error (or local discretization error), which arises because the exact solution does not satisfy the numerical scheme and represents the consistency of the numerical scheme. Moreover, η_{n-1} yields the speed of convergence of the numerical scheme.

In fact, for the forward Euler scheme in (17.3), we observe for the truncation error:

$$\begin{aligned} \eta_{n-1} &= \frac{y(t_n) - B_E y(t_{n-1})}{k} = \frac{y(t_n) - (1 + k\lambda)y(t_{n-1})}{k} \\ &= \frac{y(t_n) - y(t_{n-1})}{k} - \lambda y(t_{n-1}) \\ &= \frac{y(t_n) - y(t_{n-1})}{k} - y'(t_{n-1}). \end{aligned} \quad (17.4)$$

Thus,

$$y'(t_{n-1}) = \frac{y(t_n) - y(t_{n-1})}{k} - \eta_{n-1}, \quad (17.5)$$

which is nothing else than the approximation of the first-order derivative with the help of a difference quotient plus the truncation error. We investigate these terms further in Section 17.4.2 and concentrate first on the stability estimates in the very next section.

17.4.1 Stability

The goal of this section is to control the term $B_E = (1 + k\lambda)$. Specifically, we will justify why $|B_E| \leq 1$ should hold. The stability is often related to non-physical oscillations of the numerical solution. Otherwise speaking, a nonstable scheme shows artificial oscillations as for example illustrated in Figure 17.2.

We introduce (absolute) **stability** and **A-stability**. From the model problem

$$y'(t) = \lambda y(t), \quad y(t_0) = y_0, \quad \lambda \in \mathbb{C},$$

we know the solution $y(t) = y_0 \exp(\lambda t)$. For $t \rightarrow \infty$ the solution is characterized by the sign of $\Re \lambda$:

$$\begin{aligned} \Re \lambda < 0 &\Rightarrow |y(t)| = |y_0| \exp(t \Re \lambda) \rightarrow 0, \\ \Re \lambda = 0 &\Rightarrow |y(t)| = |y_0| \exp(t \Re \lambda) = |y_0|, \\ \Re \lambda > 0 &\Rightarrow |y(t)| = |y_0| \exp(t \Re \lambda) \rightarrow \infty. \end{aligned}$$

In the first two cases it is particularly interesting to know whether a numerical scheme can produce a bounded discrete solution, as the continuous solution is bounded.

Definition 17.5 ((Absolute) stability). *A (one-step) method is absolute stable for $\lambda k \neq 0$ if its application to the model problem produces in the case $\Re \lambda \leq 0$ a sequence of bounded discrete solutions: $\sup_{n \geq 0} |y_n| < \infty$.*

To find the stability region, we work with the stability function $R(z)$, $z = \lambda k$, such that $B_E = R(z)$. The region of absolute stability is defined as:

$$SR = \{z \in \mathbb{C} : |R(z)| \leq 1\}.$$

The stability functions to explicit Euler, implicit Euler and trapezoidal rule are given by:

Proposition 17.1. *For the simplest time-stepping schemes forward Euler, backward Euler and the trapezoidal rule, the stability functions $R(z)$ read:*

$$\begin{aligned} R(z) &= 1 + z && \text{(forward Euler),} \\ R(z) &= \frac{1}{1 - z} && \text{(backward Euler),} \\ R(z) &= \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z} && \text{(trapezoidal rule).} \end{aligned}$$

Let us look at the proof. We take again the model problem $y' = \lambda y$. We recall the discretization of this problem with the forward Euler method:

$$\frac{y_n - y_{n-1}}{k} = \lambda y_{n-1} \tag{17.6}$$

$$\Rightarrow y_n = (y_{n-1} + \lambda k) y_{n-1} \tag{17.7}$$

$$= (1 + \lambda k) y_{n-1} \tag{17.8}$$

$$= (1 + z) y_{n-1} \tag{17.9}$$

$$= R(z) y_{n-1}. \tag{17.10}$$

For the implicit Euler method we obtain:

$$\frac{y_n - y_{n-1}}{k} = \lambda y_n \tag{17.11}$$

$$\Rightarrow y_n = y_{n-1} + \lambda k y_n \tag{17.12}$$

$$\Rightarrow y_n = \frac{1}{1 - \lambda k} y_{n-1} \tag{17.13}$$

$$\Rightarrow y_n = \underbrace{\frac{1}{1 - z}}_{=: R(z)} y_{n-1}. \tag{17.14}$$

The procedure for the trapezoidal rule is again the analogous.

Definition 17.6 (A-stability). *A difference method is A-stable if its stability region is part of the absolute stability region:*

$$\{z \in \mathbb{C} : \Re z \leq 0\} \subset SR.$$

In other words: Let $\{y_n\}_n$ the sequence of solutions of a difference method for solving the ODE model problem. Then, this method is A-stable if for arbitrary $\lambda \in \mathbb{C}^- = \{\lambda : \Re(\lambda) \leq 0\}$ the approximate solutions are bounded for arbitrary, but fixed, step size k . That is to say:

$$|y_{n+1}| \leq |y_n| < \infty \quad \text{for } n = 1, 2, 3, \dots$$

Remark 17.1. *A-stability is an attractive property since in particular for stiff problems we can compute with arbitrary step sizes k and do not need any step size restriction.*

Proposition 17.2. *The explicit Euler scheme cannot be A-stable.*

Indeed, for the forward Euler scheme, we have $R(z) = 1 + z$. For $|z| \rightarrow \infty$ it holds $R(z) \rightarrow \infty$ which is a violation of the definition of A-stability.

Remark 17.2. *More generally, explicit schemes can never be A-stable.*

On the contrary:

Proposition 17.3. *The implicit Euler scheme and the trapezoidal rule are A-stable.*

Example 17.1. *We illustrate the previous statements.*

1. *In Proposition 17.1 we have seen that for the forward Euler method it holds:*

$$y_n = R(z)y_{n-1},$$

where $R(z) = 1 + z$. A necessary condition of convergence of the sequence is

$$|R(z)| \leq |1 + z| \leq 1. \tag{17.15}$$

Thus if the absolute value of λ (in $z = \lambda k$) is very big, we must choose a very small time step k in order to achieve $|1 - \lambda k| < 1$. Otherwise the sequence $\{y_n\}_n$ will increase and thus diverge (recall that stability is defined with bounded continuous solutions and consequently the numerical approximation should be bounded, too). In conclusion, the forward Euler scheme is only conditionally stable, i.e., it is stable provided that (17.15) is fulfilled.

2. *For the implicit Euler scheme, we see that a large $|\lambda|$ and large k even both help to stabilize the iteration scheme (but be careful, the implicit Euler scheme actually stabilizes sometimes too much. Because it computes contracting sequences also for case where the continuous solution would grow). Thus, no time step restriction is required. Consequently, the implicit Euler scheme is well suited for problems with large parameters/coefficients λ and also for stiff problems.*

Remark 17.3. *The previous example shows that a careful design of the appropriate discretization scheme requires some work: there is no a priori best scheme. Some schemes require time step size restrictions in case of large coefficients (explicit Euler). On the other hand, the implicit Euler scheme does not need step restrictions but may have in certain cases too much damping. Which scheme should be employed for which problem depends finally on the problem itself and must be carefully thought for each problem again.*

17.4.2 Consistency / local discretization error - convergence order

We address now the second ‘error source’ in (17.3). The consistency determines the precision of the scheme and will finally carry over the local rate of consistency to the global rate of convergence. To determine the consistency we assume sufficient regularity of the exact solution such that we can apply Taylor expansion. The idea is then that discrete schemes contain some dominant terms. The lowest order remainder term determines finally the local consistency of the scheme.

Let us analyze the forward Euler scheme. The truncation error η_{n-1} in (17.5) is given through:

$$y'(t_{n-1}) + \eta_{n-1} = \frac{y(t_n) - y(t_{n-1})}{k}.$$

We need information about the solution at the old time step t^{n-1} in order to eliminate $y(t_n)$. Thus we use Taylor-Lagrange (see section 10.4) and develop $y(t^n)$ at the time point t^{n-1} :

$$y(t^n) = y(t^{n-1}) + y'(t^{n-1})k + \frac{1}{2}y''(\tau^{n-1})k^2.$$

We obtain the difference quotient of forward Euler by the following manipulation:

$$\frac{y(t^n) - y(t^{n-1})}{k} = y'(t^{n-1}) + \frac{1}{2}y''(\tau^{n-1})k.$$

We observe that the first terms correspond to the right hand side of (17.4). Thus the remainder term is

$$\frac{1}{2}y''(\tau^{n-1})k$$

and therefore the truncation error η_{n-1} can be estimated as

$$|\eta_{n-1}| \leq \max_{t \in [0, T]} \frac{1}{2}|y''(t)|k = O(k).$$

17.4.3 Convergence

With the help of the two previous subsections, we can easily show the following error estimates:

Theorem 17.4 (Convergence of implicit/explicit Euler). *We have*

$$\max_{t_n \in I} |y_n - y(t_n)| \leq c(T, y)k = O(k),$$

where $k := \max_n k_n$.

The proof does hold for both schemes, except that when we plug-in the stability estimate one should recall that the backward Euler scheme is unconditionally stable and the forward Euler scheme

is only stable when the step size k is sufficiently small. It holds for $1 \leq n \leq N$:

$$\begin{aligned}
 |y_n - y(t_n)| &= |e_n| = k \left| \sum_{k=0}^{n-1} B_E^{n-k} \eta_k \right| \\
 &\leq k \sum_{k=0}^{n-1} |B_E^{n-k} \eta_k| \quad (\text{triangle inequality}) \\
 &\leq k \sum_{k=0}^{n-1} |B_E^{n-k}| |\eta_k| \\
 &\leq k \sum_{k=0}^{n-1} |B_E^{n-k}| Ck \quad (\text{consistency}) \\
 &\leq k \sum_{k=0}^{n-1} 1 Ck \quad (\text{stability}) \\
 &= kN Ck \\
 &= T Ck, \quad \text{where we used } k = T/N \\
 &= C(T, y)k \\
 &= O(k).
 \end{aligned}$$

In Section 17.5 we demonstrate in terms of a numerical example that the forward Euler scheme will fail when the step size restriction is not satisfied, consequently the scheme is not stable and therefore, the convergence result does not hold true.

Theorem 17.5 (Convergence of trapezoidal rule). *We have*

$$\max_{t_n \in I} |y_n(t_n) - y(t_n)| \leq c(T, y)k^2 = O(k^2).$$

The main message is that the Euler schemes both converge with order $O(k)$ (which is very slow) and the trapezoidal rule converges quadratically, i.e., $O(k^2)$.

17.5 Detailed numerical tests

In this section we demonstrate the previous concepts in terms of a numerical example. The programming code is written in octave (which is the open-source sister of MATLAB).

17.5.1 Problem statement

We solve our ODE model problem numerically. Let $a = g - m$ be $a = 0.25$ (test 1) or $a = -0.25$ (test 2) or $a = -10$ (test 3). The IVP is given by:

$$y' = ay, \quad y(t_0 = 2011) = 2.$$

The end time value is $T = 2014$. The tasks are:

- Use the forward Euler (FE), backward Euler (BE), and trapezoidal rule (CN) for the numerical approximation.
- Observe the accuracy in terms of the discretization error.
- Observe for (stiff) equations with a large negative coefficient $a = -10 \ll 1$ the behavior of the three schemes.

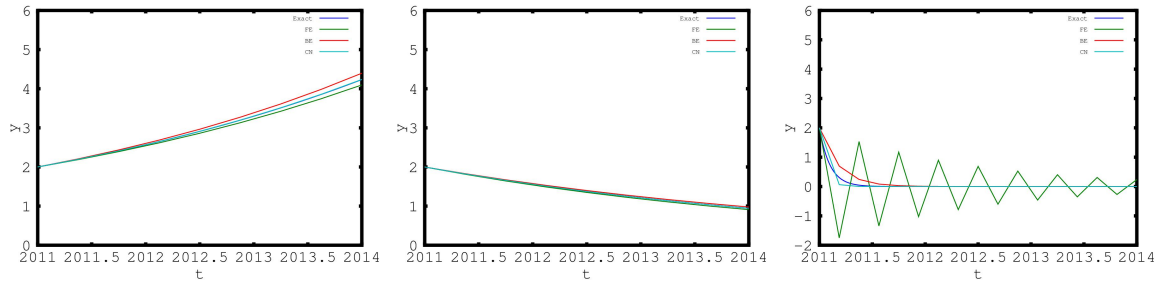


Figure 17.1: Example 17.5.1: on the left, the solution to test 1 is shown. In the middle, test 2 is plotted. On the right, the solution of test 3 with $N = 16$ (number of intervals) is shown. Here, $N = 16$ corresponds to a step size $k = 0.18$ which is slightly below the critical step size for convergence (see Section 17.5.3). Thus we observe the instable behavior of the forward Euler method, but also see slow convergence towards the continuous solution.

17.5.2 Discussion of the results for test 1 with $a=0.25$

In the following, we present our results for the end time value y_N (corresponding to $T = 2014$) for test case 1 ($a = 0.25$) on three mesh levels:

Scheme	#steps(N)	k	y_N
FE	8	0.37500	4.0961
BE	8	0.37500	4.3959
CN	8	0.37500	4.2363
FE	16	0.18750	4.1624
BE	16	0.18750	4.3115
CN	16	0.18750	4.2346
FE	32	0.093750	4.1975
BE	32	0.093750	4.2720
CN	32	0.093750	4.2341

- In the second column, i.e., 8, 16, 32, the number of steps (= number of intervals, i.e., so called mesh cells - speaking in PDE terminology) are given.
- In order to compute numerically the convergence order α with the help of formula (17.30), we work with $k = k_{max} = 0.375$. Then we identify in the third column of the above table that $P(k_{max}) = P(0.375)$, $P(k_{max}/2) = P(0.1875)$ and $P(k_{max}/4) = P(0.09375)$. Formula (17.30) yields $\alpha = 0.9175$ for FE, $\alpha = 1.0921$ for BE, $\alpha = 2.0022$ for CN. This confirms linear convergence in the first two cases, and quadratic convergence in the last one.
- A further observation is that the forward Euler scheme is unstable for $N = 16$ and $a = -10$ and has a zig-zag curve, whereas the other two schemes follow the exact solution and the decreasing exp-function. But for sufficiently small step sizes, the forward Euler scheme is also stable which we know from our A-stability calculations. These step sizes can be explicitly determined for this ODE model problem and shown below.

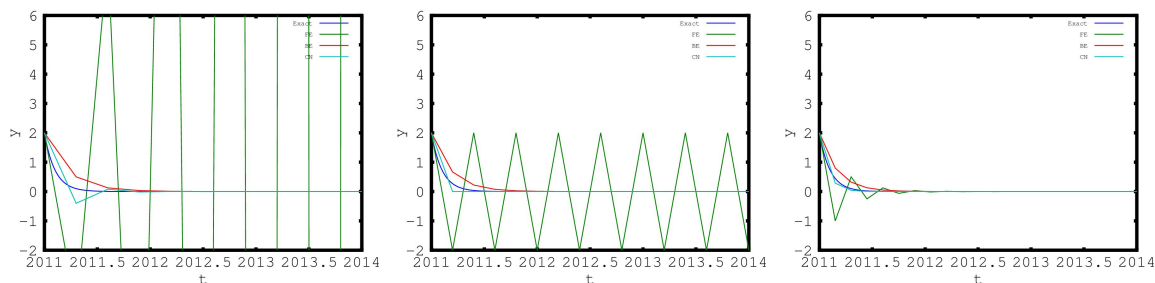


Figure 17.2: Example 17.5.1: tests 3a,3b,3d: Blow-up, constant zig-zag non-convergence, and convergence of the forward Euler method.

17.5.3 Investigating the instability of the forward Euler method: test 3 with $a=-10$

With the help of Example 17.1 let us understand how to choose stable step sizes k for the forward Euler method. The convergence interval reads:

$$|1 + z| \leq 1 \quad \Rightarrow \quad |1 + ak| \leq 1$$

In test 3, $a = -10$, which yields:

$$|1 + z| \leq 1 \quad \Rightarrow \quad |1 - 10k| \leq 1$$

Thus, we need to choose a k that fulfills the previous relation. In this case this, $k < 0.2$ is calculated. This means that for all $k < 0.2$ we should have convergence of the forward Euler method and for $k \geq 0.2$ non-convergence (and in particular no stability!). We perform the following additional tests:

- Test 3a: $N = 10$, yielding $k = 0.3$;
- Test 3b: $N = 15$, yielding $k = 0.2$; exactly the boundary of the stability interval;
- Test 3c: $N = 16$, yielding $k = 0.1875$; from before;
- Test 3d: $N = 20$, yielding $k = 0.15$.

The results of test 3a,3b,3d are provided in Figure 17.2 and visualize very nicely the theoretically predicted behavior.

17.6 A BVP: boundary value problem

17.6.1 A 1D model problem: Poisson

We study in this section a finite difference discretization for the one-dimensional Poisson problem:

$$\left\{ \begin{array}{l} \text{Find } u \in C^2([0, L]) \text{ such that} \\ -u''(x) = f \quad , \forall x \in (0, L), \\ u(0) = a, \\ u(L) = b, \end{array} \right. \quad (17.16)$$

where f is a given right hand side function $C^0([0, L])$ and $a, b \in R$.

For $f = -1$ and $a = 0$ and $L = 1$ we obtain a situation as sketched in Figure 17.3.

17.6.2 Well-posedness of the continuous problem

We investigate in the following well-posedness, the maximum principle and the regularity of the continuous problem.

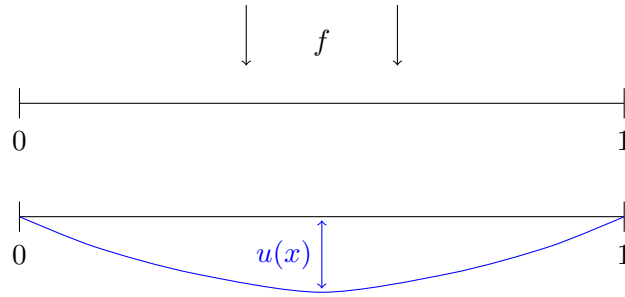


Figure 17.3: The clothesline problem: a uniform force $f = -1$ acts on a 1D line yielding a displacement $u(x)$.

Green’s function and well-posedness

We start with Green’s function that represents the integral kernel of the Laplace operator (here the second-order derivative in 1D). Using Green’s function, we obtain an explicit representation formula for the solution u . We define Green’s function as:

$$G(x, s) = \begin{cases} \frac{s(L-x)}{L} & \text{if } 0 \leq s \leq x \leq L, \\ \frac{x(L-s)}{L} & \text{if } 0 \leq x \leq s \leq L. \end{cases} \tag{17.17}$$

We now integrate two times the differential equation (17.16). For all $\forall x \in [0, L]$, we obtain

$$u(x) = - \int_0^x \left(\int_0^t f(s) ds \right) dt + C_1 x + C_2, \tag{17.18}$$

where C_1 and C_2 are integration constants. Using Fubini’s theorem (see analysis lectures), we obtain

$$\int_0^x \left(\int_0^t f(s) ds \right) dt = \int_0^x \left(\int_s^x f(s) dt \right) ds = \int_0^x (x-s) f(s) ds$$

and therefore

$$u(x) = - \int_0^x (x-s) f(s) ds + C_1 x + C_2.$$

Using the boundary conditions from (17.16), we can determine C_1 and C_2 :

$$\begin{aligned} u(x) &= - \int_0^x (x-s) f(s) ds + \frac{x}{L} \int_0^L (L-s) f(s) ds + a \frac{L-x}{L} + b \frac{x}{L} \\ &= \int_0^x \left(-\frac{L(x-s)}{L} + \frac{x}{L}(L-s) \right) f(s) ds + \frac{x}{L} \int_x^L (L-s) f(s) ds + a \frac{L-x}{L} + b \frac{x}{L} \\ &= \int_0^L G(x, s) f(s) ds + a \frac{(L-x)}{L} + b \frac{x}{L} \end{aligned} \tag{17.19}$$

in which we employed Green’s function G defined in (17.17).

Now, we are able to investigate the well-posedness (see Section 1.3.1 for the meaning of well-posedness):

Proposition 17.6. *Problem (17.16) has a unique solution and depends continuously on the problem data via the following stability estimate:*

$$\|u\|_{C^2} = \max(\|u\|_{C^0}, L\|u'\|_{C^0}, L^2\|u''\|_{C^0}) \leq C_1\|f\|_{C^0} + C_2|a| + C_3|b|,$$

where the constants may be different than before. If $a = b = 0$ (homogeneous Dirichlet conditions), the estimate simplifies to:

$$\|u\|_{C^2} \leq C_1\|f\|_{C^0}.$$

Proof. Existence and uniqueness follow immediately from the representation of the solution as

$$u(x) = \int_0^L G(x, s) f(s) ds + a \frac{(L-x)}{L} + b \frac{x}{L}. \quad (17.20)$$

The stability estimate follows via the following arguments. Green's function $G(x, s)$ is positive in $(0, L)$ and we can estimate as follows:

$$|u(x)| \leq \|f\|_{C^0} \int_0^L G(x, s) ds + |a| + |b|$$

and $\forall x \in [0, L]$, we have

$$\int_0^L G(x, s) ds = \frac{(L-x)}{L} \int_0^x s ds + \frac{x}{L} \int_x^L (L-s) ds \quad (17.21)$$

$$(17.22)$$

$$= \frac{(L-x)x^2}{2L} + \frac{x(L-x)^2}{2L} = \frac{(L-x)x}{2} \leq \frac{L^2}{8}, \quad \text{in } x = L/2, \quad (17.23)$$

yielding

$$\|u\|_{C^0} = \max_{x \in [0, L]} |u(x)| \leq \frac{L^2}{8} \|f\|_{C^0} + |a| + |b|.$$

All values on the right hand side are given by the problem data and are therefore known. However, the estimate is not yet optimal in the regularity of u . In the following we will see that we can even bound the second-order derivative of u by the right hand side f .

Indeed, by differentiating the solution (17.20) (actually better to see for the $G(x, s)$ function in (17.21)), we obtain $\forall x \in [0, L]$:

$$u'(x) = -\frac{1}{L} \int_0^x s f(s) ds + \frac{1}{L} \int_x^L (L-s) f(s) ds - \frac{a}{L} + \frac{b}{L}.$$

Since the integrand is still positive, we can again estimate:

$$|u'(x)| \leq \|f\|_{C^0} \frac{1}{L} \int_0^x s ds + \frac{1}{L} \int_x^L (L-s) ds + \frac{|a|}{L} + \frac{|b|}{L}$$

and $\forall x \in [0, L]$, the Green's function part yields:

$$\frac{1}{L} \int_0^x s ds + \frac{1}{L} \int_x^L (L-s) ds = \frac{x^2}{2L} + \frac{(L-x)^2}{2L} \leq \frac{L}{2}$$

thus

$$\|u'\|_{C^0} = \max_{x \in [0, L]} |u'(x)| \leq \frac{L}{2} \|f\|_{C^0} + \frac{|a|}{L} + \frac{|b|}{L}.$$

Differentiating a second time, we obtain (of course) $\forall x \in [0, L]$:

$$u''(x) = -f(x)$$

i.e., our original PDE. Therefore, we obtain the third estimate as:

$$\|u''\|_{C^0} = \max_{x \in [0, L]} |u''(x)| \leq \|f\|_{C^0}.$$

Consequently, the solution depends continuously on the problem data and therefore the problem is well-posed. \square

Maximum principle on the continuous level

Since $G(x, s)$ is positive in $(0, L)$, the representation (17.20) yields immediately the maximum principle. That is to say:

$$f \geq 0, a \geq 0, b \geq 0 \quad \Rightarrow \quad u \geq 0.$$

Remark 17.4. *Intuitively this is immediately clear. Just consider a parabola function $u(x) = \frac{1}{2}(-x^2 + x)$, which is $-u''(x) = 1$, where $f = 1$. The function is positive with maximum in 0.5 and has its minima on the boundaries 0 and 1.*

Moreover, the maximum principle yields also uniqueness of the solution. Assume that u_1 and u_2 are two different solutions. Their difference is $w := u_1 - u_2$. Of course the original PDE, namely $-u''(x) = f$, also holds for w . Thus: $-w''(x) = f$. We apply f, a, b and $-f, -a, -b$ and obtain $-\Delta w = f$ from which we get $w \geq 0$ and $-\Delta(-w) = -f$ from which we get $w \leq 0$, respectively. Consequently, $w = 0$ and therefore $u_1 = u_2$.

Regularity of the solution

The regularity in the case of Poisson's problem carries over from the right hand side. If $f \in \mathcal{C}^0([0, L])$ we have shown that $u \in \mathcal{C}^2([0, L])$. For higher regularity $k \geq 1$, if $f \in \mathcal{C}^k([0, L])$, therefore for $1 \leq k' \leq k$, $u^{(k'+2)} = -f^{(k')} \in \mathcal{C}^0([0, L])$ and consequently $u \in \mathcal{C}^{k+2}([0, L])$. The solution is always two orders more regular than the right hand side data f .

17.6.3 Spatial discretization

After the previous theoretical investigations we address now a finite difference discretization for Poisson's problem. As previously mentioned, using finite differences, derivatives are approximated via difference quotients. For simplicity we consider homogeneous Dirichlet boundary values. We recall the corresponding problem:

Formulation 17.7.

$$\left\{ \begin{array}{l} \text{Find } u \in \mathcal{C}^2([0, L]) \text{ such that} \\ -u''(x) = f \quad \forall x \in (0, L), \\ u(0) = 0, \\ u(L) = 0. \end{array} \right. \quad (17.24)$$

In the first step, we need to discretize the domain by introducing discrete support points. Let $n \in \mathbb{N}$ and the step size (spatial discretization parameter) h (the distance between two support points) and the support points x_j be defined as

$$h = \frac{L}{n+1}, \quad x_j = jh \quad \text{for all } 0 \leq j \leq n+1.$$

By this construction, we have

- n inner points;
- $n+2$ total points (including the two boundary points);
- $n+1$ intervals (so-called mesh or grid elements).

The difference quotient for approximating a second order derivative is given by (recall classes to the introduction of numerical methods, e.g., [15]):

$$u''(x_j) \approx \frac{U_{j+1} - 2U_j + U_{j-1}}{h^2},$$

where we used a capital letter U to distinguish the discrete solution from the exact solution u . In the following, be careful with the minus sign since we approximate the negative second order derivative of u :

$$-u''(x_j) \approx \frac{-U_{j+1} + 2U_j - U_{j-1}}{h^2}.$$

Consequently, the approximated, discrete, PDE reads:

$$\frac{-U_{j+1} + 2U_j - U_{j-1}}{h^2} = f(x_j), \quad 1 \leq j \leq n,$$

with $U_0 = U_{n+1} = 0$. Going from $j = 0, \dots, n$, we obtain a linear equation system

$$\underbrace{\frac{1}{h^2} \begin{pmatrix} 2 & -1 & 0 & \cdots & 0 \\ -1 & 2 & -1 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & -1 & 2 & -1 \\ 0 & \cdots & 0 & -1 & 2 \end{pmatrix}}_{=A \in \mathbb{R}^{n \times n}} \underbrace{\begin{pmatrix} U_1 \\ U_2 \\ \vdots \\ \vdots \\ U_{n-1} \\ U_n \end{pmatrix}}_{=U \in \mathbb{R}^n} = \underbrace{\begin{pmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ \vdots \\ f(x_{n-1}) \\ f(x_n) \end{pmatrix}}_{=F \in \mathbb{R}^n}$$

In compact form:

$$AU = F.$$

The ‘only’ remaining task consists now in solving this linear equation system. The solution U vector will then contain the discrete solution values on the corresponding support (inner) points.

Remark 17.5. *We could have formulated the above system for all points (including the two boundary points):*

$$\underbrace{\frac{1}{h^2} \begin{pmatrix} h^2 & 0 & & \cdots & \cdots & 0 \\ -1 & 2 & -1 & 0 & \cdots & \cdots & 0 \\ \vdots & -1 & 2 & -1 & \ddots & & \vdots \\ & 0 & \ddots & \ddots & \ddots & & 0 \\ \vdots & \vdots & \ddots & -1 & 2 & -1 & 0 \\ 0 & 0 & \cdots & 0 & -1 & 2 & -1 \\ 0 & 0 & \cdots & 0 & 0 & 0 & h^2 \end{pmatrix}}_{=A \in \mathbb{R}^{(n+2) \times (n+2)}} \underbrace{\begin{pmatrix} U_0 \\ U_1 \\ U_2 \\ \vdots \\ \vdots \\ U_{n-1} \\ U_n \\ U_{n+1} \end{pmatrix}}_{=U \in \mathbb{R}^{n+2}} = \underbrace{\begin{pmatrix} 0 \\ f(x_1) \\ f(x_2) \\ \vdots \\ \vdots \\ f(x_{n-1}) \\ f(x_n) \\ 0 \end{pmatrix}}_{=F \in \mathbb{R}^{n+2}}$$

But we see immediately that we can save this slightly bigger matrix since the boundary values are already known and can be computed immediately.

17.6.4 Solving the linear equation system

Depending on the fineness of the discretization (the smaller h , the smaller the discretization error as we will see later), we deal with large matrices $A \in \mathbb{R}^{n \times n}$ for $n \sim 10^6$ and larger. For moderate sizes, a **direct** solution à la LU or Cholesky is competitive. For large n , due to

- computational cost,
- and memory restrictions,

we have to approximate U with the help of iterative solution schemes such as Jacobi, Gauss-Seidel, conjugate gradient, multigrid, and so forth.

17.6.5 Well-posedness of the discrete problem

We now investigate well-posedness of the discrete problem.

Existence and uniqueness of the discrete problem

We first show that the matrix A is symmetric ($A = A^T$, where A^T is the transpose), positive definite. A direct calculation for two vectors $X, Y \in \mathbb{R}^n$ yields

$$\begin{aligned}
\langle AX, Y \rangle_n &= \sum_{i=1}^n \sum_{j=1}^n A_{ij} X_j Y_i \\
&= \frac{1}{h^2} \sum_{i=1}^n (-X_{i-1} + 2X_i - X_{i+1}) Y_i \quad (\text{apply entries of } A \text{ to } X_j) \\
&= \frac{1}{h^2} \sum_{i=1}^n (X_i - X_{i+1}) Y_i + \frac{1}{h^2} \sum_{i=1}^n (X_i - X_{i-1}) Y_i \\
&= \frac{1}{h^2} \sum_{i=0}^n (X_i - X_{i+1}) Y_i + \frac{1}{h^2} \sum_{i=1}^{n+1} (X_i - X_{i-1}) Y_i \quad (\text{index shift due to bc}) \\
&= \frac{1}{h^2} \sum_{i=1}^{n+1} (X_{i-1} - X_i) Y_{i-1} + \frac{1}{h^2} \sum_{i=1}^{n+1} (X_i - X_{i-1}) Y_i \\
&= \frac{1}{h^2} \sum_{i=1}^{n+1} (X_i - X_{i-1}) (Y_i - Y_{i-1}) \\
&= \sum_{i=1}^{n+1} \frac{X_i - X_{i-1}}{h} \frac{Y_i - Y_{i-1}}{h} \\
&= \sum_{i=1}^{n+1} \frac{Y_i - Y_{i-1}}{h} \frac{X_i - X_{i-1}}{h} \quad (\text{switch factors}) \\
&= \dots \quad (\text{do all steps backwards}) \\
&= \langle X, AY \rangle_n,
\end{aligned}$$

with $X_0 = Y_0 = X_{n+1} = Y_{n+1} = 0$. This first calculation shows the symmetry. It remains to show that

$$\langle AX, Y \rangle_n$$

is positive definite. We now consider the term

$$\langle AX, Y \rangle = \dots = \frac{1}{h^2} \sum_{i=1}^{n+1} (X_i - X_{i-1}) (Y_i - Y_{i-1})$$

from the previous calculation. Setting $Y = X$ yields:

$$\langle AX, X \rangle_n = \sum_{i=1}^{n+1} \left(\frac{X_i - X_{i-1}}{h} \right)^2 \geq 0.$$

This shows the positivity. It remains to show the definiteness: $\langle AX, X \rangle_n = 0$ when $X = 0$. For all $1 \leq i \leq n+1$, we first obtain

$$X_0 = X_1 = \dots = X_{i-1} = X_i = \dots = X_n = X_{n+1},$$

and the zero boundary conditions finally yield $X = 0$. Consequently, the matrix A is symmetric positive definite. The eigenvalues are in \mathbb{R} and strictly positive. Therefore, A is regular (recall the results from linear algebra!) and thus invertible, which finally yields:

Proposition 17.8 (Existence and uniqueness of the discrete problem). *The discrete problem*

$$\frac{-U_{j+1} + 2U_j - U_{j-1}}{h^2} = f(x_j) \quad 1 \leq j \leq n, \quad (17.25)$$

with $U_0 = U_{n+1} = 0$, has one and only one solution for all $F \in \mathbb{R}^n$.

Maximum principle on the discrete level

As previously stated, the maximum principle on the continuous level has a discrete counterpart. Let F be given such that $F_j \geq 0$ for $1 \leq j \leq n$. We perform an indirect proof and show that $\mu := \min_{1 \leq j \leq n} U_j \geq 0$. We suppose $\mu < 0$ and take $i \in \{1, \dots, n\}$ as one index (it exists at least one!) such that the minimum is taken: $U_i = \min_{1 \leq j \leq n} U_j = \mu$. If $1 < i < n$, and U is solution of $AU = F$, the i th row can be written as:

$$2U_i - U_{i+1} - U_{i-1} = h^2 F_i \geq 0$$

under the assumption on F being positive. Therefore, we split

$$0 \leq (U_i - U_{i+1}) + (U_i - U_{i-1}) = (\mu - U_{i+1}) + (\mu - U_{i-1}) \leq 0$$

because $U_i = \mu$ is the minimum value, which we assumed to be negative. To satisfy this inequality chain it must hold:

$$U_{i-1} = U_i = U_{i+1} = \mu,$$

which shows that also the neighboring indices $i - 1$ and $i + 1$ satisfy the minimum. Extending this idea to all indices yields finally, $U_1 = U_n = \mu$. Therefore, we can investigate the first or the last row of the matrix A :

$$0 \leq F_1 = 2U_1 + U_2 = U_1 + (U_1 - U_2) \leq U_1 = \mu < 0.$$

But

$$0 \leq U_1 = \mu < 0$$

is a contradiction. Consequently, $\mu > 0$. Respecting the boundary conditions, we notice that for $U_0 = U_{n+1} = 0$ and $F_j \geq 0$ for all $1 \leq j \leq n$, we have the desired result

$$\mu = \min_{0 \leq j \leq n+1} U_j \geq 0$$

showing that we have a discrete maximum principle.

Proposition 17.9 (Monotone inverse). *The matrix A corresponding to the finite difference discretization is a so-called M matrix (inverse-monotone), i.e., the inverse A^{-1} is element-wise non-negative, i.e.,*

$$A^{-1} \geq 0.$$

Proof. For $1 \leq j \leq n$, we fix $F^j = (0 \cdots 1 \cdots 0) = (\delta_{ij})_{1 \leq i \leq n}$, where δ_{ij} is the Kronecker symbol and X^j is the solution of $AX^j = F^j$. By construction, $X^j = A^{-1}F^j$ and for $1 \leq i \leq n$, the i th row can be written as

$$X_i^j = \sum_{k=1}^n A_{ik}^{-1} F_k^j = \sum_{k=1}^n A_{ik}^{-1} \delta_{kj} = A_{ij}^{-1}. \quad (17.26)$$

Using the maximum principle, we know that for $1 \leq i, j \leq n$ and $F_i^j = \delta_{ij} \geq 0$, the solution is positive, i.e.,

$$X_i^j \geq 0$$

(j is the j th solution and i is the component of that solution). Thanks to (17.26), we have $X_i^j = A_{ij}^{-1}$, and we obtain

$$A_{ij}^{-1} \geq 0,$$

which means that the inverse of the matrix A is an M matrix. \square

17.6.6 Numerical analysis: consistency, stability, and convergence

We finally investigate the convergence properties of the finite difference discretization. As we know from ODE numerics, for linear schemes, consistency plus stability will yield convergence.

Basics in error estimation

We recall the basics. Let $\bar{A}\bar{x} = \bar{b}$ the continuous problem and $Ax = b$ the discrete approximation. Then, we define:

- Approximation error: $e := \bar{x} - x$
- Truncation error: $\eta := A\bar{x} - b$ (exact solution into numerical scheme)
- Residual (defect): $\rho := \bar{A}x - \bar{b}$ (discrete solution into continuous problem)

Furthermore, when the underlying numerical scheme is linear (see ODE lectures), we have:

- Truncation error plus stability yields convergence.

In the following, it will be the same procedure that we use.

Truncation error

We first derive an expression for the **truncation error** η (the local discretization error), which is obtained by plugging in the exact (unknown) solution into the numerical scheme.

Let

$$e_j = U_j - \bar{U}_j \quad \text{for } 1 \leq j \leq n,$$

where $\bar{U}_j = (\bar{U}_i)_{1 \leq i \leq n} = (u(x_i))_{1 \leq i \leq n} \in \mathbb{R}^n$. Then, we have for $1 \leq i \leq n$:

$$\begin{aligned} (Ae)_i &= [A(U - \bar{U})]_i = (F - A\bar{U})_i = f(x_i) - (A\bar{U})_i \\ &= -u''(x_i) - \frac{-\bar{U}_{i+1} + 2\bar{U}_i - \bar{U}_{i-1}}{h^2} \\ &= \frac{u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))}{h^2} - u''(x_i) \end{aligned}$$

and thus

$$Ae = \eta \tag{17.27}$$

where $\eta \in \mathbb{R}^n$ is the local truncation error that arises because the exact solution u does not satisfy the numerical scheme (it is only the discrete solution U that is a solution of the discrete scheme!). Consequently, for all $1 \leq i \leq n$, we have

$$\eta_i := \frac{u(x_{i+1}) - 2u(x_i) + u(x_{i-1}))}{h^2} - u''(x_i).$$

Consistency

The consistency is based on estimates of the local truncation error. The usual method to work with is a Taylor expansion. We have previously established that for $f \in \mathcal{C}^2([0, L])$, we have $u \in \mathcal{C}^4([0, L])$. Consequently, we can employ a Taylor expansion up to order 4 at the point x_i , $1 \leq i \leq n$. This brings us to:

$$\begin{aligned} u(x_{i+1}) &= u(x_i) + hu'(x_i) + \frac{h^2}{2}u''(x_i) + \frac{h^3}{6}u^{(3)}(x_i) + \frac{h^4}{24}u^{(4)}(\tau_i^+) \\ u(x_{i-1}) &= u(x_i) - hu'(x_i) + \frac{h^2}{2}u''(x_i) - \frac{h^3}{6}u^{(3)}(x_i) + \frac{h^4}{24}u^{(4)}(\tau_i^-) \end{aligned}$$

with $\tau_i^+ \in [x_i, x_{i+1}]$ and $\tau_i^- \in [x_{i-1}, x_i]$. It follows that for all $1 \leq i \leq n$

$$|\eta_i| = \left| \frac{h^2}{24}u^{(4)}(\tau_i^+) + \frac{h^2}{24}u^{(4)}(\tau_i^-) \right| \leq \frac{h^2}{12}\|u^{(4)}\|_{\mathcal{C}^0} = \frac{h^2}{12}\|f''\|_{\mathcal{C}^0}$$

in which we used that $u^{(4)} = -f''$. As final result, we obtain:

Proposition 17.10. *The local truncation error of the finite difference approximation of Poisson's problem can be estimated as*

$$\|\eta\|_\infty := \max_{1 \leq i \leq n} |\eta_i| \leq C_f h^2, \quad \text{with } C_f = \frac{1}{12} \|f''\|_{C^0} \quad (17.28)$$

Therefore, the scheme (17.25) is consistent with order 2.

Remark 17.6. *We notice that the additional order of the scheme (namely order two and not one) has been obtained since we work with a central difference quotient to approximate the second order derivative. Of course, here it was necessary to assume sufficient regularity of the right hand side f and the solution u .*

Remark 17.7. *The requirements on the regularity of the solution u to show the above estimates is much higher than the corresponding estimates for the finite element method (FEM). In practice, such a high regularity can only be ensured in a few cases, which is a drawback of finite differences in comparison to finite elements. But still, we can use the method in practice, but one has to be careful whether the results are still robust and reliable.*

Stability in L^∞

We now investigate the stability of the finite difference scheme. We recapitulate the matrix norm $\|\cdot\|_\infty$ (maximum row sum) for $M \in \mathbb{R}^{n,n}$ defined by

$$\|M\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |M_{ij}|$$

Moreover, we denote by $w(x)$ the exact solution for the problem in which $f = 1$ on $[0, L]$. The choice of $f = 1$ is only for simplicity. More general f work as well, but will require more work in the proof below.

Proposition 17.11. *It holds:*

$$\|A^{-1}\|_\infty = \|\bar{W}\|_\infty,$$

and

$$\|A^{-1}\|_\infty \leq \frac{L^2}{8}.$$

Proof. For $f = 1$, the second derivative vanishes: $f'' = 0$. From Section 17.6.6 with (17.28) we know that in this case

$$\|\eta\|_\infty = 0$$

and therefore

$$0 = Ae = A(W - \bar{W})$$

where $W \in \mathbb{R}^n$ denotes the discrete solution corresponding to $f = 1$ (the full right hand side vector is denoted by F_1 in the following). Moreover, we denote $\bar{W} = (w(x_i))_{1 \leq i \leq n} \in \mathbb{R}^n$ the continuous solution. Consequently,

$$\begin{aligned} A(W - \bar{W}) = 0 &\Rightarrow AW = A\bar{W} \Rightarrow F_1 = A\bar{W} \\ &\Rightarrow \bar{W} = A^{-1}F_1 = A^{-1}. \end{aligned}$$

The last step holds because we work with $f = 1$ and therefore $F_1 = 1$. For $1 \leq i \leq n$, we obtain

$$\bar{W}_i = \sum_{j=1}^n A_{ij}^{-1}(F_1)_i = \sum_{j=1}^n A_{ij}^{-1} = \sum_{j=1}^n |A_{ij}^{-1}|$$

The last step holds because A is an M -matrix. Otherwise, we would not be allowed to take absolute signs. We finally obtain:

$$\|A^{-1}\|_{\infty} = \max_{1 \leq i \leq n} \sum_{j=1}^n |A_{ij}^{-1}| = \max_{1 \leq i \leq n} |\bar{W}_i| = \|\bar{W}\|_{\infty}.$$

The exact solution to the problem $-w''(x) = 1$ on $(0, L)$ with homogeneous Dirichlet conditions is $w(x) = \frac{x(L-x)}{2}$, which attains its maximum at $x = L/2$. This yields $\max(w) = L^2/8$ and therefore,

$$\|A^{-1}\|_{\infty} \leq \frac{L^2}{8},$$

which shows the assertion. \square

Remark 17.8. *The case $f = 1$ with $\eta = Ae = 0$ has a very practical implication. Here, the discrete solution matches the continuous solution in the support points $x_i, 1 \leq i \leq n$. This can be confirmed by numerical simulations; see e.g., Section 17.6.7. But be careful, this correspondance is in general only true for $f'' = 0$ and only in 1D cases.*

Remark 17.9. *The previous result can be extended to other right hand sides f , which are not necessarily $f = 1$.*

Convergence L^{∞}

With the previous two results, we can now proof convergence of the finite difference scheme. We assume as before $f \in C^2([0, L])$ and obtain:

Proposition 17.12. *For the convergence of the finite difference scheme (17.25), it holds:*

$$\|e\|_{\infty} = \|A^{-1}\eta\|_{\infty} \leq \|A^{-1}\|_{\infty} \|\eta\|_{\infty} \leq \frac{L^2}{8} C_f h^2 = \frac{L^2}{96} \|f''\|_{C^0} h^2 = O(h^2),$$

where we employed the stability and consistency estimates and (17.28). Consequently, we have quadratic convergence as h tends to zero.

Convergence L^2

In this final section, we proof a second convergence result, but now in the L^2 norm. We define two norms:

$$\|X\|_A = \left(\sum_{j=1}^{n+1} \frac{1}{h} |X_j - X_{j-1}|^2 \right)^{\frac{1}{2}} \quad \|X\|_2 = \left(\sum_{j=1}^n h |X_j|^2 \right)^{\frac{1}{2}}$$

and notice that $\langle \cdot, \cdot \rangle_A$ and $\langle \cdot, \cdot \rangle_2$ are the corresponding scalar products.

Proposition 17.13. *It holds for each $X \in \mathbb{R}^n$:*

$$\|X\|_2 \leq L^{1/2} \|X\|_{\infty} \leq L \|X\|_A.$$

Proof. Let $X \in \mathbb{R}^n$ for all $1 \leq i \leq n$ and we set as usually $X_0 = X_{n+1} = 0$. Then:

$$\begin{aligned} |X_i| &= \left| X_0 + \sum_{j=1}^i (X_j - X_{j-1}) \right| = \sum_{j=1}^i h \frac{X_j - X_{j-1}}{h} * 1 = \left\langle \frac{X_j - X_{j-1}}{h}, 1 \right\rangle_i \\ &\leq \left(\sum_{j=1}^i h \left| \frac{X_j - X_{j-1}}{h} \right|^2 \right)^{1/2} \left(\sum_{j=1}^i h |1|^2 \right)^{1/2} \\ &\leq \|X\|_A L^{1/2} \end{aligned}$$

where we used the Cauchy-Schwarz inequality. It follows that

$$\|X\|_\infty \leq L^{1/2} \|X\|_A.$$

Moreover,

$$\|X\|_2^2 = \left(\sum_{j=1}^n h |X_j|^2 \right) \leq \|X\|_\infty^2 \left(\sum_{j=1}^n h \right) \leq \|X\|_\infty^2 L.$$

We finally obtain

$$\|X\|_2 \leq \|X\|_\infty L^{1/2} \leq L \|X\|_A.$$

□

Remark 17.10. We notice that $\|X\|_2 \leq L \|X\|_A$ holds because $X_0 = 0$ and $\sum_{j=1}^i h |1|^2$ is bounded.

Otherwise the previous result would be not correct.

Proposition 17.14. Let $X, Y \in \mathbb{R}^n$. For a symmetric, positive, definite matrix A , we can define the so-called A -norm; see also Section 7.3. It holds

$$\langle X, Y \rangle_A = \langle AX, Y \rangle_2$$

and

$$\|e\|_A \leq L \|\eta\|_2.$$

Proof. We already have shown that for all $X, Y \in \mathbb{R}^n$, it holds

$$\langle AX, Y \rangle_n = \sum_{i=1}^{n+1} \frac{X_i - X_{i-1}}{h} \frac{Y_i - Y_{i-1}}{h}.$$

Therefore,

$$\langle AX, Y \rangle_2 = h \langle AX, Y \rangle_n = h \sum_{i=1}^{n+1} \frac{X_i - X_{i-1}}{h} \frac{Y_i - Y_{i-1}}{h} = \langle X, Y \rangle_A$$

from which follows

$$\|e\|_A^2 = \langle e, e \rangle_A = \langle Ae, e \rangle_2 = \langle \eta, e \rangle_2 \leq \|\eta\|_2 \|e\|_2 \leq L \|\eta\|_2 \|e\|_A$$

in which we used again the Cauchy-Schwarz inequality. This shows the assertion. □

Proposition 17.15. It holds:

$$\|e\|_A \leq L \|\eta\|_2 \leq L^{3/2} \|\eta\|_\infty \leq L^{3/2} C_f h^2$$

and

$$\|e\|_2 \leq L \|e\|_A \leq L^{5/2} C_f h^2.$$

In other words: the scheme (17.25) converges in the L^2 norm and the A norm with order 2 when h tends to zero.

17.6.7 Numerical test: 1D Poisson

We finish this 1D section with a numerical test implemented in octave [6].

Homogeneous Dirichlet boundary conditions**Formulation 17.16.** *Let $\Omega = (0, 1)$*

$$\begin{cases} \text{Find } u \in \mathcal{C}^2(\Omega) \text{ such that} \\ -u''(x) = f \text{ in } \Omega \\ u(0) = 0, \\ u(1) = 0. \end{cases} \quad (17.29)$$

where $f = -1$. The continuous solution can be computed here and is

$$u(x) = \frac{1}{2}(-x^2 + x).$$

For the discretization we choose $n = 4$. Thus we have five support points

$$x_0, x_1, x_2, x_3, x_4$$

and $h = 1/4 = 0.25$.

With these settings we obtain for the matrix A (see Section (17.6.3)) including boundary conditions:

$$\begin{array}{ccccc} 1 & 0 & 0 & 0 & 0 \\ -16 & 32 & -16 & 0 & 0 \\ 0 & -16 & 32 & -16 & 0 \\ 0 & 0 & -16 & 32 & -16 \\ 0 & 0 & 0 & 0 & 1 \end{array}$$

and for the right hand side vector b :

$$\begin{array}{c} 0 \\ -1 \\ -1 \\ -1 \\ 0 \end{array}$$

We use the famous backslash solution operator (in fact an LU decomposition):

$$U = A \setminus b$$

and obtain as solution $U = (U_0, U_1, U_2, U_3, U_4)$:

$$\begin{array}{c} 0.00000 \\ -0.09375 \\ -0.12500 \\ -0.09375 \\ 0.00000 \end{array}$$

A plot of the discrete solution is provided in Figure 17.4. Moreover, the continuous solution can be evaluated in the support as well. Here, we see that both solutions match. This was already suggested by our numerical analysis; see in particular Remark 17.8.

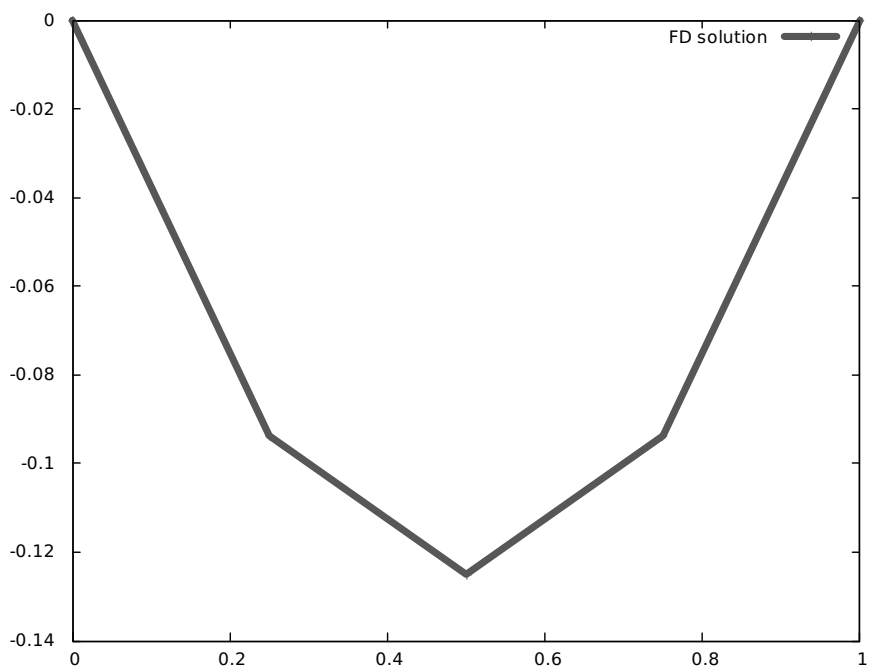


Figure 17.4: Solution of the 1D Poisson problem with $f = -1$ using five support points with $h = 0.25$. We observe that the approximation is rather coarse. A smaller h would yield more support points and more solution values and therefore a more accurate solution.

Nonhomogeneous Dirichlet boundary conditions $u(0) = u(1) = 1$

We change the boundary conditions to

$$u(0) = u(1) = 1.$$

The right hand side then reads $b = (1, -1, -1, -1, 1)^T$. We obtain as solution:

```
1.00000
0.90625
0.87500
0.90625
1.00000
```

Nonhomogeneous Dirichlet boundary conditions $u(0) = 3$ and $u(1) = 2$

We change the boundary conditions to

$$u(0) = 3, u(1) = 2.$$

The right hand side then reads $b = (3, -1, -1, -1, 2)^T$.

We obtain as solution:

```
3.0000
2.6562
2.3750
2.1562
2.0000
```

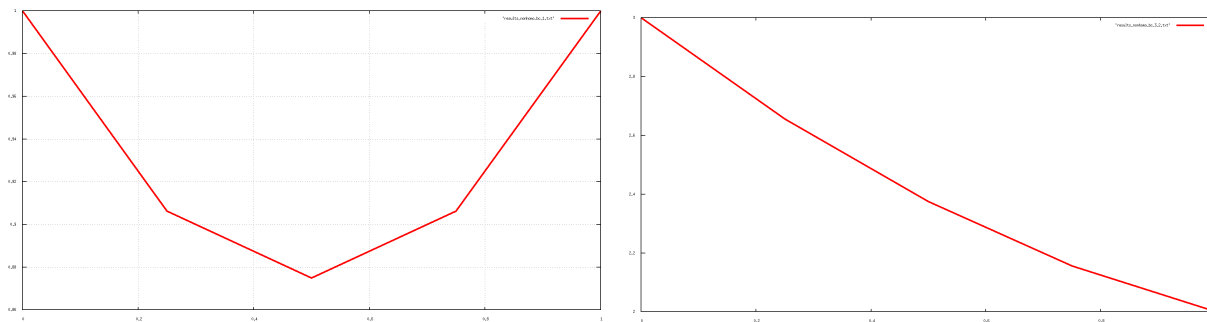


Figure 17.5: Solution of the 1D Poisson problem with $f = -1$ and nonhomogeneous boundary conditions.

How can we improve the numerical solution?

An important question is how we can improve the numerical solution? A five-point approximation as employed in Figure 17.4 is obviously not very accurate. After all previous subsections in this chapter, the answer is clear: increasing the number of support points x_n , i.e., $n \rightarrow \infty$ in the given domain Ω will yield $h \rightarrow 0$ (recall $h = x_{j+1} - x_j$). According to our findings in the numerical analysis section, the discretization error will become very small and we have therefore an accurate approximation of the given problem. On the other hand, large n will require a lot of computational resources. We have finally to find a compromise between accuracy of a numerical solution and computational cost (recall the concepts outlined at the beginning of these lecture notes).

17.7 Computational convergence analysis (complement)

We provide some tools to perform a computational convergence analysis. In these notes we faced two situations of ‘convergence’:

- **Discretization error:** Convergence of the discrete solution u_h towards the (unknown) exact solution u ;
- **Iteration error:** Convergence of an iterative scheme to approximate the discrete solution u_h through a sequence of approximate solutions $u_h^{(k)}$, $k = 1, 2, \dots$

In the following we further illustrate the terminologies ‘first order convergence’, ‘convergence of order two’, ‘quadratic convergence’, ‘linear convergence’, etc.

17.7.1 Discretization error

Before we go into detail, we discuss the relationship between the degrees of freedom (DoFs) N and the mesh size parameter h . In most cases the discretization error is measured in terms of h and all a priori and a posteriori error estimates are stated in a form

$$\|u - u_h\| = O(h^\alpha), \quad \alpha > 0.$$

In some situations it is however better to create convergence plots in terms of DoFs vs. the error. One example is when adaptive schemes are employed with different h . Then it would be not clear to which h the convergence plot should be drawn. But simply counting the total numbers of DoFs is not a problem though.

Relationship between h and N (DoFs)

The relationship of h and N depends on the basis functions (linear, quadratic), whether a Lagrange method (only nodal points) or Hermite-type method (with derivative information) is employed. Moreover, the dimension of the problem plays a role.

We illustrate the relationship for a Lagrange method with linear basis functions in 1D,2D,3D:

Proposition 17.17. *Let d be the dimension of the problem: $d = 1, 2, 3$. It holds*

$$N = \left(\frac{1}{h} + 1\right)^d$$

where h is the mesh size parameter (length of an element or diameter in higher dimensions for instance), and N the number of DoFs.

Proof. Sketch. No strict mathematical proof. We initialize as follows:

- 1D: 2 values per line;
- 2D: 4 values per quadrilaterals;
- 3D: 8 values per hexahedra.

Of course, for triangles or prisms, we have different values in 2D and 3D. We work on the unit cell with $h = 1$. All other h can be realized by just normalizing h . By simple counting the nodal values, we have in 1D

h	N
=====	
1	2
1/2	3
1/4	5
1/8	9
1/16	17
1/32	33
...	
=====	

We have in 2D

h	N
=====	
1	4
1/2	9
1/4	25
1/8	36
1/16	49
1/32	64
...	
=====	

We have in 3D

h	N
=====	
1	8
1/2	27
1/4	64
...	
=====	

□

Discretization error

With the previous considerations, we have now a relationship between h and N that we can use to display the discretization error.

Proposition 17.18. *In the approximate limit it holds:*

$$N \sim \left(\frac{1}{h}\right)^d$$

yielding

$$h \sim \frac{1}{\sqrt[d]{N}}$$

These relationships allow us to replace h in error estimates by N .

Proposition 17.19 (Linear and quadratic convergence in 1D). *When we say a scheme has a linear or quadratic convergence in 1D, (i.e., $d = 1$) respectively, we mean:*

$$O(h) = O\left(\frac{1}{N}\right)$$

or

$$O(h^2) = O\left(\frac{1}{N^2}\right)$$

In a linear scheme, the error will be divided by a factor of 2 when the mesh size h is divided by 2 and having quadratic convergence the error will decrease by a factor of 4.

Proposition 17.20 (Linear and quadratic convergence in 2D). *When we say a scheme has a linear or quadratic convergence in 2D, (i.e., $d = 2$) respectively, we mean:*

$$O(h) = O\left(\frac{1}{\sqrt{N}}\right)$$

or

$$O(h^2) = O\left(\frac{1}{N}\right)$$

17.7.2 Computationally-obtained convergence order

In order to calculate the convergence order α from numerical results, we make the following derivation. Let $P(k) \rightarrow P$ for $k \rightarrow 0$, where k is the discretization parameter (here rather the time step size $k = t_n - t_{n-1}$, but it could also be the spatial discretization parameter h) be a converging process and assume that

$$P(k) - \tilde{P} = O(k^\alpha).$$

Here \tilde{P} is either the exact limit P (in case it is known) or some ‘good’ approximation to it. Let us assume that three numerical solutions are known (this is the minimum number if the limit P is not known). That is

$$P(k), \quad P(k/2), \quad P(k/4).$$

Then, the convergence order can be calculated via the formal approach $P(k) - \tilde{P} = ck^\alpha$ with the following formula:

Proposition 17.21 (Computationally-obtained convergence order). *Given three numerically-obtained values $P(k), P(k/2)$ and $P(k/4)$, the convergence order can be estimated as:*

$$\alpha = \frac{1}{\log(2)} \log\left(\left|\frac{P(k) - P(k/2)}{P(k/2) - P(k/4)}\right|\right). \quad (17.30)$$

The order α is an estimate and heuristic because we assumed a priori a given order, which strictly speaking we have to proof first.

Proof. We assume:

$$\begin{aligned} P(k) - P(k/2) &= O(k^\alpha), \\ P(k/2) - P(k/4) &= O((k/2)^\alpha). \end{aligned}$$

First, we have

$$P(k/2) - P(k/4) = O((k/2)^\alpha) = \frac{1}{2^\alpha} O(k^\alpha)$$

We simply re-arrange:

$$\begin{aligned} P(k/2) - P(k/4) &= \frac{1}{2^\alpha} (P(k) - P(k/2)) \\ \Rightarrow 2^\alpha &= \frac{P(k) - P(k/2)}{P(k/2) - P(k/4)} \\ \Rightarrow \alpha &= \frac{1}{\log(2)} \frac{P(k) - P(k/2)}{P(k/2) - P(k/4)} \end{aligned}$$

□

Example: Temporal order for FE,BE,CN in ODEs

In the following we present results for the (absolute) end time error of an ODE problem (but it could be any other PDE problem as well) on three mesh levels (different time step sizes k) with three schemes (FE - forward Euler, BE - backward Euler, CN - Crank-Nicolson):

Scheme	#steps	k	Error at t_N
FE err.:	8	0.36	0.13786
BE err.:	8	0.36	0.16188
CN err.:	8	0.36	0.0023295
FE err.:	16	0.18	0.071567
BE err.:	16	0.18	0.077538
CN err.:	16	0.18	0.00058168
FE err.:	32	0.09	0.036483
BE err.:	32	0.09	0.037974
CN err.:	32	0.09	0.00014538

We monitor that doubling the number of intervals (i.e., halving the step size k) reduces the error in the forward and backward Euler scheme by a factor of 2. This is (almost) linear convergence, which is confirmed by using Formula (17.30) yielding $\alpha = 0.91804$. The CN scheme is much more accurate (for instance using $n = 8$ the error is 0.2% rather than 13 – 16%) and we observe that the error is reduced by a factor of 4. Thus quadratic convergence is detected. Here the ‘exact’ order on these three mesh levels is $\alpha = 1.9967$.

17.7.3 Spatial discretization error

We simply use now h and then

$$\alpha = \frac{1}{\log(2)} \log \left(\left| \frac{P(h) - P(h/2)}{P(h/2) - P(h/4)} \right| \right) \quad (17.31)$$

in order to obtain the computational convergence order α . Here,

$$P(h) := u_h, \quad P(h/2) := u_{h/2}, \quad P(h/4) := u_{h/4},$$

where u is the discrete PDE solution.

17.7.4 Temporal discretization error for fixed spatial numerical solution

If in time-dependent PDEs, no analytical solution can be constructed, the convergence can again be purely numerically determined. However, we have now influence from spatial and temporal components. Often, the spatial error is more influential than the temporal one. This requires that a numerical reference solution must be computed for sufficiently small h and k . If now ‘coarse’ temporal solutions are adopted (k coarse!), the spatial mesh must be sufficiently fine, i.e., h small!, because otherwise we see significant influences by the spatial components. Then, we perform the procedure from the previous sections: bisection of k and observing the decrease in the corresponding norm or goal functional. That is to say, observe for a discrete PDE solution u :

$$u_{k,h_{fixed},small}, \quad u_{k/2,h_{fixed},small}, \quad u_{k/4,h_{fixed},small}, \quad \dots, \quad u_{k_{ref},h_{fixed},small}$$

With this, we can set:

$$\tilde{P} := u_{k_{ref},h_{fixed},small}, \quad P(k) = u_{k,h_{fixed},small}, \quad P(k/2) = u_{k/2,h_{fixed},small}.$$

Of course, using $P(k/4)$ is also an option. But keep in mind that always h must be sufficiently small and fixed while the time step size is varied.

17.7.5 Extrapolation to the limit

In case, there is no manufactured solution, the numerical reference value can be improved by extrapolation:

Formulation 17.22. *Let h_1 and h_2 two mesh sizes with $h_2 < h_1$. The functionals of interest are denoted by J_{h_1} and J_{h_2} . Then, the extrapolation formula reads:*

$$J_{extra} = \frac{(h_1 * h_1 * J_{h_2} - h_2 * h_2 * J_{h_1})}{h_1 * h_1 - h_2 * h_2}$$

The value for J_{extra} can then be used as approximation for \tilde{P} from the previous subsection.

A specific example from [15] realized in octave is:

```
% Octave code for the extrapolation to the limit
format long

h1 = 1/2^3 % Mesh size h1
h2 = 1/2^4 % Mesh size h2

% Richter/Wick, Springer, 2017, page 383
p1 = 0.784969295 % Functional value J on h1
p2 = 0.786079344 % Functional value J on h2

% Formula MHB, page 335 (Romberg/Richardson)
% See also Quarteroni, Saleri, Ger. 2014, Springer, extrapolation

pMaHB = p2 + (p2 - p1) / ((h1*h1)/(h2*h2) - 1)

% Formula Richter/Wick for extrapolated value
pRiWi = (h1*h1 * p2 - h2*h2 * p1) / (h1*h1 - h2*h2)
```

17.7.6 Iteration error

Iterative schemes are used to approximate the discrete solution u_h . This has a priori nothing to do with the discretization error. The main interest is how fast can we get a good approximation of the discrete solution u_h . One example can be found for solving implicit methods for ODEs in which Newton's method is used to compute the discrete solutions of the backward Euler scheme.

To speak about convergence, we compare two subsequent iterations:

Proposition 17.23. *Let us assume that we have an iterative scheme to compute a root z . The iteration converges with order p when*

$$\|x_k - z\| \leq c \|x_{k-1} - z\|^p, \quad k = 1, 2, 3, \dots$$

with $p \geq 1$ and $c = \text{const.}$ In more detail:

- *Linear convergence:* $c \in (0, 1)$ and $p = 1$;
- *Superlinear convergence:* $c := c_k \rightarrow 0$, ($k \rightarrow \infty$) and $p = 1$;
- *Quadratic convergence* $c \in \mathbb{R}$ and $p = 2$.

Cubic and higher convergence are defined as quadratic convergence with the respective p .

Remark 17.11 (Other characterizations of superlinear and quadratic convergence). *Other (but equivalent) formulations for superlinear and quadratic convergence, respectively, in the case $z \neq x_k$ for all k , are:*

$$\lim_{k \rightarrow \infty} \frac{\|x_k - z\|}{\|x_{k-1} - z\|} = 0,$$

$$\limsup_{k \rightarrow \infty} \frac{\|x_k - z\|}{\|x_{k-1} - z\|^2} < \infty.$$

Corollary 17.24 (Rule of thumb). *A rule of thumb for quadratic convergence is: the number of correct digits doubles at each step. For instance, a Newton scheme to compute $f(x) = x - \sqrt{2} = 0$ yields the following results:*

Iter	x	f(x)
=====		
0	3.000000e+00	7.000000e+00
1	1.833333e+00	1.361111e+00
2	1.462121e+00	1.377984e-01
3	1.414998e+00	2.220557e-03
4	1.414214e+00	6.156754e-07
5	1.414214e+00	4.751755e-14
=====		

Chapter 18

Exercises

Basic notions on ODEs

Exercise 18.1. *Solve:*

1. $y' + y = te^t$, $y(0) = 1$;
2. $y' - y = -\frac{1}{(t+1)^2}e^t$, $y(0) = 0$.

Exercise 18.2. *Solve:*

1. $y' + y = te^t y$, $y(0) = 1$;
2. On $]0, +\infty[$, $ty' + 3y - 2t^5 = 0$, $y(2) = 1$.

Exercise 18.3. *Give the general solutions of:*

1. $y'' - 2y' - 3y = 0$;
2. $y'' + 4y' + 4y = 0$;
3. $y'' - 2y' + 5y = 0$.

Numerical schemes for ODEs

Exercise 18.4. *Given the ordinary differential equation $u'(t) = f(t, u)$, derive basic numerical schemes by taking the forward and backward difference quotient for approximating u' , respectively.*

Exercise 18.5. *Derive stability estimates for the Crank-Nicolson scheme.*

Exercise 18.6. *Recapitulate and complete the steps in the BVP:*

- *Derive step-by-step the explicit form (17.18);*
- *Determine C_1 and C_2 and convince yourself that (17.19) holds true;*
- *Convince yourself by an explicit computation that $u''(x) = -f(x)$ follows from (17.20).*

Exercise 18.7. *Given the ODE problem in Section 16.1, formulate the backward Euler (BE), forward Euler (FE), and Crank-Nicolson (trapezoidal rule) schemes and ‘compute by hand’ three steps. The initial time is $t_0 = 2021$ with $y_0 = 2$ and we are interested in $t_N = 2024$. The time step size is $k = 1$ (year). What do you obtain, how do you interpret the results?*

Part V

Projects on numerical modeling in teams

Chapter 19

Project work

The second half of this class consists of numerical projects done in pairs. The idea is to develop algorithms for a given mathematical problem statement. These algorithms then must be implemented and carefully analyzed from a computational point of view. As programming language we encourage to use an open-source package as for instance python, octave, or C++. Matlab (linked to octave as it is known) is also possible in case you have a personal license, but École Polytechnique will not provide you a license because the school supports open-source programming solutions.

19.1 Idea and formal aspects

Please have a careful look on the different projects listed below. They are related to the different topics presented previously (mandatory sections as well as complement sections). In case you have questions, before your choice, please let us know.

Once you made a choice and have built a team (2 or exceptionally 3 persons), you may start working on your project. **Please be careful that until the mid-term exam, you also work still on course materials (exercises!) of the previous chapters.**

19.1.1 Choice of your project and contact email addresses

Projects are in principle to be chosen among those listed below, including the topics from the previous years. It is also possible that you devise your own project, in this case speak with us.

At latest, please make your choice on the day of the mid-term exam such that you can immediately start working afterward. In all cases, please write us

- `samuel.amstutz@polytechnique.edu`
- `thomas.wick@polytechnique.edu`

an email to indicate:

- Subject line of this email: MAP502: choice of project
- Names of all group members;
- Your project number (Section number from below).

In the meetings No. 7,8,9, we will advise the different groups and provide further help on the specific projects. This will be individual group work and no more lectures.

19.1.2 Final exam mid December

The end-term exam (for each group) consists of:

- A report (word or latex/pdf) of your task, which contains the problem statement, the numerical approach(es), set-up of the numerical example(s), analysis/interpretation of the numerical results;
- A 20 minutes presentation (beamer/PowerPoint);
- Questions from our side to your presentation.

19.2 List of projects (to be confirmed)

19.2.1 Direct and iterative solution of linear equation systems

In this exercise, we develop numerical schemes for solving linear equation systems $Ax = b$. Let $A \in \mathbb{R}^{n \times n}$ with

$$A = \frac{1}{h^2} \begin{pmatrix} h & 0 & \dots & & \\ -1 & 2 & -1 & & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & & -1 & 2 & -1 \\ & & \dots & 0 & h \end{pmatrix} \quad (19.1)$$

and the right hand side $b \in \mathbb{R}^n$ with $b = (1, \dots, 1)^T$ with $h = 1$. This matrix arises in solving for instance the Poisson problem with finite differences or finite elements.

1. Implement Gaussian elimination and solve the above system for $n = 10$ and 100 .
2. Implement now numerical fixed-point schemes: the Richardson iteration, Jacobi method and Gauss-Seidel (see Section 7.2)
3. Compare iteration numbers and wall clock times for all four schemes. Which scheme is the most efficient one? Hint: please also recall the principal differences between direct and iterative schemes.
4. Change $h = 0.5, 0.25, 0.125$. What do you observe?
5. Extend A to $n = 10^6$ and solve again the problem. How much memory on your computer is needed and what are now the wall clock times? Hint: Consider straightforward implementations of the matrix A (each entry is stored and uses memory) as well as so-called *sparse matrices* in which only non-zero entries are stored.
6. For using Richardson's iteration, the largest eigenvalue λ_{max} must be computed. How can this be done efficiently? Please try to understand that two numerical methods (solving $Ax = b$ with Richardson and computing the eigenvalue) interact. Of course, having λ_{max} at hand, gives a very efficient Richardson approach, but obtaining λ_{max} itself might be very costly. For this reason, a trade-off in the computational cost must be envisaged.
7. Implement the LU decomposition. Can you see any difference in wall clock time / memory consumption for (i) constructing the LU decomposition and (ii) solving $Ax = b$ as $LUx = b$ with the two-step algorithm from Section 4.5.1.
8. **Optional** Implement the Cholesky decomposition and compare whether it is really only half of the computational cost in comparison to the LU decomposition.
9. **Optional** Implement the conjugate gradient (CG) scheme according to Section 7.3. This scheme is known to perform very well for the above system.

19.2.2 Numerical methods for eigenvalues

In this exercise, we develop schemes for computing numerically eigenvalue problems. Let $A \in \mathbb{R}^{n \times n}$ with

$$A = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix}, \quad h = \frac{1}{n+1} \quad (19.2)$$

with $h = 1e-1, 1e-3, 1e-6$. See e.g. [13] chap. 9 for an interpretation of this matrix in the framework of finite differences and also [21] for connections with vibrating strings / wave equation.

1. Implement the power method to compute the largest eigenvalue of A .
2. The inverse power method is the power method applied to A^{-1} , without explicitly computing A^{-1} since it is actually enough for that to solve linear systems of matrix A , see details in [13]. Implement the inverse power method to compute the smallest eigenvalue of A . It is also interesting to plot a corresponding eigenvector and compare with the notions of eigenfrequency / eigenmode you may know from physics.
3. To compute all eigenvalues, we adopt the so-called method of QR iterations, see e.g. [13]. In order to construct the orthogonal matrix Q , the Gram-Schmidt procedure (easy to implement, but numerically not stable) or Householder matrices can be used, see e.g. [21]. Further hints and materials will be given for the groups working on this project.
4. (Optional) Implement Lanczos' method for computing the eigenvalues, see e.g. [21].
5. Compare your implemented methods in terms of accuracy, efficiency (number of iterations), and give an interpretation which method you would recommend to use based on your observations.

Hint 1: The power method and QR iterations are already implemented in many scientific computing libraries. The goal **is not** to use such built-in commands, but to implement the methods by yourself.

Hint 2: It can be proven that for the given matrix A the eigenvalues admit the analytical expression

$$\lambda_k = \frac{4}{h^2} \sin^2 \frac{k\pi}{2(n+1)}, \quad k = 1, \dots, n.$$

This can be used for comparisons.

19.2.3 Image compression with SVD (singular valued decomposition) - related to eigenvalues

An interesting application of eigenvalues is the extension to non-quadratic matrices in which singular values must be computed. The first task in this project is to work through and understand the basics of SVD, e.g., with [13][Section 6.5]. The goal is then to study in detail Example 6.9 in [13]. A specific emphasis shall be on the following tasks:

1. Work out the theoretical-algorithmic details of the singular value decomposition (SVD)
2. What exactly happens when using `[U,S,V]=svd(A)` in octave/MATLAB? If you use python, please search for these commands in the python library descriptions.
3. After the implementation: What are the memory requirements?
4. Computational cost (wall clock time)?
5. Accuracy: how many singular values result into an acceptable accuracy?
6. Make a literature research and study another image.

19.2.4 Numerical optimization and application to regression problems

This project elaborates on exercise 15.21, where the steepest descent method was illustrated. Here we go further and explore Newton-type methods as well as some aspects of constrained optimization.

Unconstrained optimization: Newton-type methods

* Newton's method

Newton's method in unconstrained optimization consists in solving the first order optimality condition $\nabla J(u) = 0$. With the notation of chapter 11 (listing 11.1) and exercise 15.21 it is defined by the search direction and the stepsize:

$$\begin{cases} d^k = -\nabla^2 J(u^{k-1})^{-1} \nabla J(u^{k-1}), \\ s_k = 1, \end{cases}$$

where $\nabla^2 J(u)$ is the Hessian matrix of J at point u . Note that in practice we do not compute the inverse of $\nabla^2 J(u^{k-1})$, but 'just' solve the linear system

$$\nabla^2 J(u^{k-1}) d^k = -\nabla J(u^{k-1}).$$

Newton's method is not necessarily a descent method (this depends whether $\nabla^2 J(u^{k-1})$ is positive definite or not). Let us recall that convergence is only expected when the initial guess is sufficiently close to the solution. Of course, when J is quadratic, it converges in one iteration.

1. Define for each example from exercise 15.21 or for your own examples the function $\text{dd}J(x, y)$ which at point $u = (x, y)$ provides the matrix $\nabla^2 J(u)$.
2. Implement Newton's method, display the obtained solution and plot the trajectory.
3. Test the convergence when the initial point is varying. Show that convergence is not guaranteed even for a convex function.

* Quasi-Newton's methods

We refer to section 11.5.

1. Implement the DFP method and display the results and the trajectories, with $H^0 = I$ and Armijo's line search with $\gamma = 0$ (i.e. only check for descent).
2. Test the convergence depending on the initial point.
3. Do the same with BFGS.

* Comparisons

Compare the running times and recapitulate the pros and cons of each method.

Constrained optimization

We consider the least square approximation problem :

$$\min_{u \in K} J(u) = \frac{1}{2} \|u - p\|^2, \quad (19.3)$$

where $p \in \mathbb{R}^n$ is given, $\|\cdot\|$ is the standard Euclidean norm and K is a polyhedral convex set defined by

$$K = \{u \in \mathbb{R}^n, Au \leq b\}.$$

Therefore, our goal is to compute the projection of p onto K . Note that except in some very specific cases (we have seen the case $A = I, b = 0$), we have no formula for this projection.

We will develop a method based on the Lagrangian

$$L(u, \lambda) = J(u) + \lambda \cdot (Au - b).$$

The underlying idea is to penalize the constraint $Au - b \leq 0$ with a weight vector $\lambda \in \mathbb{R}_+^m$ called Lagrange multiplier. We expect that, by a well-chosen Lagrange multiplier λ^* , minimizing $L(u, \lambda^*)$ with respect to u in \mathbb{R}^n will provide a solution to (19.3). Intuitively, for each individual scalar constraint $(Au - b)_i \leq 0$, a two small Lagrange multiplier λ_i would lead to constraint violation, while a two large one would lead to a sub-optimal solution. The question of finding an appropriate Lagrange multiplier will be the purpose of the so-called dual problem. To introduce this we need the concept of saddle point. We say that the pair $(u^*, \lambda^*) \in \mathbb{R}^n \times \mathbb{R}_+^m$ is a saddle point of L if:

$$L(u^*, \lambda) \leq L(u^*, \lambda^*) \leq L(u, \lambda^*) \quad \forall (u, \lambda) \in \mathbb{R}^n \times \mathbb{R}_+^m.$$

* **A bit of theory**

1. Show that if $(u^*, \lambda^*) \in \mathbb{R}^n \times \mathbb{R}_+^m$ is a saddle point of L then u^* solves (19.3). Hint: see that $\lambda^* \cdot (Au^* - b) = 0$.

Note: we will admit the existence and even the uniqueness of a saddle point.

2. Given $\lambda \in \mathbb{R}_+^m$, give an expression of the unique solution u_λ to

$$\min_{u \in \mathbb{R}^n} L(u, \lambda).$$

Deduce an expression of the dual criterion

$$D(\lambda) = \min_{u \in \mathbb{R}^n} L(u, \lambda).$$

3. The dual problem is

$$\max_{\lambda \in \mathbb{R}_+^m} D(\lambda). \tag{19.4}$$

How is this related to the notion of saddle point?

4. Uzawa's method is the steepest ascent method with projection applied to the dual problem (19.4). Here the projection is known! Show that this method can be written as

$$\begin{aligned} u^k &= p - A^T \lambda^{k-1}, \\ \lambda^k &= \left(\lambda^{k-1} + \alpha (Au^k - b) \right)^+ \end{aligned}$$

(we use a fixed stepsize considering the non-normalized gradient as ascent direction).

It can be proved that if α is such that

$$0 < \alpha < \frac{2}{\|A\|_2^2}$$

then u^k converges to the solution of (19.3). Here, we use the matrix norm $\|A\|_2 = \sqrt{\varrho(A^T A)}$ where $\varrho(B)$ is the spectral radius of B .

* **Implementation**

1. Implement Uzawa's algorithm `projection(A,b,p,alpha)` with stopping criterion:

$$\|u^{k+1} - u^k\| < 10^{-6}.$$

2. Test with

$$K = \{u \in \mathbb{R}^2, \|u\|_\infty \leq 1\}.$$

Check the convergence result.

Application: regression*** Isotonic and convex regression**

Let $p = (p_1, \dots, p_n) \in \mathbb{R}^n$ and

$$K = \{(u_1, \dots, u_n) \in \mathbb{R}^n, u_{i+1} \geq u_i \forall 1 \leq i < n\}.$$

1. Construct the corresponding matrix A .
2. Use `projection(A,b,p,alpha)` in this case for some examples of vector p . This is isotonic regression. Provide a graphical output.
3. Do the same with the convex regression:

$$K = \{(u_1, \dots, u_n) \in \mathbb{R}^n, u_i \leq \frac{1}{2}(u_{i-1} + u_{i+1}) \forall 1 < i < n\}$$

*** Application to statistics**

We want to sell honey glasses and we are interested in the optimal price to maximize both the sales volume and our profit. To this end, we go to different cities in France and offer glasses for different prices and take notes how many glasses were sold. These measurement data are listed in Table 19.1.

No x of sold honey glasses	50	30	25	22	27	30	26	32	28	26	21	16	8	4
Price $p(x)$ per glass	3.5	4.0	4.5	5.0	5.5	6.0	6.5	7.0	7.5	8.0	8.5	9.0	9.5	10.0

Table 19.1: Data of 14 measurements of different prices and the corresponding number of sold glasses.

On the same graph display the linear regression (see section 12.7.1), the isotonic (decreasing) regression, the convex regression and the concave regression. Be careful that the implementation of convexity / concavity must be adapted due to non-uniform x data.

19.2.5 Neural network for image classification

Implement an artificial neural network for the classification of the handwritten digits 0 and 1, for instance taken from the MNIST database. Use a single hidden layer (see section 12.7.2) and a classical steepest descent method for training. The detailed structure of the network and the computation of the gradient (backpropagation) is to be discussed with us. Training is performed with a subset of the database and testing with another subset. Study the accuracy of the classifier with respect to the number of neurons in the hidden layer and the number of training samples. Of course, if time allows, the classification of more digits or other images could be addressed.

19.2.6 Numerical methods for root-finding problems (I)

Develop a nonlinear iteration schemes (fixed-point and Newton) in \mathbb{R}^2 to find the root of the problem:

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \quad f(x, y) = \left(2xay^2, 2(x^2 + \kappa)ay \right)^T, \quad (19.5)$$

where $\kappa = 0.001$ and $a = 3$.

1. Justify first that an antiderivative of f is $F(x, y) = (x^2 + \kappa)ay^2$. What is the relation between f and F ?
2. Compute the root of f by hand. Derive the derivative f' and study its properties.
3. Finally, design the requested Newton algorithm. As initial guess, take $(x_0, y_0) = (3, 6)$.

4. What do you observe with respect to the number of Newton iterations?
5. How could we possibly reduce the number of Newton steps?
6. Implement as a second example the system

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \quad f(x, y) = \left(\exp(x^2 + y^2) - 3, x + y - \sin(3(x + y)) \right)^T. \quad (19.6)$$

and solve the root-finding problem $f(x, y) = (0, 0)$. How does Newton's method perform here?

19.2.7 Numerical methods for root-finding problems (II)

We are given the Beale function

$$f(x, y) := \left(x(y - 1) + \frac{3}{2}\right)^2 + \left(x(y^2 - 1) + \frac{9}{4}\right)^2 + \left(x(y^3 - 1) + \frac{21}{8}\right)^2.$$

The task is to compute the minimum $\min f(x, y)$ with the help of nonlinear solvers.

1. Implement Newton's method.
2. Implement Newton's method with line search.
3. Implement a quasi-Newton method (BFGS).
4. What is the minimum? Use as initial guess $(x_0, y_0) = (0, 0)$.
5. How do both methods perform? Study the number of iterations.
6. Vary the initial guesses, e.g., $(x_0, y_0) = (1, -3)$ to start the iterations: what do you observe in terms of iterations? Do you still approach the same minimum?

19.2.8 Shape optimization for fluids using FreeFem++

This project is intended to students willing to make use of concepts of differential calculus and optimization in infinite dimensional vector spaces. The free software FreeFem++ will be used.

The goal of this project is to design optimal shapes for Navier-Stokes flows using as main building block the code

<https://github.com/flomnes/optiflow>

related to the paper [3] (an online version is freely available).

1. Explain in your own words the notion of shape derivative.
2. Explain the main aspects of the algorithm described in [3].
3. Test the algorithm on the built-in examples and study the role of some algorithmic parameters.
4. Implement your own examples, with geometries and physical parameters of your choice.

19.2.9 Discrete Fourier transform

This project deals with discrete Fourier analysis applied to signal and image processing.

1. Implement the discrete Fourier transform of a set of pairs

$$(x_j, y_j = f(x_j)), \quad x_j = \frac{2j\pi}{n+1}, \quad 0 \leq j \leq n.$$

Test on easy cases and plot the corresponding interpolating trigonometric polynomial.

2. Apply it to the function of exercise 14.18. What are your conclusions?
3. Apply it to the function of exercise 14.17 and discuss Shannon's theorem. Illustrate the aliasing phenomenon for $f(x) = \sin(x) + \sin(5x)$ with 9 nodes. How many nodes do we need for an exact representation of $f(x) = \cos(x) \cos(5x)$?
4. Improve the implementation using a simple form of Cooley-Tukey's algorithm [21]. Hint: assume that $n+1$ is a power of 2, in particular n is odd ($n = 2M - 1$), see [13], and using periodicity compute c_k for indices $k = 0, \dots, n$. Discuss the computer effort in comparison with the basic implementation.
5. Illustrate the effect of truncating the discrete Fourier transform (i.e. omitting high frequencies). Discuss some applications in signal compression.
6. Implement the discrete Fourier transform in two dimensions [21], first in a basic way then optionally in an improved way inspired from question 4.
7. Show the effect of truncating the DFT on greyscale images. Give illustrations in the context of image denoising / compression.
8. Possible continuation: FFT for solving linear PDEs with applications in image denoising. To be discussed with us.

19.2.10 Numerical solution of ODEs with finite differences

Let the following ODE initial-value problem be given:

$$y'(t) = ay(t), \quad y(t_0) = y_0 \tag{19.7}$$

on the time interval (t_0, T) where $t_0 = 1$ and $T = 4$. Let the initial value be $y_0 = 2$. Furthermore consider two different test cases with $a = -0.25$ (test 1) and $a = -10$ (test 2).

1. Implement the Euler method, backward Euler method, and the Crank-Nicolson method in octave or python.
2. Recapitulate the stability regions for these three numerical schemes and compute the critical (time) step size for the (forward) Euler scheme for test 1 and test 2.
3. Using the backward Euler method and the Crank-Nicolson method, an implicit system arises. Formulate this system as root finding problem and formulate Newton's method to solve these implicit systems.
4. Using different (time) step sizes, investigate and analyze the findings (instability of the Euler method) by evaluating the end time value $y(T)$. Compare these values for different time step sizes and determine computationally the convergence order (see Section 17.7).

5. When the above system is running, please test your solver(s) at the nonlinear equation:

$$y' = ay - by^2, \quad y(t_0) = y_0. \quad (19.8)$$

For population growth, usually $a \gg b$. For instance, one example for the growth of the world population is:

$$t_0 = 1965, \quad y_0 = 3.34 \times 10^9, \quad a = 0.029, \quad b = 2.695 \times 10^{-12}.$$

What was the expected world population in the year $t = 2000$? What is the expected world population in $t = 2021$? Are these numerical solutions reasonable?

6. **Optional** In case of interest and time, implement higher order time-stepping schemes such as Runge-Kutta schemes. Determine and compare the accuracy by again evaluating the end time value $y(T)$.

19.2.11 Predator-prey systems

In this project, we consider in more detail two coupled ODEs, namely predator-prey systems: Find prey $p(t) : [0, T] \rightarrow \mathbb{R}$ and predators $s(t) : [0, T] \rightarrow \mathbb{R}$ such that

$$\begin{aligned} p'(t) &= r_p p(t) - \alpha s(t)p(t) & \text{in } I = (0, T] \\ s'(t) &= r_s s(t) + \mu \alpha p(t)s(t) & \text{in } I = (0, T] \\ p(t_0) &= p_0 \\ s(t_0) &= s_0 \end{aligned}$$

where $r_p > 0$ is the prey growth rate, $r_s < 0$ is the predator growth rate, $\alpha \in [0, 1]$ the per-capita predation rate, and μ the conversion rate. A good and recent nice introduction can be found in [12][Chapter 6], which is strongly recommended since useful background information is provided. In Section 6.5.3 of [12], the following values are chosen:

$$r_p = 0.6, \quad r_s = -0.2, \quad \alpha = 0.001, \quad \mu = 0.01.$$

and the initial values $p(t_0) = 5000$ and $s(t_0) = 40$.

We have the following tasks in mind:

1. Implement the Euler method, backward Euler method, and the Crank-Nicolson method in octave or python.
2. Run simulations with the above given ODE system, parameters and initial values for $T = 80$ years.
3. Using the backward Euler method and the Crank-Nicolson method, an implicit system arises. Formulate this system as root finding problem and formulate Newton's method to solve these implicit systems. Hints are given in Section 12.6.
4. Using different (time) step sizes, investigate and analyze the findings (instability of the Euler method) by evaluating the end time value $p(T)$ and $s(T)$. Compare these values for different time step. sizes and determine computationally the convergence order (see Section 17.7)
5. Extend the predator-prey model and include the internal interaction rates σ_p and σ_s , [12][p. 183ff].
6. Run again numerical simulations.

19.2.12 Neural network approximation of ODEs

Implement a neural network to approximate the solution of an ODE problem. The performance shall be compared to a classical implementation using finite differences as we have studied previously in class. As direction use the work presented in [10] (available via open access).

The problem statement is:

$$\begin{aligned} u'(x) &= u(x) - u(x)^2, \quad \text{in } (0, T] \\ u(0) &= 0.5 \end{aligned}$$

with $T = 10$.

The specific tasks are:

1. Read and understand [10].
2. Formulate a root-finding problem of the form $F(x, u) = 0$ and formulate the least-squares problem

$$E(w) = \frac{1}{2} \sum_{i=1}^N \|t_i - y_i\|^2$$

where w is the neural network weight vector, t_i contains given information (initial conditions and simply zeros because of $F(x, u) = 0$) and y_i the unknown information (initial data and the discrete values of $F(x_i, y) = 0$) to be approximated (see [10], Section 3).

3. Please make clear and understand how the weight vector w is trained, how the collocation points x_1, \dots, x_{N_2} (make clear yourself why N_2 and not simply N is taken here) are chosen, and how it is possible to determine with this approach a solution of the ODE in the time interval $[0, T]$.
4. Implement a neural network approximation for the previous problem statement.
5. Implement a finite difference (e.g., forward Euler) for the previous problem statement.
6. Compare the computational cost: which one is faster?
7. Study in more detail initial weights for the neural network.
8. Study the learning rate in more detail.
9. Except the computational cost, what are other advantages and/or shortcomings of using neural networks for solving ODEs.
10. Apply the neural network and also finite differences to a second problem:

$$\begin{aligned} u'(x) &= \sin(x) - \frac{u(x)}{x}, \quad \text{in } (0, T] \\ u(\pi) &= 1 \end{aligned}$$

with $T = 4$ and $T = 12$.

19.3 Practical hints during the project work and final presentation

19.3.1 Typical guideline questions

During the projects, some typical aspects of interest are:

- Why does the algorithm converge/not converge?
- Could the algorithm be more efficient? Why/why not?

- Are there alternative algorithms? Do they converge better?
- Try to explain your observations either mathematically or in a rigorous computational fashion.
- What do you observe? What is your personal interpretation of the results?

These (and possibly more) results shall be presented in a presentation and a written report. For more guideline questions, please also see Section 1.4.2.

19.3.2 Hints for presenting results in a talk (final exam)

We collect some hints to make good presentations using, for instance, Beamer Latex or PowerPoint:

1. Structure your talk clearly: what is the goal? How did you manage to reach this goal? What are your results? What is your interpretation? Why did you choose specific methods (advantages / shortcomings in comparison to other existing techniques)?
2. Title page: title, author names, co-authors (if applicable), your institute/company, name of event you are talking.
3. Page 2: list of contents.
4. Conclusions at the end: a few take-home messages.
5. Rule of thumb: one minute per slide. For a 20 minutes presentation a good mean value is to have 20 slides.
6. A big risk of beamer/Power-Point presentations is that we rush over the slides. You are the expert and know everything! But your audience is not. Give the audience time to read the slide!
7. Taylor your talk to the audience!
8. Put a title on each slide!
9. In plots and figures: make lines, the legend, and axes big enough such that people from the back can still read your plots. A good **no go, bad** example is shown in Figure 17.1 in which the legend is by far too small. Do not do that!
10. Tables with numbers: if you create columns of numbers do not put all of them. If you have many, you may mark the most important in a different color to highlight them.
11. Do not glue to your slides. Despite you have everything there, try to explain things in a free speech.
12. Less is more:
 - (a) Do not put too much text on each slide!
 - (b) Do not write full sentences, but use bullet points!
 - (c) Do not use too fancy graphics charts, plots, etc, where everything is moving, etc.
 - (d) Do not try to pack everything you have learnt into one presentation. Everybody believes that you know many, many things. The idea of a presentation is to present in a clear way a piece of your work!
13. As always: the first sentence is the most difficult one: practice for yourself some nice welcome words to get smoothly started.
14. Practice your final presentation ahead of time either alone, with friends or colleagues.
15. Just be authentic during the presentation.
16. If you easily become nervous, avoid coffee before the presentation.

19.3.3 Some example pages of a presentation

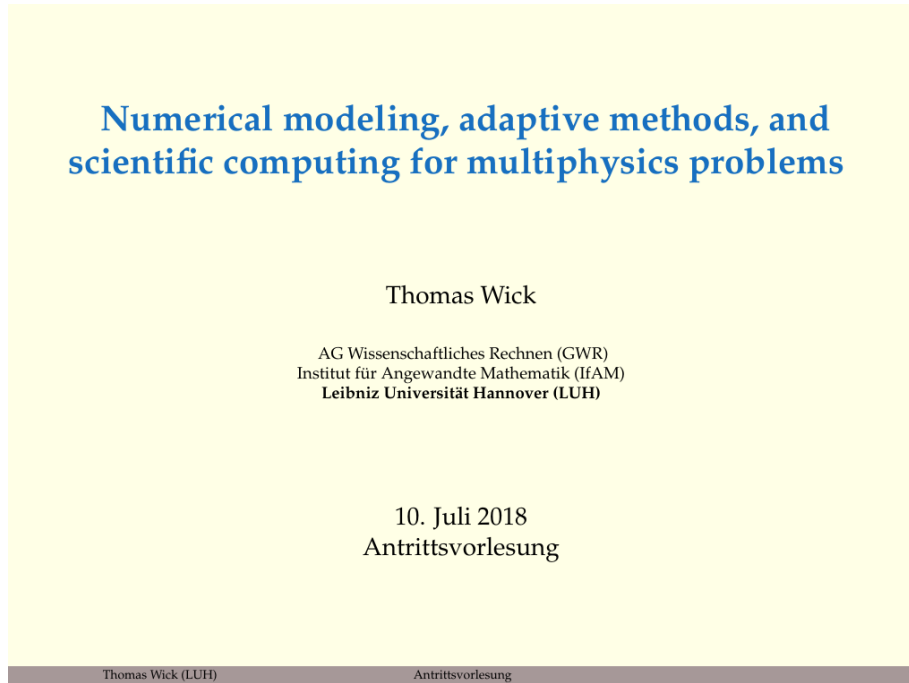


Figure 19.1: Title page. In the footline, the name, topic, and page number is displayed. Page numbers help the audience to point more precisely to a specific page when the talk is open for questions.

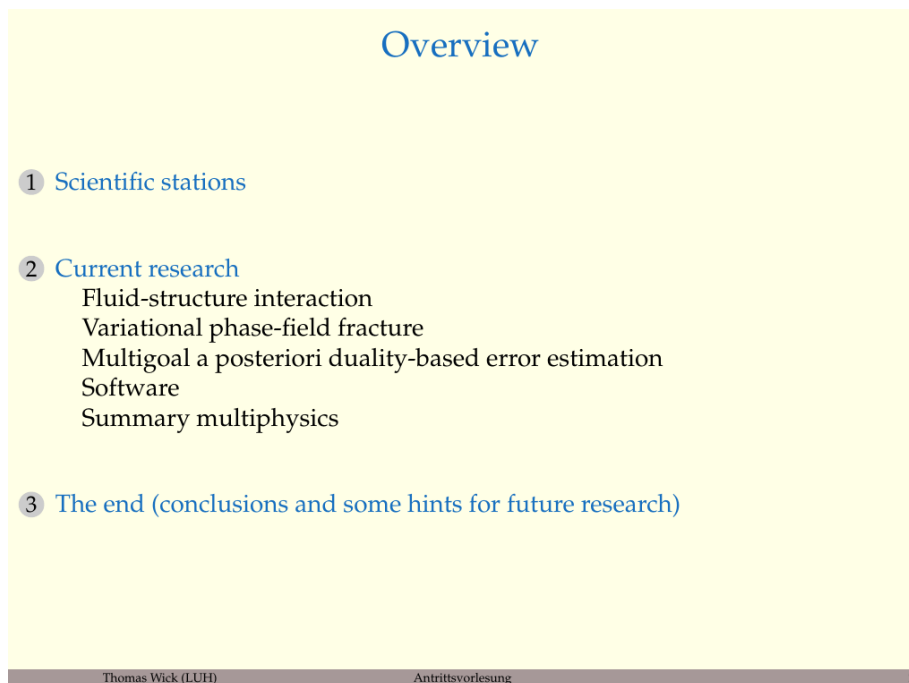


Figure 19.2: Page 2.

Explaining the key words from the title slide

- **Numerical modeling:** Taking practical applications and designing numerical (approximate) descriptions
- **Adaptivity:** Designing numerical methods that solve in a flexible way depending on the goal of the simulation.
- **Scientific computing:** mathematical modeling, numerical design and analysis, and software development.
- **Multiphysics:** Multiple physical phenomena interact with each other (e.g., solids, fluids, temperature, chemical reactions, and so forth).

→ Strong focus on **time-dependent, nonlinear, coupled** equations.

Thomas Wick (LUH) Antrittsvorlesung 4

Figure 19.3: Some page with precise statements; not overpacked, bullet points, if not too exhausting colors may be used, ...

- 1 Scientific stations
- 2 **Current research**
 - Fluid-structure interaction
 - Variational phase-field fracture
 - Multigoal a posteriori duality-based error estimation
 - Software
 - Summary multiphysics
- 3 The end (conclusions and some hints for future research)

Thomas Wick (LUH) Antrittsvorlesung 5

Figure 19.4: When switching to a new section, you may want to display the table of contents again to focus on this new section.

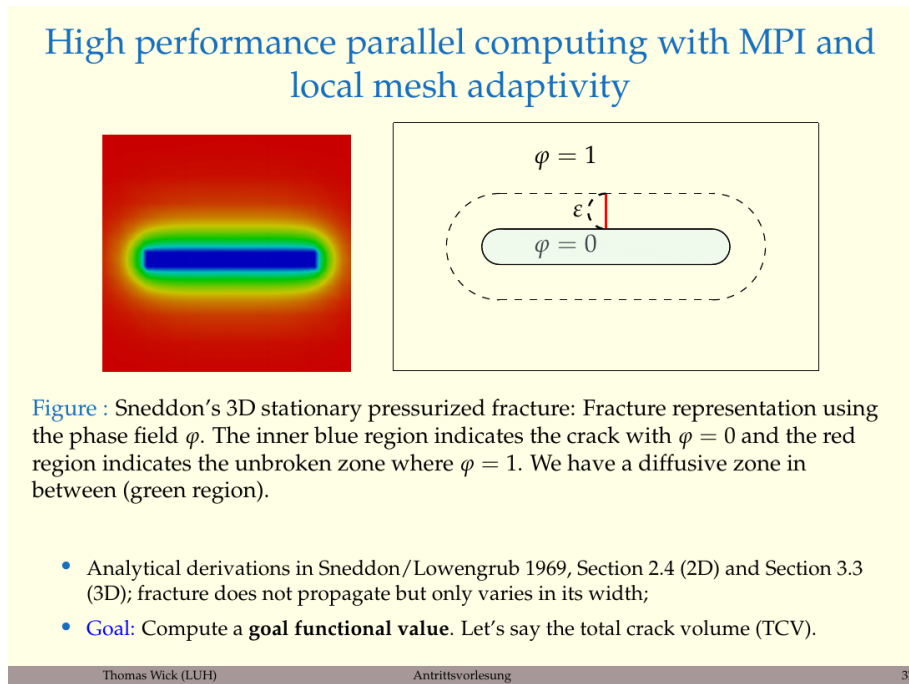


Figure 19.5: Some page representing the configuration, goals, and references of a numerical example.

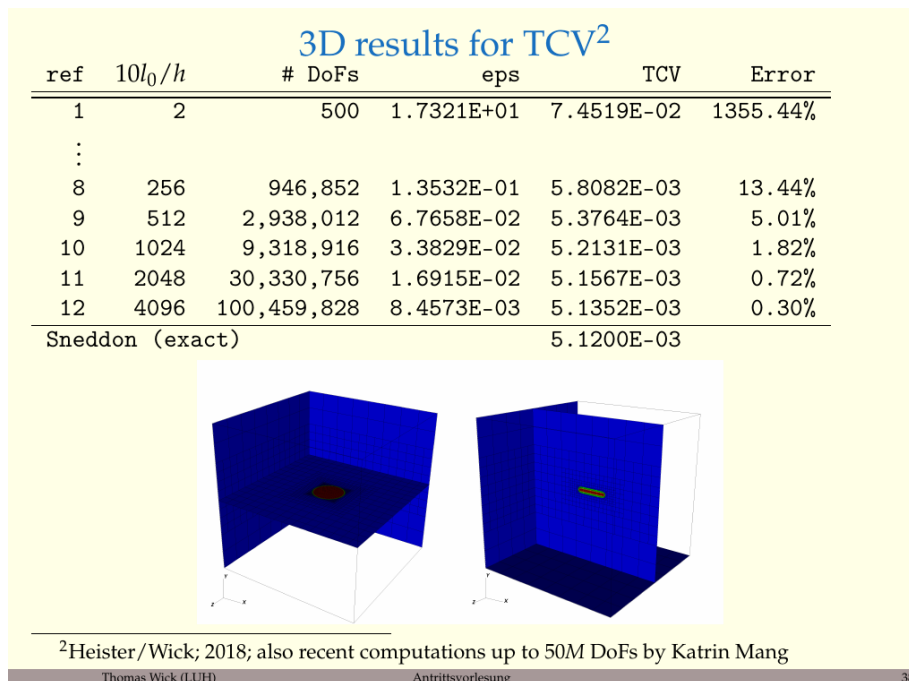




Figure 19.6: Some page with results. If too many lines in a table, please use ... (dots). A reference to published literature can be added as well. In such slides, please also interpret your findings: why did you get those numbers? How did you obtain these findings? Do they make sense? Did you expect something else? Do theory and algorithms (numerics) match?

Software based in C++ (II)

- 2015 / 2018 with T. Heister: open-source 2D / 3D code for solving variational phase-field fracture problems in [deal.II](#).
- 2013: solving fluid-structure interaction in [deal.II](#).
- since May 2009: **DOpElib** (maintainer: Wollner/Wick); Thanks to Christian Goll, Daniel Jodlbauer, Bernhard Endtmayer, Michael Geiger, Masoud Ghaderi, Uwe Köcher, Francesco Ludovici, Matthias Maier

Thomas Wick (LUH)
Antrittsvorlesung
51

Figure 19.7: Some other page with information on the software that was used to obtain the results.

Conclusions and future work

Conclusions:

- Fluid-structure interaction (FSI)
- Variational phase-field fracture
- Multiple goal-oriented error estimates

Future work:

- Adaptive parallel framework for multiphysics (first important steps by [Daniel Jodlbauer](#) and [Bernhard Endtmayer](#) in Linz)
- coupling of FSI with chemical reactions, fractures for instance
- Fractures in incompressible solids ([Katrin Mang](#))
- Nonisothermal fractures in porous media ([Nomi Noii](#))
- Parallel and adaptive fluid flow ([Philipp Thiele](#))
- Treating different temporal scales - towards homogenization in time

Thomas Wick (LUH)
Antrittsvorlesung
55

Figure 19.8: Conclusions with short and precise take-home messages and summary of the key results. If applicable (rather in a scientific research presentation and not in this class in terms of an exam), you may want to add ongoing or future work.

Bibliography

- [1] W. Arnoldi. The principle of minimized iteration in the solution of the matrix eigenvalue problem. *Quart. Appl. Math.*, 9:17–29, 1951.
- [2] M. Braun. *Differential equations and their applications*. Springer, 1993.
- [3] C. Dapogny, P. Frey, F. Omnès, and Y. Privat. Geometrical shape optimization in fluid mechanics using FreeFem++. *Struct. Multidiscip. Optim.*, 58(6):2761–2788, 2018.
- [4] P. Deuffhard. *Newton Methods for Nonlinear Problems*, volume 35 of *Springer Series in Computational Mathematics*. Springer Berlin Heidelberg, 2011.
- [5] P. Deuffhard and F. Bornemann. *Scientific computing with ordinary differential equations*. Texts in applied mathematics. Springer, New York, 2002.
- [6] J. W. Eaton, D. Bateman, S. Hauberg, and R. Wehbring. *GNU Octave version 3.8.1 manual: a high-level interactive language for numerical computations*. CreateSpace Independent Publishing Platform, 2014. ISBN 1441413006.
- [7] C. F. Higham and D. J. Higham. Deep learning: An introduction for applied mathematicians. *SIAM review*, 61(4):860–891, 2019.
- [8] L. Kantorovich and G. Akhilov. *Functional analysis in normed spaces*. Fizmatgiz, Moscow, 1959, German translation: Berlin, Akademie-Verlag, 1964.
- [9] S. Kinnewig, L. Kolditz, J. Roth, and T. Wick. *Numerical Methods for Algorithmic Systems and Neural Networks*. Hannover : Institutionelles Repositorium der Leibniz Universität Hannover, Lecture Notes. Institut für Angewandte Mathematik, Leibniz Universität Hannover, 2022.
- [10] T. Knoke and T. Wick. Solving differential equations via artificial neural networks: Findings and failures in a model problem. *Examples and Counterexamples*, 1:100035, 2021.
- [11] J. Nocedal and S. J. Wright. *Numerical optimization*. Springer Ser. Oper. Res. Financial Engrg., 2006.
- [12] A. Quarteroni and P. Gervasio. *A Primer on Mathematical Modelling*. Springer, 2020.
- [13] A. Quarteroni, F. Saleri, and P. Gervasio. *Scientific computing with MATLAB and Octave, Fourth Edition*. Texts in computational science and engineering. Springer, 2014.
- [14] R. Rannacher. Einführung in die numerische Mathematik (Numerische Mathematik 0). Vorlesungsskriptum, 2000.
- [15] T. Richter and T. Wick. *Einführung in die numerische Mathematik - Begriffe, Konzepte und zahlreiche Anwendungsbeispiele*. Springer, 2017.
- [16] Y. Saad. *Iterative methods for sparse linear systems*. SIAM, 2003.
- [17] Y. Saad and M. H. Schultz. GMRES: A generalized minimal residual algorithm for solving non-symmetric linear systems. *SIAM J. Sci. Stat. Comput.*, 7(3), 1986.

- [18] G. Strang. *Computational Science and Engineering*. Wellesley-Cambridge Press, 2007.
- [19] T. Wick. *Multiphysics Phase-Field Fracture: Modeling, Adaptive Discretizations, and Solvers*. Radon Series on Computational and Applied Mathematics, de Gruyter, Vol. 28, 2020.
- [20] T. Wick. Sergey I. Repin, Stefan A. Sauter: Accuracy of Mathematical Models. *Jahresber. Dtsch. Math. Ver.*, 122:269–274, 2020.
- [21] Wikipedia. Wikipedia. <https://www.wikipedia.org/>. Online; accessed XX XX 2021.

Index

- A norm, 186
- A-stability, 172
- Accuracy, 13, 22
- Algorithm
 - Analysis, 13
- Algorithms, 22
- Approximation, 14
- Armijo step length rule, 115

- Backward Euler method, 168
- Ball, 62
- Banach spaces, 65
- Basis
 - Canonical, 43
- basis, 43
- Big O, 20
- bisection, 122
- Boundary value problem, 176
- BVP, 176

- Canonical basis, 43
- Change of variables, 109
- closed set, 63
- Clothesline problem, 15
- compatible norms, 56
- Complex numbers, 24
- complex numbers, 28
- Computational convergence, 189
- Computational convergence analysis, 189
- Computational cost LU decomposition, 48
- Computational cost LU solution, 46
- Concepts
 - In numerics, 14
- Condition number
 - Matrix, 68, 69
- condition number, 26
- cone, 89
- Consistency, 170, 173, 183
- Convergence, 14, 173
- convergence (sequence), 64
- Convergence order
 - Computationally obtained, 191
- Cost function, 111
- Crank-Nicolson, 169
- Critical point, 112

- Cross product, 106

- Damped Newton, 134
- damping parameter, 72
- Defect, 126
- descent method, 73
- Descent methods, 72
- determinant, 43
- Diagonal dominant, 55
- Diagonal-dominant matrix, 55
- Diagonalization, 57
- Difference quotient, 143
- Differential equations, 21
- dimension, 43
- Direct solution linear systems, 35
- Discrete Fourier transform, 149
- Discretization
 - Finite differences, 179
 - Spatial, 179
- Discretization error, 17, 191
- Divergence, 106
- double precision, 25

- Efficiency, 13, 22
- Eigenspace, 57
- Eigenvalues, 57
 - Power method, 59
- Eigenvectors, 57
- equivalent norms, 64
- Error, 17
 - Discretization, 17
 - Implementation, 17
 - Interpolation, 142
 - Model, 17
 - Numerical quadrature, 17
 - Programming, 17
 - Regularization, 17
 - Round-off, 17
- Errors, 14
- Errors in numerical mathematics, 17
- Euler
 - Backward, 169
 - Implicit, 169
- Euler method, 168
- Existence

- Finite differences, 181
- Explicit schemes, 168
- Extrapolation to the limit, 193
- Fast Fourier transform, 149
- FD, 167
- FFT, 149
- Finite differences, 167
- fixed point, 123
- Fixed point methods, 70
- Floating point numbers, 24
- floating-point numbers, 24
- Forward Euler method, 168
- Fourier series, 151
 - Convergence, 151
- Fourier transform, 152
- Frobenius matrix, 47
- Galerkin orthogonality, 74
- Gauss elimination, 46
- Gauss quadrature, 147
- Gauss-Seidel, 70
- Gauss-Seidel method, 72
- Gaussian elimination, 35
- Generalized minimal residual, 80
- GMRES, 80
- Gradient descent, 72, 73
- gradient descent
 - projection, 116
- Green's function, 177
- Growth of a species, 163
- Guiding questions
 - In differential equations, 21
- Hermitian matrix, 39, 58
- Hermitian positive definite matrix, 55
- Hessian, 106
- Hilbert spaces, 65
- Implicit Euler method, 168
- Implicit schemes, 168
- induced norm, 56
- Interpolation
 - Polynomial, 141
- Interpolation error, 142
- Inverse matrix, 41
- Iteration error, 194
- Iterative solution, 70
- IVP, 167
- Jacobi, 70
- Jacobi method, 72
- Jacobian, 105
- Jacobian matrix, 129
- Krylov space, 74
- Landau symbols, 20
- Line search, 134, 135
- linear convergence, 121
- Linear regression, 137
- Linear systems, 35
 - Direct solution, 35
 - LU decomposition, 46
- Little o, 20
- Local truncation error, 183
- LU, 48
 - Computational cost, 48
- LU decomposition, 46
- M matrix, 182
- machine epsilon, 25
- Machine learning, 138
- Machine precision, 25
- machine precision, 25
- Mathematical modeling, 13
- Matrices, 35, 37
 - Invertible, 41
 - Properties, 37
- matrix norm, 56, 66
- Maximum principle
 - Continuous level, 179
 - Discrete, 182
- Metric, 18, 19
- Metric space, 18
- metric space, 18
- Minimum
 - strict, 113
- Multiindex notation, 105
- Nabla operator, 105
- Neural network, 138, 204
- Newton globalization, 135
- Newton's method, 118, 125, 136
 - Damped, 134
 - Globalization, 133
 - Line search, 134
 - Defect-correction, 126
 - overview, 126
- Newton-Kantorovich theorem, 130
- Norm, 19
- norm, 55, 62, 66
- Normal cone, 90
- Normed space, 18
- normed space, 19
- Norms, 18
- Numbers
 - Complex, 24

- Floating point, 24
 - Real, 24
- numbers, 24
- Numerical analysis
 - Finite differences, 182
 - ODE, 169
- Numerical concepts, 14, 22
 - Guiding questions, 21
- numerical differentiation, 143
 - first derivative, 143
 - second derivative, 144
- Numerical integration, 144
- Numerical methods, 13
- numerical modeling, 13
- Numerical projects, 199
- Numerical quadrature, 144
 - Interpolatory, 146
- Objective function, 111
- octave, 186
- ODE, 167
 - A-stability, 172
 - Boundary value problem, 176
 - Consistency, 173
 - Convergence, 173
 - Finite differences, 167
 - First order, 164
 - Logistic, 164
 - Model problem, 164, 169
 - Newton's method, 136
 - Nonlinear coupled system, 136
 - Second order, 165
 - System, 136
 - Well-posedness, 165
 - Model problem, 164
- One-step θ , 168
- One-step schemes, 168
- Open set, 63
- Optimality conditions, 112
- Optimization
 - Constrained, 111
 - Continuous, 111
 - Discrete, 111
 - Unconstrained, 111
- order of convergence, 121
- Overflow, 25
- Parseval equality, 152
- PCG, 79
- Pivot element, 50
- Pivot matrix, 51
- Pivoting, 50, 52
- Poisson
 - 1D simulations, 186
- Poisson problem, 176
- Polynomial interpolation, 141
- Population, 163
- Positive definite matrix, 40
- positive orthant, 89
- Power method, 59
- Preconditioning, 78
 - CG, 79
 - GMRES, 82
- Programming, 199
- quadratic convergence, 121
- Quadrature formula, 144
 - Composite, 146
 - Order, 146
- Real numbers, 24
- real numbers, 24
- Regression
 - Linear, 137
- Regularity
 - Poisson 1D, 179
- Relationship mesh size and DoFs, 190
- relaxation parameter, 72
- Research software, 13
- Richardson, 70
- Richardson extrapolation, 193
- Robustness, 13, 22
- Saddle point, 113
- Schwarz' theorem, 105
- scientific computing, 13
- Search direction, 72
- set, 23
- single precision, 25
- Software, 13
- Software development, 199
- spectral norm, 56
- spectral radius, 58, 59, 66
- Stability, 170, 171, 184
- Stationary point, 112
- Stiffness, 167
- Symmetric matrix, 39
- symmetric matrix, 58
- Taylor expansion, 183
- tolerance, 25
- Trapezoidal rule, 169
- Trigonometric interpolation, 149
- Truncation error, 183
- Underflow, 25
- Uniqueness

Finite differences, 181

vector space, 61

vectors, 37

Well-posedness

Finite differences, 181

Poisson in 1D, 177

well-posedness, 17