

ETH Lecture 401-0663-00L Numerical Methods for CSE

# Numerical Methods for Computational Science and Engineering

Prof. R. Hiptmair, SAM, ETH Zurich

(with contributions from Prof. P. Arbenz and Dr. V. Gradinaru)

Autumn Term 2020  
(C) Seminar für Angewandte Mathematik, ETH Zürich

[Link](#) to the current version of this lecture document



## Always under construction!

The online version will always be work in progress and subject to change.

(Nevertheless, structure and main contents can be expected to be stable)



**Do not print** before the end of term!

# Contents

<b>0 Introduction</b>	<b>8</b>
0.1 Course Fundamentals . . . . .	8
0.1.1 Focus of this Course . . . . .	8
0.1.2 Goals . . . . .	11
0.1.3 Literature . . . . .	11
0.2 Teaching Style and Model . . . . .	12
0.2.1 Flipped Classroom . . . . .	12
0.2.1.1 Course Videos . . . . .	12
0.2.1.2 Following the Course . . . . .	13
0.2.2 To avoid misunderstandings . . . . .	14
0.2.3 Assignments . . . . .	15
0.2.4 Information on Examinations . . . . .	16
0.3 Programming in C++ . . . . .	17
0.3.1 Function Arguments and Overloading . . . . .	18
0.3.2 Templates . . . . .	19
0.3.3 Function Objects and Lambda Functions . . . . .	20
0.3.4 Multiple Return Values . . . . .	21
0.3.5 A Vector Class . . . . .	23
<b>1 Computing with Matrices and Vectors</b>	<b>36</b>
1.1 Fundamentals . . . . .	37
1.1.1 Notations . . . . .	37
1.1.2 Classes of Matrices . . . . .	39
1.2 Software and Libraries . . . . .	41
1.2.1 EIGEN . . . . .	41
1.2.2 PYTHON . . . . .	46
1.2.3 (Dense) Matrix Storage Formats . . . . .	47
1.3 Basic Linear Algebra Operations . . . . .	52
1.3.1 Elementary Matrix-Vector Calculus . . . . .	52
1.3.2 BLAS – Basic Linear Algebra Subprograms . . . . .	59
1.4 Computational Effort . . . . .	64
1.4.1 (Asymptotic) Complexity . . . . .	65
1.4.2 Cost of Basic Linear-Algebra Operations . . . . .	67
1.4.3 Reducing Complexity in Numerical Linear Algebra: Some Tricks . . . . .	68
1.5 Machine Arithmetic and Consequences . . . . .	73
1.5.1 Experiment: Loss of Orthogonality . . . . .	73
1.5.2 Machine Numbers . . . . .	76
1.5.3 Roundoff Errors . . . . .	79
1.5.4 Cancellation . . . . .	82
1.5.5 Numerical Stability . . . . .	97

<b>2 Direct Methods for (Square) Linear Systems of Equations</b>	<b>104</b>
2.1 Preface . . . . .	105
2.2 Theory: Linear systems of equations . . . . .	107
2.2.1 Existence and uniqueness of solutions . . . . .	107
2.2.2 Sensitivity of linear systems . . . . .	108
2.3 Gaussian Elimination . . . . .	112
2.3.1 Basic algorithm . . . . .	112
2.3.2 LU-Decomposition . . . . .	120
2.3.3 Pivoting . . . . .	128
2.4 Stability of Gaussian Elimination . . . . .	134
2.5 Survey: Elimination solvers for linear systems of equations . . . . .	141
2.6 Exploiting Structure when Solving Linear Systems . . . . .	145
2.7 Sparse Linear Systems . . . . .	150
2.7.1 Sparse matrix storage formats . . . . .	152
2.7.2 Sparse matrices in EIGEN . . . . .	153
2.7.3 Direct Solution of Sparse Linear Systems of Equations . . . . .	160
2.7.4 LU-factorization of sparse matrices . . . . .	163
2.7.5 Banded matrices [DR08, Sect. 3.7] . . . . .	169
2.8 Stable Gaussian elimination without pivoting . . . . .	175
<b>3 Direct Methods for Linear Least Squares Problems</b>	<b>183</b>
3.0.1 Overdetermined Linear Systems of Equations: Examples . . . . .	184
3.1 Least Squares Solution Concepts . . . . .	186
3.1.1 Least Squares Solutions . . . . .	187
3.1.2 Normal Equations . . . . .	189
3.1.3 Moore-Penrose Pseudoinverse . . . . .	194
3.1.4 Sensitivity of Least Squares Problem . . . . .	195
3.2 Normal Equation Methods [DR08, Sect. 4.2], [Han02, Ch. 11] . . . . .	197
3.3 Orthogonal Transformation Methods [DR08, Sect. 4.4.2] . . . . .	200
3.3.1 Transformation Idea . . . . .	200
3.3.2 Orthogonal/Unitary Matrices . . . . .	201
3.3.3 QR-Decomposition [Han02, Sect. 13], [Gut09, Sect. 7.3] . . . . .	202
3.3.3.1 Theory . . . . .	202
3.3.3.2 Computation of QR-Decomposition . . . . .	205
3.3.3.3 QR-Decomposition: Stability . . . . .	211
3.3.3.4 QR-Decomposition in EIGEN . . . . .	212
3.3.4 QR-Based Solver for Linear Least Squares Problems . . . . .	214
3.3.5 Modification Techniques for QR-Decomposition . . . . .	218
3.3.5.1 Rank-1 Modifications . . . . .	219
3.3.5.2 Adding a Column . . . . .	220
3.3.5.3 Adding a Row . . . . .	223
3.4 Singular Value Decomposition (SVD) . . . . .	224
3.4.1 SVD: Definition and Theory . . . . .	224
3.4.2 SVD in EIGEN . . . . .	227
3.4.3 Solving General Least-Squares Problems by SVD . . . . .	230
3.4.4 SVD-Based Optimization and Approximation . . . . .	233
3.4.4.1 Norm-Constrained Extrema of Quadratic Forms . . . . .	233
3.4.4.2 Best Low-Rank Approximation . . . . .	235
3.4.4.3 Principal Component Data Analysis (PCA) . . . . .	239
3.5 Total Least Squares . . . . .	251
3.6 Constrained Least Squares . . . . .	252

3.6.1	Solution via Lagrangian Multipliers . . . . .	253
3.6.2	Solution via SVD . . . . .	254
<b>4</b>	<b>Filtering Algorithms</b>	<b>257</b>
4.1	Discrete Convolutions . . . . .	258
4.2	Discrete Fourier Transform (DFT) . . . . .	267
4.2.1	Discrete Convolution via DFT . . . . .	273
4.2.2	Frequency filtering via DFT . . . . .	274
4.2.3	Real DFT . . . . .	281
4.2.4	Two-dimensional DFT . . . . .	283
4.2.5	Semi-discrete Fourier Transform [QSS00, Sect. 10.11] . . . . .	288
4.3	Fast Fourier Transform (FFT) . . . . .	296
4.4	Trigonometric transformations . . . . .	304
4.4.1	Sine transform . . . . .	305
4.4.2	Cosine transform . . . . .	311
4.5	Toeplitz Matrix Techniques . . . . .	312
4.5.1	Toeplitz Matrix Arithmetic . . . . .	315
4.5.2	The Levinson Algorithm . . . . .	316
<b>5</b>	<b>Data Interpolation and Data Fitting in 1D</b>	<b>319</b>
5.1	Abstract Interpolation . . . . .	319
5.2	Global Polynomial Interpolation . . . . .	326
5.2.1	Polynomials . . . . .	326
5.2.2	Polynomial Interpolation: Theory . . . . .	327
5.2.3	Polynomial Interpolation: Algorithms . . . . .	331
5.2.3.1	Multiple evaluations . . . . .	331
5.2.3.2	Single evaluation . . . . .	334
5.2.3.3	Extrapolation to zero . . . . .	337
5.2.3.4	Newton basis and divided differences . . . . .	340
5.2.4	Polynomial Interpolation: Sensitivity . . . . .	345
5.3	Shape-Preserving Interpolation . . . . .	348
5.3.1	Shape Properties of Functions and Data . . . . .	349
5.3.2	Piecewise Linear Interpolation . . . . .	351
5.4	Cubic Hermite Interpolation . . . . .	352
5.4.1	Definition and Algorithms . . . . .	353
5.4.2	Local Monotonicity-Preserving Hermite Interpolation . . . . .	356
5.5	Splines . . . . .	359
5.5.1	Cubic Spline Interpolation . . . . .	360
5.5.2	Structural Properties of Cubic Spline Interpolants . . . . .	363
5.5.3	Shape Preserving Spline Interpolation . . . . .	366
5.6	Trigonometric Interpolation . . . . .	370
5.6.1	Trigonometric Polynomials . . . . .	371
5.6.2	Reduction to Lagrange Interpolation . . . . .	372
5.6.3	Equidistant Trigonometric Interpolation . . . . .	374
5.7	Least Squares Data Fitting . . . . .	378
<b>6</b>	<b>Approximation of Functions in 1D</b>	<b>384</b>
6.1	Approximation by Global Polynomials . . . . .	387
6.1.1	Polynomial Approximation: Theory . . . . .	387
6.1.2	Error Estimates for Polynomial Interpolation . . . . .	393
6.1.3	Chebychev Interpolation . . . . .	402
6.1.3.1	Motivation and Definition . . . . .	403

6.1.3.2	Chebychev Interpolation Error Estimates . . . . .	407
6.1.3.3	Chebychev Interpolation: Computational Aspects . . . . .	412
6.2	Mean Square Best Approximation . . . . .	416
6.2.1	Abstract Theory . . . . .	416
6.2.1.1	Mean Square Norms . . . . .	416
6.2.1.2	Normal Equations . . . . .	417
6.2.1.3	Orthonormal Bases . . . . .	419
6.2.2	Polynomial Mean Square Best Approximation . . . . .	420
6.3	Uniform Best Approximation . . . . .	426
6.4	Approximation by Trigonometric Polynomials . . . . .	430
6.4.1	Approximation by Trigonometric Interpolation . . . . .	431
6.4.2	Trigonometric Interpolation Error Estimates . . . . .	432
6.4.3	Trigonometric Interpolation of Analytic Periodic Functions . . . . .	436
6.5	Approximation by Piecewise Polynomials . . . . .	439
6.5.1	Piecewise Polynomial Lagrange Interpolation . . . . .	441
6.5.2	Cubic Hermite Interpolation: Error Estimates . . . . .	444
6.5.3	Cubic Spline Interpolation: Error Estimates [Han02, Ch. 47] . . . . .	446
<b>7</b>	<b>Numerical Quadrature</b> . . . . .	<b>449</b>
7.1	Quadrature Formulas . . . . .	451
7.2	Polynomial Quadrature Formulas . . . . .	454
7.3	Gauss Quadrature . . . . .	456
7.4	Composite Quadrature . . . . .	470
7.5	Adaptive Quadrature . . . . .	478
<b>8</b>	<b>Iterative Methods for Non-Linear Systems of Equations</b> . . . . .	<b>487</b>
8.1	Iterative Methods . . . . .	489
8.1.1	Speed of Convergence . . . . .	491
8.1.2	Termination Criteria/Stopping Rules . . . . .	495
8.2	Fixed-Point Iterations . . . . .	498
8.2.1	Consistent Fixed-Point Iterations . . . . .	499
8.2.2	Convergence of Fixed-Point Iterations . . . . .	501
8.3	Finding Zeros of Scalar Functions . . . . .	507
8.3.1	Bisection . . . . .	507
8.3.2	Model Function Methods . . . . .	509
8.3.2.1	Newton Method in the Scalar Case . . . . .	509
8.3.2.2	Special One-Point Methods . . . . .	512
8.3.2.3	Multi-Point Methods . . . . .	516
8.3.3	Asymptotic Efficiency of Iterative Methods for Zero Finding . . . . .	520
8.4	Newton's Method in $\mathbb{R}^n$ . . . . .	522
8.4.1	The Newton Iteration . . . . .	522
8.4.2	Convergence of Newton's Method . . . . .	532
8.4.3	Termination of Newton Iteration . . . . .	534
8.4.4	Damped Newton Method . . . . .	536
8.4.5	Quasi-Newton Method . . . . .	539
8.5	Non-linear Least Squares [DR08, Ch. 6] . . . . .	544
8.5.1	(Damped) Newton Method . . . . .	546
8.5.2	Gauss-Newton Method . . . . .	547
8.5.3	Trust Region Method (Levenberg-Marquardt Method) . . . . .	550
<b>9</b>	<b>Computation of Eigenvalues and Eigenvectors</b> . . . . .	<b>553</b>
9.1	Theory of eigenvalue problems . . . . .	557

9.2	"Direct" Eigensolvers . . . . .	559
9.3	Power Methods . . . . .	563
9.3.1	Direct power method . . . . .	563
9.3.2	Inverse Iteration [DR08, Sect. 7.6], [QSS00, Sect. 5.3.2] . . . . .	572
9.3.3	Preconditioned inverse iteration (PINVIT) . . . . .	584
9.3.4	Subspace iterations . . . . .	587
9.3.4.1	Orthogonalization . . . . .	592
9.3.4.2	Ritz projection . . . . .	597
9.4	Krylov Subspace Methods . . . . .	601
<b>10</b>	<b>Krylov Methods for Linear Systems of Equations</b>	<b>616</b>
10.1	Descent Methods [QSS00, Sect. 4.3.3] . . . . .	617
10.1.1	Quadratic minimization context . . . . .	617
10.1.2	Abstract steepest descent . . . . .	618
10.1.3	Gradient method for s.p.d. linear system of equations . . . . .	619
10.1.4	Convergence of the gradient method . . . . .	620
10.2	Conjugate gradient method (CG) [Han02, Ch. 9], [DR08, Sect. 13.4], [QSS00, Sect. 4.3.4] . . . . .	624
10.2.1	Krylov spaces . . . . .	625
10.2.2	Implementation of CG . . . . .	626
10.2.3	Convergence of CG . . . . .	629
10.3	Preconditioning [DR08, Sect. 13.5], [Han02, Ch. 10], [QSS00, Sect. 4.3.5] . . . . .	634
10.4	Survey of Krylov Subspace Methods . . . . .	640
10.4.1	Minimal residual methods . . . . .	640
10.4.2	Iterations with short recursions [QSS00, Sect. 4.5] . . . . .	641
<b>11</b>	<b>Numerical Integration – Single Step Methods</b>	<b>645</b>
11.1	Initial-Value Problems (IVP) for ODEs . . . . .	645
11.1.1	Modeling with ordinary differential equations: Examples . . . . .	647
11.1.2	Theory of Initial-Value-Problems (IVPs) . . . . .	650
11.1.3	Evolution Operators . . . . .	654
11.2	Introduction: Polygonal Approximation Methods . . . . .	656
11.2.1	Explicit Euler method . . . . .	657
11.2.2	Implicit Euler method . . . . .	659
11.2.3	Implicit midpoint method . . . . .	660
11.3	General single step methods . . . . .	661
11.3.1	Definition . . . . .	661
11.3.2	(Asymptotic) Convergence of Single-Step Methods . . . . .	664
11.4	Explicit Runge-Kutta Methods . . . . .	670
11.5	Adaptive Stepsize Control . . . . .	677
<b>12</b>	<b>Single Step Methods for Stiff Initial Value Problems</b>	<b>692</b>
12.1	Model Problem Analysis . . . . .	693
12.2	Stiff Initial Value Problems . . . . .	706
12.3	Implicit Runge-Kutta Single Step Methods . . . . .	710
12.3.1	The implicit Euler method for stiff IVPs . . . . .	711
12.3.2	Collocation single step methods . . . . .	712
12.3.3	General implicit RK-SSMs . . . . .	715
12.3.4	Model Problem Analysis for Implicit Runge-Kutta Singlge-Step Method (IRK-SSMs) . . . . .	717
12.4	Semi-implicit Runge-Kutta Methods . . . . .	723
12.5	Splitting methods . . . . .	725
<b>Index</b>		<b>732</b>

Symbols . . . . .	746
Examples . . . . .	748
Glossary . . . . .	756

# Chapter 0

## Introduction

### 0.1 Course Fundamentals

#### 0.1.1 Focus of this Course

Emphasis is put

- ▷ on **algorithms** (principles, computational cost, scope, and limitations),
- ▷ on (efficient and stable) **implementation** in **C++** based on the numerical linear algebra **EIGEN**, a **Domain Specific Language** (DSL) embedded into C++.
- ▷ on **numerical experiments** (design and interpretation).

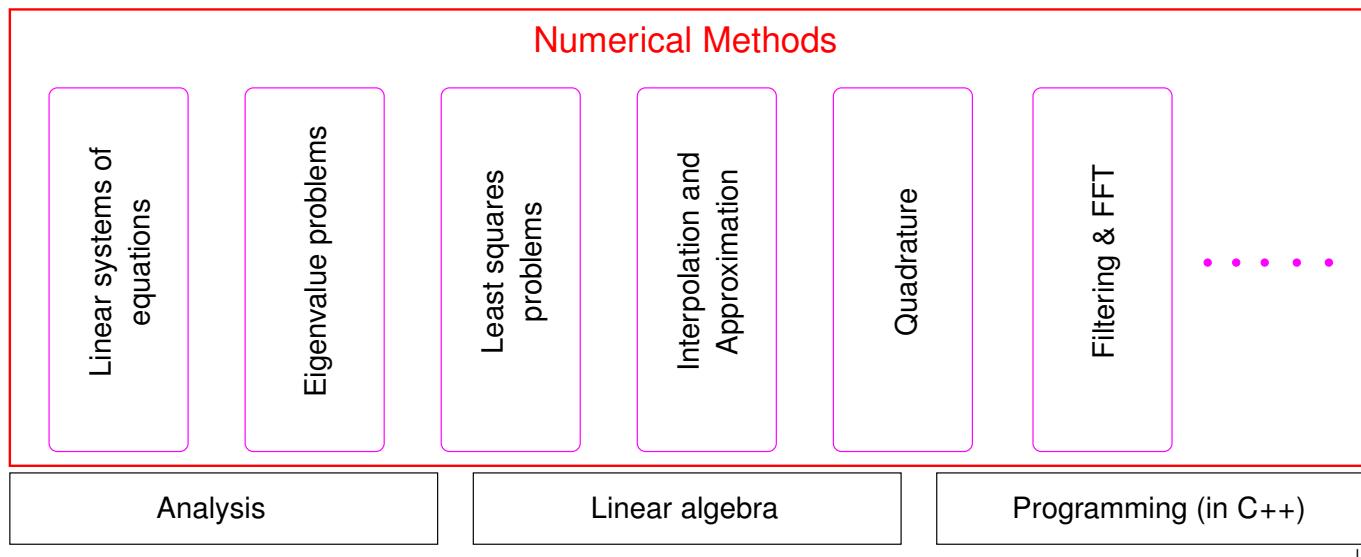
§0.1.1.1 (Aspects outside the scope of this course) No emphasis will be put on

- theory and proofs (unless essential for derivation and understanding of algorithms).
  - ☞ 401-3651-00L Numerical Methods for Elliptic and Parabolic Partial Differential Equations
  - 401-3652-00L Numerical Methods for Hyperbolic Partial Differential Equations
  - (both courses offered in BSc Mathematics)
- hardware aware implementation (cache hierarchies, CPU pipelining, vectorization, etc.)
  - ☞ 263-0007-00L Advanced System Lab (How To Write Fast Numerical Code, Prof. M. Püschel, D-INFK)
- issues of high-performance computing (HPC, shard and distributed memory parallelisation, vectorization)
  - ☞ 151-0107-20L High Performance Computing for Science and Engineering (HPCSE, Prof. P. Koumoutsakos, D-MAVT)
  - 263-2800-00L Design of Parallel and High-Performance Computing (Prof. T. Höfler, D-INFK)

However, note that these other courses partly rely on knowledge of elementary numerical methods, which is covered in this course. □

## Contents

**§0.1.1.2 (Prerequisites)** This course will take for granted basic knowledge of linear algebra, calculus, and programming, that you should have acquired during your first year at ETH.



### §0.1.1.3 (Numerical methods: A motley toolbox)

This course discusses **elementary numerical methods and techniques**

They are vastly different in terms of ideas, design, analysis, and scope of application. They are the items in a **toolbox**, some only loosely related by the common purpose of being building blocks for codes for numerical simulation.



Do not expect much coherence between the chapters of this course!

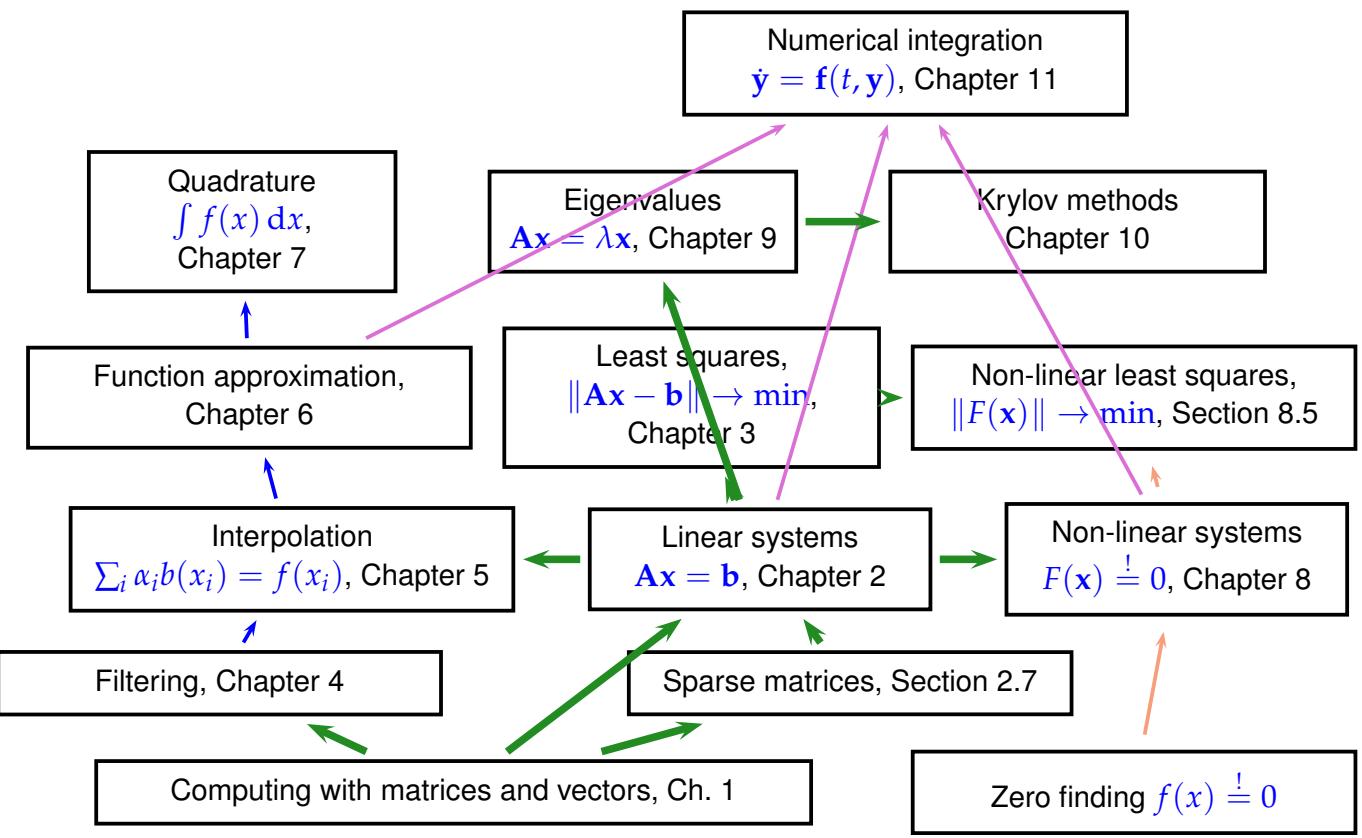


A purpose-oriented notion of “Numerical methods for CSE”:

- A: “Stop putting a hammer, a level, and duct tape in one box! They have nothing to do with each other!”
- B: “I might need any of these tools when fixing something about the house”

Fig. 1

**§0.1.1.4 (Dependencies of topics)** Despite the diverse nature of the individual topics covered in this course, some depend on others for providing essential building blocks. The following directed graph tries to capture these relationships. The arrows have to be read as “uses results or algorithms of”.



Any one-semester course “Numerical methods for CSE” will cover only selected chapters and sections of this document. Only topics addressed in **class** or in **homework problems** will be relevant for exams!

**§0.1.1.5 (Relevance of this course)** I am a student of computer science. After the exam, may I safely forget everything I have learned in this mandatory “numerical methods” course? **No**, because it is highly likely that other courses or projects will rely on the contents of this course:

- singular value decomposition  
least squares }      ➔ Computational statistics, machine learning
- function approximation  
numerical quadrature  
numerical integration }      ➔ Numerical methods for PDEs
- interpolation  
least squares }      ➔ Computer graphics
- eigensolvers  
sparse linear systems }      ➔ Graph theoretic algorithms
- numerical integration }      ➔ Computer animation

and many more applications of fundamental numerical methods . . .

Hardly anyone will need everything covered in this course, but *most of you will need something*.

### 0.1.2 Goals

This course is meant to impart

- ◆ knowledge of some fundamental algorithms forming the basis of numerical simulations,
- ◆ familiarity with essential terms in numerical mathematics and the techniques used for the analysis of numerical algorithms
- ◆ the skill to choose the appropriate numerical methods for concrete problems,
- ◆ the ability to interpret numerical results,
- ◆ proficiency in implementing numerical algorithms efficiently in C++, using numerical libraries.

Indispensable:

Learning by doing (→ exercises)

### 0.1.3 Literature

Parts of the following textbooks may be used as supplementary reading for this course. References to relevant sections will be provided in the course material.

Studying extra literature is not important for following this course!

- ◆ [AG11] U. ASCHER AND C. GREIF, *A First Course in Numerical Methods*, SIAM, Philadelphia, 2011.  
Comprehensive introduction to numerical methods with an algorithmic focus based on MATLAB.  
(Target audience: students of engineering subjects)
- ◆ [DR08] W. DAHMEN AND A. REUSKEN, *Numerik für Ingenieure und Naturwissenschaftler*, Springer, Heidelberg, 2006.

Good reference for large parts of this course; provides a lot of simple examples and lucid explanations, but also rigorous mathematical treatment.

(Target audience: undergraduate students in science and engineering)

Available for download at [PDF](#)

- ◆ [Han02] M. HANKE-BOURGEOIS, *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*, Mathematische Leitfäden, B.G. Teubner, Stuttgart, 2002.  
Gives detailed description and mathematical analysis of algorithms and relies on MATLAB. Profound treatment of theory way beyond the scope of this course. (Target audience: undergraduates in mathematics)
- ◆ [QSS00] A. QUARTERONI, R. SACCO, AND F. SALERI, *Numerical mathematics*, vol. 37 of Texts in Applied Mathematics, Springer, New York, 2000.

Classical introductory numerical analysis text with many examples and detailed discussion of algorithms. (Target audience: undergraduates in mathematics and engineering)  
Can be obtained from [website](#).

- ◆ [DH03] P. DEUFLHARD AND A. HOHMANN, *Numerische Mathematik. Eine algorithmisch orientierte Einführung*, DeGruyter, Berlin, 1 ed., 1991.

Modern discussion of numerical methods with profound treatment of theoretical aspects (Target audience: undergraduate students in mathematics).

- ◆ [GGK14]: W.. GANDER, M.J. GANDER, AND F. KWOK, *Scientific Computing*, Text in Computational Science and Engineering, Springer, 2014.

Comprehensive treatment of elementary numerical methods with an algorithmic focus.

D-INFK maintains a [webpage](#) with links to some of these books.

Essential prerequisite for this course is a solid knowledge in linear algebra and calculus. Familiarity with the topics covered in the first semester courses is taken for granted, see

- ◆ [NS02] K. NIPP AND D. STOFFER, *Lineare Algebra*, vdf Hochschulverlag, Zürich, 5 ed., 2002.
- ◆ [Gut09] M. GUTKNECHT, *Lineare algebra*, lecture notes, SAM, ETH Zürich, 2009, available online.
- ◆ [Str09] M. STRUWE, *Analysis für Informatiker*. Lecture notes, ETH Zürich, 2009, available online.

## 0.2 Teaching Style and Model

### 0.2.1 Flipped Classroom

This course will depart from the usual academic teaching arrangement centering around classes taught by a lecturer addressing an audience in a lecture hall.

#### A **flipped-classroom** course

This course will follow the **flipped-classroom** paradigm:

Learning by **self-study** guided by

instruction videos  
tablet notes

lecture notes

interactive  
Q&A sessions

tutorial classes

All the course material will be published online through the course Moodle Page. All notes jotted down by the lecturer during the creation of videos or during the Q&A sessions will be made available as PDF.

#### 0.2.1.1 Course Videos

In the flipped-classroom teaching model regular lectures will be replaced with pre-recorded videos. These videos are not commercial-grade clips, but resemble video recordings from a standard classroom setting; they convey the development of the material on a tablet accompanied by the lecturer's voice.

Fig. 2

The Videos will be published through the course **Moodle Page**.

Every video comes with a PDF containing the **tablet notes** taken during the creation of the video. However, the PDF may have been *corrected, updated, or supplemented* later.

### §0.2.1.2 (“Pause” and “fast forward”) Videos have *two big advantages*:



Fig. 3

You can stop a video at any time, whenever

- you need more time to think,
- you want to look up related information,
- you want to work for yourself.

Make use of this possibility!

The video portal also allows you to play the videos at *1.5x speed*. This can be useful, if the current topic is very clear to you. You can also *skip* entire *parts* using the scroll bar. ↳

**§0.2.1.3 (Review questions)** Most lecture units (corresponding to a video) are accompanied with a list of review questions. You should try to answer them off the top of your head *without consulting any written material* shortly after you have finished studying the unit . ↳

Failure to answer a review question indicates that you need resume studying some of the unit's topics. ↳

## 0.2.1.2 Following the Course

### Weekly study assignments

- For every week there is a list of course units and associated videos published on the course Moodle Page.
- The corresponding contents **must** be studied in that same week.

### §0.2.1.5 (How to organize your learning)

- ☛ **Develop a routine:** Plan *fixed slots*, with a total duration of *four hours*, for studying for the course material in your weekly calendar. This does not include homework.

- ☛ **Choose a stable setting**, in which you can really concentrate (quiet area, headphones, coffee, etc.)
- ☛ **Take breaks**, when concentration is declining, usually after 20 to 45 minutes, but *avoid online distractions* during breaks.



### You must not procrastinate!

Do not put off studying for this course. Dependencies between the topics will make it very hard to catch up.

**§0.2.1.6 (“Personalized learning”)** The flipped classroom model allows students to pursue their preferred ways of studying. The following approaches can be tried.

- **Traditional:** You watch the assigned videos similar to attending a conventional classroom lecture. Afterwards digest the material based on the tablet notes and/or the lecture document. Finally, answer the review questions and look up more information in the lecture document.
- **Reading-centered:** You work through the unit reading the tablet notes, and, sometimes, related sections of the lecture document. You occasionally watch parts of the videos, in case some considerations and arguments have not become clear to you already.



Fig. 4

**Collaborative studying** is encouraged:

- You may watch course videos together with classmates.
- You may meet to discuss course units.
- You may solve homework problems in a group assigning different parts to different members.
- ☛ Explaining to others is a great way to deepen understanding.
- ☛ It is easy to sustain motivation and avoid distraction in a peer study group.

## 0.2.2 To avoid misunderstandings ...

The PDF you are reading is referred to as lecture document and is an important source of information, but ...

### §0.2.2.1 (“Lecture notes”)

This course document is neither a textbook nor comprehensive lecture notes.  
They are meant to supplement and be supplemented by explanations given in class.

Some pieces of advice:

- ◆ The lecture document is only partly designed to be self-contained and can/should be studied *in parts* in addition to attending to watching the course videos and/or reading the tablet notes.
- ◆ This text is not meant for mere reading, but for working with,
- ◆ Turn pages all the time and follow the numerous cross-references,
- ◆ study the relevant section of the course material when doing homework problems,

- ◆ You may study referenced literature to refresh prerequisite knowledge and for alternative presentation of the material (from a different angle, maybe), but be careful about not getting confused or distracted by information overload.

↓

**§0.2.2.2 (Reporting errors)** As the documents for this course will always be in a state of flux, they will inevitably and invariably teem with small errors, mainly typos and omissions.

Please report errors in the lecture material through the **Course Moodle Page!**

When reporting an error, *please specify the section and the number of the paragraph, remark, equation, etc.* where it hides. You need not give a page number.

↓

### §0.2.2.3 (Comprehension is a process . . .)

- ◆ This course will require

hard work — perseverance — patience

- ◆ Do **not** expect to understand everything at once. Most students will

- understand about one third of the material when watching videos and studying the course material
- understand another third when making a *serious effort* to solve the homework problems,
- hopefully understand the remaining third when studying for the main examination after the end of the course.

Perseverance will be rewarded!

↓

## 0.2.3 Assignments

A steady and persistent effort spent on homework problems is essential for success in this course.

You should expect to spend 3-5 hours per week on trying to solve the homework problems. Since many involve small coding projects, the time it will take an individual student to arrive at a solution is hard to predict.



For the sake of efficiency:  
Avoid coding errors (bugs) in your homework coding projects!

The problems are published online together with plenty of hints. A master solution will also be made available, but it is foolish to read the master solution parallel to working on a problem sheet, because *trying to find the solution on one's own is essential for developing problem solving skills*, though it may occasionally be frustrating.

### §0.2.3.1 (Homeworks and tutors' corrections)

- ◆ The weekly assignments will be a few problems from the **NCSE Problem Collection** available online as PDF, see course Moodle page for the link. The particular problems to be solved will be communicated through that Moodle page every week.

Please note that this problem collection is being extended throughout the semester. Thus, make sure that you **obtain the most current version** every week. A polybox link will also be distributed; if you install the [Polybox Client](#) the most current version of all course documents will always be uploaded to your machine.

- ◆ Some or all of the problems of an assignment sheet will be discussed in the tutorial classes at least one week after the problems have been assigned.

Note: Homework problems can be *handed in* on paper or electronically, if you want to receive feedback about your solution. Please write your name on your homework paper clearly and hand it to your tutor or deposit it in the tray in the aisle in front of HG G 53-54.3 marked "Numerical Methods for Partial Differential Equations". Alternatively you may send your assistant. One week later the corrections will usually be returned.

#### §0.2.3.2 ([CODEEXPERT](#) C++ online IDE and testing environment)



[CODEEXPERT](#) has been developed at ETH as an online IDE for small programming assignment and coding homeworks. It will be used in this course for all C++ homework problems.

Please study the [documentation](#)!

## 0.2.4 Information on Examinations

### §0.2.4.1 (Examinations during the teaching period)

From the ETH course directory:

Computer based examination involving coding problems beside theoretical questions. Parts of the lecture documents and other materials will be made available online during the examination. A 30-minute **mid-term exam** and a 30-minute **end term exam** will be held during the teaching period on dates specified in the beginning of the semester. Points earned in these exams will be taken into account through a **bonus of up to 20%** of the full (100%) points in the final session exam.

Both will be closed book examinations on paper. The dates of the exams will be communicated in the beginning of the term and published on the course webpage.

### §0.2.4.2 (Main examination during exam session)

- ◆ Three-hour written examination involving coding problems to be done at the computer. The date of the exam will be set and communicated by the ETH exam office, and will also be published on the course webpage.
- ◆ A dry-run session for the computer based examination will be offered. Scan the course website for information, in particular concerning the mandatory **registration**.
- ◆ Subjects of examination:

- All topics, which have been addressed in class or in a homework problem (including the homework problems not labelled as “core problems”)

The lecture document contains much more material than covered in class. All these extra topics are not relevant for the exam.

- ◆ Lecture documents (as PDF), the EIGEN documentation, and the online [C++ REFERENCE PAGES](#) will be available PDF during the examination. The corresponding final version of the lecture document will be made available at least two weeks before the exam.
- ◆ No other **materials** may be used during the exam.
- ◆ The homework problem collection cannot be accessed during the exam.
- ◆ The exam questions will be asked in English.
- ◆ In case you come to the conclusion that you have too little time to prepare for the main exam a few weeks before the exam, contemplate withdrawing in order not to squander an attempt.

#### §0.2.4.3 (Repeating the main exam)

- Bonus points earned in term exams in last year's course can be taken into account for this course's main exam.
- If you are going to repeat the main exam, but also want to earn a bonus through this year's term exams, please **declare this intention** before the mid-term exam.

## 0.3 Programming in C++

C++17 is the *current* ANSI/ISO standard for the programming language C++. On the one hand, it offers a wealth of features and possibilities. On the other hand, this can be confusing and even be prone to inconsistencies. A major cause of inconsistent design is the requirement with backward compatibility with the C programming language and the earlier standard C++ 98.

However, C++ has become the main language in computational science and engineering and high performance computing. Therefore this course relies on C++ to discuss the implementation of numerical methods.

In fact C++ is a blend of different programming paradigms:

- an **object oriented** core providing classes, inheritance, and runtime polymorphism,
- a powerful **template mechanism** for parametric types and partial specialization, enabling *template meta-programming* and compile-time polymorphism,
- a collection of abstract data containers and basic algorithms provided by the **Standard Template Library** (STL).



**Supplementary literature.** A popular book for learning C++ that has been upgraded to include the C++11 standard is [[LLM12](#)].

The book [Jos12] gives a comprehensive presentation of the new features of C++11 compared to earlier versions of C++.

There are plenty of online reference pages for C++, for instance <http://en.cppreference.com> and <http://www.cplusplus.com/>.

The following sections highlight a few particular aspects of C++ that may be important for code development in this course.

### 0.3.1 Function Arguments and Overloading

**§0.3.1.1 (Function overloading, [LLM12, Sect. 6.4])** Argument types are an integral part of a function declaration in C++. Hence the following functions are different

```
int* f(int);      // use this in the case of a single numeric argument
double f(int *); // use only, if pointer to a integer is given
void f(const MyClass &); // use when called for a MyClass object
```

and the compiler selects the function to be used depending on the type of the arguments following rather sophisticated rules, refer to [overload resolution rules](#). Complications arise, because implicit type conversions have to be taken into account. In case of ambiguity a compile-time error will be triggered. Functions cannot be distinguished by return type!

For member functions (methods) of classes an additional distinction can be introduced by the **const** specifier:

```
struct MyClass {
    double f(double);           // use for a mutable object of type MyClass
    double f(double) const;     // use this version for a constant object
    ...
};
```

The second version of the method `f` is invoked for *constant objects* of type **MyClass**. □

**§0.3.1.2 (Operator overloading [LLM12, Chapter 14])** In C++ unary and binary **operators** like `=`, `==`, `+`, `-`, `*`, `/`, `+=`, `-=`, `*=`, `/=`, `%`, `&&`, `||`, `<<`, `>>`, etc. are regarded as functions with a fixed number of arguments (one or two). For built-in numeric and logic types they are defined already. They can be extended to any other type, for instance

```
MyClass operator +(const MyClass &, const MyClass &);
MyClass operator +(const MyClass &, double);
MyClass operator +(const MyClass &); // unary + !
```

The same selection rules as for function overloading apply. Of course, operators can also be introduced as class member functions.

C++ gives complete freedom to overload operators. However, the semantics of the new operators should be close to the customary use of the operator. □

**§0.3.1.3 (Passing arguments by value and by reference [LLM12, Sect. 6.2])** Consider a generic function declared as follows:

```
void f(MyClass x); // Argument x passed by value.
```

When `f` is invoked, a *temporary copy* of the argument is created through the **copy constructor** or the **move constructor** of **MyClass**. The new temporary object is a *local variable* inside the function body.

When a function is declared as follows

```
void f( MyClass &x); // Argument x passed by reference.
```

then the argument is passed to the scope of the function and can be changed inside the function. *No copies* are created. If one wants to avoid the creation of temporary objects, which may be costly, but also wants to indicate that the argument will not be modified inside `f`, then the declaration should read

```
void f(const MyClass &x); // Argument x passed by constant reference.
```

New in C++11 is **move semantics**, enabled in the following definition

```
void f(const MyClass &&x); // Optional shallow copy
```

In this case, if the scope of the object passed as the argument is merely the function or `std::move()` tags it as disposable, the **move constructor** of **MyClass** is invoked, which will usually do a *shallow copy* only. Refer to Code 0.3.5.10 for an example. ↴

## 0.3.2 Templates

**§0.3.2.1 (Function templates)** The template mechanism supports parameterization of definitions of classes and functions by type. An example of a **function templates** is

```
template <typename ScalarType, typename VectorType>
VectorType saxpy(ScalarType alpha, const VectorType &x, const
                  VectorType &y)
{ return (alpha*x+y); }
```

Depending on the concrete type of the arguments the compiler will instantiate particular versions of this function, for instance `saxpy<float, double>`, when `alpha` is of type `float` and both `x` and `y` are of type `double`. In this case the return type will be `float`.

For the above example the compiler will be able to deduce the types `ScalarType` and `VectorType` from the arguments. The programmer can also specify the types directly through the `< >`-syntax as in

```
saxpy<double , double>(a,x,y);
```

If an instantiation for all arguments of type `double` is desired. In case, the arguments do not supply enough information about the type parameters, specifying (some of) them through `< >` is mandatory. ↴

**§0.3.2.2 (Class templates)** A **class template** defines a class depending on one or more type parameters, for instance

```
template <typename T>
class MyClsTempl {
public:
    using parm_t = T; // T-dependent type
    MyClsTempl(void); // Default constructor
    MyClsTempl(const T&); // Constructor with an argument
    template <typename U>
    T memfn(const T&, const U&) const; // Templatized member function
```

```
private:
    T *ptr;      // Data member, T-pointer
};
```

Types **MyClsTempl<T>** for a concrete choice of **T** are instantiated when a corresponding object is declared, for instance via

```
double x = 3.14;
 MyClass myobj; // Default construction of an object
 MyClsTempl<double> tinstd; // Instantiation for T = double
 MyClsTempl<MyClass> mytinst(myobj); // Instantiation for T = MyClass
 MyClass ret = mytinst.memfn(myobj,x); // Instantiation of member
 function for U = double, automatic type deduction
```

The types spawned by a template for different parameter types have nothing to do with each other. □

### Requirements on parameter types

The parameter types for a template have to provide all type definitions, member functions, operators, and data to make possible the instantiation (“compilation”) of the class of function template.

### 0.3.3 Function Objects and Lambda Functions

A function object is an object of a type that provides an overloaded “function call” **operator ()**. Function objects can be implemented in two different ways:

- (I) through special classes like the following that realizes a function  $\mathbb{R} \mapsto \mathbb{R}$

```
class MyFun {
public:
    ...
    double operator (double x) const; // Evaluation operator
    ...
};
```

The evaluation operator can take more than one argument and need not be declared **const**.

- (II) through **lambda functions**, an “anonymous function” defined as

```
[<capture list>] (<arguments>) -> <return type> { body; }
```

where **<capture list>** is a list of variables from the local scope to be passed to the lambda function; an **&** indicates passing by reference,  
**<arguments>** is a comma separated list of function arguments complete with types,  
**<return type>** is an *optional* return type; often the compiler will be able to deduce the return type from the definition of the function.

Function classes should be used, when the function is needed in different places, whereas lambda functions for short functions intended for single use.

#### C++ code 0.3.3.1: Demonstration of use of lambda function → [GITLAB](#)

```
1 int main() {
2     // initialize a vector from an initializer list
```

```

3   std::vector<double> v({1.2,2.3,3.4,4.5,5.6,6.7,7.8});
4   // A vector of the same length
5   std::vector<double> w(v.size());
6   // Do cumulative summation of v and store result in w
7   double sum = 0;
8   std::transform(v.begin(),v.end(),w.begin(),
9                 [&sum] (double x) { sum += x; return sum; });
10  cout << "sum = " << sum << ", w = [ ";
11  for(auto x: w) cout << x << ' '; cout << ']' << endl;
12  return(0);
13 }

```

In this code the lambda function captures the local variable `sum` by reference, which enables the lambda function to change its value in the surrounding scope.

**§0.3.3.2 (Function type wrappers)** The special class `std::function` provides types for general polymorphic function wrappers.

```
std::function<return type(arg types)>
```

#### C++ code 0.3.3.3: Use of `std::function` → [GITLAB](#)

```

1 double binop(double arg1,double arg2) { return (arg1/arg2); }
2
3 void stdfunctiontest(void) {
4   // Vector of objects of a particular signature
5   std::vector<std::function<double(double,double)>> fnvec;
6   // Store reference to a regular function
7   fnvec.push_back(binop);
8   // Store a
9   lambda function
10  fnvec.push_back([] (double x,double y) -> double { return y/x; });
11  for (auto fn : fnvec) { std::cout << fn(3,2) << std::endl; }
12 }

```

In this example an object of type `std::function<double (double, double)>` can hold a regular function taking two `double` arguments and returning another `double` or a `lambda function` with the same signature. Guess the output of `stdfunctiontest!` ↴

### 0.3.4 Multiple Return Values

In PYTHON it is customary to return several variables from a function call, which, in fact, amounts to returning a tuple of mixed-type objects:

```

1 def f(a, b):
2     return min(a, b), max(a, b), (a+b)/2
3 x, y, z = f(1, 2)

```

In C++ this is also possible by using the `tuple` utility. For instance, the following function computes the mimimal and maximal element of a vector and also returns its cumulative sum. It returns all these values.

#### C++ code 0.3.4.1: Function with multiple return values → [GITLAB](#)

```

1 template<typename T>
2 std::tuple<T,T,std::vector<T>> extcumsum(const std::vector<T> &v) {

```

```

3 // Local summation variable captured by reference by lambda function
4 T sum{};
5 // temporary vector for returning cumulative sum
6 std::vector<T> w{};
7 // cumulative summation
8 std::transform(v.cbegin(), v.cend(), back_inserter(w),
9 [&sum] (T x) { sum += x; return(sum); });
10 return (std::make_tuple(*std::min_element(v.cbegin(), v.cend()),
11 *std::max_element(v.cbegin(), v.cend()),
12 std::move(w)));
13 }

```

This code snippet shows how to extract the individual components of the tuple returned by the previous function.

#### C++ code 0.3.4.2: Calling a function with multiple return values → [GITLAB](#)

```

1 int main () {
2     // initialize a vector from an initializer list
3     std::vector<double> v({1.2,2.3,3.4,4.5,5.6,6.7,7.8});
4     // Variables for return values
5     double minv,maxv; // Extremal elements
6     std::vector<double> cs; // Cumulative sums
7     std::tie(minv,maxv,cs) = extcumsum(v);
8     cout << "min = " << minv << ", max = " << maxv << endl;
9     cout << "CS = [ "; for(double x: cs) cout << x << ' '; cout << "]" << endl;
10    return(0);
11 }

```

Be careful: many temporary objects might be created! A demonstration of this hidden cost is given in Exp. 0.3.5.27. From C++17 a more compact syntax is available:

#### C++ code 0.3.4.3: Calling a function with multiple return values → [GITLAB](#)

```

1 int main () {
2     // initialize a vector from an initializer list
3     std::vector<double> v({1.2,2.3,3.4,4.5,5.6,6.7,7.8});
4     // Definition of variables and assignment of return values all at once
5     auto [minv, maxv, cs] = extcumsum(v);
6     cout << "min = " << minv << ", max = " << maxv << endl;
7     cout << "CS = [ "; for(double x: cs) cout << x << ' '; cout << "]" << endl;
8     return(0);
9 }

```

**Remark 0.3.4.4 (“auto” considered harmful)** C++ is a strongly typed programming language and every variable must have a precise type. However, the developer of templated classes and functions may not know the type of some variables in advance, because it can be deduced only after instantiation through the compiler. The **auto** keyword has been introduced to handle this situation.

There is a temptation to use **auto** profligately, because it is convenient, in particular when *using* templated data types. However, this denies a major benefit of types, consistency checking at compile time and, as a developer, one may eventually lose track of the types completely, which can lead to errors that are hard to detect.

Thus, the use of **auto** should be avoided, unless in the following situations:

- for variables inside templated functions or classes, whose precise type will only become clear during instantiation,
- for lambda functions, see Section 0.3.3,
- for return values of templated library (member) functions, whose type is “impossible to deduce” by the user. An example is expression templates in EIGEN, refer to Rem. 1.2.1.11 below.

### 0.3.5 A Vector Class

Since C++ is an object oriented programming language, datatypes defined by **classes** play a pivotal role in every C++ program. Here, we demonstrate the main ingredients of a class definition and other important facilities of C++ for the class **MyVector** meant for objects representing vectors from  $\mathbb{R}^n$ . The codes can be found in → [GITLAB](#). A similar vector class is presented in [Fri19, Ch. 6].

#### C++ 11 class 0.3.5.1: Definition of a simple vector class **MyVector** → [GITLAB](#)

```

1  namespace myvec {
2  class MyVector {
3  public:
4      using value_t = double;
5      // Constructor creating constant vector, also default constructor
6      explicit MyVector(std::size_t n = 0, double val = 0.0);
7      // Constructor: initialization from an STL container
8      template <typename Container> MyVector(const Container &v);
9      // Constructor: initialization from an STL iterator range
10     template <typename Iterator> MyVector(Iterator first, Iterator last);
11     // Copy constructor, computational cost O(n)
12     MyVector(const MyVector &mv);
13     // Move constructor, computational cost O(1)
14     MyVector(MyVector &&mv);
15     // Copy assignment operator, computational cost O(n)
16     MyVector &operator = (const MyVector &mv);
17     // Move assignment operator, computational cost O(1)
18     MyVector &operator = (MyVector &&mv);
19     // Destructor
20     virtual ~MyVector(void);
21     // Type conversion to STL vector
22     operator std::vector<double> () const;
23
24     // Returns length of vector
25     std::size_t size(void) const { return n; }
26     // Access operators: rvalue & lvalue, with range check
27     double operator [] (std::size_t i) const;
28     double &operator [] (std::size_t i);
29     // Comparison operators
30     bool operator == (const MyVector &mv) const;
31     bool operator != (const MyVector &mv) const;
32     // Transformation of a vector by a function  $\mathbb{R} \rightarrow \mathbb{R}$ 
33     template <typename Functor>
34     MyVector &transform(Functor &&f);
35
36     // Overloaded arithmetic operators
37     // In place vector addition:  $x += y;$ 
38     MyVector &operator +=(const MyVector &mv);
39     // In place vector subtraction:  $x -= y;$ 

```

```

40  MyVector &operator -=(const MyVector &mv);
41  // In place scalar multiplication: x *= a;
42  MyVector &operator *=(double alpha);
43  // In place scalar division: x /= a;
44  MyVector &operator /=(double alpha);
45  // Vector addition
46  MyVector operator + (MyVector mv) const;
47  // Vector subtraction
48  MyVector operator - (const MyVector &mv) const;
49  // Scalar multiplication from right and left: x = a*y; x = y*a
50  MyVector operator * (double alpha) const;
51  friend MyVector operator * (double alpha,const MyVector &);
52  // Scalar division: x = y/a;
53  MyVector operator / (double alpha) const;
54
55  // Euclidean norm
56  double norm(void) const;
57  // Euclidean inner product
58  double operator *(const MyVector &) const;
59  // Output operator
60  friend std::ostream &
61  operator << (std::ostream &,const MyVector &mv);
62
63  static bool dbg; // Flag for verbose output
64 private:
65  std::size_t n; // Length of vector
66  double *data; // data array (standard C array)
67 };
68 }

```

Note the use of a public **static** data member `dbg` in Line 63 that can be used to control debugging output by setting `MyVector::dbg = true` or `MyVector::dbg = false`.

**Remark 0.3.5.2 (Contiguous arrays in C++)** The class **MyVector** uses a C-style array and dynamic memory management with **new** and **delete** to store the vector components. This is for demonstration purposes only and not recommended.

### Arrays in C++

In C++ use the **STL container `std::vector<T>`** for storing data in contiguous memory locations.  
Exception: use **`std::array<T>`**, if the number of elements is known at compile time.

### §0.3.5.4 (Member and friend functions of **MyVector** → [GITLAB](#))

#### C++ code 0.3.5.5: Constructor for constant vector, also default constructor, see Line 28

```

1  MyVector::MyVector(std::size_t _n,double _a):n(_n),data(nullptr) {
2      if (dbg) cout << "{Constructor MyVector(" << _n
3          << ") called" << '}' << endl;
4      if (_n > 0) data = new double [_n];
5      for (std::size_t l=0;l<_n;++l) data[l] = _a;
6  }

```

This constructor can also serve as default constructor (a constructor that can be invoked without any argument), because defaults are supplied for all its arguments.

The following two constructors initialize a vector from sequential containers according to the conventions of the STL.

#### C++ code 0.3.5.6: Templated **constructors** copying vector entries from an STL container

[→ GITLAB](#)

```

1 template <typename Container>
2 MyVector::MyVector(const Container &v):n(v.size()),data(nullptr) {
3     if (dbg) cout << "{MyVector(length " << n
4                     << ") constructed from container" << '}' << endl;
5     if (n > 0) {
6         double *tmp = (data = new double [n]);
7         for(auto i: v) *tmp++ = i; // foreach loop
8     }
9 }
```

Note the use of the new C++ 11 facility of a “foreach loop” iterating through a container in Line 7.

#### C++ code 0.3.5.7: Constructor initializing vector from STL iterator range [→ GITLAB](#)

```

1 template <typename Iterator>
2 MyVector::MyVector(Iterator first,Iterator last):n(0),data(nullptr) {
3     n = std::distance(first,last);
4     if (dbg) cout << "{MyVector(length " << n
5                     << ") constructed from range" << '}' << endl;
6     if (n > 0) {
7         data = new double [n];
8         std::copy(first,last,data);
9     }
10 }
```

The use of these constructors is demonstrated in the following code

#### C++ code 0.3.5.8: Initialization of a **MyVector** object from an STL vector [→ GITLAB](#)

```

1 int main() {
2     myvec::MyVector::dbg = true;
3     std::vector<int> ivec = { 1,2,3,5,7,11,13 }; // initializer list
4     myvec::MyVector v1(ivec.cbegin(),ivec.cend());
5     myvec::MyVector v2(ivec);
6     myvec::MyVector vr(ivec.crbegin(),ivec.crend());
7     cout << "v1 = " << v1 << endl;
8     cout << "v2 = " << v2 << endl;
9     cout << "vr = " << vr << endl;
10    return(0);
11 }
```

The following output is produced:

```
{MyVector(length 7) constructed from range}
{MyVector(length 7) constructed from container}
{MyVector(length 7) constructed from range}
v1 = [ 1,2,3,5,7,11,13 ]
```

```
v2 = [ 1,2,3,5,7,11,13 ]
vr = [ 13,11,7,5,3,2,1 ]
{ Destructor for MyVector(length = 7)}
{ Destructor for MyVector(length = 7)}
{ Destructor for MyVector(length = 7)}
```

The copy constructor listed next relies on the *STL algorithm* `std::copy` to copy the elements of an existing object into a newly created object. This takes  $n$  operations.

#### C++ code 0.3.5.9: Copy constructor → GITLAB

```
1 MyVector::MyVector(const MyVector &mv) :n(mv.n),data(nullptr) {
2     if (dbg) cout << "{Copy construction of MyVector(length "
3             << n << ")" << '}' << endl;
4     if (n > 0) {
5         data = new double [n];
6         std::copy_n(mv.data,n,data);
7     }
8 }
```

An important new feature of C++11 is **move semantics** which helps avoid expensive copy operations. The following implementation just performs a shallow copy of pointers and, thus, for large  $n$  is much cheaper than a call to the copy constructor from Code 0.3.5.9. The source vector is left in an empty vector state.

#### C++ code 0.3.5.10: Move constructor → GITLAB

```
1 MyVector::MyVector(MyVector &&mv) :n(mv.n),data(mv.data) {
2     if (dbg) cout << "{Move construction of MyVector(length "
3             << n << ")" << '}' << endl;
4     mv.data = nullptr; mv.n = 0; // Reset victim of data theft
5 }
```

The following code demonstrates the use of `std::move()` to mark a vector object as disposable and allow the compiler the use of the move constructor. The code also uses left multiplication with a scalar, see Code 0.3.5.23.

#### C++ code 0.3.5.11: Invocation of copy and move constructors → GITLAB

```
1 int main() {
2     myvec::MyVector::dbg = true;
3     myvec::MyVector v1(std::vector<double>(
4         {1.2,2.3,3.4,4.5,5.6,6.7,7.8,8.9}));
5     myvec::MyVector v2(2.0*v1); // Scalar multiplication
6     myvec::MyVector v3(std::move(v1));
7     cout << "v1 = " << v1 << endl;
8     cout << "v2 = " << v2 << endl;
9     cout << "v3 = " << v3 << endl;
10    return (0);
11 }
```

This code produces the following output. We observe that `v1` is empty after its data have been “stolen” by `v2`.

```
{MyVector(length 8) constructed from container}
{operator a*, MyVector of length 8}
```

```
{Copy construction of MyVector(length 8)}
{operator *=, MyVector of length 8}
{Move construction of MyVector(length 8)}
{Destructor for MyVector(length = 0)}
{Move construction of MyVector(length 8)}
v1 = [ ]
v2 = [ 2.4,4.6,6.8,9,11.2,13.4,15.6,17.8 ]
v3 = [ 1.2,2.3,3.4,4.5,5.6,6.7,7.8,8.9 ]
{Destructor for MyVector(length = 8)}
{Destructor for MyVector(length = 8)}
{Destructor for MyVector(length = 0)}
```

We observe that the object `v1` is reset after having been moved to `v3`.



Use `std::move` only for special purposes like above and only if an object has a move constructor. Otherwise a 'move' will trigger a plain copy operation. In particular, do not use `std::move` on objects at the end of their scope, e.g., within `return` statements.

The next operator effects copy assignment of an rvalue `MyVector` object to an lvalue `MyVector`. This involves  $O(n)$  operations.

#### C++ code 0.3.5.12: Copy assignment operator → GITLAB

```
1 MyVector &MyVector::operator = (const MyVector &mv) {
2     if (dbg) cout << "{Copy assignment of MyVector(length "
3             << n << "←" << mv.n << ")"<< '}' << endl;
4     if (this == &mv) return(*this);
5     if (n != mv.n) {
6         n = mv.n;
7         if (data != nullptr) delete [] data;
8         if (n > 0) data = new double [n]; else data = nullptr;
9     }
10    if (n > 0) std::copy_n(mv.data,n,data);
11    return(*this);
12 }
```

The move semantics is realized by an assignment operator relying on shallow copying.

#### C++ code 0.3.5.13: Move assignment operator → GITLAB

```
1 MyVector &MyVector::operator = (MyVector &&mv) {
2     if (dbg) cout << "{Move assignment of MyVector(length "
3             << n << "←" << mv.n << ")"<< '}' << endl;
4     if (data != nullptr) delete [] data;
5     n = mv.n; data = mv.data;
6     mv.n = 0; mv.data = nullptr;
7     return(*this);
8 }
```

The destructor releases memory allocated by `new` during construction or assignment.

#### C++ code 0.3.5.14: Destructor: releases allocated memory → GITLAB

```
1 MyVector::~MyVector(void) {
2     if (dbg) cout << "{Destructor for MyVector(length = "
```

```

3     << n << ")" << '}' << endl;
4     if (data != nullptr) delete [] data;
5 }
```

The **operator** keyword is also used to define [implicit type conversions](#).

**C++ code 0.3.5.15: Type conversion operator: copies contents of vector into STL vector** → [GITLAB](#)

```

1 MyVector::operator std::vector<double> () const {
2     if (dbg) cout << "{Conversion to std::vector, length = " << n << '}' << endl;
3     return std::vector<double>(data, data+n);
4 }
```

The bracket operator `[]` can be used to fetch and set vector components. Note that index range checking is performed; an *exception* is thrown for invalid indices. The following code also gives an example of operator overloading as discussed in § 0.3.1.2.

**C++ code 0.3.5.16: rvalue and lvalue access operators** → [GITLAB](#)

```

1 double MyVector::operator [] (std::size_t i) const {
2     if (i >= n) throw(std::logic_error("[] out of range"));
3     return data[i];
4 }
5
6 double &MyVector::operator [] (std::size_t i) {
7     if (i >= n) throw(std::logic_error("[] out of range"));
8     return data[i];
9 }
```

Componentwise direct comparison of vectors. Can be dangerous in numerical codes, cf. Rem. 1.5.3.15.

**C++ code 0.3.5.17: Comparison operators** → [GITLAB](#)

```

1 bool MyVector::operator == (const MyVector &mv) const
2 {
3     if (dbg) cout << "{Comparison ==: " << n << " <-> " << mv.n << '}' << endl;
4     if (n != mv.n) return(false);
5     else {
6         for(std::size_t l=0;l<n;++l)
7             if (data[l] != mv.data[l]) return(false);
8     }
9     return(true);
10 }
11
12 bool MyVector::operator != (const MyVector &mv) const {
13     return !(*this == mv);
14 }
```

The `transform` method applies a function to every vector component and overwrites it with the value returned by the function. The function is passed as an object of a type providing a `()`-operator that accepts a single argument convertible to `double` and returns a value convertible to `double`.

**C++ code 0.3.5.18: Transformation of a vector through a functor  $\text{double} \rightarrow \text{double}$**  → [GITLAB](#)

```

1 template <typename Functor>
2 MyVector &MyVector::transform(Functor &&f) {
3     for(std::size_t l=0;l<n;++l) data[l] = f(data[l]);
4     return(*this);
5 }
```

The following code demonstrates the use of the `transform` method in combination with

1. a **function object** of the following type

**C++ code 0.3.5.19: A functor type**

```

1 struct SimpleFunction {
2     SimpleFunction(double _a = 1.0):cnt(0),a(_a) {}
3     double operator () (double x) { cnt++; return(x+a); }
4     int cnt;           // internal counter
5     const double a;   // increment value
6 };
```

2. a **lambda function** defined directly inside the call to `transform`.

**C++ code 0.3.5.20: transformation of a vector via a functor object**

```

1 int main() {
2     myvec::MyVector::dbg = false;
3     double a = 2.0; // increment
4     int cnt = 0;    // external counter used by lambda function
5     myvec::MyVector mv(std::vector<double>(
6         {1.2,2.3,3.4,4.5,5.6,6.7,7.8,8.9}));
7     mv.transform([a,&cnt] (double x) { cnt++; return(x+a); });
8     cout << cnt << " operations, mv transformed = " << mv << endl;
9     SimpleFunction trf(a); mv.transform(trf);
10    cout << trf.cnt << " operations, mv transformed = " << mv << endl;
11    mv.transform(SimpleFunction(-4.0));
12    cout << "Final vector = " << mv << endl;
13    return(0);
14 }
```

The output is

```

8 operations , mv transformed = [ 3.2,4.3,5.4,6.5,7.6,8.7,9.8,10.9 ]
8 operations , mv transformed = [ 5.2,6.3,7.4,8.5,9.6,10.7,11.8,12.9 ]
Final vector = [ 1.2,2.3,3.4,4.5,5.6,6.7,7.8,8.9 ]
```

Operator overloading provides the “natural” vector operations in  $\mathbb{R}^n$  both in place and with a new vector created for the result.

**C++ code 0.3.5.21: In place arithmetic operations (one argument)** → [GITLAB](#)

```

1 MyVector &MyVector::operator+=(const MyVector &mv) {
2     if (dbg) cout << "{operator +=, MyVector of length "
3         << n << '}' << endl;
4     if (n != mv.n) throw(std::logic_error("+=: vector size mismatch"));
```

```

5   for(std::size_t l=0;l<n;++l) data[l] += mv.data[l];
6   return(*this);
7 }

8 MyVector &MyVector::operator -=(const MyVector &mv) {
9   if (dbg) cout << "{operator -=, MyVector of length "
10      << n << '}' << endl;
11   if (n != mv.n) throw(std::logic_error("-=: vector size mismatch"));
12   for(std::size_t l=0;l<n;++l) data[l] -= mv.data[l];
13   return(*this);
14 }

15 MyVector &MyVector::operator *=(double alpha) {
16   if (dbg) cout << "{operator *=, MyVector of length "
17      << n << '}' << endl;
18   for(std::size_t l=0;l<n;++l) data[l] *= alpha;
19   return(*this);
20 }

21 MyVector &MyVector::operator /=(double alpha) {
22   if (dbg) cout << "{operator /=, MyVector of length "
23      << n << '}' << endl;
24   for(std::size_t l=0;l<n;++l) data[l] /= alpha;
25   return(*this);
26 }

27

```

### C++ code 0.3.5.22: Binary arithmetic operators (two arguments) → [GITLAB](#)

```

1 MyVector MyVector::operator + (MyVector mv) const {
2   if (dbg) cout << "{operator +, MyVector of length "
3      << n << '}' << endl;
4   if (n != mv.n) throw(std::logic_error("+: vector size mismatch"));
5   mv += *this;
6   return(mv);
7 }

8 MyVector MyVector::operator - (const MyVector &mv) const {
9   if (dbg) cout << "{operator -, MyVector of length "
10      << n << '}' << endl;
11   if (n != mv.n) throw(std::logic_error("-: vector size mismatch"));
12   MyVector tmp(*this); tmp -= mv;
13   return(tmp);
14 }

15 MyVector MyVector::operator * (double alpha) const {
16   if (dbg) cout << "{operator *a, MyVector of length "
17      << n << '}' << endl;
18   MyVector tmp(*this); tmp *= alpha;
19   return(tmp);
20 }

21 MyVector MyVector::operator / (double alpha) const {
22   if (dbg) cout << "{operator /, MyVector of length " << n << '}' << endl;
23   MyVector tmp(*this); tmp /= alpha;
24   return(tmp);
25 }

26

```

**C++ code 0.3.5.23: Non-member function for left multiplication with a scalar → GITLAB**

```

1 MyVector operator * (double alpha, const MyVector &mv) {
2     if (MyVector::dbg) cout << "{operator a*, MyVector of length "
3             << mv.n << '}' << endl;
4     MyVector tmp(mv); tmp *= alpha;
5     return(tmp);
6 }
```

**C++ code 0.3.5.24: Euclidean norm → GITLAB**

```

1 double MyVector::norm(void) const {
2     if (dbg) cout << "{norm: MyVector of length " << n << '}' << endl;
3     double s = 0;
4     for(std::size_t l=0;l<n;++l) s += (data[l]*data[l]);
5     return(std::sqrt(s));
6 }
```

Adopting the notation in some linear algebra texts, the operator `*` has been chosen to designate the Euclidean inner product:

**C++ code 0.3.5.25: Euclidean inner product → GITLAB**

```

1 double MyVector::operator *(const MyVector &mv) const {
2     if (dbg) cout << "{dot *, MyVector of length " << n << '}' << endl;
3     if (n != mv.n) throw(std::logic_error("dot: vector size mismatch"));
4     double s = 0;
5     for(std::size_t l=0;l<n;++l) s += (data[l]*mv.data[l]);
6     return(s);
7 }
```

At least for debugging purposes every reasonably complex class should be equipped with output functionality.

**C++ code 0.3.5.26: Non-member function output operator → GITLAB**

```

1 std::ostream &operator << (std::ostream &o,const MyVector &mv) {
2     o << "[ ";
3     for(std::size_t l=0;l<mv.n;++l)
4         o << mv.data[l] << ((l==mv.n-1)?' ':' ');
5     return(o << "]");
6 }
```

**EXPERIMENT 0.3.5.27 (“Behind the scenes” of `MyVector` arithmetic)** The following code highlights the use of operator overloading to obtain readable and compact expressions for vector arithmetic.

**C++ code 0.3.5.28:**

```

1 int main() {
2     myvec :: MyVector :: dbg = true;
3     myvec :: MyVector x (std :: vector < double > ({1.2,2.3,3.4,4.5,5.6,6.7,7.8,8.9}));
4     myvec :: MyVector y (std :: vector < double > ({2.1,3.2,4.3,5.4,6.5,7.6,8.7,9.8}));
```

```

5   auto z = x+(x*y)*x+2.0*y/(x-y).norm();
6 }
```

We run the code and trace calls. This is printed to the console:

```

{MyVector(length 8) constructed from container}
{MyVector(length 8) constructed from container}
{dot *, MyVector of length 8}
{operator a*, MyVector of length 8}
{Copy construction of MyVector(length 8)}
{operator *=, MyVector of length 8}
{operator +, MyVector of length 8}
{operator +=, MyVector of length 8}
{Move construction of MyVector(length 8)}
{operator a*, MyVector of length 8}
{Copy construction of MyVector(length 8)}
{operator *=, MyVector of length 8}
{operator +, MyVector of length 8}
{Copy construction of MyVector(length 8)}
{operator -=, MyVector of length 8}
{norm: MyVector of length 8}
{operator /, MyVector of length 8}
{Copy construction of MyVector(length 8)}
{operator *=, MyVector of length 8}
{operator +, MyVector of length 8}
{operator +=, MyVector of length 8}
{Move construction of MyVector(length 8)}
{Destructor for MyVector(length = 0)}
{Destructor for MyVector(length = 8)}
{Destructor for MyVector(length = 8)}
{Destructor for MyVector(length = 8)}
{Destructor for MyVector(length = 0)}
{Destructor for MyVector(length = 8)}
{Destructor for MyVector(length = 8)}
{Destructor for MyVector(length = 8)}
```

Several temporary objects are created and destroyed and quite a few *copy operations* take place. The situation would be worse unless move semantics was available; if we had not supplied a move constructor, a few more copy operations would have been triggered. Even worse, the frequent copying of data runs a high risk of cache misses. This is certainly not an efficient way to do elementary vector operations though it looks elegant at first glance.

**EXAMPLE 0.3.5.29 (Gram-Schmidt orthonormalization based on **MyVector** implementation)** Gram-Schmidt orthonormalization has been taught in linear algebra and its theory will be revisited in § 1.5.1.1. Here we use this simple algorithm from linear algebra to demonstrate the use of the vector class **MyVector** defined in Code 0.3.5.1.

The templated function `gramschmidt` takes a sequence of vectors stored in a **std::vector** object. The actual vector type is passed as a template parameter. It has to supply `length` and `norm` member functions as well as in place arithmetic operations `-=`, `/` and `=`. Note the use of the highlighted methods of the **std::vector** class.

**C++ code 0.3.5.30: templated function for Gram-Schmidt orthonormalization → GITLAB**

```

1 template <typename Vec>
2 std::vector<Vec> gramschmidt(const std::vector<Vec> &A, double eps=1E-14) {
3     const int k = A.size(); // no. of vectors to orthogonalize
4     const int n = A[0].size(); // length of vectors
5     cout << "gramschmidt orthogonalization for " << k << ' ' << n << "-vectors" <<
6         endl;
7     std::vector<Vec> Q({A[0]/A[0].norm()}); // output vectors
8     for(int j=1;(j<k) && (j<n);++j) {
9         Q.push_back(A[j]);
10        for(int l=0;l<j;++l) Q.back() -= (A[j]*Q[l])*Q[l];
11        if (Q.back().norm() < eps*A[j].norm()) { // premature termination ?
12            Q.pop_back(); break;
13        }
14        Q.back() /= Q.back().norm(); // normalization
15    }
16    return (Q); // return at end of local scope
}

```

This driver program calls a function that initializes a sequence of vectors and then orthonormalizes them by means of the Gram-Schmidt algorithm. Eventually orthonormality of the computed vectors is tested. Please pay attention to

- the use of auto to avoid cumbersome type declarations,
- the for loops following the “foreach” syntax.
- automatic indirect template type deduction for the templated function `gramschmidt` from its argument. In Line 6 the function `gramschmidt<MyVector>` is instantiated.

**C++ code 0.3.5.31: Driver code for Gram-Schmidt orthonormalization**

```

1 int main() {
2     myvec::MyVector::dbg = false;
3     const int n = 7; const int k = 7;
4     auto A(initvectors(n,k,[] (int i,int j)
5         { return std::min(i+1,j+1); }));
6     auto Q(gramschmidt(A)); // instantiate template for MyVector
7     cout << "Set of vectors to be orthonormalized:" << endl;
8     for (const auto &a: A) { cout << a << endl; }
9     cout << "Output of Gram-Schmidt orthonormalization: " << endl;
10    for (const auto &q: Q) { cout << q << endl; }
11    cout << "Testing orthogonality:" << endl;
12    for (const auto &qi: Q) {
13        for (const auto &qj: Q)
14            cout << std::setprecision(3) << std::setw(9) << qi*qj << ' ';
15        cout << endl; }
16    return(0);
}

```

This initialization function takes a functor argument as discussed in Section 0.3.3.

**C++ code 0.3.5.32: Initialization of a set of vectors through a functor with two arguments**

```

1 template<typename Functor>
2 std::vector<myvec::MyVector>
3     initvectors(std::size_t n,std::size_t k,Functor &&f) {

```

```
4  std :: vector<MyVector> A{};  
5  for( int j=0;j<k;++j ) {  
6      A.push_back( MyVector(n) );  
7      for( int i=0;i<n;++i )  
8          (A.back())[i] = f(i,j);  
9     }  
10    return(A);  
11 }
```

□

# Bibliography

- [AG11] Uri M. Ascher and Chen Greif. *A first course in numerical methods*. Vol. 7. Computational Science & Engineering. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2011, pp. xxii+552. DOI: [10.1137/1.9780898719987](https://doi.org/10.1137/1.9780898719987) (cit. on p. 11).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on p. 11).
- [DH03] P. Deuflhard and A. Hohmann. *Numerical Analysis in Modern Scientific Computing*. Vol. 43. Texts in Applied Mathematics. Springer, 2003 (cit. on p. 12).
- [Fri19] F. Friedrich. *Datenstrukturen und Algorithmen*. Lecture slides. 2019 (cit. on p. 23).
- [GGK14] W. Gander, M.J. Gander, and F. Kwok. *Scientific Computing*. Vol. 11. Texts in Computational Science and Engineering. Heidelberg: Springer, 2014 (cit. on p. 12).
- [Gut09] M.H. Gutknecht. *Lineare Algebra*. Lecture Notes. SAM, ETH Zürich, 2009 (cit. on p. 12).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on p. 11).
- [Jos12] N.M. Josuttis. *The C++ Standard Library*. Boston, MA: Addison-Wesley, 2012 (cit. on p. 18).
- [LLM12] S. Lippman, J. Lajoie, and B. Moo. *C++ Primer*. 5th. Boston: Addison-Wesley, 2012 (cit. on pp. 17, 18).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on p. 12).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on p. 11).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on p. 12).

# Chapter 1

## Computing with Matrices and Vectors

**§1.0.0.1 (Prerequisite knowledge for Chapter 1)** The reader must master the basics of linear vector and matrix calculus as covered in every introductory course on linear algebra [NS02, Ch. 2].

On a few occasions we will also need results of 1D real calculus like Tylor's formula [Str09, Sect. 5.5]. ↴

**§1.0.0.2 (Levels of operations in simulation codes)** The lowest level of real arithmetic available on computers are the **elementary operations** “+”, “−”, “\*”, “\”, “^”, usually implemented in hardware. The next level comprises computations on finite arrays of real numbers, the **elementary linear algebra operations** (BLAS). On top of them we build complex algorithms involving iterations and approximations.

Complex iterative/recursive/approximative algorithms
Linear algebra operations on arrays (BLAS)
Elementary operations in $\mathbb{R}$

Hardly ever anyone will contemplate implementing elementary operations on binary data formats; similarly, well tested and optimised code libraries should be used for all elementary linear algebra operations in simulation codes. This chapter will introduce you to such libraries and how to use them smartly. ↴

### Contents

<b>1.1 Fundamentals</b> . . . . .	<b>37</b>
1.1.1 Notations . . . . .	37
1.1.2 Classes of Matrices . . . . .	39
<b>1.2 Software and Libraries</b> . . . . .	<b>41</b>
1.2.1 EIGEN . . . . .	41
1.2.2 PYTHON . . . . .	46
1.2.3 (Dense) Matrix Storage Formats . . . . .	47
<b>1.3 Basic Linear Algebra Operations</b> . . . . .	<b>52</b>
1.3.1 Elementary Matrix-Vector Calculus . . . . .	52
1.3.2 BLAS – Basic Linear Algebra Subprograms . . . . .	59
<b>1.4 Computational Effort</b> . . . . .	<b>64</b>
1.4.1 (Asymptotic) Complexity . . . . .	65
1.4.2 Cost of Basic Linear-Algebra Operations . . . . .	67
1.4.3 Reducing Complexity in Numerical Linear Algebra: Some Tricks . . . . .	68
<b>1.5 Machine Arithmetic and Consequences</b> . . . . .	<b>73</b>
1.5.1 Experiment: Loss of Orthogonality . . . . .	73
1.5.2 Machine Numbers . . . . .	76
1.5.3 Roundoff Errors . . . . .	79
1.5.4 Cancellation . . . . .	82

## 1.1 Fundamentals

### 1.1.1 Notations

The notations in this course try to adhere to established conventions. Since these may not be universal, idiosyncrasies cannot be avoided completely. Notations in textbooks may be different, beware!

Many considerations apply to real (field  $\mathbb{R}$ ) and complex (field  $\mathbb{C}$ ) numbers alike. Therefore we adopt the notation  $\mathbb{K}$  for a generic field of numbers. Thus, in this course,  $\mathbb{K}$  will designate either  $\mathbb{R}$  (real numbers) or  $\mathbb{C}$  (complex numbers); complex arithmetic [Str09, Sect. 2.5] plays a crucial role in many applications, for instance in signal processing.

### §1.1.1.1 (Notations for vectors)

- ◆ Vectors = are  $n$ -tuples ( $n \in \mathbb{N}$ ) with components  $\in \mathbb{K}$ .

**vector** = one-dimensional array (of real/complex numbers)

- ◆ Default in this lecture: vectors = column vectors

$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{K}^n$ <p>column vector</p>	$[x_1 \cdots x_n] \in \mathbb{K}^{1,n}$ <p>row vector</p>
--	---

$\mathbb{K}^n$   $\hat{=}$  vector space of *column vectors* with  $n$  components in  $\mathbb{K}$ .

“Linear algebra convention”: Unless stated otherwise, in mathematical formulas vector components are indexed from 1!

- notation for column vectors: **bold** small roman letters, e.g.  $\mathbf{x}, \mathbf{y}, \mathbf{z}$
  - Transposing:**

$$\begin{cases} \text{column vector} & \mapsto \text{row vector} \\ \text{row vector} & \mapsto \text{column vector} \end{cases}$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}^T = [x_1 \cdots x_n], \quad [x_1 \cdots x_n]^T = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$
  - Notation for row vectors:  $\mathbf{x}^\top, \mathbf{y}^\top, \mathbf{z}^\top$
  - Addressing vector components:**
    - two notations:  $\mathbf{x} = [x_1, \dots, x_n]^\top \rightarrow x_i, \quad i = 1, \dots, n$   
 $\mathbf{x} \in \mathbb{K}^n \rightarrow (\mathbf{x})_i, \quad i = 1, \dots, n$
  - Selecting sub-vectors:**
    - notation:  $\mathbf{x} = [x_1 \dots x_n]^\top \geq (\mathbf{x})_{k:l} = [x_k, \dots, x_l]^\top, 1 \leq k \leq l$

- ◆  **$j$ -th unit vector:**  $\mathbf{e}_j = [0, \dots, 1, \dots, 0]^\top, (\mathbf{e}_j)_i = \delta_{ij}, i, j = 1, \dots, n.$
- ☞ notation: Kronecker symbol  $\delta_{ij} := 1$ , if  $i = j$ ,  $\delta_{ij} := 0$ , if  $i \neq j$ .

### §1.1.1.2 (Notations and notions for matrices)

- ◆ **Matrices** = two-dimensional arrays of real/complex numbers

$$\mathbf{A} := \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \in \mathbb{K}^{n,m}, \quad n, m \in \mathbb{N}.$$

vector space of  $n \times m$ -matrices: ( $n \hat{=} \text{number of rows}$ ,  $m \hat{=} \text{number of columns}$ )

☞ notation for matrices: **bold CAPITAL roman letters**, e.g.,  $\mathbf{A}, \mathbf{S}, \mathbf{Y}$

Special cases:  $\mathbb{K}^{n,1} \leftrightarrow \text{column vectors}, \mathbb{K}^{1,n} \leftrightarrow \text{row vectors}$

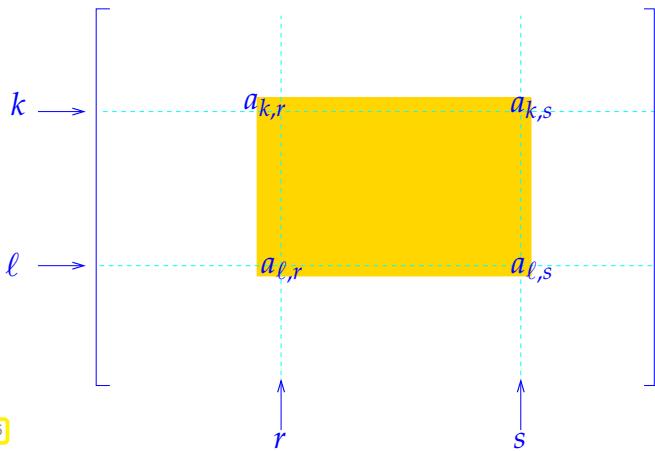
- ◆ Writing a matrix as a tuple of its columns or rows

$$\mathbf{c}_i \in \mathbb{K}^n, \quad i = 1, \dots, m \quad \Rightarrow \quad \mathbf{A} := [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_m] \in \mathbb{K}^{n,m},$$

$$\mathbf{r}_i \in \mathbb{K}^m, \quad i = 1, \dots, n \quad \Rightarrow \quad \mathbf{A} := \begin{bmatrix} \mathbf{r}_1^\top \\ \vdots \\ \mathbf{r}_n^\top \end{bmatrix} \in \mathbb{K}^{n,m}.$$

- ◆ Addressing matrix entries & sub-matrices (☞ notations):

$$\mathbf{A} := \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \quad \begin{array}{l} \rightarrow \text{matrix entry/matrix element } (\mathbf{A})_{i,j} := a_{ij}, \quad 1 \leq i \leq n, 1 \leq j \leq m, \\ \rightarrow i\text{-th row, } 1 \leq i \leq n: (\mathbf{A})_{i,:} := [a_{i,1}, \dots, a_{i,m}], \\ \rightarrow j\text{-th column, } 1 \leq j \leq m: (\mathbf{A})_{:,j} := [a_{1,j}, \dots, a_{n,j}]^\top, \\ \rightarrow \text{matrix block } (\mathbf{A})_{k:l,r:s} := [a_{ij}]_{\substack{i=k,\dots,l \\ j=r,\dots,s}}, \quad 1 \leq k \leq l \leq n, \\ \text{(sub-matrix)} \quad \quad \quad 1 \leq r \leq s \leq m. \end{array}$$



The colon (:) range notation is inspired by MATLAB's/PYTHON's matrix addressing conventions.  $(\mathbf{A})_{k:l,r:s}$  is a matrix of size  $(l - k + 1) \times (s - r + 1)$ .

- ◆ **Transposed matrix:**

$$\mathbf{A}^\top = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix}^\top := \begin{bmatrix} a_{11} & \dots & a_{n1} \\ \vdots & & \vdots \\ a_{1m} & \dots & a_{mn} \end{bmatrix} \in \mathbb{K}^{m,n}.$$

◆ Adjoint matrix (Hermitian transposed):

$$\mathbf{A}^H := \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix}^H := \begin{bmatrix} \bar{a}_{11} & \dots & \bar{a}_{n1} \\ \vdots & & \vdots \\ \bar{a}_{1m} & \dots & \bar{a}_{mn} \end{bmatrix} \in \mathbb{K}^{m,n}.$$

☞ notation:  $\bar{a}_{ij} = \Re(a_{ij}) - i\Im(a_{ij})$  complex conjugate of  $a_{ij}$ . Of course, for  $\mathbf{A} \in \mathbb{R}^{n,m}$  we have  $\mathbf{A}^H = \mathbf{A}^\top$ .

□

## 1.1.2 Classes of Matrices

Most matrices occurring in mathematical modelling have a special structure. This section presents a few of these. More will come up throughout the remainder of this chapter; see also [AG11, Sect. 4.3].

**§1.1.2.1 (Special matrices)** Terminology and notations for a few very special matrices:

Identity matrix:  $\mathbf{I} := \mathbf{I}_n := \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix} \in \mathbb{K}^{n,n}$ ,

Zero matrix:  $\mathbf{O} := \mathbf{O}_{n,m} := \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} \in \mathbb{K}^{n,m}$ ,

Diagonal matrix:  $\text{diag}(d_1, \dots, d_n) := \begin{bmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_n \end{bmatrix} \in \mathbb{K}^{n,n}$ ,  $d_j \in \mathbb{K}$ ,  $j = 1, \dots, n$ .

The creation of special matrices can usually be done by special commands or functions in the various languages or libraries dedicated to numerical linear algebra, see ??, § 1.2.1.3. □

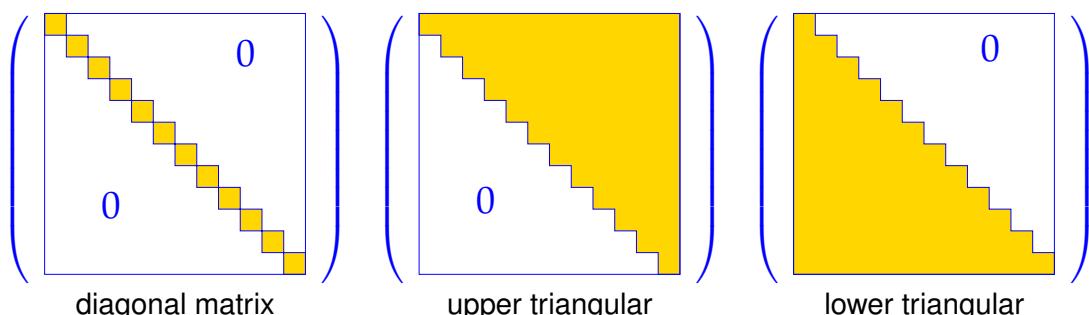
**§1.1.2.2 (Diagonal and triangular matrices)** A little terminology to quickly refer to matrices whose non-zero entries occupy special locations:

### Definition 1.1.2.3. Types of matrices

A matrix  $\mathbf{A} = [a_{ij}] \in \mathbb{K}^{m,n}$  is a

- diagonal matrix, if  $a_{ij} = 0$  for  $i \neq j$ ,
- upper triangular matrix, if  $a_{ij} = 0$  for  $i > j$ ,
- lower triangular matrix, if  $a_{ij} = 0$  for  $i < j$ .

A triangular matrix is normalized, if  $a_{ii} = 1$ ,  $i = 1, \dots, \min\{m, n\}$ .



### §1.1.2.4 (Symmetric matrices)

#### Definition 1.1.2.5. Hermitian/symmetric matrices

A matrix  $\mathbf{M} \in \mathbb{K}^{n,n}$ ,  $n \in \mathbb{N}$ , is **Hermitian**, if  $\mathbf{M}^H = \mathbf{M}$ . If  $\mathbb{K} = \mathbb{R}$ , the matrix is called **symmetric**.

**Definition 1.1.2.6. Symmetric positive definite (s.p.d.) matrices** → [DR08, Def. 3.31], [QSS00, Def. 1.22]

$\mathbf{M} \in \mathbb{K}^{n,n}$ ,  $n \in \mathbb{N}$ , is **symmetric (Hermitian) positive definite (s.p.d.)**, if

$$\mathbf{M} = \mathbf{M}^H \quad \text{and} \quad \forall \mathbf{x} \in \mathbb{K}^n: \mathbf{x}^H \mathbf{M} \mathbf{x} > 0 \Leftrightarrow \mathbf{x} \neq 0.$$

If  $\mathbf{x}^H \mathbf{M} \mathbf{x} \geq 0$  for all  $\mathbf{x} \in \mathbb{K}^n$  ▷  $\mathbf{M}$  **positive semi-definite**.

**Lemma 1.1.2.7. Necessary conditions for s.p.d.** → [DR08, Satz 3.33], [QSS00, Prop. 1.18]

For a symmetric/Hermitian positive definite matrix  $\mathbf{M} = \mathbf{M}^H \in \mathbb{K}^{n,n}$  holds true:

1.  $m_{ii} > 0$ ,  $i = 1, \dots, n$ ,
2.  $m_{ii} m_{jj} - |m_{ij}|^2 > 0 \quad \forall 1 \leq i < j \leq n$ ,
3. all eigenvalues of  $\mathbf{M}$  are positive. ( $\leftarrow$  also sufficient for symmetric/Hermitian  $\mathbf{M}$ )

**Remark 1.1.2.8 (S.p.d. Hessians)** Recall from analysis: in an isolated local minimum  $\mathbf{x}^*$  of a  $C^2$ -function  $f: \mathbb{R}^n \mapsto \mathbb{R}$  ► Hessian  $D^2 f(\mathbf{x}^*)$  s.p.d. (see Def. 8.4.1.11 for the definition of the Hessian)

To compute the minimum of a  $C^2$ -function iteratively by means of Newton's method (→ Sect. 8.4) a linear system of equations with the s.p.d. Hessian as system matrix has to be solved in each step.

The solutions of many equations in science and engineering boils down to finding the minimum of some (energy, entropy, etc.) function, which accounts for the prominent role of s.p.d. linear systems in applications.

**Review question(s) 1.1.2.9 (Notations, matrix-vector calculus, and special matrices)**

**(Q1.1.2.9.A)** Give a compact notation for the row vector containing the diagonal entries of a square matrix  $\mathbf{S} \in \mathbb{R}^{n,n}$ ,  $n \in \mathbb{N}$ .

**(Q1.1.2.9.B)** How can you write down the  $s \times s$ -submatrix,  $s \in \mathbb{N}$ , in the upper right corner of  $\mathbf{C} \in \mathbb{R}^{n,m}$ ,  $n, m \geq s$ .

**(Q1.1.2.9.C)** We consider two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n,m}$ , both with at most  $N \in \mathbb{N}$  non-zero entries. What is the maximal number of non-zero entries of  $\mathbf{A} + \mathbf{B}$ ?

**(Q1.1.2.9.D)** A matrix  $\mathbf{A} \in \mathbb{R}^{n,m}$  enjoys the following property (**banded matrix**):

$$i \in \{1, \dots, n\}, j \in \{1, \dots, m\}, i - j \notin \{-B_-, \dots, B_+\} \Rightarrow (\mathbf{A})_{ij} = 0,$$

for given  $B_-, B_+ \in \mathbb{N}_0$ . What is the maximal number of non-zero entries of  $\mathbf{A}$ .

**(Q1.1.2.9.E)** A matrix  $\mathbf{A}$  with real entries is known to be **skew-symmetric**:  $\mathbf{A}^\top = -\mathbf{A}$  What does this tell us about  $\mathbf{A}$  and its entries?

## 1.2 Software and Libraries

Whenever algorithms involve matrices and vectors (in the sense of linear algebra) it is advisable to rely on suitable code libraries or numerical programming environments.

### 1.2.1 EIGEN

Currently, the most widely used programming language for the development of new simulation software in scientific and industrial high-performance computing is C++. In this course we are going to use and discuss **EIGEN** as an example for a C++ library for numerical linear algebra (“embedded” domain specific language: DSL).

**EIGEN** is a **header-only** C++ template library designed to enable easy, natural and efficient numerical linear algebra: it provides data structures and a wide range of operations for matrices and vectors, see below. **EIGEN** also implements many more fundamental algorithms [documentation page](#) or the discussion below).

**EIGEN** relies on **expression templates** to allow the efficient evaluation of complex expressions involving matrices and vectors. Refer to the [example](#) given in the **EIGEN** documentation for details.

→ [Link](#) to an “**EIGEN** Cheat Sheet” (quick reference relating to MATLAB commands)

**§1.2.1.1 (Matrix and vector data types in EIGEN)** A generic matrix data type is given by the templated class

```
Eigen::Matrix<typename Scalar,
             int RowsAtCompileTime, int ColsAtCompileTime>
```

Here **Scalar** is the underlying scalar type of the matrix entries, which must support the usual operations '+', '-', '\*', '/', and '+=' , '\*'= , '-' etc. Usually the scalar type will be either double, float, or complex<>. The cardinal template arguments **RowsAtCompileTime** and **ColsAtCompileTime** can pass a **fixed** size of the matrix, if it is known at compile time. There is a specialization selected by the template argument **Eigen::Dynamic** supporting variable size “dynamic” matrices.

#### C++-code 1.2.1.2: Vector type and their use in EIGEN → [GITLAB](#)

```

1 #include <Eigen/Dense >
2
3 template<typename Scalar>
4 void eigenTypeDemo(unsigned int dim)
5 {
6     // General dynamic (variable size) matrices
7     using dynMat_t = Eigen::Matrix<Scalar,Eigen::Dynamic,Eigen::Dynamic>;
8     // Dynamic (variable size) column vectors
9     using dynColVec_t = Eigen::Matrix<Scalar,Eigen::Dynamic,1>;
10    // Dynamic (variable size) row vectors
11    using dynRowVec_t = Eigen::Matrix<Scalar,1,Eigen::Dynamic>;
12    using index_t = typename dynMat_t::Index;
13    using entry_t = typename dynMat_t::Scalar;
14
15    // Declare vectors of size 'dim'; not yet initialized
16    dynColVec_t colvec(dim);
17    dynRowVec_t rowvec(dim);
18    // Initialisation through component access

```

```

19   for(index_t i=0; i< colvec.size(); ++i) colvec[i] = (Scalar)i;
20   for(index_t i=0; i< rowvec.size(); ++i) rowvec[i] = (Scalar)1/(i+1);
21   colvec[0] = (Scalar)3.14; rowvec[dim-1] = (Scalar)2.718;
22   // Form tensor product, a matrix, see Section 1.3.1
23   dynMat_t vecprod = colvec*rowvec;
24   const int nrows = vecprod.rows();
25   const int ncols = vecprod.cols();
26 }
```

Note that in Line 23 we could have relied on automatic type deduction via **auto** `vectprod = ...`. However, as argued in Rem. 0.3.4.4 often it is safer to forgo this option and specify the type directly

The following convenience data types are provided by EIGEN, see  [EIGEN documentation](#):

- **MatrixXd**  $\hat{=}$  generic variable size matrix with **double** precision entries
- **VectorXd, RowVectorXd**  $\hat{=}$  dynamic column and row vectors  
(= dynamic matrices with one dimension equal to 1)
- **MatrixNd** with  $N = 2, 3, 4$  for small fixed size square  $N \times N$ -matrices (type **double**)
- **VectorNd** with  $N = 2, 3, 4$  for small column vectors with fixed length  $N$ .

The **d** in the type name may be replaced with **i** (for **int**), **f** (for **float**), and **cd** (for **complex<double>**) to select another basic scalar type.

All matrix type feature the methods **cols()**, **rows()**, and **size()** telling the number of columns, rows, and total number of entries.

Access to individual matrix entries and vector components, both as Rvalue and Lvalue, is possible through the `()`-operator taking two arguments of type **index\_t**. If only one argument is supplied, the matrix is accessed as a linear array according to its memory layout. For vectors, that is, matrices where one dimension is fixed to 1, the `[]`-operator can replace `()` with one argument, see Line 21 of Code 1.2.1.2.

→

**§1.2.1.3 (Initialization of dense matrices in EIGEN,  [EIGEN documentation](#))** The entry access operator `(int i, int j)` allows the most direct setting of matrix entries; there is hardly any runtime penalty.

Of course, in EIGEN dedicated functions take care of the initialization of the special matrices introduced in § 1.1.2.1:

```

Eigen::MatrixXd I = Eigen::MatrixXd::Identity(n,n);
Eigen::MatrixXd O = Eigen::MatrixXd::Zero(n,m);
Eigen::MatrixXd D = d_vector.asDiagonal();
```

#### C++-code 1.2.1.4: Initializing special matrices in EIGEN, → [GITLAB](#)

```

1 #include <Eigen/Dense>
2 // Just allocate space for matrix, no initialisation
3 Eigen::MatrixXd A(rows,cols);
4 // Zero matrix. Similar to matlab command zeros(rows,cols);
5 Eigen::MatrixXd B = MatrixXd::Zero(rows, cols);
6 // Ones matrix. Similar to matlab command ones(rows,cols);
7 Eigen::MatrixXd C = MatrixXd::Ones(rows, cols);
8 // Matrix with all entries same as value.
9 Eigen::MatrixXd D = MatrixXd::Constant(rows, cols, value);
10 // Random matrix, entries uniformly distributed in [0,1]
```

```

11 Eigen::MatrixXd E = MatrixXd::Random(rows, cols);
12 // (Generalized) identity matrix, 1 on main diagonal
13 Eigen::MatrixXd I = MatrixXd::Identity(rows,cols);
14 std::cout << "size of A = (" << A.rows() << ',' << A.cols() << ')' << std::endl;

```

A versatile way to initialize a matrix relies on a combination of the operators `<<` and `,`, which allows the construction of a matrix from blocks, see → [GITLAB](#), function `blockinit()`.

```

MatrixXd mat3(6,6);
mat3 <<
    MatrixXd::Constant(4,2,1.5), // top row, first block
    MatrixXd::Constant(4,3,3.5), // top row, second block
    MatrixXd::Constant(4,1,7.5), // top row, third block
    MatrixXd::Constant(2,4,2.5), // bottom row, left block
    MatrixXd::Constant(2,2,4.5); // bottom row, right block

```

The matrix is filled top to bottom left to right, block dimensions have to match (like in MATLAB).

**§1.2.1.5 (Access to submatrices in EIGEN,  [EIGEN documentation](#))** The method `block(int i, int j, int p, int q)` returns a reference to the submatrix with upper left corner at position  $(i,j)$  and size  $p \times q$ .

The methods `row(int i)` and `col(int j)` provide a reference to the corresponding row and column of the matrix. Even more specialised access methods are

```

topLeftCorner(p,q), bottomLeftCorner(p,q),
topRightCorner(p,q), bottomRightCorner(p,q),
topRows(q), bottomRows(q),
leftCols(p), and rightCols(q),

```

with obvious purposes.

#### C++ code 1.2.1.6: Demonstration code for access to matrix blocks in EIGEN → [GITLAB](#)

```

2 template<typename MatType>
3 void blockAccess(Eigen::MatrixBase<MatType> &M)
4 {
5     using index_t = typename Eigen::MatrixBase<MatType>::Index;
6     using entry_t = typename Eigen::MatrixBase<MatType>::Scalar;
7     const index_t nrows(M.rows()); // No. of rows
8     const index_t ncols(M.cols()); // No. of columns
9
10    cout << "Matrix M = " << endl << M << endl; // Print matrix
11    // Block size half the size of the matrix
12    index_t p = nrows/2,q = ncols/2;
13    // Output submatrix with left upper entry at position (i,i)
14    for(index_t i=0; i < min(p,q); i++)
15        cout << "Block (" << i << ',' << i << ',' << p << ',' << q
16        << ") = " << M.block(i,i,p,q) << endl;
17    // l-value access: modify sub-matrix by adding a constant
18    M.block(1,1,p,q) += Eigen::MatrixBase<MatType>::Constant(p,q,1.0);
19    cout << "M = " << endl << M << endl;
20    // r-value access: extract sub-matrix
21    MatrixXd B = M.block(1,1,p,q);
22    cout << "Isolated modified block = " << endl << B << endl;
23    // Special sub-matrices

```

```

24   cout << p << " top rows of m = " << M.topRows(p) << endl;
25   cout << p << " bottom rows of m = " << M.bottomRows(p) << endl;
26   cout << q << " left cols of m = " << M.leftCols(q) << endl;
27   cout << q << " right cols of m = " << M.rightCols(p) << endl;
28   // r-value access to upper triangular part
29   const MatrixXd T = M.template triangularView<Upper>(); //
30   cout << "Upper triangular part = " << endl << T << endl;
31   // l-value access to upper triangular part
32   M.template triangularView<Lower>() *= -1.5; //
33   cout << "Matrix M = " << endl << M << endl;
34 }
```

- Note that the function `blockAccess()` is templated and that the matrix argument passed through `M` has a type derived from `Eigen::MatrixBase`. The deeper reason for this alien looking signature of `blockAccess()` is explained in  [EIGEN documentation](#).
- EIGEN offers `views` for access to triangular parts of a matrix, see Line 29 and Line 32, according to

`M.triangularView<XX>()`

where `XX` can stand for one of the following: `Upper`, `Lower`, `StrictlyUpper`, `StrictlyLower`, `UnitUpper`, `UnitLower`, see  [EIGEN documentation](#).

- For column and row vectors references to sub-vectors can be obtained by the methods `head(int length)`, `tail(int length)`, and `segment(int pos, int length)`.

Note: Unless the preprocessor switch `NDEBUG` is set, EIGEN performs range checks on all indices. 

**§1.2.1.7 (Componentwise operations in EIGEN)** Running out of overloadable operators, EIGEN uses the `Array concept` to furnish entry-wise operations on matrices. An EIGEN-Array contains the same data as a matrix, supports the same methods for initialisation and access, but replaces the operators of matrix arithmetic with entry-wise actions. Matrices and arrays can be converted into each other by the `array()` and `matrix()` methods, see  [EIGEN documentation](#) for details. Information about functions that enable entry-wise operation is available in the  [EIGEN documentation](#).

### C++ code 1.2.1.8: Using Array in EIGEN → GITLAB

```

2 void matArray( int nrows , int ncols ){
3     Eigen::MatrixXd m1(nrows , ncols ),m2(nrows , ncols );
4     for( int i = 0; i < m1.rows(); i++)
5         for( int j=0; j < m1.cols(); j++) {
6             m1(i , j) = (double)(i+1)/( j+1);
7             m2(i , j) = (double)( j+1)/( i+1);
8         }
9     // Entry-wise product, not a matrix product
10    Eigen::MatrixXd m3 = (m1.array() * m2.array()).matrix();
11    // Explicit entry-wise operations on matrices are possible
12    Eigen::MatrixXd m4(m1.cwiseProduct(m2));
13    // Entry-wise logarithm
14    cout << "Log(m1) = " << endl << log(m1.array()) << endl;
15    // Entry-wise boolean expression, true cases counted
16    cout << (m1.array() > 3).count() << " entries of m1 > 3" << endl;
17 }
```

The application of a `functor` (→ Section 0.3.3) to all entries of a matrix can also been done via the `unaryExpr()` method of a matrix:

```
// Apply a lambda function to all entries of a matrix
auto fnct = [] (double x) { return (x+1.0/x); };
cout << "f(m1) = " << endl << m1.unaryExpr(fnct) << endl;
```

**§1.2.1.9 (Reduction operations in EIGEN)** According to EIGEN's terminology, **reductions** are operations that access all entries of a matrix and accumulate some information in the process

of EIGEN documentation. A typical example is the summation of the entries.

#### C++ code 1.2.1.10: Summation reduction in EIGEN → GITLAB

```
template <class Matrix> void sumEntries(Eigen::MatrixBase<Matrix> &M) {
    using Scalar = typename Eigen::MatrixBase<Matrix>::Scalar;
    // Compute sum of all entries
    const Scalar s = M.sum();
    // Row-wise and column-wise sum of entries: results are vectors
    Eigen::Matrix<Scalar, 1, Eigen::Dynamic> colsums{M.colwise().sum()};
    Eigen::Matrix<Scalar, Eigen::Dynamic, 1> rowsums{M.colwise().sum()};
    std::cout << M.rows() << 'x' << M.cols() << "—matrix: " << colsums.sum()
           << " = " << rowsums.sum() << " = " << s << std::endl;
}
```

**Remark 1.2.1.11 ('auto' in EIGEN codes)** The expression template programming model relies on complex intermediary data types hidden from the user. They support the efficient evaluation of complex expressions

of EIGEN documentation. Let us look at the following two code snippets that assume that both `M` and `R` are of type `Eigen::MatrixXd`.

Code I:

```
auto D = M.diagonal().asDiagonal(); R = D.inverse();
```

Code II:

```
Eigen::MatrixXd D = M.diagonal().asDiagonal(); R = D.inverse();
```

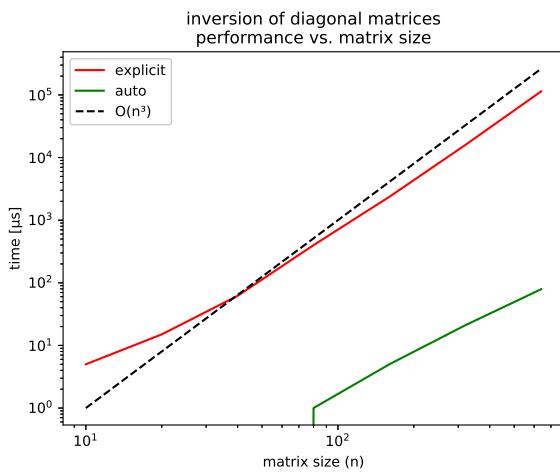


Fig. 6

We observe that for large matrices Code I ("auto", —, values  $\leq 0.5\mu\text{s}$  suppressed) runs much faster than Code II ("explicit", —) though they are "algebraically equivalent".

The reason is that in Code I `D` is of a complex type that preserves the information that the matrix is diagonal. Of course, inverting a diagonal matrix is cheap. Conversely forcing `D` to be of type `Eigen::MatrixXd` loses this information and the expensive `invert()` method for a generic densely populated matrix is invoked.

This is one of the exceptions to Rem. 0.3.4.4: for variables holding the result of EIGEN expressions `auto` is recommended.

**Remark 1.2.1.12 (EIGEN-based code: debug mode and release mode)** If you want a C++ code built using the EIGEN library run fast, for instance, for large computations or runtime measurements, you should compile in **release mode**, that is, with the compiler switches `-O2 -DNDEBUG` (for gcc or clang). In a `cmake`-based build system you can achieve this by setting the flag `CMAKE_BUILD_TYPE` to “Release”.

The default setting for EIGEN is debug mode, which makes EIGEN do a lot of consistency checking and considerably slows down execution of a code. □

### Remark 1.2.1.13 (EIGEN in use)

- ☞ EIGEN is used as one of the base libraries for the [Robot Operating System](#) (ROS), an open source project with strong ETH participation.
- ☞ The geometry processing library [libigl](#) uses EIGEN as its basic linear algebra engine. At ETH it is being used and developed at ETH Zurich, at the [Interactive Geometry Lab](#) and [Advanced Technologies Lab](#).

□

### Review question(s) 1.2.1.14 (EIGEN)

For the following questions you may consult the EIGEN documentation.

**(Q1.2.1.14.A)** Outline a C++ function

```
template <typename Matrix>
void replaceWithId(Eigen::DenseBase<Matrix> &M);
```

that checks whether the matrix is an  $n \times n$ -matrix with even  $n \in \mathbb{N}$  and then replaces its upper right  $n/2 \times n/2$ -block with an identity matrix. Do not use any C++ loops.

**(Q1.2.1.14.B)** Given an [Eigen::VectorXd](#) object  $\mathbf{v}$  ( $\leftrightarrow \mathbf{v} \in \mathbb{R}^n$ ), sketch a C++ code snippet that replaces it with a vector  $\tilde{\mathbf{v}}$  defined by

$$(\tilde{\mathbf{v}})_i := \begin{cases} (\mathbf{v})_n & \text{for } i = 1, \\ (\mathbf{v})_{i-1} & \text{for } i = 2, \dots, n. \end{cases}$$

Do not use C++ loops. Can you see a problem?

**(Q1.2.1.14.C)** Given a matrix  $\mathbf{M} \in \mathbb{R}^{m,n}$  stored in an [Eigen::MatrixXd](#) object  $\mathbf{M}$  write down a C++ code snippet that initializes another variable  $\mathbf{M}_{\text{ext}}$  of type [Eigen::MatrixXd](#) corresponding to

$$\widetilde{\mathbf{M}} := \begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0}^\top & 1 \end{bmatrix} \in \mathbb{R}^{m+1,n+1}$$

using EIGEN's `<<` matrix construction operator  [EIGEN documentation](#).

**(Q1.2.1.14.D)** Learn about the methods `head()` and `tail()` from the EIGEN documentation and express them by means of the `block()` method applied to the same variable.

△

## 1.2.2 PYTHON

**PYTHON** is a widely used general-purpose and open source programming language. Together with the packages like **NUMPY** and **MATPLOTLIB** it delivers similar functionality like MATLAB for free. For interactive computing **IPYTHON** can be used. All those packages belong to the **SCIPY** ecosystem.

PYTHON features a good documentation and several scientific distributions are available (e.g. [Anaconda](#), [Enthought](#)) which contain the most important packages. On most Linux-distributions the SciPY ecosystem is also available in the software repository, as well as many other packages including for example the Spyder IDE delivered with Anaconda.

A good introduction tutorial to numerical PYTHON are the [SciPy-lectures](#). The full documentation of NUMPY and SciPY can be found [here](#). For former MATLAB-users there's also a [guide](#). The scripts in this lecture notes follow the official [PYTHON style guide](#).

Note that in PYTHON we have to import the numerical packages explicitly before use. This is normally done at the beginning of the file with lines like `import numpy as np` and `from matplotlib import pyplot as plt`. Those import statements are often skipped in this lecture notes to focus on the actual computations. But you can always assume the import statements as given here, e.g. `np.ravel(A)` is a call to a NUMPY function and `plt.loglog(x, y)` is a call to a MATPLOTLIB pyplot function.

PYTHON is not used in the current version of the lecture. Nevertheless a few PYTHON codes are supplied in order to convey similarities and differences to implementations in MATLAB and C++.

**§1.2.2.1 (Matrices and Vectors in PYTHON)** The basic numeric data type in PYTHON are NUMPY's n-dimensional arrays. Vectors are normally implemented as 1D arrays and no distinction is made between row and column vectors. Matrices are represented as 2D arrays.

- ☞ `v = np.array([1, 2, 3])` creates a 1D array with the three elements 1, 2 and 3.
- ☞ `A = np.array([[1, 2], [3, 4]])` creates a 2D array.
- ☞ `A.shape` gives the n-dimensional size of an array.
- ☞ `A.size` gives the total number of entries in an array.

Note: There's also a matrix class in NUMPY with different semantics but its use is officially discouraged and it might even be removed in future release.

**§1.2.2.2 (Manipulating arrays in PYTHON)** There are many possibilities listed in the documentation how to [create](#), [index](#) and [manipulate](#) arrays.

An important difference to MATLAB is, that all arithmetic operations are normally performed element-wise, e.g. `A * B` is not the matrix-matrix product but element-wise multiplication (in MATLAB: `A.*A`). Also `A * v` does a [broadcasted](#) element-wise product. For the matrix product one has to use `np.dot(A, B)` or `A.dot(B)` explicitly.

## 1.2.3 (Dense) Matrix Storage Formats

All numerical libraries store the entries of a (generic = *dense*) matrix  $\mathbf{A} \in \mathbb{K}^{m,n}$  in a [linear](#) array of length *mn* (or longer). Accessing entries entails suitable index computations.

Two natural options for “vectorisation” of a matrix: *row major*, *column major*

		Row major (C-arrays, bitmaps, Python):																		
		Column major (Fortran, MATLAB, EIGEN):																		
$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$		<table border="1"> <tr> <td>A_arr</td> <td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td><td>9</td> </tr> </table>									A_arr	1	2	3	4	5	6	7	8	9
A_arr	1	2	3	4	5	6	7	8	9											
		<table border="1"> <tr> <td>A_arr</td> <td>1</td><td>4</td><td>7</td><td>2</td><td>5</td><td>8</td><td>3</td><td>6</td><td>9</td> </tr> </table>									A_arr	1	4	7	2	5	8	3	6	9
A_arr	1	4	7	2	5	8	3	6	9											

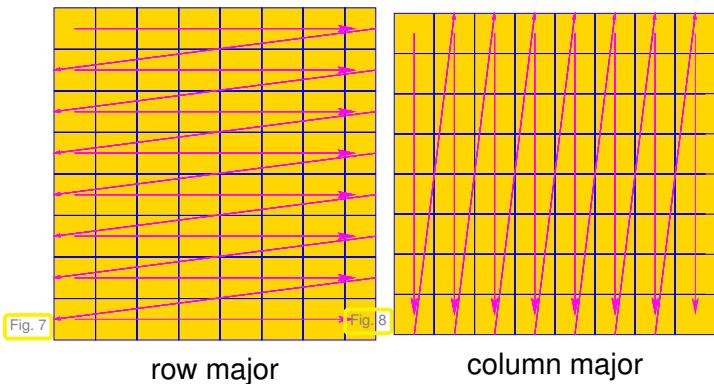
Access to entry  $(\mathbf{A})_{ij}$  of  $\mathbf{A} \in \mathbb{K}^{n,m}$ ,  
 $i = 1, \dots, n, j = 1, \dots, m$ :

row major:

$$(\mathbf{A})_{ij} \leftrightarrow \text{A\_arr}(m^*(i-1)+(j-1))$$

column major:

$$(\mathbf{A})_{ij} \leftrightarrow \text{A\_arr}(n^*(j-1)+(i-1))$$



**EXAMPLE 1.2.3.1 (Accessing matrix data as a vector)** In EIGEN the single index access operator relies on the linear data layout:

In EIGEN the data layout can be controlled by a template argument; default is column major.

#### C++ code 1.2.3.2: Single index access of matrix entries in EIGEN → GITLAB

```

2 void storageOrder(int nrows=6,int ncols=7)
3 {
4     cout << "Different matrix storage layouts in Eigen" << endl;
5     // Template parameter ColMajor selects column major data layout
6     Matrix<double, Dynamic, Dynamic, ColMajor> mcm(nrows,ncols);
7     // Template parameter RowMajor selects row major data layout
8     Matrix<double, Dynamic, Dynamic, RowMajor> mmr(nrows,ncols);
9     // Direct initialization; lazy option: use int as index type
10    for (int l=1,i= 0; i< nrows; i++)
11        for (int j= 0; j< ncols; j++,l++)
12            mcm(i,j) = mmr(i,j) = l;
13
14    cout << "Matrix mmr = " << endl << mmr << endl;
15    cout << "mcm linear = ";
16    for (int l=0;l < mcm.size(); l++) cout << mcm(l) << ',';
17    cout << endl;
18
19    cout << "mmr linear = ";
20    for (int l=0;l < mmr.size(); l++) cout << mmr(l) << ',';
21    cout << endl;
22 }
```

The function call `storageOrder(3, 3)`, cf. Code 1.2.3.2 yields the output

```

1 Different matrix storage layouts in Eigen
2 Matrix mmr =
3 1 2 3
4 4 5 6
5 7 8 9
6 mcm linear = 1,4,7,2,5,8,3,6,9,
7 mmr linear = 1,2,3,4,5,6,7,8,9,
```

In PYTHON the default data layout is row major, but it can be explicitly set. Further, array transposition does not change any data, but only the memory order and array shape.

**PYTHON-code 1.2.3.3: Storage order in PYTHON**

```

1 # array creation
2 A = np.array([[1, 2], [3, 4]]) # default (row major) storage
3 B = np.array([[1, 2], [3, 4]], order='F') # column major storage
4
5 # show internal storage
6 np.ravel(A, 'K') # array elements as stored in memory: [1, 2, 3, 4]
7 np.ravel(B, 'K') # array elements as stored in memory: [1, 3, 2, 4]
8
9 # nothing happens to the data on transpose, just the storage order
10 changes
11 np.ravel(A.T, 'K') # array elements as stored in memory: [1, 2, 3, 4]
12 np.ravel(B.T, 'K') # array elements as stored in memory: [1, 3, 2, 4]
13
14 # storage order can be accessed by checking the array's flags
15 A.flags['C_CONTIGUOUS'] # True
16 B.flags['F_CONTIGUOUS'] # True
17 A.T.flags['F_CONTIGUOUS'] # True
18 B.T.flags['C_CONTIGUOUS'] # True

```

↓

**Remark 1.2.3.4 (Vectorisation of a matrix)** Mapping a column-major matrix to a column vector with the same number of entries is called **vectorization** or linearization in numerical linear algebra, in symbols

$$\text{vec} : \mathbb{K}^{n,m} \rightarrow \mathbb{K}^{n \cdot m}, \quad \text{vec}(\mathbf{A}) := \begin{bmatrix} (\mathbf{A})_{:,1} \\ (\mathbf{A})_{:,2} \\ \vdots \\ (\mathbf{A})_{:,m} \end{bmatrix} \in \mathbb{R}^{n \cdot m}. \quad (1.2.3.5)$$

↓

**Remark 1.2.3.6 (Reshaping matrices in EIGEN)** If you need a reshaped view of a matrix' data in EIGEN you can obtain it via the raw data vector belonging to the matrix. Then use this information to create a matrix view by means of [Map](#) → [documentation](#).

**C++ code 1.2.3.7: Demonstration on how reshape a matrix in EIGEN ➔ GITLAB**

```

2 template<typename MatType>
3 void reshapeltest(MatType &M)
4 {
5     using index_t = typename MatType::Index;
6     using entry_t = typename MatType::Scalar;
7     const index_t nsize(M.size());
8
9     // reshaping possible only for matrices with non-prime dimensions
10    if ((nsize % 2) == 0) {
11        entry_t *Mdat = M.data(); // raw data array for M
12        // Reinterpretation of data of M
13        Map<Eigen::Matrix<entry_t, Dynamic, Dynamic>> R(Mdat, 2, nsize / 2);
14        // (Deep) copy data of M into matrix of different size
15        Eigen::Matrix<entry_t, Dynamic, Dynamic> S =
16            Map<Eigen::Matrix<entry_t, Dynamic, Dynamic>>(Mdat, 2, nsize / 2);
17
18        cout << "Matrix M = " << endl << M << endl;
19        cout << "reshaped to " << R.rows() << 'x' << R.cols()
20            << " = " << endl << R << endl;

```

↓

```

21 // Modifying R affects M, because they share the data space !
22 R *= -1.5;
23 cout << "Scaled (!) matrix M = " << endl << M << endl;
24 // Matrix S is not affected, because of deep copy
25 cout << "Matrix S = " << endl << S << endl;
26 }
27 }
```

This function has to be called with a mutable (l-value) matrix type object. A sample output is printed next:

```

1 Matrix M =
2 0 -1 -2 -3 -4 -5 -6
3 1 0 -1 -2 -3 -4 -5
4 2 1 0 -1 -2 -3 -4
5 3 2 1 0 -1 -2 -3
6 4 3 2 1 0 -1 -2
7 5 4 3 2 1 0 -1
8 reshaped to 2x21 =
9 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1 -6 -4 -2
10 1 3 5 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1
11 Scaled (!) matrix M =
12 -0 1.5 3 4.5 6 7.5 9
13 -1.5 -0 1.5 3 4.5 6 7.5
14 -3 -1.5 -0 1.5 3 4.5 6
15 -4.5 -3 -1.5 -0 1.5 3 4.5
16 -6 -4.5 -3 -1.5 -0 1.5 3
17 -7.5 -6 -4.5 -3 -1.5 -0 1.5
18 Matrix S =
19 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1 -6 -4 -2
20 1 3 5 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1
```

**Remark 1.2.3.8 (NumPy function `reshape`)** NumPy offers the function `np.reshape` for changing the dimensions of a matrix  $A \in \mathbb{K}^{m,n}$ :

```
# read elements of A in row major order (default)
B = np.reshape(A, (k, l)) # error, in case kl ≠ mn
B = np.reshape(A, (k, l), order='C') # same as above
# read elements of A in column major order
B = np.reshape(A, (k, l), order='F')
# read elements of A as stored in memory
B = np.reshape(A, (k, l), order='A')
```

This command will create an  $k \times l$ -array by reinterpreting the array of entries of  $A$  as data for an array with  $k$  rows and  $l$  columns. The order in which the elements of  $A$  are read can be set by the `order` argument to row major (default, '`C`'), column major ('`F`') or  $A$ 's internal storage order, i.e. row major if  $A$  is row major or column major if  $A$  is column major ('`A`').

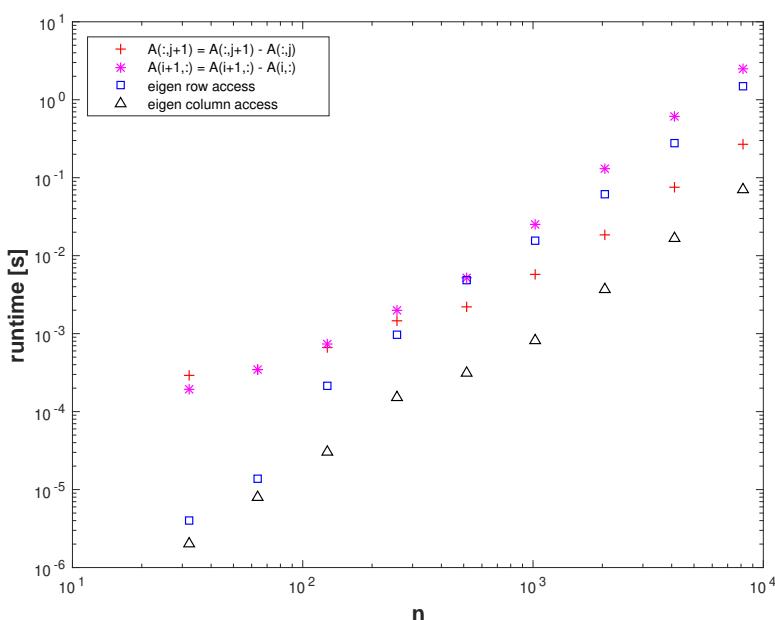
**EXPERIMENT 1.2.3.9 (Impact of matrix data access patterns on runtime)** Modern CPU feature several levels of memories (registers, L1 cache, L2 cache, ..., main memory) of different latency, bandwidth, and size. Frequently accessing memory locations with widely different addresses results in many cache misses and will considerably slow down the CPU.

The following C++ code sequentially runs through the entries of a *column major* matrix (EIGEN's default) in two ways and measures the (average) time required for the loops to complete. It relies on the `std::chrono` library [C++ reference](#).

#### C++-code 1.2.3.10: Timing for row and column oriented matrix access for EIGEN → [GITLAB](#)

```

2 void rowcolaccesstiming(void)
3 {
4     const int K = 3; // Number of repetitions
5     const int N_min = 5; // Smalles matrix size 32
6     const int N_max = 13; // Scan until matrix size of 8192
7     unsigned long n = (1L << N_min);
8     Eigen::MatrixXd times(N_max-N_min+1,3);
9
10    for(int l=N_min; l<= N_max; l++, n*=2) {
11        Eigen::MatrixXd A = Eigen::MatrixXd::Random(n,n);
12        double t1 = 1000.0;
13        for(int k=0;k<K;k++) {
14            auto tic = high_resolution_clock::now();
15            for(int j=0; j < n-1; j++) A.row(j+1) -= A.row(j); // row access
16            auto toc = high_resolution_clock::now();
17            double t = (double)duration_cast<microseconds>(toc-tic).count()/1E6;
18            t1 = std::min(t1,t);
19        }
20        double t2 = 1000.0;
21        for(int k=0;k<K;k++) {
22            auto tic = high_resolution_clock::now();
23            for(int j=0; j < n-1; j++) A.col(j+1) -= A.col(j); //column access
24            auto toc = high_resolution_clock::now();
25            double t = (double)duration_cast<microseconds>(toc-tic).count()/1E6;
26            t2 = std::min(t2,t);
27        }
28        times(l-N_min,0) = n;      times(l-N_min,1) = t1;      times(l-N_min,2) = t2;
29    }
30    std::cout << times << std::endl;
31 }
```



▷ Plot of average runtimes as measured with code Code 1.2.3.10.

Platform:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz — 4
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB, Mem 8 GB
- ◆ gcc 4.8.4, -O3, -DNDEBUG

The compiler flags `-O3` and `-DNDEBUG` are essential. The C++ code would be significantly slower if the default compiler options were used!

We observe a blatant discrepancy of CPU time required for accessing entries of a matrix in rowwise or columnwise fashion. This reflects the impact of features of the underlying hardware architecture, like cache

size and memory bandwidth:

Interpretation of timings: Since standard matrices in EIGEN are stored column major all the matrix elements in a column occupy contiguous memory locations, which will all reside in the cache together. Hence, column oriented access will mainly operate on data in the cache even for large matrices. Conversely, row oriented access addresses matrix entries that are stored in distant memory locations, which incurs frequent cache misses (**cache thrashing**).

The impact of hardware architecture on the performance of algorithms will **not** be taken into account in this course, because hardware features tend to be both intricate and ephemeral. However, for modern high performance computing it is essential to adapt implementations to the hardware on which the code is supposed to run.

#### Review question(s) 1.2.3.11 (Dense matrix storage formats)

**(Q1.2.3.11.A)** Write *efficient* elementary C++ loops that realize the matrix×vector product  $\mathbf{M}\mathbf{x}$ ,  $\mathbf{m} \in \mathbb{R}^{m,n}$ ,  $\mathbf{x} \in \mathbb{R}^n$ , where  $\mathbf{M}$  is stored in an **Eigen::MatrixXd** object  $\mathbf{M}$  and  $\mathbf{x}$  given as a **Eigen::VectorXd** object  $\mathbf{x}$ . Assume the default (Column major) memory layout for  $\mathbf{M}$ . Discuss the memory access pattern.

**(Q1.2.3.11.B)** A black-box function has the following signature:

```
template <typename Vector>
double processVector(const Eigen::DenseBase<Vector> &v);
```

It is known that it accesses each vector entry only once.

Recall that the vectorisation of a matrix  $\mathbf{A} \in \mathbb{K}^{n,m}$  is defined as

$$\text{vec} : \mathbb{K}^{n,m} \rightarrow \mathbb{K}^{n \cdot m}, \quad \text{vec}(\mathbf{A}) := \begin{bmatrix} (\mathbf{A})_{:,1} \\ (\mathbf{A})_{:,2} \\ \vdots \\ (\mathbf{A})_{:,m} \end{bmatrix} \in \mathbb{R}^{n \cdot m}. \quad (1.2.3.5)$$

Given a matrix  $\mathbf{A} \in \mathbb{R}^{n,m}$  stored in an **Eigen::MatrixXd** object  $\mathbf{A}$  (*column major* memory layout), how can you efficiently realize the following function calls in C++:

- `processVector(vec(A))`,
- `processVector(vec(AT))` ?

△

## 1.3 Basic Linear Algebra Operations

First we refresh the basic rules of vector and matrix calculus. Then we will learn about a very old programming interface for simple dense linear algebra operations.

### 1.3.1 Elementary Matrix-Vector Calculus

What you should know from linear algebra [NS02, Sect. 2.2]:

- ◆ vector space operations in matrix space  $\mathbb{K}^{m,n}$  (addition, multiplication with scalars)



dot product:  $\mathbf{x}, \mathbf{y} \in \mathbb{K}^n, n \in \mathbb{N}$ :  $\mathbf{x} \cdot \mathbf{y} := \mathbf{x}^H \mathbf{y} = \sum_{i=1}^n \bar{x}_i y_i \in \mathbb{K}$   
 (in EIGEN:  $\mathbf{x}.dot(\mathbf{y})$  or  $\mathbf{x}.adjoint() * \mathbf{y}$ ,  $\mathbf{x}, \mathbf{y} \hat{=} \text{column vectors}$ )



tensor product:  $\mathbf{x} \in \mathbb{K}^m, \mathbf{y} \in \mathbb{K}^n, n \in \mathbb{N}$ :  $\mathbf{x}\mathbf{y}^H = (x_i \bar{y}_j)_{\substack{i=1, \dots, m \\ j=1, \dots, n}} \in \mathbb{K}^{m,n}$   
 (in EIGEN:  $\mathbf{x} * \mathbf{y}.adjoint()$ ,  $\mathbf{x}, \mathbf{y} \hat{=} \text{column vectors}$ )

◆ All are special cases of the matrix product:

$$\mathbf{A} \in \mathbb{K}^{m,n}, \quad \mathbf{B} \in \mathbb{K}^{n,k}: \quad \mathbf{AB} = \left[ \sum_{j=1}^n a_{ij} b_{jl} \right]_{\substack{i=1, \dots, m \\ l=1, \dots, k}} \in \mathbb{K}^{m,k}. \quad (1.3.1.1)$$

Recall from linear algebra basic properties of the matrix product: for all  $\mathbb{K}$ -matrices  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  (of suitable sizes),  $\alpha, \beta \in \mathbb{K}$

associative:  $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$ ,

bi-linear:  $(\alpha\mathbf{A} + \beta\mathbf{B})\mathbf{C} = \alpha(\mathbf{AC}) + \beta(\mathbf{BC})$ ,  $\mathbf{C}(\alpha\mathbf{A} + \beta\mathbf{B}) = \alpha(\mathbf{CA}) + \beta(\mathbf{CB})$ ,

non-commutative:  $\mathbf{AB} \neq \mathbf{BA}$  in general .

### §1.3.1.2 (Visualisation of (special) matrix products)

Dependency of an entry of a product matrix:

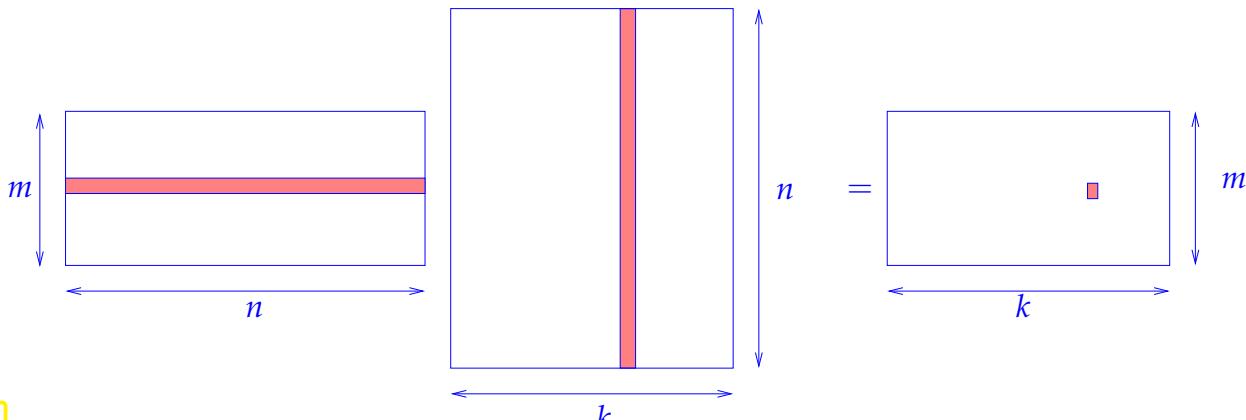
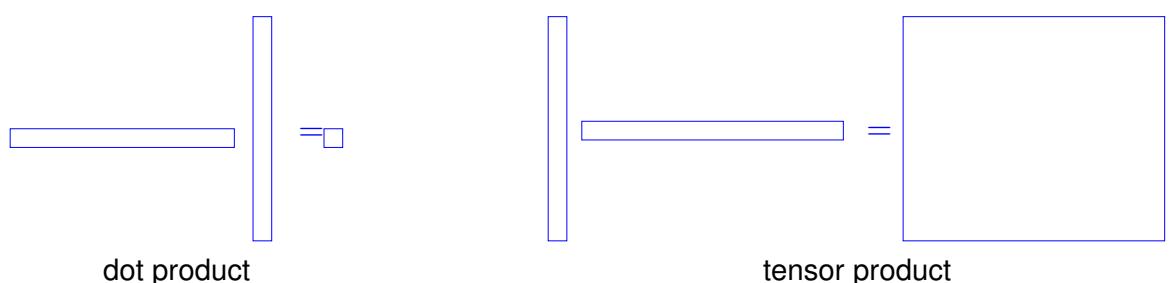


Fig. 10



**Remark 1.3.1.3 (Row-wise & column-wise view of matrix product)** To understand what is going on when forming a matrix product, it is often useful to decompose it into matrix  $\times$  vector operations in one of the following two ways:

$\mathbf{A} \in \mathbb{K}^{m,n}$ ,  $\mathbf{B} \in \mathbb{K}^{n,k}$ :

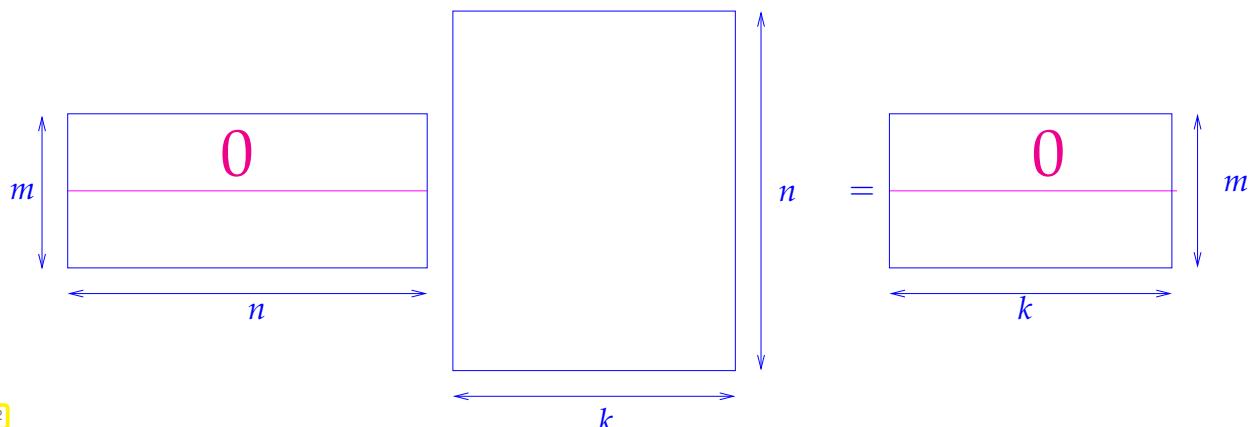
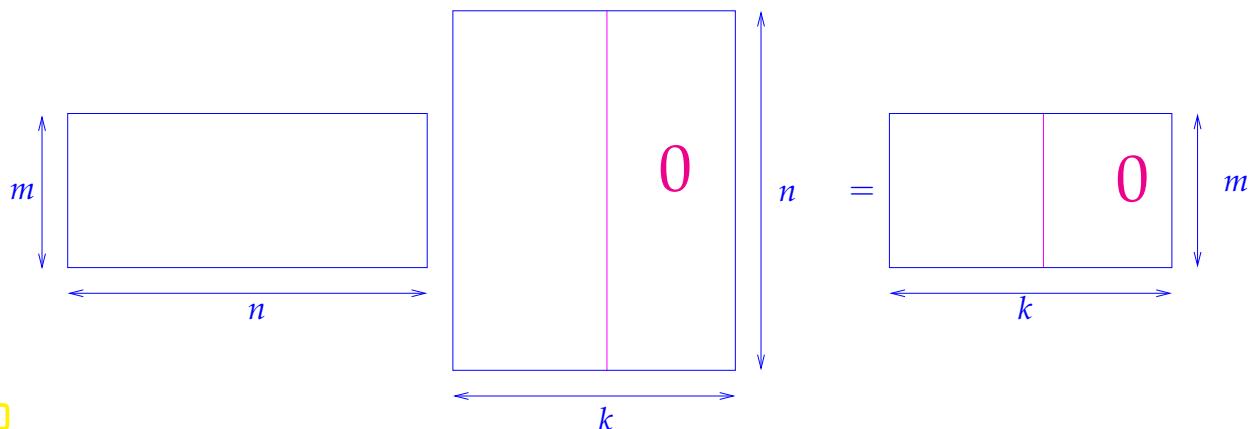
$$\mathbf{AB} = \begin{bmatrix} \mathbf{A}(\mathbf{B})_{:,1} & \dots & \mathbf{A}(\mathbf{B})_{:,k} \end{bmatrix}, \quad \mathbf{AB} = \begin{bmatrix} (\mathbf{A})_{1,:}\mathbf{B} \\ \vdots \\ (\mathbf{A})_{m,:}\mathbf{B} \end{bmatrix}. \quad (1.3.1.4)$$

↓    ↓  
matrix assembled from columns      matrix assembled from rows

For notations refer to Sect. 1.1.1. □

**Remark 1.3.1.5 (Understanding the structure of product matrices)** A “mental image” of matrix multiplication is useful for telling special properties of product matrices.

For instance, zero blocks of the product matrix can be predicted easily in the following situations using the idea explained in Rem. 1.3.1.3 (try to understand how):



A clear understanding of matrix multiplication enables you to “see”, which parts of a matrix factor matter in a product:

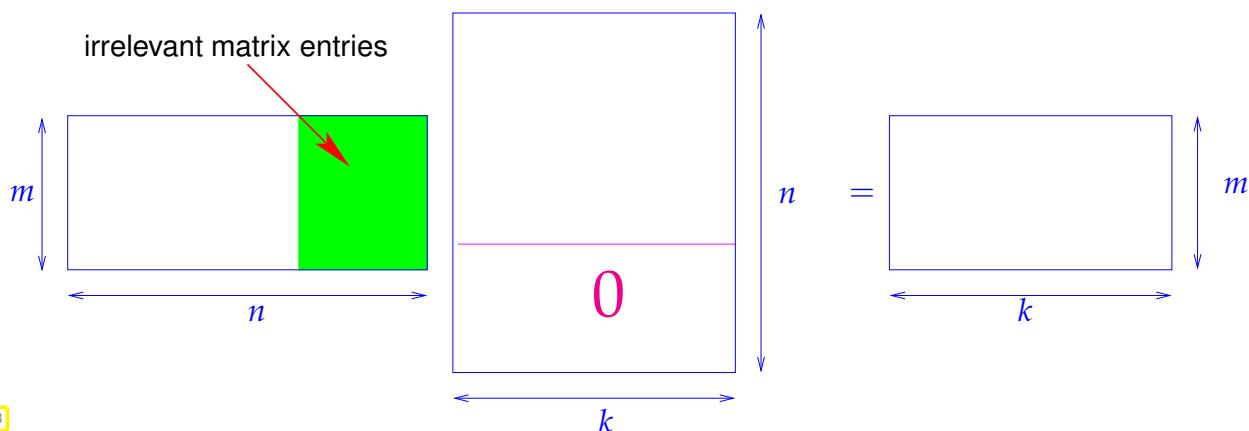
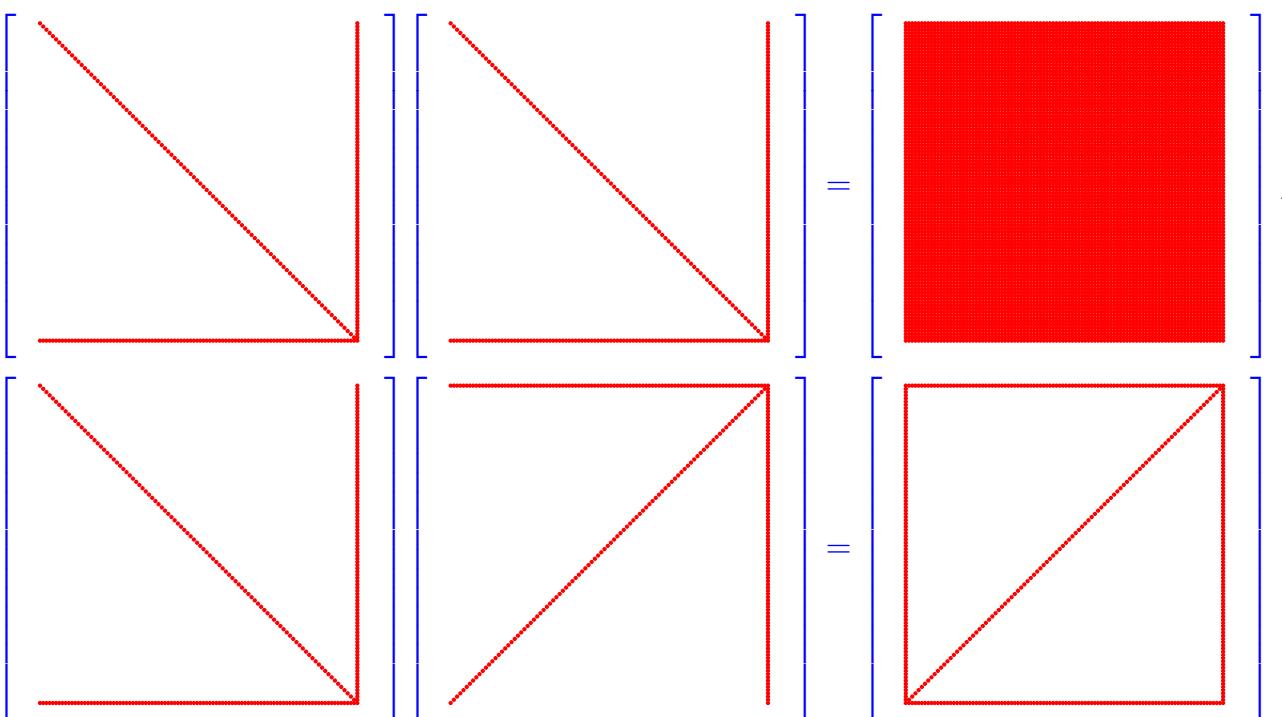


Fig. 13

“Seeing” the structure/pattern of a matrix product:



These nice renderings of the so-called patterns of matrices, that is, the distribution of their non-zero entries have been created by a special EIGEN/**Figure**-command for visualizing the structure of a matrix: `fig.spy(M)`

#### C++ code 1.3.1.6: Visualizing the structure of matrices in EIGEN → GITLAB

```

2 #include "matplotlibcpp.h" // Tools for plotting, see
3   https://github.com/lava/matplotlib-cpp
4 #include <Eigen/Dense>
5 #include <string>
6 namespace plt = matplotlibcpp;
7 using namespace Eigen;
8
9 // Produce spy-plot of a dense Eigen matrix.
10 void spy(const Eigen::MatrixXd &M, const std::string &fname) {
11     plt::figure();
12     plt::spy(M, {"marker", "o"}, {"markersize", "2"}, {"color", "b"});
13     plt::title("nnz = " + std::to_string(M.nonZeros()));
14     plt::savefig(fname);
15 }
```

```

16 int main() {
17     int n = 100;
18     MatrixXd A(n, n), B(n, n);
19     A.setZero();
20     B.setZero();
21     // Initialize matrices, see Fig. 13
22     A.diagonal() = VectorXd::LinSpaced(n, 1, n);
23     A.col(n - 1) = VectorXd::LinSpaced(n, 1, n);
24     A.row(n - 1) = RowVectorXd::LinSpaced(n, 1, n);
25     B = A.colwise().reverse();
26     // Matrix products
27     MatrixXd C = A * A, D = A * B;
28     spy(A, "Aspy_cpp.eps"); // Sparse arrow matrix
29     spy(B, "Bspy_cpp.eps"); // Sparse arrow matrix
30     spy(C, "Cspy_cpp.eps"); // Fully populated matrix
31     spy(D, "Dspy_cpp.eps"); // Sparse "framed" matrix
32     return 0;
33 }
```

This code also demonstrates the use of `diagonal()`, `col()`, `row()` for L-value access to parts of a matrix.

PYTHON/MATPLOTLIB-command for visualizing the structure of a matrix: `plt.spy(M)`

#### PYTHON-code 1.3.1.7: Visualizing the structure of matrices in PYTHON

```

1 n = 100
2 A = np.diag(np.mgrid [:n])
3 A[:, -1] = A[-1, :] = np.mgrid [:n]
4 plt.spy(A)
5 plt.spy(A[::-1, :])
6 plt.spy(np.dot(A, A))
7 plt.spy(np.dot(A, B))
```

**Remark 1.3.1.8 (Multiplying triangular matrices)** The following result is useful when dealing with matrix decompositions that often involve triangular matrices.

#### Lemma 1.3.1.9. Group of regular diagonal/triangular matrices

$$\mathbf{A}, \mathbf{B} \left\{ \begin{array}{l} \text{diagonal} \\ \text{upper triangular} \\ \text{lower triangular} \end{array} \right. \Rightarrow \mathbf{AB} \text{ and } \mathbf{A}^{-1} \left\{ \begin{array}{l} \text{diagonal} \\ \text{upper triangular} \\ \text{lower triangular} \end{array} \right. .$$

(assumes that  $\mathbf{A}$  is regular)

“Proof by visualization” → Rem. 1.3.1.5

**EXPERIMENT 1.3.1.10 (Scaling a matrix)** Scaling = multiplication with diagonal matrices (with non-zero diagonal entries):

It is important to know the different effect of multiplying with a diagonal matrix from left or right:

- ◆ multiplication with diagonal matrix *from left* ► row scaling

$$\begin{bmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & \dots & 0 \\ & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & d_n \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{bmatrix} = \begin{bmatrix} d_1 a_{11} & d_1 a_{12} & \dots & d_1 a_{1m} \\ d_2 a_{21} & d_2 a_{22} & \dots & d_2 a_{2m} \\ \vdots & & \ddots & \vdots \\ d_n a_{n1} & d_n a_{n2} & \dots & d_n a_{nm} \end{bmatrix} = \begin{bmatrix} d_1 (\mathbf{A})_{1,:} \\ \vdots \\ d_n (\mathbf{A})_{n,:} \end{bmatrix}.$$

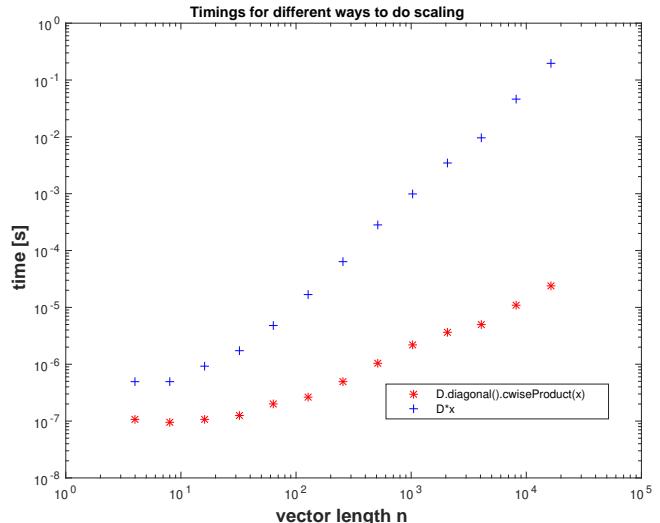
- ◆ multiplication with diagonal matrix *from right* ► column scaling

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{bmatrix} \begin{bmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & \dots & 0 \\ & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & d_m \end{bmatrix} = \begin{bmatrix} d_1 a_{11} & d_2 a_{12} & \dots & d_m a_{1m} \\ d_1 a_{21} & d_2 a_{22} & \dots & d_m a_{2m} \\ \vdots & & \ddots & \vdots \\ d_1 a_{n1} & d_2 a_{n2} & \dots & d_m a_{nm} \end{bmatrix} = \begin{bmatrix} d_1 (\mathbf{A})_{:,1} & \dots & d_m (\mathbf{A})_{:,m} \end{bmatrix}.$$

Multiplication with a scaling matrix  $\mathbf{D} := \text{diag}(d_1, \dots, d_n) \in \mathbb{R}^{n,n}$  in EIGEN can be realised in three ways, see Code 1.3.1.11, Line 9-Line 1 Line 13, and Line 15.

Measured runtimes

The code will be slowed down massively in case a temporary dense matrix is created inadvertently. Notice that EIGEN's expression templates avoid the pointless effort, see Line 15.



### C++ code 1.3.1.11: Timing multiplication with scaling matrix in EIGEN → GITLAB

```

2 int nruns = 3, minExp = 2, maxExp = 14;
3 MatrixXd tms(maxExp-minExp+1,4);
4 for(int i = 0; i <= maxExp-minExp; ++i){
5     Timer tbad, tgood, topt; // timer class
6     int n = std::pow(2, minExp + i);
7     VectorXd d = VectorXd::Random(n,1), x = VectorXd::Random(n,1), y(n);
8     for(int j = 0; j < nruns; ++j) {
9         MatrixXd D = d.asDiagonal(); //
10        // matrix vector multiplication
11        tbad.start(); y = D*x; tbad.stop(); //
12        // componentwise multiplication
13        tgood.start(); y= d.cwiseProduct(x); tgood.stop(); //

```

```

14     // matrix multiplication optimized by Eigen
15     topt.start(); y = d.asDiagonal()*x; topt.stop(); //
16 }
17 tms(i,0)=n;
18 tms(i,1)=tgood.min(); tms(i,2)=tbad.min(); tms(i,3)= topt.min();
19 }

```

Hardly surprising, the component-wise multiplication of the two vectors is way faster than the intermittent initialisation of a diagonal matrix (main populated by zeros) and the computation of a matrix  $\times$  vector product. Nevertheless, such blunders keep on haunting numerical codes. Do not rely solely on EIGEN optimizations!

**Remark 1.3.1.12 (Row and column transformations)** Simple operations on rows/columns of matrices, cf. what was done in Exp. 1.2.3.9, can often be expressed as multiplication with special matrices: For instance, given  $\mathbf{A} \in \mathbb{K}^{n,m}$  we obtain  $\mathbf{B}$  by adding row  $(\mathbf{A})_{j,:}$  to row  $(\mathbf{A})_{j+1,:}$ ,  $1 \leq j < n$ .

Realisation through matrix product      ➡      
$$\mathbf{B} = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & 1 & 1 \\ & & & \ddots \\ & & & & 1 \end{bmatrix} \mathbf{A}.$$

The matrix multiplying  $\mathbf{A}$  from the left is a specimen of a **transformation matrix**, a matrix that coincides with the identity matrix  $\mathbf{I}$  except for a single off-diagonal entry.

left-multiplication right-multiplication	with transformation matrices	→	row transformations column transformations
---	------------------------------	---	---

Row/column transformations will play a central role in Section 2.3.

**§1.3.1.13 (Block matrix product)** Given matrix dimensions  $M, N, K \in \mathbb{N}$  block sizes  $1 \leq n < N$  ( $n' := N - n$ ),  $1 \leq m < M$  ( $m' := M - m$ ),  $1 \leq k < K$  ( $k' := K - k$ ) we start from the following matrices:

$$\begin{array}{lll} \mathbf{A}_{11} \in \mathbb{K}^{m,n} & \mathbf{A}_{12} \in \mathbb{K}^{m,n'} & \mathbf{B}_{11} \in \mathbb{K}^{n,k} & \mathbf{B}_{12} \in \mathbb{K}^{n,k'} \\ \mathbf{A}_{21} \in \mathbb{K}^{m',n} & \mathbf{A}_{22} \in \mathbb{K}^{m',n'} & \mathbf{B}_{21} \in \mathbb{K}^{n',k} & \mathbf{B}_{22} \in \mathbb{K}^{n',k'} \end{array},$$

This matrices serve as sub-matrices or matrix blocks and are assembled into larger matrices

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \in \mathbb{K}^{M,N}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \in \mathbb{K}^{N,K}.$$

It turns out that the matrix product  $\mathbf{AB}$  can be computed by the same formula as the product of simple  $2 \times 2$ -matrices:

➡ 
$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11}\mathbf{B}_{11} + \mathbf{A}_{12}\mathbf{B}_{21} & \mathbf{A}_{11}\mathbf{B}_{12} + \mathbf{A}_{12}\mathbf{B}_{22} \\ \mathbf{A}_{21}\mathbf{B}_{11} + \mathbf{A}_{22}\mathbf{B}_{21} & \mathbf{A}_{21}\mathbf{B}_{12} + \mathbf{A}_{22}\mathbf{B}_{22} \end{bmatrix}. \quad (1.3.1.14)$$

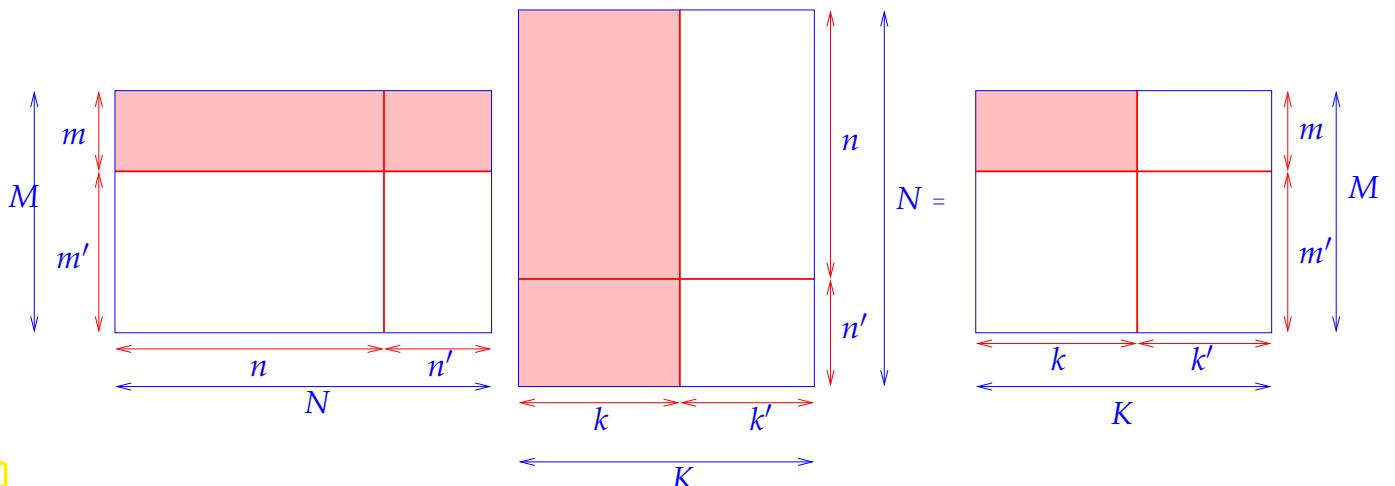


Fig. 15

Bottom line: one can compute with block-structured matrices in *almost* (\*) the same ways as with matrices with real/complex entries, see [QSS00, Sect. 1.3.3].



(\*): you must not use the commutativity of multiplication (because matrix multiplication is not commutative).

□

### 1.3.2 BLAS – Basic Linear Algebra Subprograms

**BLAS** (Basic Linear Algebra Subprograms) is a specification (API) that prescribes a set of low-level routines for performing common linear algebra operations such as vector addition, scalar multiplication, dot products, linear combinations, and matrix multiplication. They are the de facto low-level routines for linear algebra libraries ([Wikipedia](#)).

The BLAS API is standardised by the [BLAS technical forum](#) and, due to its history dating back to the 70s, follows conventions of FORTRAN 77, see the [Quick Reference Guide](#) for examples. However, wrappers for other programming languages are available. CPU manufacturers and/or developers of operating systems usually supply highly optimised implementations:

- **OpenBLAS**: open source implementation with some general optimisations, available under BSD license.
- **ATLAS** (Automatically Tuned Linear Algebra Software): open source BLAS implementation with auto-tuning capabilities. Comes with C and FORTRAN interfaces and is included in Linux distributions.
- **Intel MKL** (Math Kernel Library): commercial highly optimised BLAS implementation available for all Intel CPUs. Used by most proprietary simulation software and also MATLAB.

#### EXPERIMENT 1.3.2.1 (Multiplying matrices in EIGEN)

The following EIGEN-based C++ code performs a multiplication of densely populated matrices in three different ways:

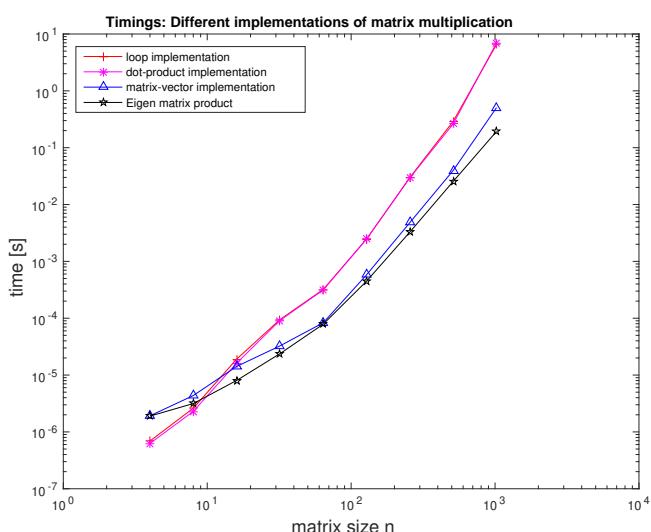
1. Direct implementation of three nested loops
2. Realization by matrix  $\times$  vector products
3. Use of built-in matrix multiplication of EIGEN

### C++ code 1.3.2.2: Timing different implementations of matrix multiplication in EIGEN → GITLAB

```

2 void mmtiming() {
3     int nruns = 3, minExp = 2, maxExp = 10;
4     MatrixXd timings(maxExp - minExp + 1, 5);
5     for (int p = 0; p <= maxExp - minExp; ++p) {
6         Timer t1, t2, t3, t4; // timer class
7         int n = std::pow(2, minExp + p);
8         MatrixXd A = MatrixXd::Random(n, n);
9         MatrixXd B = MatrixXd::Random(n, n);
10        MatrixXd C = MatrixXd::Zero(n, n);
11        for (int q = 0; q < nruns; ++q) {
12            // Loop based implementation no template magic
13            t1.start();
14            for (int i = 0; i < n; ++i)
15                for (int j = 0; j < n; ++j)
16                    for (int k = 0; k < n; ++k)
17                        C(i, j) += A(i, k) * B(k, j);
18            t1.stop();
19            // dot product based implementation little template magic
20            t2.start();
21            for (int i = 0; i < n; ++i)
22                for (int j = 0; j < n; ++j)
23                    C(i, j) = A.row(i).dot(B.col(j));
24            t2.stop();
25            // matrix-vector based implementation middle template magic
26            t3.start();
27            for (int j = 0; j < n; ++j)
28                C.col(j) = A * B.col(j);
29            t3.stop();
30            // Eigen matrix multiplication template magic optimized
31            t4.start();
32            C = A * B;
33            t4.stop();
34        }
35        timings(p, 0) = n;
36        timings(p, 1) = t1.min();
37        timings(p, 2) = t2.min();
38        timings(p, 3) = t3.min();
39        timings(p, 4) = t4.min();
40    }
41    std::cout << std::scientific << std::setprecision(3) << timings << std::endl;

```



Platform:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB, Mem 8 GB
- ◆ gcc 4.8.4, -O3

In EIGEN we can achieve some gain in execution speed by relying on compact matrix/vector operations that invoke efficient EIGEN built-in functions. However, compiler optimizations make plain and simple loops almost competitive.

BLAS routines are grouped into “levels” according to the amount of data and computation involved (asymptotic complexity, see Section 1.4.1 and [GV89, Sect. 1.1.12]):

- **Level 1:** vector operations such as scalar products and vector norms.  
asymptotic complexity  $O(n)$ , (with  $n \triangleq$  vector length),  
e.g.: dot product:  $\rho = \mathbf{x}^\top \mathbf{y}$
- **Level 2:** vector-matrix operations such as matrix-vector multiplications.  
asymptotic complexity  $O(mn)$ , (with  $(m, n) \triangleq$  matrix size),  
e.g.: matrix  $\times$  vector multiplication:  $\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$
- **Level 3:** matrix-matrix operations such as matrix additions or multiplications.  
asymptotic complexity often  $O(nmk)$ , (with  $(n, m, k) \triangleq$  matrix sizes),  
e.g.: matrix product:  $\mathbf{C} = \mathbf{A} \mathbf{B}$

### Syntax of BLAS calls:

The functions have been implemented for different types, and are distinguished by the first letter of the function name. E.g. *sdot* is the dot product implementation for single precision and *ddot* for double precision.

- ◆ **BLAS LEVEL 1:** vector operations, asymptotic complexity  $O(n)$ ,  $n \triangleq$  vector length

- dot product  $\rho = \mathbf{x}^\top \mathbf{y}$

**$\text{xDOT}(N, X, INCX, Y, INCY)$**

- $x \in \{S, D\}$ , scalar type: S  $\triangleq$  type float, D  $\triangleq$  type double
- $N \triangleq$  length of vector (modulo stride INCX)
- $X \triangleq$  vector  $\mathbf{x}$ : array of type  $x$
- INCX  $\triangleq$  stride for traversing vector  $X$
- $Y \triangleq$  vector  $\mathbf{y}$ : array of type  $x$
- INCY  $\triangleq$  stride for traversing vector  $Y$

- vector operations  $\mathbf{y} = \alpha \mathbf{x} + \mathbf{y}$

**$\text{xAXPY}(N, ALPHA, X, INCX, Y, INCY)$**

- $x \in \{S, D, C, Z\}$ , S  $\triangleq$  type float, D  $\triangleq$  type double, C  $\triangleq$  type complex
- $N \triangleq$  length of vector (modulo stride INCX)
- ALPHA  $\triangleq$  scalar  $\alpha$
- $X \triangleq$  vector  $\mathbf{x}$ : array of type  $x$
- INCX  $\triangleq$  stride for traversing vector  $X$
- $Y \triangleq$  vector  $\mathbf{y}$ : array of type  $x$
- INCY  $\triangleq$  stride for traversing vector  $Y$

- ◆ **BLAS LEVEL 2:** matrix-vector operations, asymptotic complexity  $O(mn)$ ,  $(m, n) \triangleq$  matrix size

- matrix  $\times$  vector multiplication  $\mathbf{y} = \alpha \mathbf{Ax} + \beta \mathbf{y}$

**XGEMV (TRANS, M, N, ALPHA, A, LDA, X,**

**INCX, BETA, Y, INCY)**

- $\mathbf{x} \in \{\mathbf{S}, \mathbf{D}, \mathbf{C}, \mathbf{Z}\}$ , scalar type:  $\mathbf{S} \hat{=} \text{type float}$ ,  $\mathbf{D} \hat{=} \text{type double}$ ,  $\mathbf{C} \hat{=} \text{type complex}$
- $M, N \hat{=} \text{size of matrix } \mathbf{A}$
- $\text{ALPHA} \hat{=} \text{scalar parameter } \alpha$
- $\mathbf{A} \hat{=} \text{matrix } \mathbf{A} \text{ stored in } \text{linear array} \text{ of length } M \cdot N \text{ (column major arrangement)}$   

$$(\mathbf{A})_{i,j} = \mathbf{A}[N * (j - 1) + i] .$$
- $\text{LDA} \hat{=} \text{"leading dimension" of } \mathbf{A} \in \mathbb{K}^{n,m}$ , that is, the number  $n$  of rows.
- $\mathbf{X} \hat{=} \text{vector } \mathbf{x}: \text{array of type } \mathbf{x}$
- $\text{INCX} \hat{=} \text{stride for traversing vector } \mathbf{X}$
- $\text{BETA} \hat{=} \text{scalar parameter } \beta$
- $\mathbf{Y} \hat{=} \text{vector } \mathbf{y}: \text{array of type } \mathbf{x}$
- $\text{INCY} \hat{=} \text{stride for traversing vector } \mathbf{Y}$

- BLAS LEVEL 3:** matrix-matrix operations, asymptotic complexity  $O(mnk)$ ,  $(m, n, k) \hat{=} \text{matrix sizes}$

- matrix  $\times$  matrix multiplication  $\mathbf{C} = \alpha \mathbf{AB} + \beta \mathbf{C}$

**XGEMM (TRANSA, TRANSB, M, N, K,  
ALPHA, A, LDA, X, B, LDB,  
BETA, C, LDC)**

(☞ meaning of arguments as above)

**Remark 1.3.2.3 (BLAS calling conventions)** The BLAS calling syntax seems queer in light of modern object oriented programming paradigms, but it is a legacy of FORTRAN77, which was (and partly still is) the programming language, in which the BLAS routines were coded.

It is a very common situation in scientific computing that one has to rely on old codes and libraries implemented in an old-fashioned style. □

**EXAMPLE 1.3.2.4 (Calling BLAS routines from C/C++)** When calling BLAS library functions from C, all arguments have to be passed by reference (as pointers), in order to comply with the argument passing mechanism of FORTRAN77, which is the model followed by BLAS.

#### C++-code 1.3.2.5: BLAS-based SAXPY operation in C++

```

1 #define daxpy_ daxpy
2 #include <iostream>
3

```

```

4 // Definition of the required BLAS function. This is usually done
5 // in a header file like blas.h that is included in the EIGEN3
6 // distribution
7 extern "C" {
8     int daxpy_(const int* n, const double* da, const double* dx,
9                 const int* incx, double* dy, const int* incy);
10 }
11
12 using namespace std;
13
14 int main(){
15     const int      n      = 5; // length of vector
16     const int      incx   = 1; // stride
17     const int      incy   = 1; // stride
18     double alpha = 2.5;      // scaling factor
19
20     // Allocated raw arrays of doubles
21     double* x  = new double [n];
22     double* y  = new double [n];
23
24     for (size_t i=0; i<n; i++){
25         x[i] = 3.1415 * i;
26         y[i] = 1.0 / (double)(i+1);
27     }
28
29     cout << "x=[";
30     for (size_t i=0; i<n; i++) cout << x[i] << ' ';
31     cout << "]" << endl;
32     cout << "y=[";
33     for (size_t i=0; i<n; i++) cout << y[i] << ' ';
34     cout << "]" << endl;
35
36     // Call the BLAS library function passing pointers to all arguments
37     // (Necessary when calling FORTRAN routines from C
38     daxpy_(&n, &alpha, x, &incx, y, &incy);
39
40     cout << "y = " << alpha << " * x + y = ]";
41     for (int i=0; i<n; i++) cout << y[i] << ' ';
42     cout << "]" << endl;
43     return(0);
44 }
```

When using EIGEN in a mode that includes an external BLAS library, all these calls are wrapped into EIGEN methods.

**EXAMPLE 1.3.2.6 (Using Intel Math Kernel Library (Intel MKL) from EIGEN)** The Intel Math Kernel Library is a highly optimized math library for Intel processors and can be called directly from EIGEN, see  [EIGEN documentation](#) on “Using Intel® Math Kernel Library from Eigen”.

#### C++-code 1.3.2.7: Timing of matrix multiplication in EIGEN for MKL comparison → [GITLAB](#)

```

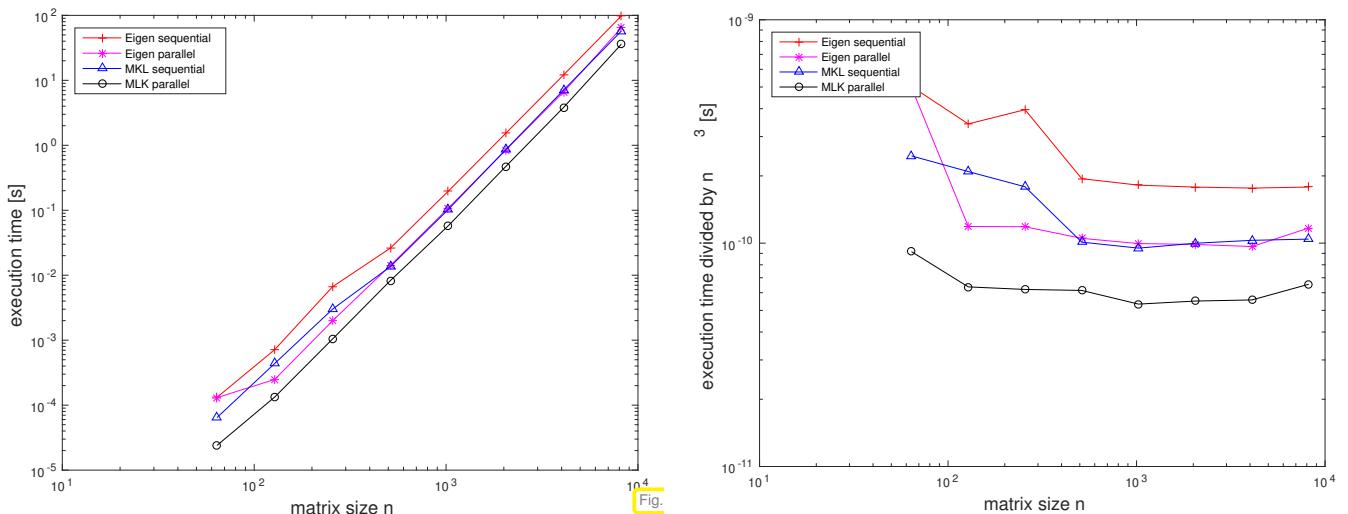
2 //! script for timing different implementations of matrix
3 // multiplications
4 void mmeigenmkl(){
5     int nruns = 3, minExp = 6, maxExp = 13;
6     MatrixXd timings(maxExp-minExp+1,2);
7     for(int p = 0; p <= maxExp-minExp; ++p){
8         Timer t1; // timer class
9         int n = std::pow(2, minExp + p);
10        MatrixXd A = MatrixXd::Random(n,n);
11        MatrixXd B = MatrixXd::Random(n,n);
12        MatrixXd C = MatrixXd::Zero(n,n);
```

```

12     for(int q = 0; q < nruns; ++q){
13         t1.start();
14         C = A * B;
15         t1.stop();
16     }
17     timings(p,0)=n; timings(p,1)=t1.min();
18 }
19 std::cout << std::scientific << std::setprecision(3) << timings << std::endl;
20 }
```

Timing results:

<i>n</i>	EIGEN sequential [s]	EIGEN parallel [s]	MKL sequential [s]	MKL parallel [s]
64	1.318e-04	1.304e-04	6.442e-05	2.401e-05
128	7.168e-04	2.490e-04	4.386e-04	1.336e-04
256	6.641e-03	1.987e-03	3.000e-03	1.041e-03
512	2.609e-02	1.410e-02	1.356e-02	8.243e-03
1024	1.952e-01	1.069e-01	1.020e-01	5.728e-02
2048	1.531e+00	8.477e-01	8.581e-01	4.729e-01
4096	1.212e+01	6.635e+00	7.075e+00	3.827e+00
8192	9.801e+01	6.426e+01	5.731e+01	3.598e+01



Timing environment:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB, Mem 8 GB
- ◆ gcc 4.8.4, -O3

## 1.4 Computational Effort

Large scale numerical computations require immense resources and execution time of numerical codes often becomes a central concern. Therefore, much emphasis has to be put on

1. designing algorithms that produce a desired result with (nearly) minimal computational effort (defined precisely below),
2. exploit possibilities for parallel and vectorised execution,

3. organising algorithms in order to make them fit memory hierarchies,
4. implementing codes that make optimal use of hardware resources and capabilities,

While Item 2–Item 4 are out of the scope of this course and will be treated in more advanced lectures, Item 1 will be a recurring theme.

The following definition encapsulates what is regarded as a measure for the “cost” of an algorithm in computational mathematics.

#### Definition 1.4.0.1. Computational effort

The **computational effort/computational cost** required/incurred by a numerical code amounts to the number of **elementary operations** (additions, subtractions, multiplications, divisions, square roots) executed in a run.

**§1.4.0.2 (What computational effort does not tell us)** Fifty years ago counting elementary operations provided good predictions of runtimes, but nowadays this is no longer true.

#### “Computational effort $\not\propto$ runtime”



The computational effort involved in a run of a numerical code is only loosely related to overall execution time on modern computers.

This is conspicuous in Exp. 1.2.3.9, where algorithms incurring exactly the same computational effort took different times to execute.

The reason is that on today’s computers a key bottleneck for fast execution is latency and bandwidth of memory, *cf.* the discussion at the end of Exp. 1.2.3.9 and [KW03]. Thus, concepts like **I/O-complexity** [AV88; GJ10] might be more appropriate for gauging the efficiency of a code, because they take into account the pattern of memory access.  $\square$

### 1.4.1 (Asymptotic) Complexity

The concept of computational effort from Def. 1.4.0.1 is still useful in a particular context:

#### Definition 1.4.1.1. (Asymptotic) complexity

The **asymptotic complexity** of an algorithm characterises the worst-case dependence of its computational effort ( $\rightarrow$  Def. 1.4.0.1) on one or more **problem size parameter(s)** when these tend to  $\infty$ .

- *Problem size parameters* in numerical linear algebra usually are the lengths and dimensions of the vectors and matrices that an algorithm takes as inputs.
- *Worst case* indicates that the maximum effort over a set of admissible data is taken into account.

When dealing with asymptotic complexities a mathematical formalism comes handy:

### Definition 1.4.1.2. Landau symbol [AG11, p. 7]

We write  $F(n) = O(G(n))$  for two functions  $F, G : \mathbb{N} \rightarrow \mathbb{R}$ , if there exists a constant  $C > 0$  and  $n_* \in \mathbb{N}$  such that

$$F(n) \leq C G(n) \quad \forall n \geq n_* .$$

More generally,  $F(n_1, \dots, n_k) = O(G(n_1, \dots, n_k))$  for two functions  $F, G : \mathbb{N}^k \rightarrow \mathbb{R}$  implies the existence of a constant  $C > 0$  and a threshold value  $n_* \in \mathbb{N}$  such that

$$F(n_1, \dots, n_k) \leq C G(n_1, \dots, n_k) \quad \forall n_1, \dots, n_k \in \mathbb{N}, \quad n_\ell \geq n_*, \quad \ell = 1, \dots, k .$$

**Remark 1.4.1.3 (Meaningful “ $O$ -bounds” for complexity)** Of course, the definition of the Landau symbol leaves ample freedom for stating meaningless bounds; an algorithm that runs with linear complexity  $O(n)$  can be correctly labelled as possessing  $O(\exp(n))$  complexity.

Yet, whenever the Landau notation is used to describe asymptotic complexities, the bounds have to be **sharp** in the sense that no function with slower asymptotic growth will be possible inside the  $O$ . To make this precise we stipulate the following.

### Sharpness of a complexity bound

Whenever the asymptotic complexity of an algorithm is stated as  $O(n^\alpha \log^\beta n \exp(\gamma n^\delta))$  with non-negative parameters  $\alpha, \beta, \gamma, \delta \geq 0$  in terms of the problem size parameter  $n$ , we take for granted that choosing a smaller value for any of the parameters will no longer yield a valid (or provable) asymptotic bound.

In particular

- ◆ complexity  $O(n)$  means that the complexity is not  $O(n^\alpha)$  for any  $\alpha < 1$ ,
- ◆ complexity  $O(\exp(n))$  excludes asymptotic complexity  $O(n^p)$  for any  $p \in \mathbb{R}$ .

Terminology: If the asymptotic complexity of an algorithm is  $O(n^p)$  with  $p = 1, 2, 3$  we say that it is of “linear”, “quadratic”, and “cubic” complexity, respectively. □

**Remark 1.4.1.5 (Relevance of asymptotic complexity)** § 1.4.0.2 warned us that computational effort and, thus, asymptotic complexity, of an algorithm for a concrete problem on a particular platform may not have much to do with the actual runtime (the blame goes to memory hierarchies, internal pipelining, vectorisation, etc.).

Then, why do we pay so much attention to asymptotic complexity in this course?

To a certain extent, the asymptotic complexity allows to predict the *dependence of the runtime* of a particular implementation of an algorithm *on the problem size* (for large problems).

For instance, an algorithm with asymptotic complexity  $O(n^2)$  is likely to take  $4\times$  as much time when the problem size is doubled. □

### §1.4.1.6 (Concluding polynomial complexity from runtime measurements)

Available: “Measured runtimes”  $t_i = t_i(n_i)$  for different values  $n_1, n_2, \dots, n_N, n_i \in \mathbb{N}$ , of the problem size parameter

Conjectured: power law dependence  $t_i \approx Cn_i^\alpha$  (also “algebraic dependence”),  $\alpha \in \mathbb{R}$

How can we glean evidence that supports or refutes our conjecture from the data? Look at the data in **doubly logarithmic scale!**

$$t_i = Cn_i^\alpha \Rightarrow \log(t_i) \approx \log C + \alpha \log(n_i), \quad i = 1, \dots, N.$$

► If the conjecture holds true, then the points  $(n_i, t_i)$  will approximately lie on a *straight line* with **slope**  $\alpha$  in a doubly logarithmic plot (which can be created in PYTHON by the `matplotlib.pyplot.loglog` plotting command).

➤ Offers a quick “visual test” of conjectured asymptotic complexity

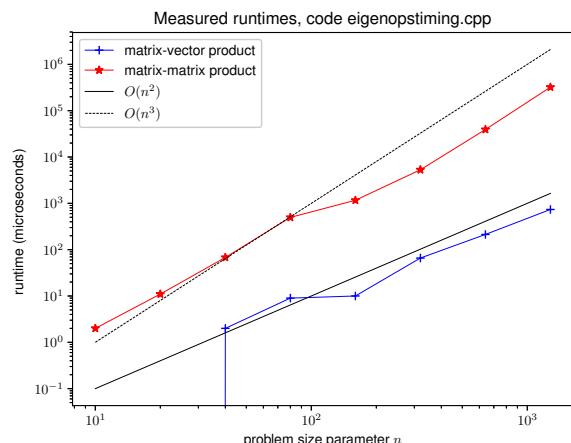
More rigorous: Perform linear regression on  $(\log n_i, \log t_i), i = 1, \dots, N$  (→ Chapter 3)

## 1.4.2 Cost of Basic Linear-Algebra Operations

Performing elementary BLAS-type operations through simple (nested) loops, we arrive at the following obvious complexity bounds:

operation	description	#mul/div	#add/sub	asymp. complexity
dot product	$(\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^n) \mapsto \mathbf{x}^H \mathbf{y}$	$n$	$n - 1$	$O(n)$
tensor product	$(\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n) \mapsto \mathbf{x}\mathbf{y}^H$	$nm$	0	$O(mn)$
Matrix × vector	$(\mathbf{x} \in \mathbb{R}^n, \mathbf{A} \in \mathbb{R}^{m,n}) \mapsto \mathbf{Ax}$	$mn$	$(n - 1)m$	$O(mn)$
matrix product <sup>(*)</sup>	$(\mathbf{A} \in \mathbb{R}^{m,n}, \mathbf{B} \in \mathbb{R}^{n,k}) \mapsto \mathbf{AB}$	$mnk$	$mk(n - 1)$	$O(mnk)$

### EXPERIMENT 1.4.2.1 (Runtimes of elementary linear algebra operations in EIGEN)



- Timing code → GITLAB,  $m = n$
- Linux kernel 5.5.9-100.fc30.x86\_64, Fedora 30
- gcc 9.3.1 -O2, Release mode
- Intel(R) Core(TM) i7-7600U CPU @ 2.80GHz

Runtime data points approximately on lines in **doubly logarithmic plot**, see § 1.4.1.6: Asymptotic behavior of measured runtimes match predictions from above table for large  $n$ .

### Remark 1.4.2.2 (“Fast” matrix multiplication)

(\*) The  $O(mnk)$  complexity bound applies to “straightforward” matrix multiplication according to (1.3.1.1).

For  $m = n = k$  there are (sophisticated) variants with better asymptotic complexity, e.g., the **divide-and-conquer Strassen algorithm** [Str69] with asymptotic complexity  $O(n^{\log_2 7})$ :

Start from  $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{n,n}$  with  $n = 2\ell$ ,  $\ell \in \mathbb{N}$ . The idea relies on the block matrix product (1.3.1.14) with  $\mathbf{A}_{ij}, \mathbf{B}_{ij} \in \mathbb{K}^{\ell,\ell}$ ,  $i, j \in \{1, 2\}$ . Let  $\mathbf{C} := \mathbf{AB}$  be partitioned accordingly:  $\mathbf{C} = \begin{pmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{pmatrix}$ . Then tedious elementary computations reveal

$$\begin{aligned}\mathbf{C}_{11} &= \mathbf{Q}_0 + \mathbf{Q}_3 - \mathbf{Q}_4 + \mathbf{Q}_6, \\ \mathbf{C}_{21} &= \mathbf{Q}_1 + \mathbf{Q}_3, \\ \mathbf{C}_{12} &= \mathbf{Q}_2 + \mathbf{Q}_4, \\ \mathbf{C}_{22} &= \mathbf{Q}_0 + \mathbf{Q}_2 - \mathbf{Q}_1 + \mathbf{Q}_5,\end{aligned}$$

where the  $\mathbf{Q}_k \in \mathbb{K}^{\ell,\ell}$ ,  $k = 1, \dots, 7$  are obtained from

$$\begin{aligned}\mathbf{Q}_0 &= (\mathbf{A}_{11} + \mathbf{A}_{22}) * (\mathbf{B}_{11} + \mathbf{B}_{22}), \\ \mathbf{Q}_1 &= (\mathbf{A}_{21} + \mathbf{A}_{22}) * \mathbf{B}_{11}, \\ \mathbf{Q}_2 &= \mathbf{A}_{11} * (\mathbf{B}_{12} - \mathbf{B}_{22}), \\ \mathbf{Q}_3 &= \mathbf{A}_{22} * (-\mathbf{B}_{11} + \mathbf{B}_{21}), \\ \mathbf{Q}_4 &= (\mathbf{A}_{11} + \mathbf{A}_{12}) * \mathbf{B}_{22}, \\ \mathbf{Q}_5 &= (-\mathbf{A}_{11} + \mathbf{A}_{21}) * (\mathbf{B}_{11} + \mathbf{B}_{12}), \\ \mathbf{Q}_6 &= (\mathbf{A}_{12} - \mathbf{A}_{22}) * (\mathbf{B}_{21} + \mathbf{B}_{22}).\end{aligned}$$

Beside a considerable number of matrix additions (computational effort  $O(n^2)$ ) it takes only 7 multiplications of matrices of size  $n/2$  to compute  $\mathbf{C}$ ! Strassen's algorithm boils down to the *recursive application* of these formulas for  $n = 2^k$ ,  $k \in \mathbb{N}$ .

A refined algorithm of this type can achieve complexity  $O(n^{2.36})$ , see [CW90].  $\square$

### 1.4.3 Reducing Complexity in Numerical Linear Algebra: Some Tricks

In computations involving matrices and vectors complexity of algorithms can often be reduced by performing the operations in a particular order:

**EXAMPLE 1.4.3.1 (Efficient associative matrix multiplication)** We consider the multiplication with a **rank-1-matrix**. Matrices with rank 1 can always be obtained as the tensor product of two vectors, that is, the matrix product of a column vector and a row vector. Given  $\mathbf{a} \in \mathbb{K}^m$ ,  $\mathbf{b} \in \mathbb{K}^n$ ,  $\mathbf{x} \in \mathbb{K}^n$  we may compute the vector  $\mathbf{y} = \mathbf{ab}^\top \mathbf{x}$  in two ways:

$$\mathbf{y} = (\mathbf{ab}^\top) \mathbf{x}. \quad (1.4.3.2) \qquad \qquad \qquad \mathbf{y} = \mathbf{a}(\mathbf{b}^\top \mathbf{x}). \quad (1.4.3.3)$$

$\text{T} = (\mathbf{a} * \mathbf{b}.\text{transpose}()) * \mathbf{x};$

$\text{t} = \mathbf{a} * \mathbf{b}.\text{dot}(\mathbf{x});$

► complexity  $O(mn)$

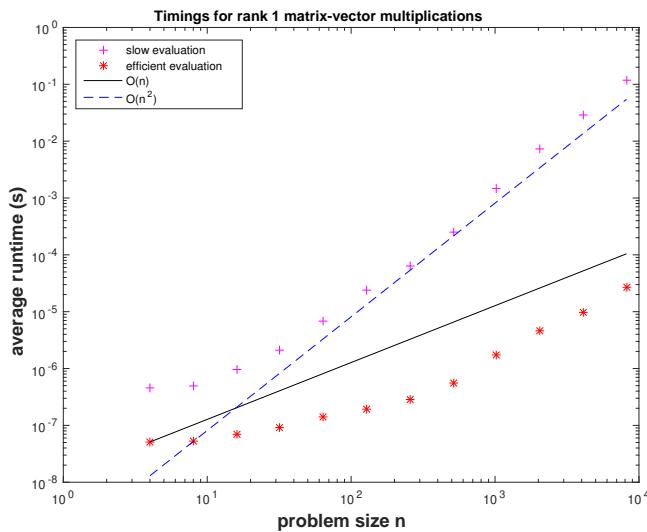
► complexity  $O(n + m)$  ("linear complexity")

Visualization of evaluation according to (1.4.3.2):

$$\left[ \begin{array}{|c|} \hline \end{array} \right] \cdot \left[ \begin{array}{|c|} \hline \end{array} \right] = \left[ \begin{array}{|c|} \hline \end{array} \right] \left[ \begin{array}{|c|} \hline \end{array} \right]$$

Visualization of evaluation according to (1.4.3.3):

$$\begin{bmatrix} \text{slow eval} \\ \text{efficient eval} \end{bmatrix} \begin{bmatrix} \text{slow eval} \\ \text{efficient eval} \end{bmatrix} \cdot \begin{bmatrix} \text{slow eval} \\ \text{efficient eval} \end{bmatrix} = \begin{bmatrix} \text{slow eval} \\ \text{efficient eval} \end{bmatrix} \begin{bmatrix} \text{slow eval} \\ \text{efficient eval} \end{bmatrix}$$



▷ average runtimes for efficient/inefficient matrix  $\times$  vector multiplication with rank-1 matrices , see § 1.4.1.6 for the rationale behind choosing a *doubly logarithmic plot*.

Platform:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz — 4
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB,
- ◆ 8 GB main memory
- ◆ gcc 4.8.4, -O3

#### C++ code 1.4.3.4: EIGEN code for Ex. 1.4.3.1 → GITLAB

```

2  ///! This function compares the runtimes for the multiplication
3  ///! of a vector with a rank-1 matrix  $\mathbf{ab}^\top$ ,  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ 
4  ///! using different associative evaluations.
5  ///! Runtime measurements consider minimal time for
6  ///! several (nrunc) runs
7  MatrixXd dottenstming () {
8      const int nrunc = 3, minExp = 2, maxExp = 13;
9      // Matrix for storing recorded runtimes
10     MatrixXd timings(maxExp-minExp+1,3);
11     for(int i = 0; i <= maxExp-minExp; ++i){
12         Timer tfool, tsmart; // Timer objects
13         const int n = std::pow(2, minExp + i);
14         VectorXd a = VectorXd::LinSpaced(n,1,n);
15         VectorXd b = VectorXd::LinSpaced(n,1,n).reverse();
16         VectorXd x = VectorXd::Random(n,1), y(n);
17         for(int j = 0; j < nrunc; ++j){
18             // Grossly wasteful evaluation
19             tfool.start(); y = (a*b.transpose())*x; tfool.stop();
20             // Efficient implementation
21             tsmart.start(); y = a * b.dot(x); tsmart.stop();
22         }
23         timings(i,0)=n;
24         timings(i,1)=tsmart.min(); timings(i,2)=tfool.min();
25     }
26     return timings;
27 }
```

Complexity can sometimes be reduced by reusing intermediate results.

**EXAMPLE 1.4.3.5 (Hidden summation)** The asymptotic complexity of the EIGEN code

```
Eigen::MatrixXd AB = A*B.transpose();
y = AB.triangularView<Eigen::Upper> () *x;
```

when supplied with two **low-rank** matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{n,p}$ ,  $p \ll n$ , in terms of  $n \rightarrow \infty$  obviously is  $O(n^2)$ , because an intermediate  $n \times n$ -matrix  $\mathbf{AB}^T$  is built.

First, consider the case of a tensor product (= rank-1) matrix, that is,  $p = 1$ ,  $\mathbf{A} \leftrightarrow \mathbf{a} = [a_1, \dots, a_n]^\top \in \mathbb{K}^n$ ,  $\mathbf{B} \leftrightarrow \mathbf{b} = [b_1, \dots, b_n] \in \mathbb{K}^n$ . Then

$$\begin{aligned} \mathbf{y} &= \text{triu}(\mathbf{ab}^T)\mathbf{x} = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \dots & \dots & \dots & a_1 b_n \\ 0 & a_2 b_2 & a_2 b_3 & \dots & \dots & a_2 b_n \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & a_n b_n & \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &= \begin{bmatrix} a_1 \\ \vdots \\ \ddots \\ \ddots \\ a_n \end{bmatrix} \left( \underbrace{\begin{bmatrix} 1 & 1 & \dots & \dots & 1 \\ 0 & 1 & 1 & \dots & \dots & 1 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & 1 \end{bmatrix}}_{\mathbf{T}} \left( \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \right) \right) \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}. \end{aligned}$$

The brackets indicate the order of the matrix×vector multiplications. Thus, the core problem is the fast multiplication of a vector with an upper triangular matrix  $\mathbf{T}$  described in EIGEN syntax by `Eigen::MatrixXd::Ones(n,n).triangularView<Eigen::Upper>()`. Note that multiplication of a vector  $\mathbf{x}$  with  $\mathbf{T}$  yields a vector of partial sums of components of  $\mathbf{x}$  starting from last component:

$$\begin{bmatrix} 1 & 1 & \dots & \dots & 1 \\ 0 & 1 & 1 & \dots & \dots & 1 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & 1 \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} s_n \\ s_{n-1} \\ \vdots \\ s_1 \end{bmatrix}, \quad s_j := \sum_{k=1}^j v_k.$$

This can be achieved by invoking the special C++ command `std::partial_sum` from the C++ standard library ([documentation](#)). We also observe that

$$\mathbf{AB}^T = \sum_{\ell=1}^p (\mathbf{A})_{:, \ell} ((\mathbf{B})_{:, \ell})^\top,$$

so that the computations for the special case  $p = 1$  discussed above can simply be reused  $p$  times!

**C++ code 1.4.3.6: Efficient multiplication with the upper diagonal part of a rank- $p$ -matrix in EIGEN** → [GITLAB](#)

```

2  ///! Computation of  $\mathbf{y} = \text{triu}(\mathbf{AB}^T)\mathbf{x}$ 
3  ///! Efficient implementation with backward cumulative sum
4  ///! (partial_sum)
5  template<class Vec, class Mat>
6  void ltrimulteff(const Mat& A, const Mat& B, const Vec& x, Vec& y){
7      const int n = A.rows(), p = A.cols();
8      assert( n == B.rows() && p == B.cols()); // size mismatch
9      for(int l = 0; l < p; ++l){
10          Vec tmp = (B.col(l).array() * x.array()).matrix().reverse();
11          std::partial_sum(tmp.data(), tmp.data() + n, tmp.data());
12          y += (A.col(l).array() * tmp.reverse().array()).matrix();
13      }
14  }

```

This code enjoys the obvious complexity of  $O(pn)$  for  $p, n \rightarrow \infty$ ,  $p < n$ . The code offers an example of a function templated with its argument types, see § 0.3.2.1. The types **Vec** and **Mat** must fit the concept of EIGEN vectors/matrices.

The next concept from linear algebra is important in the context of computing with multi-dimensional arrays.

**Definition 1.4.3.7. Kronecker product**

The **Kronecker product**  $\mathbf{A} \otimes \mathbf{B}$  of two matrices  $\mathbf{A} \in \mathbb{K}^{m,n}$  and  $\mathbf{B} \in \mathbb{K}^{l,k}$ ,  $m, n, l, k \in \mathbb{N}$ , is the  $(ml) \times (nk)$ -matrix

$$\mathbf{A} \otimes \mathbf{B} := \begin{bmatrix} (\mathbf{A})_{1,1}\mathbf{B} & (\mathbf{A})_{1,2}\mathbf{B} & \dots & \dots & (\mathbf{A})_{1,n}\mathbf{B} \\ (\mathbf{A})_{2,1}\mathbf{B} & (\mathbf{A})_{2,2}\mathbf{B} & & & \vdots \\ \vdots & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ (\mathbf{A})_{m,1}\mathbf{B} & (\mathbf{A})_{m,2}\mathbf{B} & \dots & \dots & (\mathbf{A})_{m,n}\mathbf{B} \end{bmatrix} \in \mathbb{K}^{ml,nk}.$$

**EXAMPLE 1.4.3.8 (Multiplication of Kronecker product with vector)** The function  $(\mathbf{A} \otimes \mathbf{B})\mathbf{x}$  when invoked with two matrices  $\mathbf{A} \in \mathbb{K}^{m,n}$  and  $\mathbf{B} \in \mathbb{K}^{l,k}$  and a vector  $\mathbf{x} \in \mathbb{K}^{nk}$ , will suffer an asymptotic complexity of  $O(m \cdot n \cdot l \cdot k)$ , determined by the size of the intermediate dense matrix  $\mathbf{A} \otimes \mathbf{B} \in \mathbb{K}^{ml,nk}$ .

Using the partitioning of the vector  $\mathbf{x}$  into  $n$  equally long sub-vectors

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}^1 \\ \mathbf{x}^2 \\ \vdots \\ \mathbf{x}^n \end{bmatrix}, \quad \mathbf{x}^j \in \mathbb{K}^k,$$

we find the representation

$$(\mathbf{A} \otimes \mathbf{B})\mathbf{x} = \begin{bmatrix} (\mathbf{A})_{1,1}\mathbf{Bx}^1 + (\mathbf{A})_{1,2}\mathbf{Bx}^2 + \dots + (\mathbf{A})_{1,n}\mathbf{Bx}^n \\ (\mathbf{A})_{2,1}\mathbf{Bx}^1 + (\mathbf{A})_{2,2}\mathbf{Bx}^2 + \dots + (\mathbf{A})_{2,n}\mathbf{Bx}^n \\ \vdots \\ \vdots \\ (\mathbf{A})_{m,1}\mathbf{Bx}^1 + (\mathbf{A})_{m,2}\mathbf{Bx}^2 + \dots + (\mathbf{A})_{m,n}\mathbf{Bx}^n \end{bmatrix}.$$

The idea is to form the products  $\mathbf{Bx}^j$ ,  $j = 1, \dots, n$ , once, and then combine them linearly with coefficients given by the entries in the rows of  $\mathbf{A}$ :

**C++ code 1.4.3.9: Efficient multiplication of Kronecker product with vector in EIGEN**  
[→ GITLAB](#)

```

2 template <class Matrix, class Vector>
3 void kronmultv(const Matrix &A, const Matrix &B, const Vector &x, Vector &y){
4     unsigned int m = A.rows(); unsigned int n = A.cols();
5     unsigned int l = B.rows(); unsigned int k = B.cols();
6     // 1st matrix mult. computes the products  $\mathbf{Bx}^j$ 
7     // 2nd matrix mult. combines them linearly with the coefficients of
8     A
9     Matrix t = B * Matrix::Map(x.data(), k, n) * A.transpose(); //
10    y = Matrix::Map(t.data(), m*l, 1);
11 }
```

Recall the reshaping of a matrix in EIGEN in order to understand this code: Rem. 1.2.3.6.

The asymptotic complexity of this code is determined by the two matrix multiplications in Line 8. This yields the asymptotic complexity  $O(lkn + mnl)$  for  $l, k, m, n \rightarrow \infty$ .

**PYTHON-code 1.4.3.10: Efficient multiplication of Kronecker product with vector in PYTHON**

```

1 def kronmultv(A, B, x):
2     n, k = A.shape[1], B.shape[1]
3     assert x.size == n * k, 'size mismatch'
4     xx = np.reshape(x, (n, k))
5     Z = np.dot(xx, B.T)
6     yy = np.dot(A, Z)
7     return np.ravel(yy)
```

Note that different reshaping is used in the PYTHON code due to the default row major storage order. ↗

**Review question(s) 1.4.3.11 (Computational effort)**

**(Q1.4.3.11.A)** Explain why the classical concept of “computational effort” (= computational cost) is only loosely related to the runtime of a concrete implementation of an algorithm.

**(Q1.4.3.11.B)** We are given two dense matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n,p}$ ,  $n, p \in \mathbb{N}$ ,  $p < n$  fixed. What is the asymptotic complexity of each of the following two lines of code in terms of  $n, p \rightarrow \infty$ ?

```
Eigen::MatrixXd AB = A*B.transpose();
y = AB.triangularView<Eigen::Upper>()*x;
```

**(Q1.4.3.11.C)** Given a vector  $\mathbf{u} \in \mathbb{R}^n$ ,  $n \in \mathbb{N}$ , we consider the matrix

$$\mathbf{A} \in \mathbb{R}^{n,n}: \quad (\mathbf{A})_{ij} = u_i + u_j + u_i u_j, \quad i, j \in \{1, \dots, n\}.$$

Outline an efficient algorithm for computing the matrix-vector product  $\mathbf{Ax}$ ,  $\mathbf{x} \in \mathbb{R}^n$ . △

## 1.5 Machine Arithmetic and Consequences

### 1.5.1 Experiment: Loss of Orthogonality

**§1.5.1.1 (Gram-Schmidt orthogonalisation)** From linear algebra [NS02, Sect. 4.4] or Ex. 0.3.5.29 we recall the fundamental algorithm of **Gram-Schmidt orthogonalisation** of an ordered finite set  $\{\mathbf{a}^1, \dots, \mathbf{a}^k\}$ ,  $k \in \mathbb{N}$ , of vectors  $\mathbf{a}^\ell \in \mathbb{K}^n$ :

Input:  $\{\mathbf{a}^1, \dots, \mathbf{a}^k\} \subset \mathbb{K}^n$

```

1:  $\mathbf{q}^1 := \frac{\mathbf{a}^1}{\|\mathbf{a}^1\|_2}$  % 1st output vector
2: for  $j = 2, \dots, k$  do
   { % Orthogonal projection
3:    $\mathbf{q}^j := \mathbf{a}^j$ 
4:   for  $\ell = 1, 2, \dots, j-1$  do (GS)
5:     {  $\mathbf{q}^j \leftarrow \mathbf{q}^j - \mathbf{a}^\ell \cdot \mathbf{q}^\ell \mathbf{q}^\ell$  }
6:     if ( $\mathbf{q}^j = \mathbf{0}$ ) then STOP
7:     else {  $\mathbf{q}^j \leftarrow \frac{\mathbf{q}^j}{\|\mathbf{q}^j\|_2}$  }
8:   }

```

Output:  $\{\mathbf{q}^1, \dots, \mathbf{q}^j\}$

☞ Notation:  $\|\cdot\|_2$   $\hat{=}$  Euclidean norm of a vector  $\in \mathbb{K}^n$

The following code implements the Gram-Schmidt orthonormalization of a set of vectors passed as the columns of a matrix  $\mathbf{A} \in \mathbb{R}^{n,k}$ . The template parameter **Matrix** should match the concept of a matrix type in EIGEN like **Eigen::MatrixXd** or

In linear algebra we have learnt that, if it does not **STOP** prematurely, this algorithm will compute **orthonormal vectors**  $\mathbf{q}^1, \dots, \mathbf{q}^k$  satisfying

$$\text{Span}\{\mathbf{q}^1, \dots, \mathbf{q}^k\} = \text{Span}\{\mathbf{a}^1, \dots, \mathbf{a}^k\}, \quad (1.5.1.2)$$

for all  $\ell \in \{1, \dots, k\}$ .

More precisely, if  $\mathbf{a}^1, \dots, \mathbf{a}^\ell, \ell \leq k$ , are linearly independent, then the Gram-Schmidt algorithm will not terminate before the  $\ell+1$ -th step.

#### C++ code 1.5.1.3: Gram-Schmidt orthogonalisation in EIGEN → GITLAB

```

2 template <class Matrix> Matrix gramschmidt(const Matrix & $\mathbf{A}$ ) {
3   Matrix  $\mathbf{Q} = \mathbf{A}$ ;
4   // First vector just gets normalized, Line 1 of (GS)
5    $\mathbf{Q}.\text{col}(0).\text{normalize}()$ ;
6   for (unsigned int  $j = 1$ ;  $j < \mathbf{A}.\text{cols}()$ ;  $++j$ ) {
7     // Replace inner loop over each previous vector in Q with fast
8     // matrix-vector multiplication (Lines 4, 5 of (GS))
9      $\mathbf{Q}.\text{col}(j) -= \mathbf{Q}.\text{leftCols}(j) *$ 
10    ( $\mathbf{Q}.\text{leftCols}(j).\text{adjoint}() * \mathbf{A}.\text{col}(j)$ ); //
11    // Normalize vector, if possible.
12    // Otherwise columns of A must have been linearly dependent
13    if ( $\mathbf{Q}.\text{col}(j).\text{norm}() \leq 10e-9 * \mathbf{A}.\text{col}(j).\text{norm}()$ ) { //
14      std::cerr << "Gram-Schmidt failed: A has lin. dep columns." << std::endl;
15      break;
16    } else {
17       $\mathbf{Q}.\text{col}(j).\text{normalize}()$ ;
18    } // Line 7 of (GS)
19  }
20  return  $\mathbf{Q}$ ;
21}

```

We will soon learn the rationale behind the odd test in Line 13.

In PYTHON the same algorithm can be implemented as follows:

**PYTHON-code 1.5.1.4: Gram-Schmidt orthogonalisation in PYTHON**

```

1 def gramschmidt(A):
2     _, k = A.shape
3     Q = A[:, [0]] / np.linalg.norm(A[:, 0])
4     for j in range(1, k):
5         q = A[:, j] - np.dot(Q, np.dot(Q.T, A[:, j]))
6         nq = np.linalg.norm(q)
7         if nq < 1e-9 * np.linalg.norm(A[:, j]):
8             break
9     Q = np.column_stack([Q, q / nq])
10    return Q

```

Note the different loop range due to the zero-based indexing in PYTHON.

**EXPERIMENT 1.5.1.5 (Unstable Gram-Schmidt orthonormalization)** If  $\{a^1, \dots, a^k\}$  are linearly independent we expect the output vectors  $q^1, \dots, q^k$  to be orthonormal:

$$(q^\ell)^\top q^m = \delta_{\ell,m}, \quad \ell, m \in \{1, \dots, k\}. \quad (1.5.1.6)$$

This property can be easily tested numerically, for instance by computing  $Q^\top Q$  for a matrix  $Q = [q^1, \dots, q^k] \in \mathbb{R}^{n,k}$ .

**C++ code 1.5.1.7: Wrong result from Gram-Schmidt orthogonalisation EIGEN → GITLAB**

```

2 void gsroundoff(MatrixXd& A){
3     // Gram-Schmidt orthogonalization of columns of A, see Code 1.5.1.3
4     MatrixXd Q = gramschmidt(A);
5     // Test orthonormality of columns of Q, which should be an
6     // orthogonal matrix according to theory
7     cout << setprecision(4) << fixed << "I = "
8     << endl << Q.transpose()*Q << endl;
9     // EIGEN'S stable internal Gram-Schmidt orthogonalization by
10    // QR-decomposition, see Rem. 1.5.1.9 below
11    HouseholderQR<MatrixXd> qr(A.rows(), A.cols()); //
12    qr.compute(A); MatrixXd Q1 = qr.householderQ(); //
13    // Test orthonormality
14    cout << "I1 = " << endl << Q1.transpose()*Q1 << endl;
15    // Check orthonormality and span property (1.5.1.2)
16    MatrixXd R1 = qr.matrixQR().triangularView<Upper>();
17    cout << scientific << "A-Q1*R1 = " << endl << A-Q1*R1 << endl;
18 }

```

We test the orthonormality of the output vectors of Gram-Schmidt orthogonalization for a special matrix  $A \in \mathbb{R}^{10,10}$ , a so-called Hilbert matrix, defined by  $(A)_{i,j} = (i+j-1)^{-1}$ . Then Code 1.5.1.7 produces the following output:

```

I =
 1.0000   0.0000  -0.0000   0.0000  -0.0000   0.0000  -0.0000  -0.0000  -0.0000  -0.0000
 0.0000   1.0000  -0.0000   0.0000  -0.0000   0.0000  -0.0000  -0.0000  -0.0000  -0.0000
 -0.0000  -0.0000   1.0000   0.0000  -0.0000   0.0000  -0.0000  -0.0000  -0.0000  -0.0000
 0.0000   0.0000   0.0000   1.0000  -0.0000   0.0000  -0.0000  -0.0000  -0.0000  -0.0000
 -0.0000  -0.0000  -0.0000  -0.0000   1.0000   0.0000  -0.0008  -0.0007  -0.0007  -0.0006
 0.0000   0.0000   0.0000   0.0000   0.0000   1.0000  -0.0540  -0.0430  -0.0360  -0.0289
 -0.0000  -0.0000  -0.0000  -0.0000  -0.0008   -0.0540   1.0000  0.9999  0.9998  0.9996
 -0.0000  -0.0000  -0.0000  -0.0000  -0.0007  -0.0430   0.9999  1.0000  1.0000  0.9999
 -0.0000  -0.0000  -0.0000  -0.0000  -0.0007  -0.0360   0.9998  1.0000  1.0000  1.0000
 -0.0000  -0.0000  -0.0000  -0.0000  -0.0006  -0.0289   0.9996  0.9999  1.0000  1.0000

```

Obviously, the vectors produced by the function `gramschmidt` fail to be orthonormal, contrary to the predictions of rigorous results from linear algebra!

However, Line 11, Line 12 of Code 1.5.1.7 demonstrate another way to orthonormalize the columns of a matrix using EIGEN's built-in class template `HouseholderQR` (more details in Section 3.3.3).

```
I1 =
1.0000 -0.0000 0.0000 -0.0000 -0.0000 -0.0000 -0.0000 -0.0000 -0.0000 0.0000
-0.0000 1.0000 -0.0000 0.0000 0.0000 -0.0000 -0.0000 0.0000 0.0000 -0.0000
0.0000 -0.0000 1.0000 -0.0000 -0.0000 0.0000 0.0000 0.0000 0.0000 0.0000
-0.0000 0.0000 -0.0000 1.0000 0.0000 -0.0000 -0.0000 0.0000 0.0000 0.0000
-0.0000 0.0000 -0.0000 0.0000 1.0000 -0.0000 0.0000 -0.0000 -0.0000 0.0000
-0.0000 -0.0000 0.0000 -0.0000 -0.0000 1.0000 -0.0000 -0.0000 0.0000 0.0000
-0.0000 -0.0000 0.0000 -0.0000 0.0000 -0.0000 1.0000 0.0000 0.0000 -0.0000
-0.0000 -0.0000 0.0000 0.0000 -0.0000 0.0000 -0.0000 1.0000 0.0000 0.0000
-0.0000 -0.0000 0.0000 0.0000 -0.0000 0.0000 1.0000 -0.0000 0.0000 0.0000
0.0000 -0.0000 0.0000 0.0000 -0.0000 -0.0000 -0.0000 0.0000 1.0000 0.0000
```

Now we observe apparently perfect orthogonality (1.5.1.6) of the columns of the matrix `Q1` in Code 1.5.1.7. Obviously, there is another algorithm that reliably yields the theoretical output of Gram-Schmidt orthogonalization. There is no denying that it is possible to compute Gram-Schmidt orthonormalization in a “clean” way.  $\square$

### “Computers cannot compute”

Computers cannot compute “properly” in  $\mathbb{R}$ : numerical computations may not respect the laws of analysis and linear algebra!

This introduces an important new aspect in the study of numerical algorithms.

**Remark 1.5.1.9 (Stable orthonormalization by QR-decomposition)** In Code 1.5.1.7 we saw the use of the EIGEN class `HouseholderQR<MatrixType>` for the purpose of Gram-Schmidt orthogonalisation. The underlying theory and algorithms will be explained later in Section 3.3.3. There we will have the following insight:

- Up to signs the columns of the matrix `Q` available from the QR-decomposition of `A` are the same vectors as produced by the Gram-Schmidt orthogonalisation of the columns of `A`.

! Code 1.5.1.7 demonstrates a case where a desired result can be obtained by two *algebraically equivalent* computations, that is, they yield the same result in a mathematical sense. Yet, when implemented on a computer, the results can be *vastly different*. One algorithm may produce junk (“*unstable algorithm*”), whereas the other lives up to the expectations (“*stable algorithm*”)

Supplement to Exp. 1.5.1.5: despite its ability to produce orthonormal vectors, we get as output for `D=A-Q1*R1` in Code 1.5.1.7:

```
D =
2.2204e-16 3.3307e-16 3.3307e-16 1.9429e-16 1.9429e-16 5.5511e-17 1.3878e-16 6.9389e-17 8.3267e-17 9.7145e-17
0.0000e+00 1.1102e-16 8.3267e-17 5.5511e-17 0.0000e+00 5.5511e-17 -2.7756e-17 0.0000e+00 0.0000e+00 4.1633e-17
-5.5511e-17 5.5511e-17 2.7756e-17 5.5511e-17 0.0000e+00 0.0000e+00 0.0000e+00 -1.3878e-17 1.3878e-17 1.3878e-17
0.0000e+00 5.5511e-17 2.7756e-17 2.7756e-17 0.0000e+00 1.3878e-17 -1.3878e-17 0.0000e+00 1.3878e-17 2.7756e-17
0.0000e+00 2.7756e-17 0.0000e+00 1.3878e-17 1.3878e-17 0.0000e+00 1.3878e-17 -1.3878e-17 1.3878e-17 4.1633e-17
-2.7756e-17 2.7756e-17 1.3878e-17 4.1633e-17 2.7756e-17 1.3878e-17 0.0000e+00 -1.3878e-17 2.7756e-17 2.7756e-17
0.0000e+00 2.7756e-17 0.0000e+00 2.7756e-17 2.7756e-17 1.3878e-17 0.0000e+00 1.3878e-17 2.7756e-17 2.0817e-17
0.0000e+00 2.7756e-17 2.7756e-17 1.3878e-17 1.3878e-17 0.0000e+00 1.3878e-17 2.0817e-17 2.7756e-17
1.3878e-17 1.3878e-17 1.3878e-17 2.7756e-17 1.3878e-17 0.0000e+00 -1.3878e-17 6.9389e-18 -6.9389e-18 1.3878e-17
0.0000e+00 2.7756e-17 1.3878e-17 1.3878e-17 0.0000e+00 0.0000e+00 0.0000e+00 1.3878e-17 1.3878e-17
```

- The computed QR-decomposition apparently fails to meet the exact algebraic requirements stipulated by Thm. 3.3.3.4. However, note the tiny size of the “defect”.  $\square$

## 1.5.2 Machine Numbers

**§1.5.2.1 (The finite and discrete set of machine numbers)** The reason, why computers must fail to execute exact computations with real numbers is clear:

Computer = finite automaton



can handle only *finitely many* numbers, not  $\mathbb{R}$

machine numbers, set  $\mathbb{M}$

Essential property:  $\mathbb{M}$  is a *finite, discrete* subset of  $\mathbb{R}$  (its numbers separated by gaps)

The set of machine numbers  $\mathbb{M}$  cannot be closed under elementary arithmetic operations  $+, -, \cdot, /$ , that is, when adding, multiplying, etc., two machine numbers the result may not belong to  $\mathbb{M}$ .

The results of elementary operations with operands in  $\mathbb{M}$  have to be mapped back to  $\mathbb{M}$ , an operation called **rounding**.

► roundoff errors (*ger.*: Rundungsfehler) are inevitable

The impact of roundoff means that mathematical identities may not carry over to the computational realm. As we have seen above in Exp. 1.5.1.5

Computers cannot compute “properly” !



**numerical computations**  $\neq$  **analysis  
linear algebra**

This introduces a *new and important* aspect in the study of numerical algorithms!

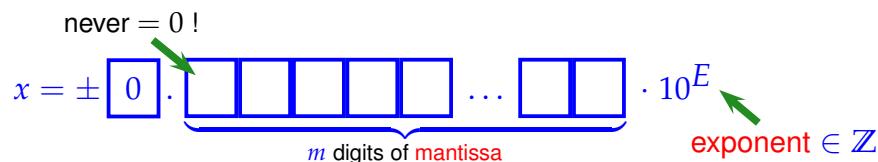
**§1.5.2.2 (Internal representation of machine numbers)** Now we give a brief sketch of the internal structure of machine numbers  $\in \mathbb{M}$ . The main insight will be that

“Computers use floating point numbers (scientific notation)”

**EXAMPLE 1.5.2.3 (Decimal floating point numbers)** Some 3-digit **normalized** decimal floating point numbers:

valid:  $0.723 \cdot 10^2$ ,  $0.100 \cdot 10^{-20}$ ,  $-0.801 \cdot 10^5$   
 invalid:  $0.033 \cdot 10^2$ ,  $1.333 \cdot 10^{-4}$ ,  $-0.002 \cdot 10^3$

General form of an  $m$ -digit **normalized decimal floating point number**:



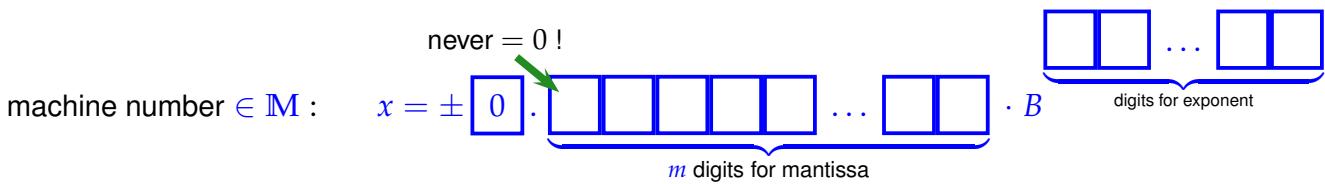
Of course, computers are restricted to a *finite range* of exponents:

**Definition 1.5.2.4. Machine numbers/floating point numbers** → [AG11, Sect. 2.2]

- Given
- ↳ basis  $B \in \mathbb{N} \setminus \{1\}$ ,
  - ↳ exponent range  $\{e_{\min}, \dots, e_{\max}\}$ ,  $e_{\min}, e_{\max} \in \mathbb{Z}$ ,  $e_{\min} < e_{\max}$ ,
  - ↳ number  $m \in \mathbb{N}$  of digits (for mantissa),

the corresponding set of machine numbers is

$$\mathbb{M} := \{d \cdot B^E : d = i \cdot B^{-m}, i = B^{m-1}, \dots, B^m - 1, E \in \{e_{\min}, \dots, e_{\max}\}\}$$



**Remark 1.5.2.5 (Extremal numbers in  $\mathbb{M}$ )** Clearly, there is a largest element of  $\mathbb{M}$  and two that are closest to zero. These are mainly determined by the range for the exponent  $E$ , cf. Def. 1.5.2.4.

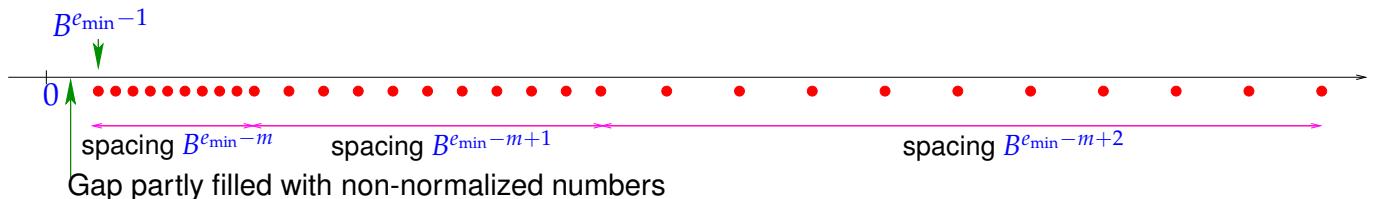
- Largest machine number (in modulus) :  $x_{\max} = \max |\mathbb{M}| = (1 - B^{-m}) \cdot B^{e_{\max}}$   
Smallest machine number (in modulus) :  $x_{\min} = \min |\mathbb{M}| = B^{-1} \cdot B^{e_{\min}}$

In C++ these extremal machine numbers are accessible through the

```
std::numeric_limits<double>::max()
and std::numeric_limits<double>::min()
```

functions. Other properties of arithmetic types can be queried accordingly from the `numeric_limits` header.

**Remark 1.5.2.6 (Distribution of machine numbers)** From Def. 1.5.2.4 it is clear that there are equi-spaced sections of  $\mathbb{M}$  and that the gaps between machine numbers are bigger for larger numbers, see also [AG11, Fig. 2.3].



Non-normalized numbers violate the lower bound for the mantissa  $i$  in Def. 1.5.2.4.

**§1.5.2.7 (IEEE standard 754 for machine numbers)** → [Ove01], [AG11, Sect. 2.4], → link) No surprise: for modern computers  $B = 2$  (binary system), the other parameters of the universally implemented machine number system are

single precision :  $m = 24^*, E \in \{-125, \dots, 128\} \geq 4$  bytes

double precision :  $m = 53^*, E \in \{-1021, \dots, 1024\} \geq 8$  bytes

\*: including bit indicating sign

The standardisation of machine numbers is important, because it ensures that the same numerical algorithm, executed on different computers will nevertheless produce the same result.  $\square$

#### Remark 1.5.2.8 (Special cases in IEEE standard)

The IEEE standard makes provisions for exceptions triggered by overflow or invalid floating point operations. The following C++ code snippet shows cases when these “flags” are raised.

```
double x = exp(1000), y = 3/x, z = x*sin(M_PI), w = x*log(1);
cout << x << endl << y << endl << z << endl << w << endl;
```

Output:

1	inf
2	0
3	inf
4	-nan



$E = e_{\max}, M \neq 0 \hat{=} \text{NaN} = \text{Not a number} \rightarrow \text{exception}$

$E = e_{\max}, M = 0 \hat{=} \text{Inf} = \text{Infinity} \rightarrow \text{overflow}$

$E = 0 \hat{=} \text{Non-normalized numbers} \rightarrow \text{underflow}$

$E = 0, M = 0 \hat{=} \text{number 0}$

In C++ these flags can be tested with the functions `std::isnan()` [C++ reference](#) and `std::isinf()` [C++ reference](#).  $\square$

#### §1.5.2.9 (Characteristic parameters of IEEE floating point numbers (double precision))

☞ C++ does not always fulfill the requirements of the IEEE 754 standard and it needs to be checked with `std::numeric_limits<T>::is_iec559`.

#### C++-code 1.5.2.10: Querying characteristics of double numbers → [GITLAB](#)

```
2 #include <climits>
3 #include <iostream>
4 #include <iomanip>
5
6 using namespace std;
7
8 int main() {
9     cout << std::numeric_limits<double>::is_iec559 << endl
10    << std::defaultfloat << numeric_limits<double>::min() << endl
11    << std::hexfloat << numeric_limits<double>::min() << endl
12    << std::defaultfloat << numeric_limits<double>::max() << endl
13    << std::hexfloat << numeric_limits<double>::max() << endl;
14 }
```

Output:

1	true
2	2.22507e-308
3	0010000000000000
4	1.79769e+308
5	7 f e f f f f f f f f f f f f

### 1.5.3 Roundoff Errors

**EXPERIMENT 1.5.3.1 (Input errors and roundoff errors)** The following computations would always result in 0, if done in exact arithmetic.

#### C++-code 1.5.3.2: Demonstration of roundoff errors → GITLAB

```

2 #include <iostream>
3 int main(){
4     std::cout.precision(15);
5     double a = 4.0/3.0, b = a-1, c = 3*b, e = 1-c;
6     std::cout << e << std::endl;
7     a = 1012.0/113.0; b = a-9; c = 113*b; e = 5+c;
8     std::cout << e << std::endl;
9     a = 83810206.0/6789.0; b = a-12345; c = 6789*b; e = c-1;
10    std::cout << e << std::endl;
11 }
```

Output:

---

1	2.22044604925031e-16
2	6.75015598972095e-14
3	-1.60798663273454e-09

---

Can you devise a similar calculation, whose result is even farther off zero? Apparently the rounding that inevitably accompanies arithmetic operations in  $\mathbb{M}$  can lead to results that are far away from the true result. □

For the discussion of errors introduced by rounding we need important notions.

#### Definition 1.5.3.3. Absolute and relative error → [AG11, Sect. 1.2]

Let  $\tilde{x} \in \mathbb{K}$  be an approximation of  $x \in \mathbb{K}$ . Then its **absolute error** is given by

$$\epsilon_{\text{abs}} := |x - \tilde{x}|,$$

and its **relative error** is defined as

$$\epsilon_{\text{rel}} := \frac{|x - \tilde{x}|}{|x|}.$$

**Remark 1.5.3.4 (Relative error and number of correct digits)** The **number of correct** (significant, valid) **digits** of an approximation  $\tilde{x}$  of  $x \in \mathbb{K}$  is defined through the relative error:

$$\text{If } \epsilon_{\text{rel}} := \frac{|x - \tilde{x}|}{|x|} \leq 10^{-\ell}, \text{ then } \tilde{x} \text{ has } \ell \text{ correct digits, } \ell \in \mathbb{N}_0$$

**§1.5.3.5 (Floating point operations)** We may think of the elementary binary operations  $+, -, *, /$  in  $\mathbb{M}$  comprising two steps:

- ① Compute the exact result of the operation.
- ② Perform rounding of the result of ① to map it back to  $\mathbb{M}$ .

**Definition 1.5.3.6. Correct rounding**

Correct rounding (“rounding up”) is given by the function

$$\text{rd} : \begin{cases} \mathbb{R} & \rightarrow \mathbb{M} \\ x & \mapsto \max \arg\min_{\tilde{x} \in \mathbb{M}} |x - \tilde{x}| \end{cases}$$

(Recall that  $\arg\min_x F(x)$  is the set of arguments of a real valued function  $F$  that makes it attain its (global) minimum.)

Of course, ① above is not possible in a strict sense, but the effect of both steps can be realised and yields a **floating point realization of  $\star \in \{+, -, \cdot, /\}$** .

☞ Notation: write  $\tilde{\star}$  for the floating point realization of  $\star \in \{+, -, \cdot, /\}$ :

Then ① and ② may be summed up into

$$\text{For } \star \in \{+, -, \cdot, /\}: \quad x \tilde{\star} y := \text{rd}(x \star y).$$

↓

**Remark 1.5.3.7 (Breakdown of associativity)** As a consequence of rounding addition  $\tilde{+}$  and multiplication  $\tilde{\cdot}$  as implemented on computers fail to be associative. They will usually be commutative, though this is not guaranteed. ↓

**§1.5.3.8 (Estimating roundoff errors → [AG11, p. 23])** Let us denote by **EPS** the largest **relative error** ( $\rightarrow$  Def. 1.5.3.3) incurred through rounding:

$$\text{EPS} := \max_{x \in I \setminus \{0\}} \frac{|\text{rd}(x) - x|}{|x|}, \quad (1.5.3.9)$$

where  $I = [\min |\mathbb{M}|, \max |\mathbb{M}|]$  is the range of positive machine numbers.

For machine numbers according to Def. 1.5.2.4 **EPS** can be computed from the defining parameters **B** (base) and **m** (length of mantissa) [AG11, p. 24]:

$$\text{EPS} = \frac{1}{2} B^{1-m}. \quad (1.5.3.10)$$

↓

However, when studying roundoff errors, we do not want to delve into the intricacies of the internal representation of machine numbers. This can be avoided by just using a single bound for the relative error due to rounding, and, thus, also for the relative error potentially suffered in each elementary operation.

**Assumption 1.5.3.11. “Axiom” of roundoff analysis**

There is a small positive number **EPS**, the **machine precision**, such that for the elementary arithmetic operations  $\star \in \{+, -, \cdot, /\}$  and “hard-wired” functions\*  $f \in \{\exp, \sin, \cos, \log, \dots\}$  holds

$$x \tilde{\star} y = (x \star y)(1 + \delta) \quad , \quad \tilde{f}(x) = f(x)(1 + \delta) \quad \forall x, y \in \mathbb{M},$$

with  $|\delta| < \text{EPS}$ .

\*: this is an ideal, which may not be accomplished even by modern CPUs.

- relative roundoff errors of elementary steps in a program bounded by machine precision !

**EXAMPLE 1.5.3.12 (Machine precision for IEEE standard)** C++ tells the machine precision as following:

**C++ code 1.5.3.13: Finding out EPS in C++ → GITLAB**

```
2 #include <iostream>
3 #include <climits> // get various properties of arithmetic types
4 int main() {
5     std::cout.precision(15);
6     std::cout << std::numeric_limits<double>::epsilon() << std::endl;
7 }
```

Output:

1 2.22044604925031e-16

Knowing the machine precision can be important for checking the validity of computations or coding termination conditions for iterative approximations.

**EXPERIMENT 1.5.3.14 (Adding EPS to 1)**

```
cout.precision(25);
double eps = numeric_limits<double>::epsilon();
cout << fixed << 1.0 + 0.5*eps << endl
    << 1.0 - 0.5*eps << endl
    << (1.0 + 2/eps) - 2/eps << endl;
```

Output:

1 1.0000000000000000000000000000000
2 0.999999999999998889776975
3 0.0000000000000000000000000000000

In fact, the following “definition” of EPS is sometimes used:

EPS is the smallest positive number  $\in \mathbb{M}$  for which  $1 + \tilde{\text{EPS}} \neq 1$  (in  $\mathbb{M}$ ):

We find that  $1 + \tilde{\text{EPS}} = 1$  actually complies with the “axiom” of roundoff error analysis, Ass. 1.5.3.11:

$$1 = (1 + \text{EPS})(1 + \delta) \Rightarrow |\delta| = \left| \frac{\text{EPS}}{1 + \text{EPS}} \right| < \text{EPS},$$

$$\frac{2}{\text{EPS}} = (1 + \frac{2}{\text{EPS}})(1 + \delta) \Rightarrow |\delta| = \left| \frac{\text{EPS}}{2 + \text{EPS}} \right| < \text{EPS}.$$

Do we have to worry about these tiny roundoff errors ?

**YES** ( $\rightarrow$  Exp. 1.5.1.5):

- accumulation of roundoff errors
- amplification of roundoff errors

**Remark 1.5.3.15 (Testing equality with zero)**



Since results of numerical computations are almost always polluted by roundoff errors:  
Tests like `if (x == 0)` are pointless and even dangerous, if  $x$  contains the result  
of a numerical computation.  
► Remedy: Test `if (abs(x) < eps*s) ...`,  
 $s \triangleq$  positive number, compared to which  $|x|$  should be small.

We saw a first example of this practise in Code 1.5.1.3, Line 13. □

**Remark 1.5.3.16 (Overflow and underflow)** Since the set of machine numbers  $\mathbb{M}$  is a finite set, the result of an arithmetic operation can lie outside the range covered by it. In this case we have to deal with

$$\begin{aligned}\text{overflow } &\triangleq |\text{result of an elementary operation}| > \max\{\mathbb{M}\} \\ &\triangleq \text{IEEE standard} \Rightarrow \text{Inf} \\ \text{underflow } &\triangleq 0 < |\text{result of an elementary operation}| < \min\{|\mathbb{M} \setminus \{0\}|\} \\ &\triangleq \text{IEEE standard} \Rightarrow \text{use non-normalized numbers (!)}\end{aligned}$$

The Axiom of roundoff analysis Ass. 1.5.3.11 does *not hold* once non-normalized numbers are encountered:

#### C++-code 1.5.3.17: Demonstration of over-/underflow → GITLAB

```

2 #include <iostream>
3 #define _USE_MATH_DEFINES
4 #include <cmath>
5 #include <limits>
6 using namespace std;
7 int main() {
8     cout.precision(15);
9     double min = numeric_limits<double>::min();
10    double res1 = M_PI*min/123456789101112;
11    double res2 = res1*123456789101112/min;
12    cout << res1 << endl << res2 << endl;
13 }
```

Output: 1 5.68175492717434e-322  
2 3.15248510554597



Try to avoid underflow and overflow

A simple example teaching how to avoid overflow during the computation of the norm of a 2D vector [AG11, Ex. 2.9]:

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

straightforward evaluation: overflow, when  $|x| > \sqrt{\max |\mathbb{M}|}$  or  $|y| > \sqrt{\max |\mathbb{M}|}$ .

$$r = \begin{cases} |x| \sqrt{1 + (y/x)^2} & , \text{if } |x| \geq |y| , \\ |y| \sqrt{1 + (x/y)^2} & , \text{if } |y| > |x| . \end{cases}$$

► no overflow!

#### 1.5.4 Cancellation

In general, predicting the impact of roundoff errors on the result of a multi-stage computation is very difficult, if possible at all. However, there is a constellation that is particularly prone to dangerous amplification of roundoff and still can be detected easily.

**EXAMPLE 1.5.4.1 (Computing the zeros of a quadratic polynomial)** The following simple EIGEN code computes the real roots of a quadratic polynomial  $p(\xi) = \xi^2 + \alpha\xi + \beta$  by the discriminant formula

$$p(\xi_1) = p(\xi_2) = 0, \quad \xi_{1/2} = \frac{1}{2}(\alpha \pm \sqrt{D}), \text{ if } D := \alpha^2 - 4\beta \geq 0. \quad (1.5.4.2)$$

**C++ code 1.5.4.3: Discriminant formula for the real roots of  $p(\xi) = \xi^2 + \alpha\xi + \beta \rightarrow \text{GITLAB}$**

```

2  ///! C++ function computing the zeros of a quadratic polynomial
3  ///!  $\xi \rightarrow \xi^2 + \alpha\xi + \beta$  by means of the familiar discriminant
4  ///! formula  $\xi_{1,2} = \frac{1}{2}(-\alpha \pm \sqrt{\alpha^2 - 4\beta})$ . However
5  ///! this implementation is vulnerable to round-off! The zeros are
6  ///! returned in a column vector
7  Vector2d zerosquadpol(double alpha, double beta) {
8      Vector2d z;
9      double D = std::pow(alpha, 2) - 4 * beta; // discriminant
10     if (D < 0)
11         throw "no real zeros";
12     else {
13         // The famous discriminant formula
14         double wD = std::sqrt(D);
15         z << (-alpha - wD) / 2, (-alpha + wD) / 2; //
16     }
17     return z;
18 }
```

This formula is applied to the quadratic polynomial  $p(\xi) = (\xi - \gamma)(\xi - \frac{1}{\gamma})$  after its coefficients  $\alpha, \beta$  have been computed from  $\gamma$ , which will have introduced *small* relative roundoff errors (of size  $\text{EPS}$ ).

**C++ code 1.5.4.4: Testing the accuracy of computed roots of a quadratic polynomial  
→ GITLAB**

```

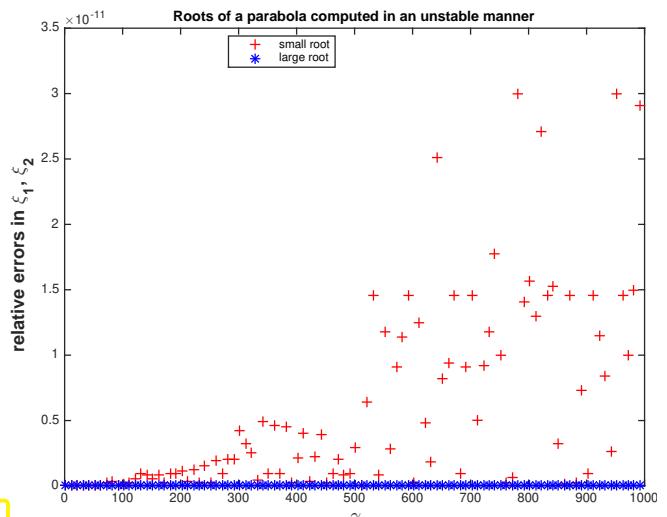
2  ///! Eigen Function for testing the computation of the zeros of a
3  ///! parabola
4  void compzeros() {
5      int n = 100;
6      MatrixXd res(n, 4);
7      VectorXd gamma = VectorXd::LinSpaced(n, 2, 992);
8      for (int i = 0; i < n; ++i) {
9          double alpha = -(gamma(i) + 1. / gamma(i));
10         double beta = 1.;
11         Vector2d z1 = zerosquadpol(alpha, beta);
12         Vector2d z2 = zerosquadpolstab(alpha, beta);
13         double ztrue = 1. / gamma(i), z2true = gamma(i);
14         res(i, 0) = gamma(i);
15         res(i, 1) = std::abs((z1(0) - ztrue) / ztrue);
16         res(i, 2) = std::abs((z2(0) - ztrue) / ztrue);
17         res(i, 3) = std::abs((z1(1) - z2true) / z2true);
    }
```

Plot of relative errors Def. 1.5.3.3  $\triangleright$

We observe that roundoff incurred during the computation of  $\alpha$  and  $\beta$  leads to “wrong” roots.

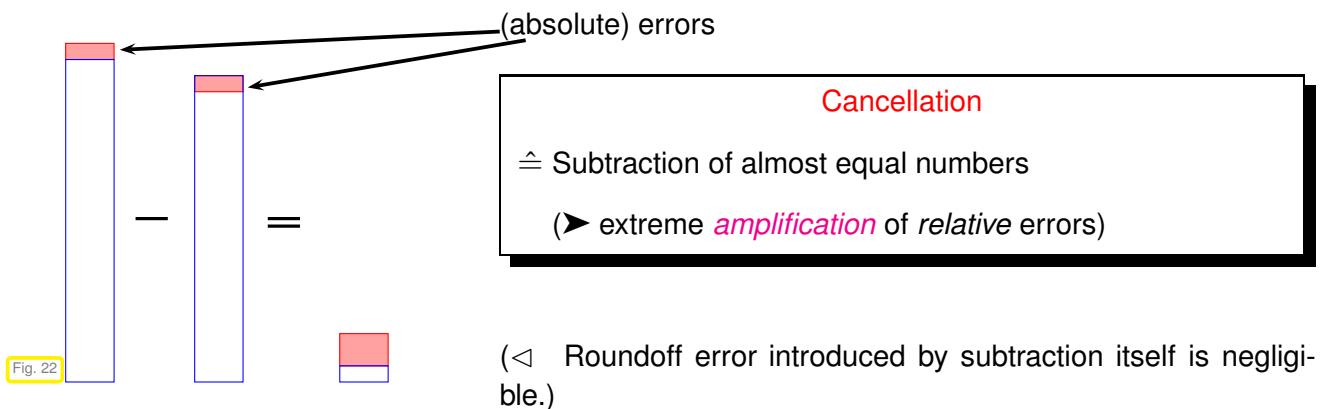
For large  $\gamma$  the computed small root may be fairly inaccurate as regards its *relative error*, which can be several orders of magnitude larger than machine precision  $\text{EPS}$ .

The large root always enjoys a small relative error about the size of  $\text{EPS}$ .



In order to understand why the small root is much more severely affected by roundoff, note that its computation involves the subtraction of two large numbers, if  $\gamma$  is large. This is the typical situation, in which **cancellation** occurs.  $\lrcorner$

**§1.5.4.5 (Visualisation of cancellation effect)** We look at the *exact* subtraction of two almost equal positive numbers both of which have small relative errors (red boxes) with respect to some desired exact value (indicated by blue boxes). The result of the subtraction will be small, but the errors may add up during the subtraction, ultimately constituting a large fraction of the result.



**EXAMPLE 1.5.4.6 (Cancellation in decimal system)** We consider two positive numbers  $x, y$  of about the same size afflicted with relative errors  $\approx 10^{-7}$ . This means that their seventh decimal digits are perturbed, here indicated by \*. When we subtract the two numbers the perturbed digits are shifted to the left, resulting in a possible relative error of  $\approx 10^{-3}$ :

$$\begin{array}{r}
 x = 0.123467* \\
 y = 0.123456* \\
 \hline
 x - y = 0.000011* = 0.11*000 \cdot 10^{-4}
 \end{array}$$

← 7th digit perturbed  
 ← 7th digit perturbed  
 ← 3rd digit perturbed  
 ↑  
 padded zeroes

Again, this example demonstrates that cancellation wreaks havoc through **error amplification**, not through the roundoff error due to the subtraction.  $\lrcorner$

**EXAMPLE 1.5.4.7 (Cancellation when evaluating difference quotients → [DR08, Sect. 8.2.6], [AG11, Ex. 1.3])** From analysis we know that the derivative of a differentiable function  $f : I \subset \mathbb{R} \rightarrow \mathbb{R}$  at a point

$x \in I$  is the limit of a difference quotient

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

This suggests the following approximation of the derivative by a difference quotient with small but finite  $h > 0$

$$f'(x) \approx \frac{f(x+h) - f(x)}{h} \quad \text{for } |h| \ll 1.$$

Results from analysis tell us that the approximation error should tend to zero for  $h \rightarrow 0$ . More precise quantitative information is provided by the Taylor formula for a twice continuously differentiable function [AG11, p. 5]

$$f(x+h) = f(x) + f'(x)h + \frac{1}{2}f''(\xi)h^2 \quad \text{for some } \xi = \xi(x, h) \in [\min\{x, x+h\}, \max\{x, x+h\}], \quad (1.5.4.8)$$

from which we infer

$$\frac{f(x+h) - f(x)}{h} - f'(x) = \frac{1}{2}hf''(\xi) \quad \text{for some } \xi = \xi(x, h) \in [\min\{x, x+h\}, \max\{x, x+h\}]. \quad (1.5.4.9)$$

We investigate the approximation of the derivative by difference quotients for  $f = \exp$ ,  $x = 0$ , and different values of  $h > 0$ :

**C++ code 1.5.4.10: Difference quotient approximation of the derivative of exp → GITLAB**

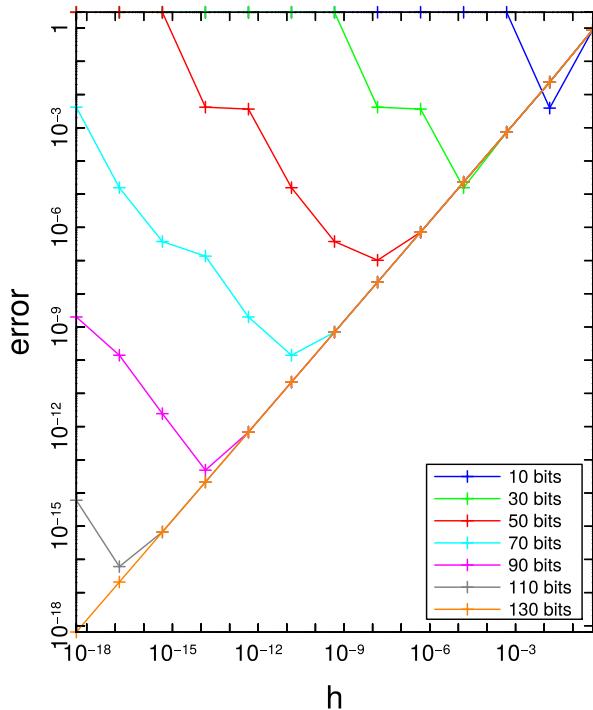
```

2 //! Difference quotient approximation
3 //! of the derivative of exp
4 void diffq () {
5     double h = 0.1, x = 0.0;
6     for (int i = 1; i <= 16; ++i) {
7         double df = (exp(x+h)-exp(x))/h;
8         cout << setprecision(14) << fixed ;
9         cout << setw(5) << -i
10        << setw(20) << df-1 << endl;
11        h /= 10;
12    }
13 }
```

Measured relative errors ▷

We observe an initial decrease of the relative approximation error followed by a steep increase when  $h$  drops below  $10^{-8}$ .

$\log_{10}(h)$	relative error
-1	0.05170918075648
-2	0.00501670841679
-3	0.00050016670838
-4	0.00005000166714
-5	0.00000500000696
-6	0.00000049996218
-7	0.00000004943368
-8	-0.00000000607747
-9	0.00000008274037
-10	0.00000008274037
-11	0.00000008274037
-12	0.00008890058234
-13	-0.00079927783736
-14	-0.00079927783736
-15	0.11022302462516
-16	-1.000000000000000



That the observed errors are really due to round-off errors is confirmed by the numerical results reported besides, using a variable precision floating point module of EIGEN, the [MPFR++ Support module](#).

The C++ used to generate these results can be found in [GITLAB](#).

Fig. 23

- Obvious culprit for what we see in Fig. 23: **cancellation** when computing the numerator of the difference quotient for small  $|h|$  leads to a strong amplification of inevitable errors introduced by the evaluation of the transcent exponential function.

We witness the competition of two opposite effects: Smaller  $h$  results in a better approximation of the derivative by the difference quotient, but the impact of cancellation is the stronger the smaller  $|h|$ .

$$\left. \begin{array}{l} \text{Approximation error } f'(x) - \frac{f(x+h) - f(x)}{h} \rightarrow 0 \\ \text{Impact of roundoff } \rightarrow \infty \end{array} \right\} \text{as } h \rightarrow 0 .$$

In order to provide a rigorous underpinning for our conjecture, in this example we embark on our first **roundoff error analysis** merely based on the “Axiom of roundoff analysis” Ass. 1.5.3.11: As in the computational example above we study the approximation of  $f'(x) = e^x$  for  $f = \exp, x \in \mathbb{R}$ .

$$\begin{aligned} df &= \frac{e^{x+h} (1 + \delta_1) - e^x (1 + \delta_2)}{h} && \text{correction factors take into account roundoff:} \\ &= e^x \left( \frac{e^h - 1}{h} + \frac{\delta_1 e^h - \delta_2}{h} \right) && (\rightarrow \text{"axiom of roundoff analysis"}, \text{ Ass. 1.5.3.11}) \\ \Rightarrow |df| &\leq e^x \left( \frac{e^h - 1}{h} + \text{eps} \frac{1+e^h}{h} \right) && |\delta_1|, |\delta_2| \leq \text{eps} . \end{aligned}$$

$1 + O(h)$        $O(h^{-1})$     for  $h \rightarrow 0$

(Note that the estimate for the term  $(e^h - 1)/h$  is a particular case of (1.5.4.9).)

► relative error:  $\left| \frac{e^x - df}{e^x} \right| \approx h + \frac{2\text{eps}}{h} \rightarrow \min \text{ for } h = \sqrt{2\text{eps}} .$

In double precision:  $\sqrt{2\text{eps}} = 2.107342425544702 \cdot 10^{-8}$

**Remark 1.5.4.11 (Cancellation during the computation of relative errors)** In the numerical experiment of Ex. 1.5.4.7 we computed the relative error of the result by subtraction, see Code 1.5.4.10. Of course, massive cancellation will occur! Do we have to worry?

In this case cancellation can be tolerated, because we are *interested only* in the **magnitude** of the relative error. Even if it was affected itself by a large relative error, this information is still not compromised.

For example, if the relative error has the exact value  $10^{-8}$ , but can be computed only with a huge relative error of 10%, then the perturbed value would still be in the range  $[0.9 \cdot 10^{-8}, 1.1 \cdot 10^{-8}]$ . Therefore it will still have the correct magnitude and still permit us to conclude the number of valid digits correctly.

**Remark 1.5.4.12 (Cancellation in Gram-Schmidt orthogonalisation of Exp. 1.5.1.5)** The **Hilbert matrix**  $\mathbf{A} \in \mathbb{R}^{10,10}$ ,  $(\mathbf{A})_{i,j} = (i+j-1)^{-1}$ , considered in Exp. 1.5.1.5 has columns that are almost linearly dependent.

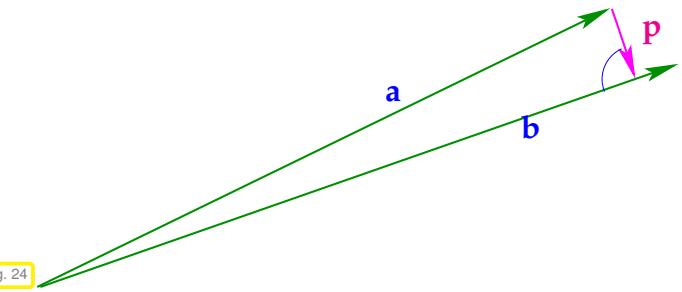
Cancellation when computing orthogonal projection of vector  $\mathbf{a}$  onto space spanned by vector  $\mathbf{b}$



$$\mathbf{p} = \mathbf{a} - \frac{\mathbf{a} \cdot \mathbf{b}}{\mathbf{b} \cdot \mathbf{b}} \mathbf{b}.$$

If  $\mathbf{a}$ ,  $\mathbf{b}$  point in almost the same direction,  $\|\mathbf{p}\| \ll \|\mathbf{a}\|, \|\mathbf{b}\|$ , so that a “tiny” vector  $\mathbf{p}$  is obtained by subtracting two “long” vectors, which implies **cancellation**.

Fig. 24



This can happen in Line 10 of Code 1.5.1.3.



**EXAMPLE 1.5.4.13 (Cancellation: roundoff error analysis)** We consider a simple arithmetic expression written in two ways:

$$a^2 - b^2 = (a+b)(a-b), \quad a, b \in \mathbb{R}.$$

We evaluate this term by means of two **algebraically equivalent** algorithms for the input data  $a = 1.3$ ,  $b = 1.2$  in 2-digit decimal arithmetic with standard rounding. (“Algebraically equivalent” means that two algorithms will produce the same results in the absence of roundoff errors.

Algorithm A	Algorithm B
$x := \tilde{a \cdot a} = 1.7$ (rounded)	$x := \tilde{a+b} = 2.5$ (exact)
$y := \tilde{b \cdot b} = 1.4$ (rounded)	$y := \tilde{a-b} = 0.1$ (exact)
$\tilde{x-y} = 0.30$ (exact)	$x * y = 0.25$ (exact)

Algorithm B produces the exact result, whereas Algorithm A fails to do so. Is this pure coincidence or an indication of the superiority of algorithm B? This question can be answered by **roundoff error analysis**. We demonstrate the approach for the two algorithms A & B and general input  $a, b \in \mathbb{R}$ .

Roundoff error analysis heavily relies on Ass. 1.5.3.11 and dropping terms of “higher order” in the machine precision, that is terms that behave like  $O(\text{EPS}^q)$ ,  $q > 1$ . It involves introducing the relative roundoff error for every elementary operation through a factor  $(1 + \delta)$ ,  $|\delta| \leq \text{EPS}$ .

**Algorithm A:**

$$x = a^2(1 + \delta_1), \quad y = b^2(1 + \delta_2)$$

$$\tilde{f} = (a^2(1 + \delta_1) - b^2(1 + \delta_2))(1 + \delta_3) = f + a^2\delta_1 - b^2\delta_2 + (a^2 - b^2)\delta_3 + O(\text{EPS}^2)$$

$$\frac{|\tilde{f} - f|}{|f|} \leq \text{EPS} \frac{a^2 + b^2 + |a^2 - b^2|}{|a^2 - b^2|} + O(\text{EPS}^2) = \text{EPS} \left( 1 + \frac{|a^2 + b^2|}{|a^2 - b^2|} \right) + O(\text{EPS}^2). \quad (1.5.4.14)$$

will be neglected

- For  $a \approx b$  the relative error of the result of Algorithm A will be much larger than the machine precision  $\text{EPS}$ . This reflects cancellation in the last subtraction step.

### Algorithm B:

$$x = (a + b)(1 + \delta_1), y = (a - b)(1 + \delta_2)$$

$$\tilde{f} = (a + b)(a - b)(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) = f + (a^2 - b^2)(\delta_1 + \delta_2 + \delta_3) + O(\text{EPS}^2)$$

$$\frac{|\tilde{f} - f|}{|f|} \leq |\delta_1 + \delta_2 + \delta_3| + O(\text{EPS}^2) \leq 3\text{EPS} + O(\text{EPS}^2). \quad (1.5.4.15)$$

- Relative error of the result of Algorithm B is always  $\approx \text{EPS}$  !

In this example we see a general guideline at work:

If inevitable, subtractions prone to cancellation should be done as early as possible.

The reason is that input data and initial intermediate results are usually not as much tainted by roundoff errors as numbers computed after many steps. □

**§1.5.4.16 (Avoiding disastrous cancellation)** The following examples demonstrate a few fundamental techniques for steering clear of cancellation by using alternative formulas that yield the same value (in exact arithmetic), but do not entail subtracting two numbers of almost equal size.

**EXAMPLE 1.5.4.17 (Stable discriminant formula → Ex. 1.5.4.1, [AG11, Ex. 2.10])** If  $\xi_1$  and  $\xi_2$  are the two roots of the quadratic polynomial  $p(\xi) = \xi^2 + \alpha\xi + \beta$ , then  $\xi_1 \cdot \xi_2 = \beta$  (Vieta's formula). Thus once we have computed a root, we can obtain the other by simple division.



Idea:

- ① Depending on the sign of  $\alpha$  compute “stable root” without cancellation.
- ② Compute other root from Vieta’s formula (avoiding subtraction)

### C++ code 1.5.4.18: Stable computation of real root of a quadratic polynomial → GITLAB

```

2 //! C++ function computing the zeros of a quadratic polynomial
3 //!  $\xi \rightarrow \xi^2 + \alpha\xi + \beta$  by means of the familiar discriminant
4 //! formula  $\xi_{1,2} = \frac{1}{2}(-\alpha \pm \sqrt{\alpha^2 - 4\beta})$ .
5 //! This is a stable implementation based on Vieta's theorem.
6 //! The zeros are returned in a column vector
7 VectorXd zerosquadrabpolstab(double alpha, double beta) {
8     Vector2d z(2);

```

```

9  double D = std::pow(alpha, 2) - 4 * beta; // discriminant
10 if (D < 0)
11     throw "no real zeros";
12 else {
13     double wD = std::sqrt(D);
14     // Use discriminant formula only for zero far away from 0
15     // in order to avoid cancellation. For the other zero
16     // use Vieta's formula.
17     if (alpha >= 0) {
18         double t = 0.5 * (-alpha - wD); //
19         z << t, beta / t;
20     } else {
21         double t = 0.5 * (-alpha + wD); //
22         z << beta / t, t;
23     }
24 }
25 return z;
26 }
```

→ Invariably, we add numbers with the same sign in Line 18 and Line 21.

Numerical experiment based on the driver code  
Code 1.5.4.4.

Observation:

The new code can also compute the small root of the polynomial  $p(\xi) = (\xi - \gamma)(\xi - \frac{1}{\gamma})$  (expanded in monomials) with a relative error  $\approx \text{EPS}$ .

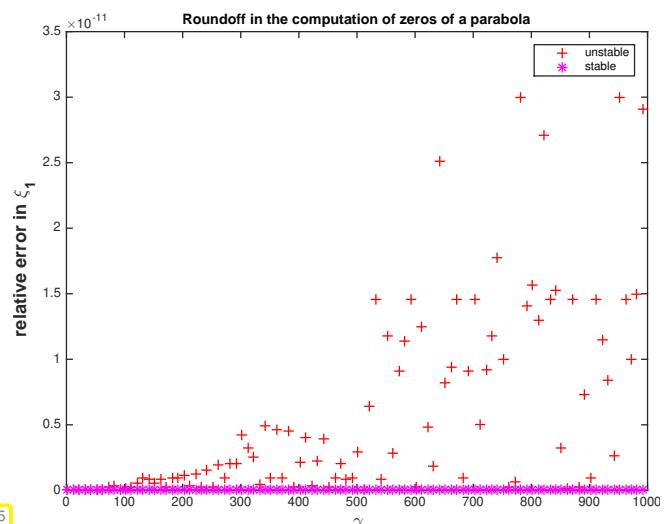


Fig. 25

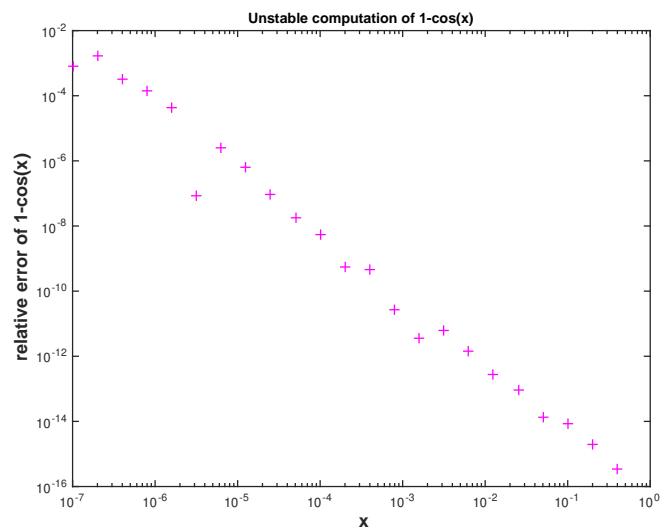
**EXAMPLE 1.5.4.19 (Exploiting trigonometric identities to avoid cancellation)** The task is to evaluate the integral

$$\int_0^x \sin t dt = \underbrace{1 - \cos x}_I = \underbrace{2 \sin^2(x/2)}_{II} \quad \text{for } 0 < x \ll 1, \quad (1.5.4.20)$$

and this can be done by the two different formulas  $I$  and  $II$ .

Relative error of expression  $I$  ( $1 - \cos(x)$ ) with respect to equivalent expression  $II$  ( $2 * \sin(x/2)^{2}$ )

▷ Expression  $I$  is affected by cancellation for  $|x| \ll 1$ , since then  $\cos x \approx 1$ , whereas expression  $II$  can be evaluated with a relative error  $\approx \text{EPS}$  for all  $x$ .



Analytic manipulations offer ample opportunity to rewrite expressions in equivalent form immune to cancellation.

#### EXAMPLE 1.5.4.21 (Switching to equivalent formulas to avoid cancellation)

Now we see an example of a computation allegedly dating back to Archimedes, who tried to approximate the area of a circle by the areas of inscribed regular polygons.

Approximation of circle by regular  $n$ -gon,  $n \in \mathbb{N}$

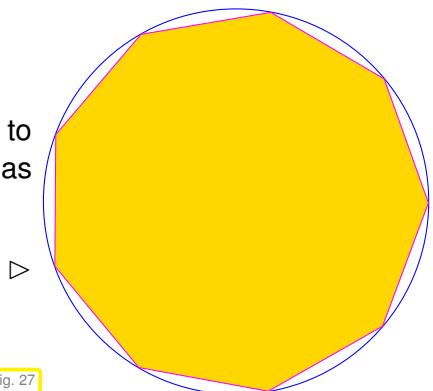


Fig. 27

Area of  $n$ -gon:

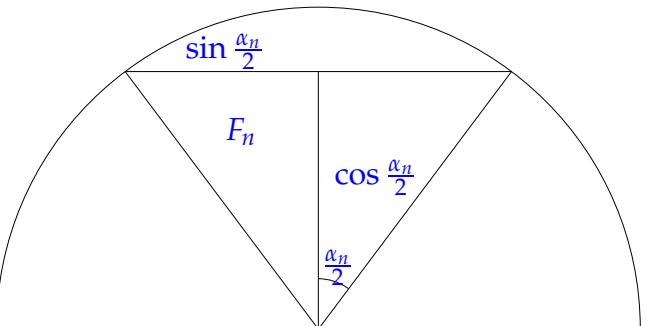
$$A_n = n \cos \frac{\alpha_n}{2} \sin \frac{\alpha_n}{2} = \frac{n}{2} \sin \alpha_n = \frac{n}{2} \sin\left(\frac{2\pi}{n}\right).$$

Recursion formula for  $A_n$  derived from

Fig. 28

$$\sin \frac{\alpha_n}{2} = \sqrt{\frac{1 - \cos \alpha_n}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}},$$

$$\text{Initial approximation: } A_6 = \frac{3}{2}\sqrt{3}.$$



#### C++ code 1.5.4.22: Tentative computation of circumference of regular polygon → [GITLAB](#)

```

2 //! Approximation of Pi by approximating the circumference of a
3 //! regular polygon
4 MatrixXd ApproxPiInstable(double tol = 1e-8, int maxIt = 50){
5     double s=sqrt(3)/2.; double An=3.*s;// initialization (hexagon case)

```

```

6  unsigned int n = 6, it = 0;
7  MatrixXd res(maxIt,4); // matrix for storing results
8  res(it,0) = n; res(it,1) = An;
9  res(it,2) = An - M_PI; res(it,3)=s;
10  while( it < maxIt && s > tol ){// terminate when s is 'small enough'
11      s = sqrt((1. - sqrt(1.-s*s))/2.); // recursion for area
12      n *= 2; An = n/2.*s;           // new estimate for circumference
13      ++it;
14      res(it,0) =n; res(it,1) =An;// store results and (absolute) error
15      res(it,2) = An - M_PI; res(it,3)=s;
16  }
17  return res.topRows(it);
18 }

```

The approximation deteriorates after applying the recursion formula many times:

$n$	$A_n$	$A_n - \pi$	$\sin \alpha_n$
6	2.598076211353316	-0.543516442236477	0.866025403784439
12	3.000000000000000	-0.141592653589794	0.500000000000000
24	3.105828541230250	-0.035764112359543	0.258819045102521
48	3.132628613281237	-0.008964040308556	0.130526192220052
96	3.139350203046872	-0.002242450542921	0.065403129230143
192	3.141031950890530	-0.000560702699263	0.032719082821776
384	3.141452472285344	-0.000140181304449	0.016361731626486
768	3.141557607911622	-0.000035045678171	0.008181139603937
1536	3.141583892148936	-0.000008761440857	0.004090604026236
3072	3.141590463236762	-0.000002190353031	0.002045306291170
6144	3.141592106043048	-0.000000547546745	0.001022653680353
12288	3.141592516588155	-0.000000137001638	0.000511326906997
24576	3.141592618640789	-0.000000034949004	0.000255663461803
49152	3.141592645321216	-0.000000008268577	0.000127831731987
98304	3.141592645321216	-0.000000008268577	0.000063915865994
196608	3.141592645321216	-0.000000008268577	0.000031957932997
393216	3.141592645321216	-0.000000008268577	0.000015978966498
786432	3.141593669849427	0.0000001016259634	0.000007989485855
1572864	3.141592303811738	-0.000000349778055	0.000003994741190
3145728	3.141608696224804	0.000016042635011	0.000001997381017
6291456	3.141586839655041	-0.0000005813934752	0.000000998683561
12582912	3.141674265021758	0.000081611431964	0.000000499355676
25165824	3.141674265021758	0.000081611431964	0.000000249677838
50331648	3.143072740170040	0.001480086580246	0.000000124894489
100663296	3.159806164941135	0.018213511351342	0.00000062779708
201326592	3.181980515339464	0.040387861749671	0.000000031610136
402653184	3.354101966249685	0.212509312659892	0.000000016660005
805306368	4.242640687119286	1.101048033529493	0.000000010536712
1610612736	6.000000000000000	2.858407346410207	0.000000007450581

Where does cancellation occur in Line 11 of Code 1.5.4.22? Since  $s \ll 1$ , computing  $1 - s$  will not trigger cancellation. However, the subtraction  $1 - \sqrt{1 - s^2}$  will, because  $\sqrt{1 - s^2} \approx 1$  for  $s \ll 1$ :

$$\sin \frac{\alpha_n}{2} = \sqrt{\frac{1 - \cos \alpha_n}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}}$$

For  $\alpha_n \ll 1$ :  $\sqrt{1 - \sin^2 \alpha_n} \approx 1$

Cancellation here!

We arrive at an *equivalent* formula not vulnerable to cancellation essentially using the identity  $(a + b)(a - b) = a^2 - b^2$  in order to eliminate the difference of square roots in the numerator.

$$\begin{aligned}\sin \frac{\alpha_n}{2} &= \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2} \cdot \frac{1 + \sqrt{1 - \sin^2 \alpha_n}}{1 + \sqrt{1 - \sin^2 \alpha_n}}} \\ &= \sqrt{\frac{1 - (1 - \sin^2 \alpha_n)}{2(1 + \sqrt{1 - \sin^2 \alpha_n})}} = \frac{\sin \alpha_n}{\sqrt{2(1 + \sqrt{1 - \sin^2 \alpha_n})}}.\end{aligned}$$

#### C++ code 1.5.4.23: Stable recursion for area of regular $n$ -gon [→ GITLAB](#)

```

2  //! Approximation of Pi by approximating the circumference of a
3  //! regular polygon
4  MatrixXd apprpistable(double tol = 1e-8, int maxIt = 50){
5      double s=sqrt(3)/2.; double An=3.*s;// initialization (hexagon case)
6      unsigned int n = 6, it = 0;
7      MatrixXd res(maxIt,4); // matrix for storing results
8      res(it,0) = n; res(it,1) = An;
9      res(it,2) = An - M_PI; res(it,3)=s;
10     while( it < maxIt && s > tol ){// terminate when s is 'small enough'
11         s = s/sqrt(2*(1+sqrt((1+s)*(1-s))));// Stable recursion without
12         cancellation
13         n *= 2; An = n/2.*s;           // new estimate for circumference
14         ++it;
15         res(it,0) =n; res(it,1) =An;// store results and (absolute) error
16         res(it,2) = An - M_PI; res(it,3)=s;
17     }
18     return res.topRows(it);
}

```

Using the stable recursion, we observe better approximation for polygons with more corners:

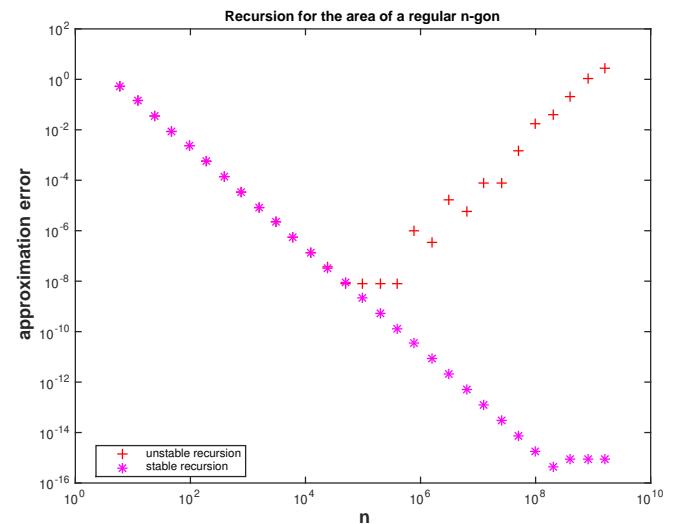
$n$	$A_n$	$A_n - \pi$	$\sin \alpha_n$
6	2.598076211353316	-0.543516442236477	0.866025403784439
12	3.000000000000000	-0.141592653589793	0.500000000000000
24	3.105828541230249	-0.035764112359544	0.258819045102521
48	3.132628613281238	-0.008964040308555	0.130526192220052
96	3.139350203046867	-0.002242450542926	0.065403129230143
192	3.141031950890509	-0.000560702699284	0.032719082821776
384	3.141452472285462	-0.000140181304332	0.016361731626487
768	3.141557607911857	-0.000035045677936	0.008181139603937
1536	3.141583892148318	-0.000008761441475	0.004090604026235
3072	3.141590463228050	-0.000002190361744	0.002045306291164
6144	3.141592105999271	-0.000000547590522	0.001022653680338
12288	3.141592516692156	-0.000000136897637	0.000511326907014
24576	3.141592619365383	-0.000000034224410	0.000255663461862
49152	3.141592645033690	-0.000000008556103	0.000127831731976
98304	3.141592651450766	-0.000000002139027	0.000063915866118
196608	3.141592653055036	-0.000000000534757	0.000031957933076
393216	3.141592653456104	-0.000000000133690	0.000015978966540
786432	3.141592653556371	-0.000000000033422	0.000007989483270
1572864	3.141592653581438	-0.000000000008355	0.000003994741635
3145728	3.141592653587705	-0.0000000000002089	0.000001997370818
6291456	3.141592653589271	-0.0000000000000522	0.000000998685409
12582912	3.141592653589663	-0.0000000000000130	0.000000499342704
25165824	3.141592653589761	-0.0000000000000032	0.000000249671352
50331648	3.141592653589786	-0.0000000000000008	0.000000124835676
100663296	3.141592653589791	-0.0000000000000002	0.000000062417838
201326592	3.141592653589794	0.0000000000000000	0.000000031208919
402653184	3.141592653589794	0.0000000000000001	0.000000015604460
805306368	3.141592653589794	0.0000000000000001	0.000000007802230
1610612736	3.141592653589794	0.0000000000000001	0.000000003901115

Plot of errors for approximations of  $\pi$  as computed by the two algebraically equivalent recursion formulas▷

Observation, cf. Ex. 1.5.4.7

Amplified roundoff errors due to cancellation supersedes approximation error for  $n \geq 10^5$ .

Roundoff errors merely of magnitude  $\text{EPS}$  in the case of stable recursion



### EXAMPLE 1.5.4.24 (Summation of exponential series)

In principle, the function value  $\exp(x)$  can be approximated up to any accuracy by summing sufficiently many terms of the globally convergent exponential series.

$$\begin{aligned}\exp(x) &= \sum_{k=0}^{\infty} \frac{x^k}{k!} \\ &= 1 + x + \frac{x^2}{2} + \frac{x^3}{6} + \frac{x^4}{24} + \dots\end{aligned}$$

#### C++ code 1.5.4.25: Summation of exponential series → [GITLAB](#)

```

2  double expeval(double x,
3      double tol=1e-8){
4      // Initialization
5      double y = 1.0, term = 1.0;
6      long int k = 1;
7      // Termination criterion
8      while(abs(term) > tol*y) {
9          term *= x/k;      // next summand
10         y += term; // Summation
11         ++k;
12     }
13     return y;
14 }
```

Results for  $\text{tol} = 10^{-8}$ ,  $\widetilde{\exp}$  designates the approximate value for  $\exp(x)$  returned by the function from Code 1.5.4.25. Rightmost column lists relative errors, which tells us the number of valid digits in the approximate result.

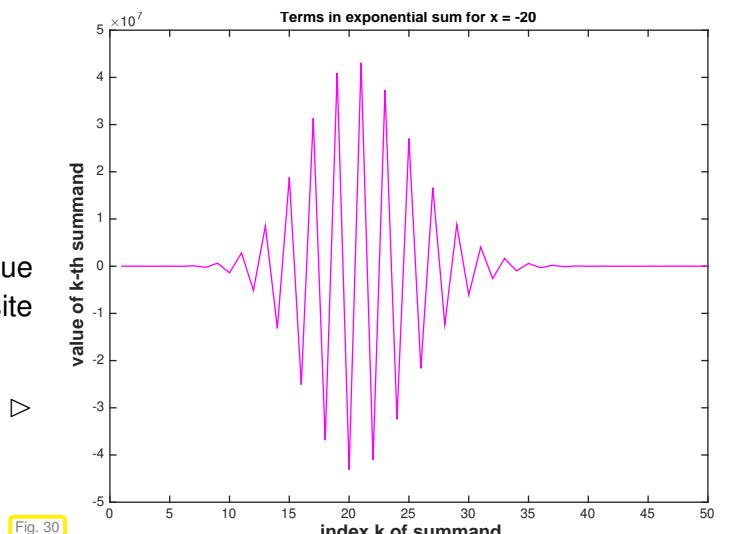
$x$	Approximation $\widetilde{\exp}(x)$	$\exp(x)$	$\frac{ \exp(x)-\widetilde{\exp}(x) }{\exp(x)}$
-20	6.1475618242e-09	2.0611536224e-09	1.982583033727893
-18	1.5983720359e-08	1.5229979745e-08	0.049490585500089
-16	1.1247503300e-07	1.1253517472e-07	0.000534425951530
-14	8.3154417874e-07	8.3152871910e-07	0.000018591829627
-12	6.1442105142e-06	6.1442123533e-06	0.000000299321453
-10	4.5399929604e-05	4.5399929762e-05	0.000000003501044
-8	3.3546262812e-04	3.3546262790e-04	0.000000000662004
-6	2.4787521758e-03	2.4787521767e-03	0.000000000332519
-4	1.8315638879e-02	1.8315638889e-02	0.000000000530724
-2	1.3533528320e-01	1.3533528324e-01	0.000000000273603
0	1.0000000000e+00	1.0000000000e+00	0.000000000000000
2	7.3890560954e+00	7.3890560989e+00	0.000000000479969
4	5.4598149928e+01	5.4598150033e+01	0.0000000001923058
6	4.0342879295e+02	4.0342879349e+02	0.0000000001344248
8	2.9809579808e+03	2.9809579870e+03	0.0000000002102584
10	2.2026465748e+04	2.2026465795e+04	0.0000000002143799
12	1.6275479114e+05	1.6275479142e+05	0.0000000001723845
14	1.2026042798e+06	1.2026042842e+06	0.0000000003634135
16	8.8861105010e+06	8.8861105205e+06	0.0000000002197990
18	6.5659968911e+07	6.5659969137e+07	0.0000000003450972
20	4.8516519307e+08	4.8516519541e+08	0.0000000004828737

Observation:

Large relative approximation errors for  $x \ll 0$ .

For  $x \ll 0$  we have  $\exp(x)| \ll 1$ , but this value is computed by summing large numbers of opposite sign.

Terms summed up for  $x = -20$



Remedy: Cancellation can be avoided by using identity



$$\exp(x) = \frac{1}{\exp(-x)} , \text{ if } x < 0 .$$

**EXAMPLE 1.5.4.26 (Trade cancellation for approximation)** In a computer code we have to provide a routine for the evaluation of the “hidden difference quotient”

$$I(a) := \int_0^1 e^{at} dt = \frac{\exp(a) - 1}{a} \quad \text{for any } a > 0 , \quad (1.5.4.27)$$

cf. the discussion of cancellation in the context of numerical differentiation in Ex. 1.5.4.7. There we observed massive cancellation.

**Trick.** Recall the **Taylor expansion** formula in one dimension for a function that is  $m + 1$  times continuously differentiable in a neighborhood of  $x$  [Str09, Satz 5.5.1]

$$f(x_0 + h) = \sum_{k=0}^m \frac{1}{k!} f^{(k)}(x_0) h^k + R_m(x_0, h) , \quad R_m(x_0, h) = \frac{1}{(m+1)!} f^{(m+1)}(\xi) h^{m+1} , \quad (1.5.4.28)$$

for some  $\xi \in [\min\{x_0, x_0 + h\}, \max\{x_0, x_0 + h\}]$ , and for all sufficiently small  $|h|$ . Here  $R(x_0, h)$  is called the **remainder** term and  $f^{(k)}$  denotes the  $k$ -th derivative of  $f$ .

Cancellation in (1.5.4.27) can be avoided by replacing  $\exp(a)$ ,  $a > 0$ , with a suitable Taylor expansion of  $a \mapsto e^a$  around  $a = 0$  and then dividing by  $a$ :

$$\frac{\exp(a) - 1}{a} = \sum_{k=0}^m \frac{1}{(k+1)!} a^k + R_m(a) , \quad R_m(a) = \frac{1}{(m+1)!} \exp(\xi) a^{m+1} \text{ for some } 0 \leq \xi \leq a .$$

Then use as an approximation the point value of the Taylor polynomial

$$I(a) \approx \tilde{I}_m(a) := \sum_{k=0}^m \frac{1}{(k+1)!} a^k .$$

For a similar discussion see [AG11, Ex. 2.12].

Issue: A finite Taylor sum usually offers only an **approximation** and we incur an approximation error. This begs the question how to choose the number  $m$  of terms to be retained in the Taylor expansion. We have to pick  $m$  large enough such that the relative approximation error remains below a prescribed threshold

`tol`. To estimate the relative approximation error, we use the expression for the remainder together with the simple estimate  $(\exp(a) - 1)/a > 1$  for all  $a > 0$ :

$$\begin{aligned} \text{rel. err.} &= \frac{|I(a) - \tilde{I}_m(a)|}{|I(a)|} = \frac{(e^a - 1)/a - \sum_{k=0}^m \frac{1}{(k+1)!} a^k}{(e^a - 1)/a} \\ &\leq \frac{1}{(m+1)!} \exp(\xi) a^m \leq \frac{1}{(m+1)!} \exp(a) a^m . \end{aligned}$$

For  $a = 10^{-3}$  we get

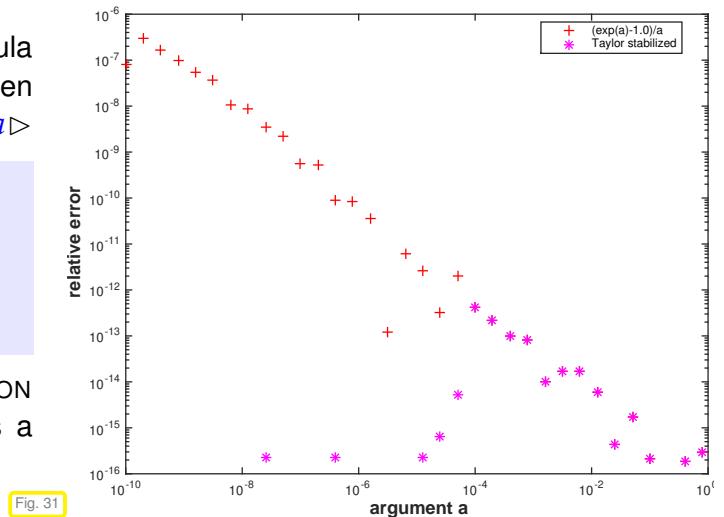
<i>m</i>	1	2	3	4	5
	1.0010e-03	5.0050e-07	1.6683e-10	4.1708e-14	8.3417e-18

Hence, keeping  $m = 3$  terms is enough for achieving about 10 valid digits.

Relative error of unstable formula  $(\exp(a) - 1.0)/a$  and relative error, when using a Taylor expansion approximation for small  $a \triangleright$

```
if (abs(a) < 1E-3)
    v = 1.0 + (1.0/2 + 1.0/6*a)*a;
else
    v = (exp(a)-1.0)/a;
end
```

Error computed by comparison with the PYTHON library function `numpy.expm1()` that provides a stable implementation of  $\exp(x) - 1$ .



**EXAMPLE 1.5.4.29 (Complex step differentiation [LM67])** This is a technique from complex analysis that can be applied to real-valued analytic functions. Let  $f : I \rightarrow \mathbb{R}$ ,  $I \subset \mathbb{R}$  an interval, be analytic in a neighborhood of  $x_0 \in I$ , which means that it can be written as a convergent power series there, see also Def. 6.1.2.32:

$$f(x) = \sum_{j=0}^{\infty} a_j (x - x_0)^j \quad \forall x : |x - x_0| < \rho \quad \text{and some } \rho > 0, a_j \in \mathbb{R} . \quad (1.5.4.30)$$

Note that  $f$  is infinitely many times differentiable in a neighborhood of  $x_0$  and that its derivatives satisfy  $f^{(n)}(x_0) = n! a_n \in \mathbb{R}$ ,  $n \in \mathbb{N}_0$ .

A power series like in (1.5.4.30) also makes sense for  $x \in \mathbb{C}$ ,  $|x - x_0| < \rho$ ! Thus, for  $0 < h < \rho$  we can approximate  $f$  in a neighborhood of  $x_0$  by means of a complex Taylor polynomial

$$f(x_0 + ih) = f(x_0) + f'(x_0)ih - f''(x_0)h^2 + O(h^3) \quad \text{for } h \in \mathbb{R} \rightarrow 0 . \quad (1.5.4.31)$$

**Trick.** Take the imaginary part on both sides of (1.5.4.31) using that all derivatives are real:

$$\Rightarrow f(x_0 + ih) = hf'(x_0) + O(h^3) \quad \text{for } h \in \mathbb{R} \rightarrow 0 .$$

As a consequence we obtain the approximation

$$f'(x_0) = \frac{\operatorname{Im} f(x_0 + ih)}{h} + O(h^2) \quad \text{for } h \in \mathbb{R} \rightarrow 0 ,$$

which suggests that we may rely on the *cancellation-free* expression

$$f'(x_0) \approx \frac{\operatorname{Im} f(x_0 + ih)}{h}$$

for  $h \approx \sqrt{\text{EPS}}$  to compute the derivative of  $f$  in  $x_0$ . □

### Review question(s) 1.5.4.32.

(Q1.5.4.32.A) Give an expression for

$$I(a, b) := \int_a^b \frac{1}{x} dx$$

that allows cancellation-free evaluation for integration bounds  $1 \ll a \approx b$ .

(Q1.5.4.32.B) For integration bounds  $a, b > 1$ ,  $a \approx b$ , propose a numerically sound way of computing

$$I(a, b) := \int_a^b \frac{1}{1+x^2} dx .$$

(Q1.5.4.32.C) What is the problem with the C++ expression

```
y = std::log(std::acosh(x));
```

$y$  where  $x$  is of type **double**? Rewrite this line of code into an algebraically equivalent one so that that problem does no longer occur.

△

## 1.5.5 Numerical Stability

We have seen that a particular “problem” can be tackled by different “algorithms”, which produce different results due to roundoff errors. This section will clarify what distinguishes a “good” algorithm from a rather abstract point of view.

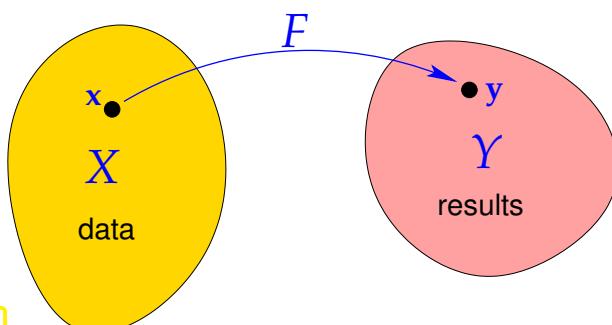
### §1.5.5.1 (The “problem”)

A mathematical notion of “problem”:

- ◆ data space  $X$ , usually  $X \subset \mathbb{R}^n$
- ◆ result space  $Y$ , usually  $Y \subset \mathbb{R}^m$
- ◆ mapping (problem function)  $F : X \mapsto Y$

A problem is a well defined *function* that assigns to each datum a result.

Fig. 32



Note: In this course, both the data space  $X$  and the result space  $Y$  will always be subsets of finite dimensional **vector spaces**.

**EXAMPLE 1.5.5.2 (The “matrix  $\times$  vector-multiplication problem”)** We consider the “problem” of computing the product  $\mathbf{Ax}$  for a given matrix  $\mathbf{A} \in \mathbb{K}^{m,n}$  and a given vector  $\mathbf{x} \in \mathbb{K}^n$ .

- • Data space  $X = \mathbb{K}^{m,n} \times \mathbb{K}^n$  (input is a matrix and a vector)
- Result space  $Y = \mathbb{K}^m$  (space of column vectors)

- Problem function  $F : X \rightarrow Y$ ,  $F(\mathbf{a}, \mathbf{x}) := \mathbf{Ax}$

↓

**§1.5.5.3 (Norms on spaces of vectors and matrices)** Norms provide tools for measuring errors. Recall from linear algebra and calculus [NS02, Sect. 4.3], [Gut09, Sect. 6.1]:

#### Definition 1.5.5.4. Norm

$X$  = vector space over field  $\mathbb{K}$ ,  $\mathbb{K} = \mathbb{C}, \mathbb{R}$ . A map  $\|\cdot\| : X \mapsto \mathbb{R}_0^+$  is a **norm** on  $X$ , if it satisfies

- $\forall \mathbf{x} \in X: \mathbf{x} \neq 0 \Leftrightarrow \|\mathbf{x}\| > 0$  (definite),
- $\|\lambda \mathbf{x}\| = |\lambda| \|\mathbf{x}\| \quad \forall \mathbf{x} \in X, \lambda \in \mathbb{K}$  (homogeneous),
- $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in X$  (triangle inequality).

Examples: (for vector space  $\mathbb{K}^n$ , vector  $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{K}^n$ )

name	:	definition	EIGEN function
Euclidean norm	:	$\ \mathbf{x}\ _2 := \sqrt{ x_1 ^2 + \dots +  x_n ^2}$	<code>x.norm()</code>
1-norm	:	$\ \mathbf{x}\ _1 :=  x_1  + \dots +  x_n $	<code>x.lpNorm&lt;1&gt;()</code>
$\infty$ -norm, max norm	:	$\ \mathbf{x}\ _\infty := \max\{ x_1 , \dots,  x_n \}$	<code>x.lpNorm&lt;Eigen::Infinity&gt;()</code>

**Remark 1.5.5.5 (Inequalities between vector norms)** All norms on the vector space  $\mathbb{K}^n$ ,  $n \in \mathbb{N}$ , are **equivalent** in the sense that for arbitrary two norms  $\|\cdot\|_1$  and  $\|\cdot\|_2$  we can always find a constant  $C > 0$  such that

$$\|\mathbf{v}\|_1 \leq C \|\mathbf{v}\|_2 \quad \forall \mathbf{v} \in \mathbb{K}^n. \quad (1.5.5.6)$$

Of course, the constant  $C$  will usually depend on  $n$  and the norms under consideration.

For the vector norms introduced above, explicit expressions for the constants “ $C$ ” are available: for all  $\mathbf{x} \in \mathbb{K}^n$

$$\|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1 \leq \sqrt{n} \|\mathbf{x}\|_2, \quad (1.5.5.7)$$

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \sqrt{n} \|\mathbf{x}\|_\infty, \quad (1.5.5.8)$$

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_1 \leq n \|\mathbf{x}\|_\infty. \quad (1.5.5.9)$$

The matrix space  $\mathbb{K}^{m,n}$  is a vector space, of course, and can also be equipped with various norms. Of particular importance are norms *induced by vector norms* on  $\mathbb{K}^n$  and  $\mathbb{K}^m$ .

#### Definition 1.5.5.10. Matrix norm

Given vector norms  $\|\cdot\|_x$  and  $\|\cdot\|_y$  on  $\mathbb{K}^n$  and  $\mathbb{K}^m$ , respectively, the associated **matrix norm** is defined by

$$\mathbf{M} \in \mathbb{R}^{m,n}: \quad \|\mathbf{M}\| := \sup_{\mathbf{x} \in \mathbb{R}^n \setminus \{0\}} \frac{\|\mathbf{M}\mathbf{x}\|_y}{\|\mathbf{x}\|_x}.$$

By virtue of definition the matrix norms enjoy an important property, they are **sub-multiplicative**:

$$\forall \mathbf{A} \in \mathbb{K}^{n,m}, \mathbf{B} \in \mathbb{K}^{m,k}: \quad \|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|. \quad (1.5.5.11)$$

☞ notations for matrix norms for *quadratic matrices* associated with standard vector norms:

$$\|\mathbf{x}\|_2 \rightarrow \|\mathbf{M}\|_2, \quad \|\mathbf{x}\|_1 \rightarrow \|\mathbf{M}\|_1, \quad \|\mathbf{x}\|_\infty \rightarrow \|\mathbf{M}\|_\infty$$

**EXAMPLE 1.5.5.12 (Matrix norm associated with  $\infty$ -norm and 1-norm)** Rather simple formulas are available for the matrix norms induced by the vector norms  $\|\cdot\|_\infty$  and  $\|\cdot\|_1$

$$\begin{aligned} \text{e.g. for } \mathbf{M} \in \mathbb{K}^{2,2}: \quad \|\mathbf{M}\|_\infty &= \max\{|m_{11}x_1 + m_{12}x_2|, |m_{21}x_1 + m_{22}x_2|\} \\ &\leq \max\{|m_{11}| + |m_{12}|, |m_{21}| + |m_{22}|\} \|\mathbf{x}\|_\infty, \\ \|\mathbf{M}\|_1 &= |m_{11}x_1 + m_{12}x_2| + |m_{21}x_1 + m_{22}x_2| \\ &\leq \max\{|m_{11}| + |m_{21}|, |m_{12}| + |m_{22}|\}(|x_1| + |x_2|). \end{aligned}$$

For general  $\mathbf{M} \in \mathbb{K}^{m,n}$

$$\geq \text{matrix norm} \leftrightarrow \|\cdot\|_\infty = \text{row sum norm} \quad \|\mathbf{M}\|_\infty := \max_{i=1,\dots,m} \sum_{j=1}^n |m_{ij}|, \quad (1.5.5.13)$$

$$\geq \text{matrix norm} \leftrightarrow \|\cdot\|_1 = \text{column sum norm} \quad \|\mathbf{M}\|_1 := \max_{j=1,\dots,n} \sum_{i=1}^m |m_{ij}|. \quad (1.5.5.14)$$

□

Sometimes special formulas for the Euclidean matrix norm come handy [GV89, Sect. 2.3.3]:

**Lemma 1.5.5.15. Formula for Euclidean norm of a Hermitian matrix**

$$\mathbf{A} \in \mathbb{K}^{n,n}, \mathbf{A} = \mathbf{A}^H \Rightarrow \|\mathbf{A}\|_2 = \max_{\mathbf{x} \neq 0} \frac{|\mathbf{x}^H \mathbf{A} \mathbf{x}|}{\|\mathbf{x}\|_2^2}.$$

*Proof.* Recall from linear algebra: Hermitian matrices (a special class of normal matrices) enjoy unitary similarity to diagonal matrices:

$$\exists \mathbf{U} \in \mathbb{K}^{n,n}, \text{ diagonal } \mathbf{D} \in \mathbb{R}^{n,n}: \quad \mathbf{U}^{-1} = \mathbf{U}^H \quad \text{and} \quad \mathbf{A} = \mathbf{U}^H \mathbf{D} \mathbf{U}.$$

Since multiplication with an unitary matrix preserves the 2-norm of a vector, we conclude

$$\|\mathbf{A}\|_2 = \left\| \mathbf{U}^H \mathbf{D} \mathbf{U} \right\|_2 = \|\mathbf{D}\|_2 = \max_{i=1,\dots,n} |d_i|, \quad \mathbf{D} = \text{diag}(d_1, \dots, d_n).$$

On the other hand, for the same reason:

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^H \mathbf{A} \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} (\mathbf{U} \mathbf{x})^H \mathbf{D} (\mathbf{U} \mathbf{x}) = \max_{\|\mathbf{y}\|_2=1} \mathbf{y}^H \mathbf{D} \mathbf{y} = \max_{i=1,\dots,n} |d_i|.$$

Hence, both expressions in the statement of the lemma agree with the largest modulus of eigenvalues of  $\mathbf{A}$ . □

**Corollary 1.5.5.16. Euclidean matrix norm and eigenvalues**

For  $\mathbf{A} \in \mathbb{K}^{m,n}$  the Euclidean matrix norm  $\|\mathbf{A}\|_2$  is the square root of the largest (in modulus) eigenvalue of  $\mathbf{A}^H \mathbf{A}$ .

For a *normal* matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$  (that is,  $\mathbf{A}$  satisfies  $\mathbf{A}^H \mathbf{A} = \mathbf{A} \mathbf{A}^H$ ) the Euclidean matrix norm agrees with the modulus of the largest eigenvalue.

**§1.5.5.17 ((Numerical) algorithm)** When we talk about an “algorithm” we have in mind a **concrete code function** in MATLAB or C++; the only way to describe an algorithm is through a piece of code. We assume that this function defines another mapping  $\tilde{F} : X \rightarrow Y$  on the data space of the problem. Of course, we can only feed data to the MATLAB/C++-function, if they can be represented in the set  $\mathbb{M}$  of machine numbers. Hence, implicit in the definition of  $\tilde{F}$  is the assumption that input data are subject to *rounding* before passing them to the code function proper.

$$\begin{array}{c|c} \text{Problem} & \text{Algorithm} \\ F : X \subset \mathbb{R}^n \rightarrow Y \subset \mathbb{R}^m & \tilde{F} : X \rightarrow \tilde{Y} \subset \mathbb{M} \end{array}$$

### §1.5.5.18 (Stable algorithm) → [AG11, Sect. 1.3]) [Stable algorithm]

- ◆ We study a problem ( $\rightarrow$  § 1.5.5.1)  $F : X \rightarrow Y$  on data space  $X$  into result space  $Y$ .
- ◆ We assume that both  $X$  and  $Y$  are equipped with norms  $\|\cdot\|_X$  and  $\|\cdot\|_Y$ , respectively ( $\rightarrow$  Def. 1.5.5.4).
- ◆ We consider a concrete algorithm  $\tilde{F} : X \rightarrow Y$  according to § 1.5.5.17.

We write  $w(\mathbf{x})$ ,  $\mathbf{x} \in X$ , for the **computational effort** ( $\rightarrow$  Def. 1.4.0.1, “number of elementary operations”) required by the algorithm for input  $\mathbf{x}$ .

#### Definition 1.5.5.19. Stable algorithm

An algorithm  $\tilde{F}$  for solving a problem  $F : X \mapsto Y$  is **numerically stable** if for all  $\mathbf{x} \in X$  its result  $\tilde{F}(\mathbf{x})$  (possibly affected by roundoff) is the exact result for “slightly perturbed” data:

$$\exists C \approx 1: \forall \mathbf{x} \in X: \exists \tilde{\mathbf{x}} \in X: \|\mathbf{x} - \tilde{\mathbf{x}}\|_X \leq C w(\mathbf{x}) \text{ EPS} \|\mathbf{x}\|_X \quad \wedge \quad \tilde{F}(\mathbf{x}) = F(\tilde{\mathbf{x}}).$$

Here **EPS** should be read as machine precision according to the “Axiom” of roundoff analysis Ass. 1.5.3.11.

Illustration of Def. 1.5.5.19  
( $\mathbf{y} \hat{=} \text{exact result for exact data } \mathbf{x}$ )

Terminology:

Def. 1.5.5.19 introduces stability in the sense of  
**backward error analysis**

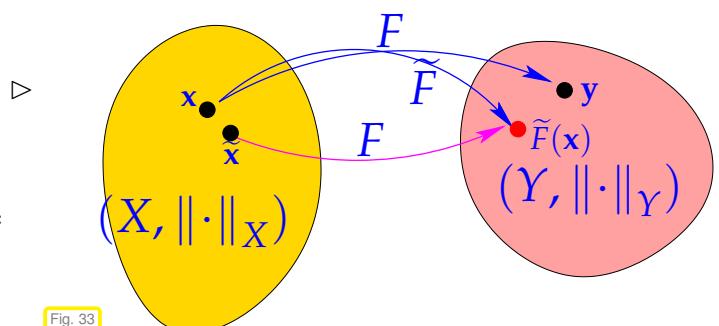


Fig. 33

Sloppily speaking, the impact of roundoff (\*) on a *stable algorithm* is of the same order of magnitude as the effect of the inevitable perturbations due to rounding the input data.

>

For stable algorithms roundoff errors are “harmless”.

(\*) In some cases the definition of  $\tilde{F}$  will also involve some approximations as in Ex. 1.5.4.26. Then the above statement also includes approximation errors.

**EXAMPLE 1.5.5.20 (Testing stability of matrix×vector multiplication)** Assume you are given a black box implementation of a function

```
VectorXd mvmult(const MatrixX &A, const VectorXd &x)
```

that purports to provide a stable implementation of  $\mathbf{Ax}$  for  $\mathbf{A} \in \mathbb{K}^{m,n}$ ,  $\mathbf{x} \in \mathbb{K}^n$ , cf. Ex. 1.5.5.2. How can we verify this claim for particular data. Both,  $\mathbb{K}^{m,n}$  and  $\mathbb{K}^n$  are equipped with the Euclidean norm.

The task is, given  $\mathbf{y} \in \mathbb{K}^n$  as returned by the function, to find conditions on  $\mathbf{y}$  that ensure the existence of a  $\tilde{\mathbf{A}} \in \mathbb{K}^{m,n}$  such that

$$\tilde{\mathbf{A}}\mathbf{x} = \mathbf{y} \quad \text{and} \quad \|\tilde{\mathbf{A}} - \mathbf{A}\|_2 \leq C_{mn} \text{EPS} \|\mathbf{A}\|_2, \quad (1.5.5.21)$$

for a small constant  $\approx 1$ .

In fact we can choose (easy computation)

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{z}\mathbf{x}^T, \quad \mathbf{z} := \frac{\mathbf{y} - \mathbf{Ax}}{\|\mathbf{x}\|_2^2} \in \mathbb{K}^m,$$

and we find

$$\|\tilde{\mathbf{A}} - \mathbf{A}\|_2 = \|\mathbf{z}\mathbf{x}^T\|_2 = \sup_{\mathbf{w} \in \mathbb{K}^n \setminus \{\mathbf{0}\}} \frac{\mathbf{x} \cdot \mathbf{w} \|\mathbf{z}\|_2}{\|\mathbf{w}\|_2} \leq \|\mathbf{x}\|_2 \|\mathbf{z}\|_2 = \frac{\|\mathbf{y} - \mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}.$$

Hence, in principle stability of an algorithm for computing  $\mathbf{Ax}$  is confirmed, if for every  $\mathbf{x} \in \mathbb{R}^n$  the computed result  $\mathbf{y} = \text{mvmult}(\mathbf{A}, \mathbf{x})$  satisfies

$$\|\mathbf{y} - \mathbf{Ax}\|_2 \leq C mn \text{EPS} \|\mathbf{x}\|_2 \|\mathbf{A}\|_2,$$

with a small constant  $C > 0$  independent of data and problem size.

**Remark 1.5.5.22 (Numerical stability and sensitive dependence on data)** A problem shows **sensitive dependence** on the data, if small perturbations of input data lead to large perturbations of the output. Such problems are also called **ill-conditioned**. For such problems stability of an algorithm is easily accomplished.

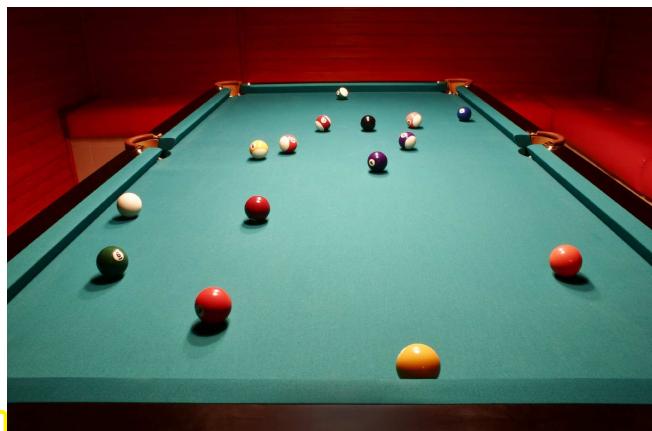


Fig. 34

Example: The problem is the prediction of the position of the billiard ball after ten bounces given the initial position, velocity, and spin.

It is well known, that tiny changes of the initial conditions can shift the final location of the ball to virtually any point on the table: the billiard problem is **chaotic**.

Hence, a stable algorithm for its solution may just output a *fixed* or *random* position without even using the initial conditions!

## Learning Outcomes

Principal take-home knowledge and skills from this chapter:

- Learning by doing: Knowledge about the syntax of fundamental operations on matrices and vectors in EIGEN.
- Understanding of the concepts of computational effort/cost and asymptotic complexity in numerics.
- Awareness of the asymptotic complexity of basic linear algebra operations
- Ability to determine the (asymptotic) computational effort for a concrete (numerical linear algebra) algorithm.
- Ability to manipulate simple expressions involving matrices and vectors in order to reduce the computational cost for their evaluation.
- Knowledge about round-off and machine precision.
- Familiarity with the phenomenon of “cancellation”: cause, effect, remedies, and tricks

# Bibliography

- [AV88] A. Aggarwal and J.S. Vitter. “The input/output complexity of sorting and related problems”. In: *Communications of the ACM* 31.9 (1988), pp. 1116–1127 (cit. on p. 65).
- [AG11] Uri M. Ascher and Chen Greif. *A first course in numerical methods*. Vol. 7. Computational Science & Engineering. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2011, pp. xxii+552. DOI: [10.1137/1.9780898719987](https://doi.org/10.1137/1.9780898719987) (cit. on pp. 39, 66, 77, 79, 80, 82, 84, 85, 88, 95, 100).
- [CW90] D. Coppersmith and S. Winograd. “Matrix multiplication via arithmetic progression”. In: *J. Symbolic Computing* 9.3 (1990), pp. 251–280 (cit. on p. 68).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 40, 84).
- [GV89] G.H. Golub and C.F. Van Loan. *Matrix computations*. 2nd. Baltimore, London: John Hopkins University Press, 1989 (cit. on pp. 61, 99).
- [GJ10] Gero Greiner and Riko Jacob. “The I/O Complexity of Sparse Matrix Dense Matrix Multiplication”. In: *LATIN 2010: THEORETICAL INFORMATICS*. Ed. by LopezOrtiz, A. Vol. 6034. Lecture Notes in Computer Science. Microsoft Res; Yahoo Res; Univ Waterloo. 2010, 143–156. DOI: [10.1007/978-3-642-12200-2\\_14](https://doi.org/10.1007/978-3-642-12200-2_14) (cit. on p. 65).
- [Gut09] M.H. Gutknecht. *Lineare Algebra*. Lecture Notes. SAM, ETH Zürich, 2009 (cit. on p. 98).
- [KW03] M. Kowarschik and C. Weiss. “An Overview of Cache Optimization Techniques and Cache-Aware Numerical Algorithms”. In: *Algorithms for Memory Hierarchies*. Vol. 2625. Lecture Notes in Computer Science. Heidelberg: Springer, 2003, pp. 213–232 (cit. on p. 65).
- [LM67] J. N. Lyness and C. B. Moler. “Numerical differentiation of analytic functions”. In: *SIAM J. Numer. Anal.* 4 (1967), pp. 202–210 (cit. on p. 96).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on pp. 36, 52, 73, 98).
- [Ove01] M.L. Overton. *Numerical Computing with IEEE Floating Point Arithmetic*. Philadelphia, PA: SIAM, 2001 (cit. on p. 77).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 40, 59).
- [Str69] V. Strassen. “Gaussian elimination is not optimal”. In: *Numer. Math.* 13 (1969), pp. 354–356 (cit. on p. 67).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on pp. 36, 37, 95).
- [Van00] Charles F. Van Loan. “The ubiquitous Kronecker product”. In: *J. Comput. Appl. Math.* 123.1-2 (2000), pp. 85–100. DOI: [10.1016/S0377-0427\(00\)00393-9](https://doi.org/10.1016/S0377-0427(00)00393-9).

# Chapter 2

## Direct Methods for (Square) Linear Systems of Equations

**§2.0.0.1 (Required prior knowledge for Chapter 2)** Also this chapter heavily relies on concepts and techniques from [linear algebra](#) as taught in the 1st semester introductory course. Knowledge of the following topics from linear algebra will be taken for granted and they should be refreshed in case of gaps:

- Operations involving matrices and vectors [[NS02](#), Ch. 2], already covered in Chapter 1
- Computations with block-structured matrices, *cf.* § 1.3.1.13
- Linear systems of equations: existence and uniqueness of solutions [[NS02](#), Sects. 1.2, 3.3]
- Gaussian elimination [[NS02](#), Ch. 2]
- LU-decomposition and its connection with Gaussian elimination [[NS02](#), Sect. 2.4]

### Contents

---

2.1	Preface . . . . .	105
2.2	Theory: Linear systems of equations . . . . .	107
2.2.1	Existence and uniqueness of solutions . . . . .	107
2.2.2	Sensitivity of linear systems . . . . .	108
2.3	Gaussian Elimination . . . . .	112
2.3.1	Basic algorithm . . . . .	112
2.3.2	LU-Decomposition . . . . .	120
2.3.3	Pivoting . . . . .	128
2.4	Stability of Gaussian Elimination . . . . .	134
2.5	Survey: Elimination solvers for linear systems of equations . . . . .	141
2.6	Exploiting Structure when Solving Linear Systems . . . . .	145
2.7	Sparse Linear Systems . . . . .	150
2.7.1	Sparse matrix storage formats . . . . .	152
2.7.2	Sparse matrices in EIGEN . . . . .	153
2.7.3	Direct Solution of Sparse Linear Systems of Equations . . . . .	160
2.7.4	LU-factorization of sparse matrices . . . . .	163
2.7.5	Banded matrices [ <a href="#">DR08</a> , Sect. 3.7] . . . . .	169
2.8	Stable Gaussian elimination without pivoting . . . . .	175

---

## 2.1 Preface

**§2.1.0.1 (The problem: solving a linear system)** What is “the problem” considered in this chapter, when we apply the notion of “Problem” introduced in § 1.5.5.1, that is, which functions do “Direct Methods for Linear Systems of Equations (LSE)” attempt to evaluate and what suitable data spaces  $X$  and result spaces  $Y$ ?

Input/data : square matrix  $A \in \mathbb{K}^{n,n}$ , vector  $b \in \mathbb{K}^n$ ,  $n \in \mathbb{N}$   $\geq$  data space  $X = \mathbb{K}^{n,n} \times \mathbb{K}^n$   
 Output/result : solution vector  $x \in \mathbb{K}^n$ :  $\boxed{Ax = b} \leftarrow$  (square) linear system of equations (LSE)  
 $\geq$  result space  $Y = \mathbb{K}^n$

(Terminology:  $A \hat{=} \text{system matrix/coefficient matrix}$ ,  $b \hat{=} \text{right hand side vector}$  )

Linear systems with rectangular system matrices  $A \in \mathbb{K}^{m,n}$ , called “overdetermined” for  $m > n$ , and “underdetermined” for  $m < n$  will be treated in Chapter 3.  $\square$

**Remark 2.1.0.2 (LSE: key components of mathematical models in many fields)** Linear systems of equations are ubiquitous in computational science: they are encountered

- with discrete linear models in network theory (see Ex. 2.1.0.3), control, statistics;
- in the case of *discretized* boundary value problems for ordinary and partial differential equations ( $\rightarrow$  course “Numerical methods for partial differential equations”, 4th semester);
- as a result of linearization (e.g., “Newton’s method”  $\rightarrow$  Sect. 8.4).

### EXAMPLE 2.1.0.3 (Nodal analysis of (linear) electric circuit [QSS00, Sect. 4.7.1])

Now we study a very important application of numerical simulation, where (large, sparse) linear systems of equations play a central role: **Numerical circuit analysis**. We begin with *linear circuits* in the *frequency domain*, which are directly modelled by complex linear systems of equations. In later chapters we will tackle circuits with non-linear elements, see Ex. 8.0.0.1, and, finally, will learn about numerical methods for computing the transient (time-dependent) behavior of circuits, see Ex. 11.1.1.9.

Modeling of simple linear circuits takes only elementary physical laws as covered in any introductory course of physics (or even in secondary school physics). There is no sophisticated physics or mathematics involved. Circuits are composed of so-called circuit elements connected by (ideal) wires.

A circuit diagram  $\Rightarrow$

•: Nodes (ger.: Knoten) that is, junctions of wires

We number the nodes  $1, \dots, n$  and write  $I_{kj}$  (physical units,  $[I_{kj}] = 1\text{A}$ ) for the electric current flowing from node  $k$  to node  $j$ . Currents have a sign:

$$I_{kj} = -I_{jk}$$

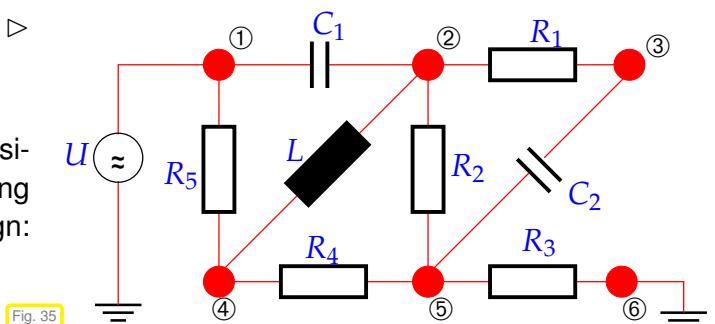


Fig. 35

The most fundamental relationship is the **Kirchhoff current law (KCL)** that demands that the sum of node currents vanishes:

$$\forall k \in \{1, \dots, n\}: \quad \sum_{j=1}^n I_{kj} = 0 . \quad (2.1.0.4)$$

The unknowns of the model are the **nodal potentials**  $U_k, k = 1, \dots, n$ . (Some of them may be known, for instance those for grounded nodes: ⑥ in Fig. 35, and nodes connected to voltage sources: ① in Fig. 35.)

The circuit elements are characterized by current-voltage relationships, so-called constitutive relations, here given in **frequency domain** for angular frequency  $\omega > 0$  (physical units  $[\omega] = 1\text{s}^{-1}$ ). We consider only the following simple circuit elements:

- Ohmic resistor:  $I = \frac{U}{R}, [R] = 1\text{VA}^{-1}$
- capacitor:  $I = \imath\omega C U$ , capacitance  $[C] = 1\text{AsV}^{-1}$  ➤  $I_{kj} = \begin{cases} R^{-1}(U_k - U_j) , \\ \imath\omega C(U_k - U_j) , \\ -\imath\omega^{-1}L^{-1}(U_k - U_j) . \end{cases}$
- coil/inductor :  $I = \frac{U}{\imath\omega L}$ , inductance  $[L] = 1\text{VsA}^{-1}$

☞ notation:  $\imath \triangleq \text{imaginary unit } \imath := \sqrt{-1}$ ,  $\imath = \exp(\imath\pi/2)$ ,  $\imath^2 = -1$

Here we face the special case of a **linear circuit**: all relationships between branch currents and voltages are of the form

$$I_{kj} = \alpha_{kj}(U_k - U_j) \quad \text{with} \quad \alpha_{kj} \in \mathbb{C} . \quad (2.1.0.5)$$

The concrete value of  $\alpha_{kj}$  is determined by the circuit element connecting node  $k$  and node  $j$ .

These constitutive relations are derived by assuming a harmonic time-dependence of all quantities, which is termed circuit analysis in the **frequency domain** (AC-mode).

$$\text{voltage: } u(t) = \text{Re}\{U \exp(\imath\omega t)\} , \quad \text{current: } i(t) = \text{Re}\{I \exp(\imath\omega t)\} . \quad (2.1.0.6)$$

Here  $U, I \in \mathbb{C}$  are called complex amplitudes. This implies for temporal derivatives (denoted by a dot):

$$\frac{du}{dt}(t) = \text{Re}\{\imath\omega U \exp(\imath\omega t)\} , \quad \frac{di}{dt}(t) = \text{Re}\{\imath\omega I \exp(\imath\omega t)\} . \quad (2.1.0.7)$$

For a capacitor the total charge is proportional to the applied voltage:

$$q(t) = Cu(t) \quad \Rightarrow \quad i(t) = \frac{dq}{dt}(t) = C \frac{du}{dt}(t) .$$

For a coil the voltage is proportional to the rate of change of current:  $u(t) = L \frac{di}{dt}(t)$ . Combined with (2.1.0.6) and (2.1.0.7) this leads to the above constitutive relations.

Now we combine the constitutive relations with the Kirchhoff current law (2.1.0.4). We end up with a linear system of equations!

$$\begin{aligned} ②: \quad \imath\omega C_1(U_2 - U_1) + R_1^{-1}(U_2 - U_3) - \imath\omega^{-1}L^{-1}(U_2 - U_4) + R_2^{-1}(U_2 - U_5) &= 0 , \\ ③: \quad &R_1^{-1}(U_3 - U_2) + \imath\omega C_2(U_3 - U_5) = 0 , \\ ④: \quad &R_5^{-1}(U_4 - U_1) - \imath\omega^{-1}L^{-1}(U_4 - U_2) + R_4^{-1}(U_4 - U_5) = 0 , \\ ⑤: \quad &R_2^{-1}(U_5 - U_2) + \imath\omega C_2(U_5 - U_3) + R_4^{-1}(U_5 - U_4) + R_3^{-1}(U_5 - U_6) = 0 , \end{aligned}$$

$$U_1 = U \quad , \quad U_6 = 0 .$$

We do not get equations for the nodes ① and ⑥, because these nodes are connected to the “outside world” so that the Kirchhoff current law (2.1.0.4) does not hold (from a local perspective). This is fitting, because the voltages in these nodes are known anyway.

$$\begin{pmatrix} i\omega C_1 + \frac{1}{R_1} - \frac{i}{\omega L} + \frac{1}{R_2} & -\frac{1}{R_1} & \frac{i}{\omega L} & -\frac{1}{R_2} \\ -\frac{1}{R_1} & \frac{1}{R_1} + i\omega C_2 & 0 & -i\omega C_2 \\ \frac{i}{\omega L} & 0 & \frac{1}{R_5} - \frac{i}{\omega L} + \frac{1}{R_4} & -\frac{1}{R_4} \\ -\frac{1}{R_2} & -i\omega C_2 & -\frac{1}{R_4} & \frac{1}{R_2} + i\omega C_2 + \frac{1}{R_4} + R_3^{-1} \end{pmatrix} \begin{pmatrix} U_2 \\ U_3 \\ U_4 \\ U_5 \end{pmatrix} = \begin{pmatrix} i\omega C_1 U \\ 0 \\ \frac{1}{R_5} U \\ 0 \end{pmatrix}$$

This is a linear system of equations with *complex* coefficients:  $\mathbf{A} \in \mathbb{C}^{4,4}$ ,  $\mathbf{b} \in \mathbb{C}^4$ . For the algorithms to be discussed below this does not matter, because they work alike for real and complex numbers.  $\square$

## 2.2 Theory: Linear systems of equations

### 2.2.1 Existence and uniqueness of solutions

The following concepts and results are known from linear algebra [NS02, Sect. 1.2], [Gut09, Sect. 1.3]:

**Definition 2.2.1.1. Invertible matrix** → [NS02, Sect. 2.3]

$$\mathbf{A} \in \mathbb{K}^{n,n} \quad \text{invertible/regular} \quad :\Leftrightarrow \quad \exists_1 \mathbf{B} \in \mathbb{K}^{n,n}: \quad \mathbf{AB} = \mathbf{BA} = \mathbf{I} .$$

$\mathbf{B}$  is called the **inverse** of  $\mathbf{A}$ , (☞ notation  $\mathbf{B} = \mathbf{A}^{-1}$ )

Now, recall a few notions from linear algebra needed to state criteria for the invertibility of a matrix.

**Definition 2.2.1.2. Image space and kernel of a matrix**

Given  $\mathbf{A} \in \mathbb{K}^{m,n}$ , the **range/image** (space) of  $\mathbf{A}$  is the subspace of  $\mathbb{K}^m$  spanned by the columns of  $\mathbf{A}$

$$\mathcal{R}(\mathbf{A}) := \{\mathbf{Ax}, \mathbf{x} \in \mathbb{K}^n\} \subset \mathbb{R}^m .$$

The **kernel/nullspace** of  $\mathbf{A}$  is

$$\mathcal{N}(\mathbf{A}) := \{\mathbf{z} \in \mathbb{R}^n : \mathbf{Az} = \mathbf{0}\} .$$

**Definition 2.2.1.3. Rank of a matrix** → [NS02, Sect. 2.4], [QSS00, Sect. 1.5]

The **rank** of a matrix  $\mathbf{M} \in \mathbb{K}^{m,n}$ , denoted by  $\text{rank}(\mathbf{M})$ , is the maximal number of linearly independent rows/columns of  $\mathbf{M}$ . Equivalently,  $\text{rank}(\mathbf{A}) = \dim \mathcal{R}(\mathbf{A})$ .

**Theorem 2.2.1.4. Criteria for invertibility of matrix** → [NS02, Sect. 2.3 & Cor. 3.8]

A square matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$  is *invertible/regular* if one of the following equivalent conditions is satisfied:

1.  $\exists \mathbf{B} \in \mathbb{K}^{n,n}: \mathbf{BA} = \mathbf{AB} = \mathbf{I}$ ,
2.  $\mathbf{x} \mapsto \mathbf{Ax}$  defines an endomorphism of  $\mathbb{K}^n$ ,
3. the columns of  $\mathbf{A}$  are linearly independent (full column rank),
4. the rows of  $\mathbf{A}$  are linearly independent (full row rank),
5.  $\det \mathbf{A} \neq 0$  (non-vanishing determinant),
6.  $\text{rank}(\mathbf{A}) = n$  (full rank).

**§2.2.1.5 (Solution of a LSE as a “problem”, recall § 2.1.0.1)** Linear algebra give us a *formal* way to denote solution of LSE:

$$\mathbf{A} \in \mathbb{K}^{n,n} \text{ regular} \quad \& \quad \mathbf{Ax} = \mathbf{b} \Rightarrow \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}.$$

↗  
inverse matrix

Now recall our notion of “problem” from § 1.5.5.1 as a function  $F$  mapping data in a data space  $X$  to a result in a result space  $Y$ . Concretely, for  $n \times n$  linear systems of equations:

$$F : \begin{cases} X := \mathbb{K}_{*}^{n,n} \times \mathbb{K}^n & \rightarrow Y := \mathbb{K}^n \\ (\mathbf{A}, \mathbf{b}) & \mapsto \mathbf{A}^{-1}\mathbf{b} \end{cases}$$

☞ notation: (open) set of regular matrices  $\subset \mathbb{K}^{n,n}$ :

$$\mathbb{K}_{*}^{n,n} := \{\mathbf{A} \in \mathbb{K}^{n,n} : \mathbf{A} \text{ regular/invertible} \rightarrow \text{Def. 2.2.1.1}\}.$$

**Remark 2.2.1.6 (The inverse matrix and solution of a LSE)** In principle, in EIGEN the inverse of a matrix  $\mathbf{A}$  is available through the member function `inverse()` of matrix type, see  [EIGEN documentation](#). However, there are only a few case that always involve fixed-size small matrices, where the actual computation of the inverse of a matrix is warranted. The general advice is the following:

Avoid computing the inverse of a matrix (which can almost always be avoided)!

 In particular, never ever even contemplate using  $\mathbf{x} = \mathbf{A}.\text{inverse}() * \mathbf{b}$  to solve the linear system of equations  $\mathbf{Ax} = \mathbf{b}$ , cf. Exp. 2.4.0.13. The next sections present a sound way to do this.  
Another reason for this advice is given in Exp. 2.4.0.13.

## 2.2.2 Sensitivity of linear systems

The **sensitivity** of a problem (for given data) gauges the impact of small perturbations of the data on the result.

Before we examine sensitivity for linear systems of equations, we look at the simpler problem of matrix  $\times$  vector multiplication.

**EXAMPLE 2.2.2.1 (Sensitivity of linear mappings)** For a *fixed* given **regular**  $\mathbf{A} \in \mathbb{K}^{n,n}$  we study the problem map

$$F : \mathbb{K}^n \rightarrow \mathbb{K}^n , \quad \mathbf{x} \mapsto \mathbf{Ax} ,$$

that is, now we consider only the vector  $\mathbf{x}$  as data.

Goal: Estimate relative perturbations in  $F(\mathbf{x})$  due to relative perturbations in  $\mathbf{x}$ .

We assume that  $\mathbb{K}^n$  is equipped with *some* vector norm ( $\rightarrow$  Def. 1.5.5.4) and we use the induced matrix norm ( $\rightarrow$  Def. 1.5.5.10) on  $\mathbb{K}^{n,n}$ . Using linearity and the elementary estimate  $\|\mathbf{M}\mathbf{x}\| \leq \|\mathbf{M}\| \|\mathbf{x}\|$ , which is a direct consequence of the definition of an induced matrix norm, we obtain

$$\begin{aligned} \mathbf{Ax} = \mathbf{y} &\Rightarrow \|\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{y}\| \\ \mathbf{A}(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{y} + \Delta\mathbf{y} &\Rightarrow \mathbf{A}\Delta\mathbf{x} = \Delta\mathbf{y} \Rightarrow \|\Delta\mathbf{y}\| \leq \|\mathbf{A}\| \|\Delta\mathbf{x}\| \\ \Rightarrow \frac{\|\Delta\mathbf{y}\|}{\|\mathbf{y}\|} &\leq \frac{\|\mathbf{A}\| \|\Delta\mathbf{x}\|}{\|\mathbf{A}^{-1}\|^{-1} \|\mathbf{x}\|} = \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|}. \end{aligned} \quad (2.2.2.2)$$

relative perturbation in result    relative perturbation in data

We have found that the quantity  $\|\mathbf{A}\| \|\mathbf{A}^{-1}\|$  bounds amplification of relative errors in the argument vector in a matrix  $\times$  vector-multiplication with the matrix  $\mathbf{A}$ . □

Now we study the sensitivity of the problem of finding the solution of a linear system of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{n,n}$  regular,  $\mathbf{b} \in \mathbb{R}$ , see § 2.1.0.1. We write  $\tilde{\mathbf{x}}$  for the solution of the perturbed linear system.

Question: To what extent do perturbations in the data  $\mathbf{A}$ ,  $\mathbf{b}$  cause a

$$(\text{normwise}) \quad \text{relative error:} \quad \epsilon_r := \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|}.$$

( $\|\cdot\| \triangleq$  suitable vector norm, e.g., maximum norm  $\|\cdot\|_\infty$ )

Perturbed linear system:

$$\mathbf{Ax} = \mathbf{b} \leftrightarrow (\mathbf{A} + \Delta\mathbf{A})\tilde{\mathbf{x}} = \mathbf{b} + \Delta\mathbf{b} \quad \blacktriangleright \quad (\mathbf{A} + \Delta\mathbf{A})(\tilde{\mathbf{x}} - \mathbf{x}) = \Delta\mathbf{b} - \Delta\mathbf{Ax}. \quad (2.2.2.3)$$



#### Theorem 2.2.2.4. Conditioning of LSEs $\rightarrow$ [QSS00, Thm. 3.1], [GGK14, Thm 3.5]

If  $\mathbf{A}$  regular,  $\|\Delta\mathbf{A}\| < \|\mathbf{A}^{-1}\|^{-1}$  and (2.2.2.3), then

- (i)  $\mathbf{A} + \Delta\mathbf{A}$  is regular/invertible,
- (ii)

$$\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{A}^{-1}\| \|\mathbf{A}\|}{1 - \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \|\Delta\mathbf{A}\| / \|\mathbf{A}\|} \left( \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} \right).$$

relative error of data    relative perturbations

The proof is based on the following fundamental result:

**Lemma 2.2.2.5** (Perturbation lemma).  $\rightarrow$  [QSS00, Thm. 1.5]

$$\mathbf{B} \in \mathbb{R}^{n,n}, \|\mathbf{B}\| < 1 \Rightarrow \mathbf{I} + \mathbf{B} \text{ regular} \wedge \|(\mathbf{I} + \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}.$$

*Proof.* We start with the  $\triangle$ -inequality

$$\|(\mathbf{I} + \mathbf{B})\mathbf{x}\| \geq (\|\mathbf{x}\| - \|\mathbf{B}\mathbf{x}\|) \geq \underbrace{(1 - \|\mathbf{B}\|)}_{\neq 0} \|\mathbf{x}\| \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (2.2.2.6)$$

We conclude that  $\mathbf{I} + \mathbf{B}$  must have trivial kernel  $\mathcal{N}(\mathbf{I} + \mathbf{B}) = \{\mathbf{0}\}$ , which implies that the square matrix  $\mathbf{I} + \mathbf{B}$  is regular. We continue using this fact, the definition of the matrix norm, and (2.2.2.6):

$$\|(\mathbf{I} + \mathbf{B})^{-1}\| = \sup_{\mathbf{x} \in \mathbb{R}^n \setminus \{0\}} \frac{\|(\mathbf{I} + \mathbf{B})^{-1}\mathbf{x}\|}{\|\mathbf{x}\|} = \sup_{\mathbf{y} \in \mathbb{R}^n \setminus \{0\}} \frac{\|\mathbf{y}\|}{\|(\mathbf{I} + \mathbf{B})\mathbf{y}\|} \leq \frac{1}{1 - \|\mathbf{B}\|}$$

□

*Proof. (of Thm. 2.2.2.4)* We use a slightly generalized version of Lemma 2.2.2.5, which gives us

$$\|(\mathbf{A} + \Delta\mathbf{A})^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\mathbf{A}^{-1}\|\|\Delta\mathbf{A}\|}.$$

We combine this estimate with (2.2.2.3):

$$\|\Delta\mathbf{x}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\mathbf{A}^{-1}\|\|\Delta\mathbf{A}\|} (\|\Delta\mathbf{b}\| + \|\Delta\mathbf{A}\mathbf{x}\|) \leq \frac{\|\mathbf{A}^{-1}\|\|\mathbf{A}\|}{1 - \|\mathbf{A}^{-1}\|\|\Delta\mathbf{A}\|} \left( \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{A}\|\|\mathbf{x}\|} + \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} \right) \|\mathbf{x}\|.$$

□

Note that the term  $\|\mathbf{A}\|\|\mathbf{A}^{-1}\|$  occurs frequently. Therefore it has been given a special name:

**Definition 2.2.2.7. Condition (number) of a matrix**

Condition (number) of a matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ :	$\text{cond}(\mathbf{A}) := \ \mathbf{A}^{-1}\ \ \mathbf{A}\ $
--	--

Note:

$\text{cond}(\mathbf{A})$  depends on the matrix norm  $\|\cdot\|$ !

Rewriting estimate of Thm. 2.2.2.4 with  $\Delta\mathbf{b} = 0$ ,

$$\epsilon_r := \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \frac{\text{cond}(\mathbf{A})\delta_A}{1 - \text{cond}(\mathbf{A})\delta_A}, \quad \delta_A := \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|}. \quad (2.2.2.8)$$

From (2.2.2.8) we conclude important messages of  $\text{cond}(\mathbf{A})$ :

- ◆ If  $\text{cond}(\mathbf{A}) \gg 1$ , small perturbations in  $\mathbf{A}$  can lead to large relative errors in the solution of the LSE.
- ◆ If  $\text{cond}(\mathbf{A}) \gg 1$ , a stable algorithm ( $\rightarrow$  Def. 1.5.5.19) can produce solutions with large relative error !

Recall Thm. 2.2.2.4: for regular  $\mathbf{A} \in \mathbb{K}^{n,n}$ , small  $\Delta\mathbf{A}$ , generic vector/matrix norm  $\|\cdot\|$

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b} \\ (\mathbf{A} + \Delta\mathbf{A})\tilde{\mathbf{x}} &= \mathbf{b} + \Delta\mathbf{b} \end{aligned} \Rightarrow \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \frac{\text{cond}(\mathbf{A})}{1 - \text{cond}(\mathbf{A})\|\Delta\mathbf{A}\|/\|\mathbf{A}\|} \left( \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} \right). \quad (2.2.2.9)$$

- $\text{cond}(\mathbf{A}) \gg 1 \Rightarrow$  small relative changes of data  $\mathbf{A}, \mathbf{b}$  may effect huge relative changes in solution.
- $\text{cond}(\mathbf{A})$  indicates sensitivity of “LSE problem”  $(\mathbf{A}, \mathbf{b}) \mapsto \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$   
(as “amplification factor” of (worst-case) relative perturbations in the data  $\mathbf{A}, \mathbf{b}$ ).

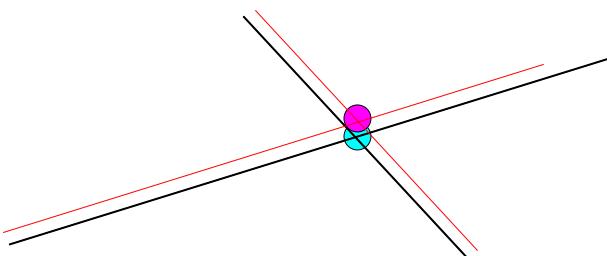
Terminology:

Small changes of data  $\Rightarrow$  small perturbations of result : well-conditioned problem  
 Small changes of data  $\Rightarrow$  large perturbations of result : ill-conditioned problem

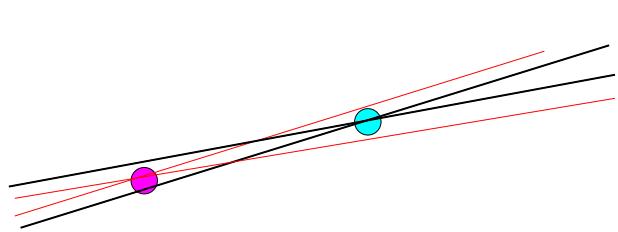
Note: sensitivity gauge depends on the chosen norm !

**EXAMPLE 2.2.2.10 (Intersection of lines in 2D)** Solving a  $2 \times 2$  linear system of equations amounts to finding the intersection of two lines in the coordinate plane. This relationship allows a geometric view of “sensitivity of a linear system”:

In distance metric (Euclidean vector norm):



nearly orthogonal intersection: well-conditioned



glancing intersection: ill-conditioned

Hessian normal form of line # $i$ ,  $i = 1, 2$ :

$$L_i = \{\mathbf{x} \in \mathbb{R}^2 : \mathbf{x}^T \mathbf{n}_i = d_i\}, \quad \mathbf{n}_i \in \mathbb{R}^2, d_i \in \mathbb{R}.$$

► LSE for finding intersection:  

$$\underbrace{\begin{bmatrix} \mathbf{n}_1^T \\ \mathbf{n}_2^T \end{bmatrix}}_{=: \mathbf{A}} \mathbf{x} = \underbrace{\begin{bmatrix} d_1 \\ d_2 \end{bmatrix}}_{=: \mathbf{b}},$$

$\mathbf{n}_i \hat{=} (\text{unit}) \text{ direction vectors}$ ,  $d_i \hat{=} \text{distance to origin}$ .

The following code investigates condition numbers for the matrix  $\begin{bmatrix} 1 & \cos \varphi \\ 0 & \sin \varphi \end{bmatrix}$  that can arise when computing the intersection of two lines enclosing the angle  $\varphi$ . As usual the directive `using namespace Eigen;` was given in the beginning of the file.

#### C++-code 2.2.2.11: condition numbers of $2 \times 2$ matrices → [GITLAB](#)

```

2  VectorXd phi = VectorXd::LinSpaced(50, M_PI / 200, M_PI / 2);
3  MatrixXd res(phi.size(), 3);
4  Matrix2d A;
5  A(0, 0) = 1;
6  A(1, 0) = 0;
7  for (int i = 0; i < phi.size(); ++i) {
8      A(0, 1) = std::cos(phi(i));
9      A(1, 1) = std::sin(phi(i));
10     // L2 condition number is the quotient of the maximal
11     // and minimal singular value of A
12     JacobiSVD<MatrixXd> svd(A);

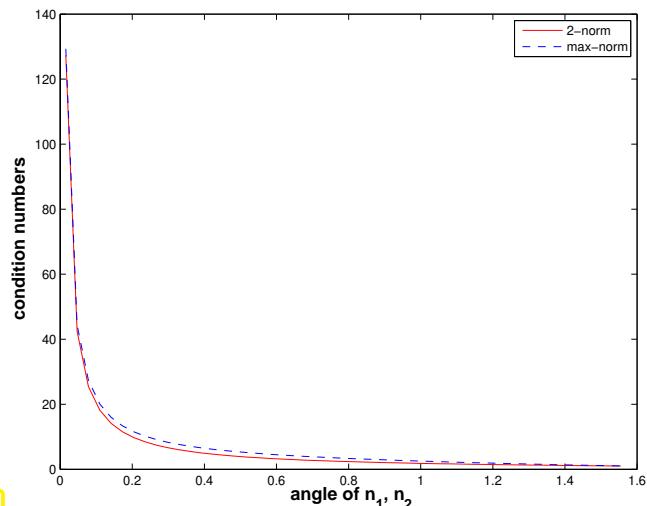
```

```

13     double C2 = svd.singularValues()(0) / //
14         svd.singularValues()(svd.singularValues().size() - 1);
15 // L-infinity condition number
16     double Cinf =
17         A.inverse().cwiseAbs().rowwise().sum().maxCoeff() *
18             A.cwiseAbs().rowwise().sum().maxCoeff(); //
19     res(i, 0) = phi(i);
20     res(i, 1) = C2;
21     res(i, 2) = Cinf;
22 }
```

In Line 13 we compute the condition number of  $\mathbf{A}$  with respect to the Euclidean vector norm using special EIGEN built-in functions.

Line 18 evaluated the condition number of a matrix for the maximum norm, recall Ex. 1.5.5.12.



We clearly observe a blow-up of  $\text{cond}(\mathbf{A})$  (with respect to the Euclidean vector norms) as the angle enclosed by the two lines shrinks.

This corresponds to a large sensitivity of the location of the intersection point in the case of glancing incidence.

Heuristics for predicting large  $\text{cond}(\mathbf{A})$

$\text{cond}(\mathbf{A}) \gg 1 \leftrightarrow \text{columns/rows of } \mathbf{A} \text{ “almost linearly dependent”}$

## 2.3 Gaussian Elimination

### 2.3.1 Basic algorithm

The problem of solving a linear system of equations is rather special compared to many other numerical tasks:

! An exceptional feature of linear systems of equations (LSE) is that its “exact” solution computable with finitely many elementary operations.

The algorithm is **Gaussian elimination (GE)** ( $\rightarrow$  secondary school, linear algebra,)

Familiarity with the algorithm of Gaussian elimination for a square linear system of equations will be taken for granted.

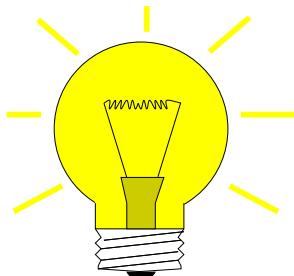


*Supplementary literature.* In case you cannot remember the main facts about Gaussian elimination, very detailed accounts and examples can be found in

- M. Gutknecht's lecture notes [[Gut09](#), Ch. 1],
- the textbook by Nipp & Stoffer [[NS02](#), Ch. 1],
- the numerical analysis text by Quarteroni et al. [[QSS00](#), Sects. 3.2 & 3.3],
- the textbook by Ascher & Greif [[AG11](#), Sect. 5.1],

and, to some extend, below, see Ex. 2.3.1.1.

[Wikipedia:](#) Although the method is named after mathematician **Carl Friedrich Gauss**, the earliest presentation of it can be found in the important Chinese mathematical text *Jiuzhang suanshu* or The Nine Chapters on the Mathematical Art, dated approximately 150 B.C., and commented on by **Liu Hui** in the 3rd century.



The idea of Gaussian elimination is the transformation of a linear system of equations into a “simpler”, but equivalent LSE by means of successive (invertible) *row transformations*.

Rem. 1.3.1.12: row transformations  $\leftrightarrow$  left-multiplication with transformation matrix

Obviously, left multiplication with a regular matrix does not affect the solution of an LSE: for any *regular*  $T \in \mathbb{K}^{n,n}$

$$\mathbf{Ax} = \mathbf{b} \Rightarrow \mathbf{A}'\mathbf{x} = \mathbf{b}' , \text{ if } \mathbf{A}' = \mathbf{T}\mathbf{A}, \mathbf{b}' = \mathbf{T}\mathbf{b} .$$

So we may try to convert the linear system of equations to a form that can be solved more easily by multiplying with regular matrices from left, which boils down to applying row transformations. A suitable target format is a diagonal linear system of equations, for which all equations are completely decoupled. This is the gist of Gaussian elimination.

### EXAMPLE 2.3.1.1 (Gaussian elimination)

Stage ① (Forward) elimination:

$$\begin{aligned}
 \left[ \begin{array}{ccc|c} 1 & 1 & 0 & 4 \\ 2 & 1 & -1 & 1 \\ 3 & -1 & -1 & -3 \end{array} \right] \xleftrightarrow{\quad} \left[ \begin{array}{ccc|c} x_1 & x_2 & & 4 \\ 2x_1 & x_2 & -x_3 & 1 \\ 3x_1 & -x_2 & -x_3 & -3 \end{array} \right] & \quad \left[ \begin{array}{ccc|c} 1 & 1 & 0 & 4 \\ 0 & -1 & -1 & 1 \\ 3 & -1 & -1 & -3 \end{array} \right] \\
 \left[ \begin{array}{ccc|c} 1 & 1 & 0 & 4 \\ 2 & 1 & -1 & 1 \\ 3 & -1 & -1 & -3 \end{array} \right] \xrightarrow{\quad} \left[ \begin{array}{ccc|c} 1 & 1 & 0 & 4 \\ 0 & -1 & -1 & 1 \\ 3 & -1 & -1 & -3 \end{array} \right] & \quad \left[ \begin{array}{ccc|c} 1 & 1 & 0 & 4 \\ 0 & -1 & -1 & 1 \\ 0 & -4 & -1 & -15 \end{array} \right] \\
 \xrightarrow{\quad} \left[ \begin{array}{ccc|c} 1 & 1 & 0 & 4 \\ 0 & -1 & -1 & 1 \\ 0 & 0 & 3 & 13 \end{array} \right] & \quad \underbrace{\left[ \begin{array}{ccc|c} 1 & 1 & 0 & 4 \\ 0 & -1 & -1 & 1 \\ 0 & 0 & 3 & 13 \end{array} \right]}_{=U}
 \end{aligned}$$

 = pivot row, pivot element **bold**.



We have transformed the LSE to **upper triangular form**

Stage ② Solve by **back substitution**:

$$\begin{aligned} x_1 + x_2 &= 4 \\ -x_2 - x_3 &= -7 \Rightarrow x_2 = 7 - \frac{13}{3} = \frac{8}{3} \\ 3x_3 &= 13 \quad x_1 = 4 - \frac{8}{3} = \frac{4}{3}. \end{aligned}$$

More detailed examples: [Gut09, Sect. 1.1], [NS02, Sect. 1.1]. 

More general:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2 \\ \vdots &\vdots \vdots \vdots \vdots \vdots \vdots \vdots \vdots \\ \vdots &\vdots \vdots \vdots \vdots \vdots \vdots \vdots \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n &= b_n \end{aligned}$$

- $i$ -th row -  $l_{i1} \cdot$  1st row (**pivot row**),  $l_{i1} := a_{i1}/a_{11}$ ,  $i = 2, \dots, n$

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1 \\ a_{22}^{(1)}x_2 + \cdots + a_{2n}^{(1)}x_n &= b_2^{(1)} \quad \text{with} \\ \vdots &\vdots \vdots \vdots \vdots \vdots \vdots \vdots \vdots \quad a_{ij}^{(1)} = a_{ij} - a_{1j}l_{i1}, \quad i, j = 2, \dots, n, \\ \vdots &\vdots \vdots \vdots \vdots \vdots \vdots \vdots \vdots \quad b_i^{(1)} = b_i - b_1l_{i1}, \quad i = 2, \dots, n. \\ a_{n2}^{(1)}x_2 + \cdots + a_{nn}^{(1)}x_n &= b_n^{(1)} \end{aligned}$$

- $i$ -th row -  $l_{i1} \cdot$  2nd row (**pivot row**),  $l_{i2} := a_{i2}^{(1)}/a_{22}^{(1)}$ ,  $i = 3, \dots, n$ .

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \cdots + a_{1n}x_n &= b_1 \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \cdots + a_{2n}^{(1)}x_n &= b_2^{(1)} \\ a_{33}^{(2)}x_3 + \cdots + a_{3n}^{(2)}x_n &= b_3^{(2)} \\ \vdots &\vdots \vdots \vdots \vdots \vdots \vdots \vdots \vdots \quad \vdots \vdots \vdots \vdots \vdots \vdots \vdots \vdots \\ a_{n3}^{(2)}x_3 + \cdots + a_{nn}^{(2)}x_n &= b_n^{(2)} \end{aligned}$$



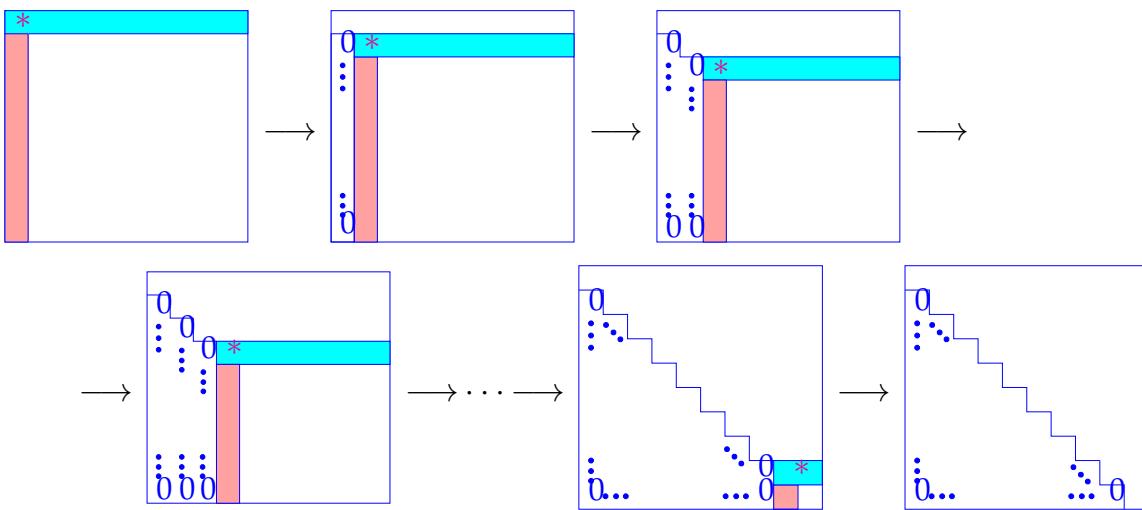
After  $n - 1$  steps:

linear systems of equations in **upper triangular form**

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \cdots + a_{1n}x_n &= b_1 \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \cdots + a_{2n}^{(1)}x_n &= b_2^{(1)} \\ a_{33}^{(2)}x_3 + \cdots + a_{3n}^{(2)}x_n &= b_3^{(2)} \\ \ddots &\ddots \ddots \vdots \vdots \vdots \vdots \vdots \vdots \\ \ddots &\ddots \ddots \vdots \vdots \vdots \vdots \vdots \vdots \\ a_{nn}^{(n-1)}x_n &= b_n^{(n-1)} \end{aligned}$$

Terminology:  $a_{11}, a_{22}^{(1)}, a_{33}^{(2)}, \dots, a_{n-1,n-1}^{(n-2)}$  = **pivots/pivot elements**

Graphical depiction:



\*  $\hat{=}$  the pivot entry (necessarily  $\neq 0$ , which we assume here), pivot row

In  $k$ -th step (starting from  $\mathbf{A} \in \mathbb{K}^{n,n}$ ,  $1 \leq k < n$ , pivot row  $\mathbf{a}_k^\top$ ):

transformation:  $\mathbf{Ax} = \mathbf{b} \rightarrow \mathbf{A}'\mathbf{x} = \mathbf{b}'$ .

with

$$a'_{ij} := \begin{cases} a_{ij} - \frac{a_{ik}}{a_{kk}} a_{kj} & \text{for } k < i, j \leq n, \\ 0 & \text{for } k < i \leq n, j = k, \\ a_{ij} & \text{else,} \end{cases} \quad b'_i := \begin{cases} b_i - \frac{a_{ik}}{a_{kk}} b_k & \text{for } k < i \leq n, \\ b_i & \text{else.} \end{cases} \quad (2.3.1.2)$$

multipliers  $l_{ik}$

**§2.3.1.3 (Gaussian elimination: algorithm)** Here we give a direct EIGEN implementation of Gaussian elimination for LSE  $\mathbf{Ax} = \mathbf{b}$  (grossly inefficient!).

#### C++ code 2.3.1.4: Solving LSE $\mathbf{Ax} = \mathbf{b}$ with Gaussian elimination → GITLAB

```

2 //! Gauss elimination without pivoting, x = A\b
3 //! A must be an n×n-matrix, b an n-vector
4 //! The result is returned in x
5 void gausselimsolve(const MatrixXd &A, const VectorXd& b,
6     VectorXd& x) {
7     int n = A.rows();
8     MatrixXd Ab(n,n+1); // Augmented matrix [A,b]
9     Ab << A, b; //
10    // Forward elimination (cf. step ① in Ex. 2.3.1.1)
11    for(int i = 0; i < n-1; ++i) {
12        double pivot = Ab(i,i);
13        for(int k = i+1; k < n; ++k) {
14            double fac = Ab(k, i) / pivot; // the multiplier
15            Ab.block(k, i+1, 1, n-i) -= fac * Ab.block(i, i+1, 1, n-i); //
16        }
17    }
18    // Back substitution (cf. step ② in Ex. 2.3.1.1)
19    Ab(n-1,n) = Ab(n-1,n) / Ab(n-1,n-1);
20    for(int i = n-2; i >= 0; --i) {
21        for(int l = i+1; l < n; ++l) Ab(i,n) -= Ab(l,n)*Ab(i,l);
22        Ab(i,n) /= Ab(i,i);
23    }

```

```

24   x = Ab.rightCols(1); // Solution in rightmost column!
25 }
```

- In Line 9 the right hand side vector set as last column of matrix, which facilitates simultaneous row transformations of matrix and r.h.s.
- In Line 14 the variable `fac` is the multiplier from (2.3.1.2).
- In Line 24 we extract solution from last column of the transformed matrix.

]

### §2.3.1.5 (Computational effort of Gaussian elimination)

- We examine Code 2.3.1.4.
- Forward elimination involves three nested loops (note that the compact vector operation in line 15 involves another loop from  $i+1$  to  $m$ ).
  - Back substitution can be done with two nested loops.

computational cost ( $\leftrightarrow$  number of elementary operations) of Gaussian elimination [NS02, Sect. 1.3]:

$$\begin{aligned} \text{forward elimination : } & \sum_{i=1}^{n-1} (n-i)(2(n-i)+3) = n(n-1)(\frac{2}{3}n + \frac{7}{6}) \text{ Ops. ,} \\ \text{back substitution : } & \sum_{i=1}^n 2(n-i) + 1 = n^2 \text{ Ops. .} \end{aligned} \quad (2.3.1.6)$$

asymptotic complexity ( $\rightarrow$  Sect. 1.4) of Gaussian elimination  
 (without pivoting) for generic LSE  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{n,n}$   $= \frac{2}{3}n^3 + O(n^2) = O(n^3)$

**EXPERIMENT 2.3.1.7 (Runtime of Gaussian elimination)** In this experiment we compare the efficiency of our hand-coded Gaussian elimination with that of library functions.

**C++ code 2.3.1.8: Measuring runtimes of Code 2.3.1.4 vs. EIGEN lu()-operator vs. MKL**  
**→ GITLAB**

```

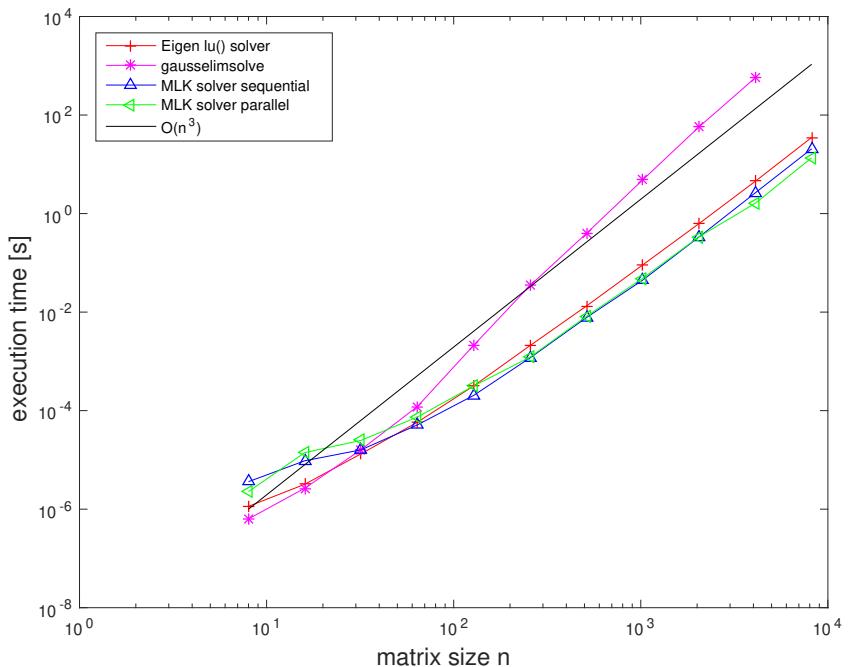
2  ///! Eigen code for timing numerical solution of linear systems
3  MatrixXd gausstiming(){
4      std::vector<int> n = {8,16,32,64,128,256,512,1024,2048,4096,8192};
5      int nruns = 3;
6      MatrixXd times(n.size(),3);
7      for(int i = 0; i < n.size(); ++i){
8          Timer t1, t2; // timer class
9          MatrixXd A = MatrixXd::Random(n[i],n[i]) + n[i]*MatrixXd::Identity(n[i],n[i]);
10         VectorXd b = VectorXd::Random(n[i]);
11         VectorXd x(n[i]);
12         for(int j = 0; j < nruns; ++j){
13             t1.start(); x = A.lu().solve(b); t1.stop(); // Eigen implementation
14             #ifndef EIGEN_USE_MKL_ALL // only test own algorithm without MKL
15             if(n[i] <= 4096) // Prevent long runs
16                 t2.start(); gausselimsolve(A,b,x); t2.stop(); // own gauss
17                 elimination
18             #endif
19         }
20         times(i,0) = n[i]; times(i,1) = t1.min(); times(i,2) = t2.min();
21     }
22     return times;
}

```

Platform:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz  
Ã— 4
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB, Mem 8 GB
- ◆ gcc 4.8.4, -O3

► EIGEN is about two orders of magnitude faster than a direct implementation, MKL is even faster.



<i>n</i>	Code 2.3.1.4 [s]	EIGEN lu() [s]	MKL sequential [s]	MKL parallel [s]
8	6.340e-07	1.140e-06	3.615e-06	2.273e-06
16	2.662e-06	3.203e-06	9.603e-06	1.408e-05
32	1.617e-05	1.331e-05	1.603e-05	2.495e-05
64	1.214e-04	5.836e-05	5.142e-05	7.416e-05
128	2.126e-03	3.180e-04	2.041e-04	3.176e-04
256	3.464e-02	2.093e-03	1.178e-03	1.221e-03
512	3.954e-01	1.326e-02	7.724e-03	8.175e-03
1024	4.822e+00	9.073e-02	4.457e-02	4.864e-02
2048	5.741e+01	6.260e-01	3.347e-01	3.378e-01
4096	5.727e+02	4.531e+00	2.644e+00	1.619e+00
8192	-	3.510e+01	2.064e+01	1.360e+01

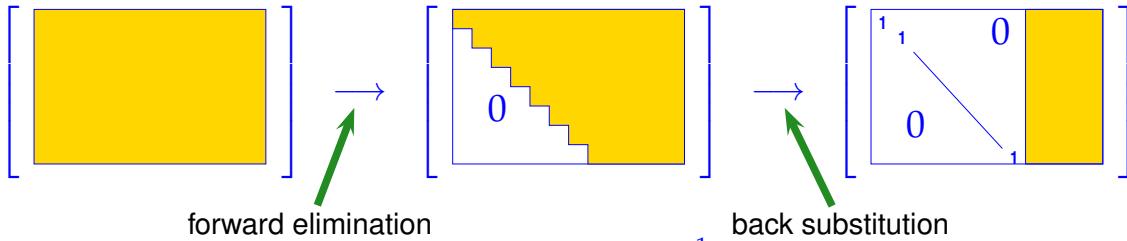
Never implement Gaussian elimination yourself !

use numerical libraries (LAPACK/MKL) or EIGEN !

A concise list of libraries for numerical linear algebra and related problems can be found [here](#).

**Remark 2.3.1.9 (Gaussian elimination for non-square matrices)** In Code 2.3.1.4: the right hand side vector  $\mathbf{b}$  was first appended to matrix  $\mathbf{A}$  as rightmost column, and then forward elimination and back substitution were carried out on the resulting matrix. This can be generalized to a Gaussian elimination for rectangular matrices  $\mathbf{A} \in \mathbb{K}^{n,n+1}$ !

Consider a “fat matrix”  $\mathbf{A} \in \mathbb{K}^{n,m}$ ,  $m > n$ :



Recall Code 2.3.1.4 ( $m = n + 1$ ): the solution vector  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$  was recovered as the rightmost column of the augmented matrix  $(\mathbf{A}, \mathbf{b})$  after forward elimination and back substitution. In the above cartoon it would be contained in the yellow part of the matrix on the right.

With this technique we have an efficient way for *simultaneously* solving of LSEs with multiple right hand sides. These multiple right-hand side can be passed as the column of a matrix  $\mathbf{B}$  and the problem of solving an LSE for several right-hand-side vectors can be stated as follows:

Given regular  $\mathbf{A} \in \mathbb{K}^{n,n}$ ,  $\mathbf{B} \in \mathbb{K}^{n,k}$ , seek  $\mathbf{X} \in \mathbb{K}^{n,k}$  such that

$$\mathbf{AX} = \mathbf{B} \Leftrightarrow \mathbf{X} = \mathbf{A}^{-1}\mathbf{B}$$

Usually library functions meant to solve LSEs also accept a matrix instead of a right-hand-side vector and then return a matrix of solution vectors. For instance, in EIGEN the following function call accomplishes this:

```
Eigen::MatrixXd X = A.lu().solve(B);
```

Its asymptotic complexity is

$O(n^2(n+k))$  for  $n, k \rightarrow \infty$ .

**C++ code 2.3.1.10: Gaussian elimination with multiple r.h.s. → Code 2.3.1.4 → GITLAB**

```

2  ///! Gauss elimination without pivoting, X = A-1B
3  ///! A must be an n × n-matrix, B an n × m-matrix
4  ///! Result is returned in matrix X
5  void gausselimsolvemult(const MatrixXd &A, const MatrixXd& B,
6   MatrixXd& X) {
7      int n = A.rows(), m = B.cols();
8      MatrixXd AB(n, n+m); // Augmented matrix [A,B]
9      AB << A, B;
10     // Forward elimination, do not forget the B part of the Matrix
11     for(int i = 0; i < n-1; ++i){
12         double pivot = AB(i, i);
13         for(int k = i+1; k < n; ++k){
14             double fac = AB(k, i)/pivot;
15             AB.block(k, i+1, 1,m+n-i-1) -= fac * AB.block(i, i+1, 1,m+n-i-1);
16         }
17     }
18     // Back substitution
19     AB.block(n-1, n, 1, m) /= AB(n-1,n-1);
20     for(int i = n-2; i >= 0; --i) {
21         for(int l = i+1; l < n; ++l) {
22             AB.block(i, n, 1, m) -= AB.block(l, n, 1, m)*AB(i, l);
23         }
24         AB.block(i, n, 1, m) /= AB(i, i);
25     }
26     X = AB.rightCols(m);
27 }
```

Concerning the next two remarks: For understanding or analyzing special variants of Gaussian elimination, it is useful to be aware of

- the effects of elimination steps on the level of *matrix blocks*, cf. Rem. 1.3.1.13,
- and of the *recursive nature* of Gaussian elimination.

**Remark 2.3.1.11 (Gaussian elimination via rank-1 modifications)** We can view Gaussian elimination from the perspective of matrix block operations: Then the first step of Gaussian elimination with pivot  $\alpha \neq 0$ , cf. (2.3.1.2), can be expressed as

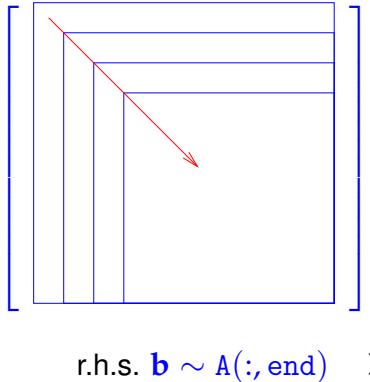
$$\mathbf{A} := \begin{bmatrix} \alpha & \mathbf{c}^T \\ \mathbf{d} & \mathbf{C} \end{bmatrix} \rightarrow \mathbf{A}' := \begin{bmatrix} \alpha & \mathbf{c}^T \\ \mathbf{0} & \mathbf{C}' := \mathbf{C} - \frac{\mathbf{d}\mathbf{c}^T}{\alpha} \end{bmatrix}. \quad (2.3.1.12)$$

rank-1 modification of  $\mathbf{C}$

Terminology: Adding a tensor product of two vectors to a matrix is called a **rank-1 modification** of that matrix, see also § 2.6.0.12 below.

Notice that the transformation (2.3.1.12) is applied to the resulting lower-right block  $\mathbf{C}'$  in the next elimination step. Thus Gaussian elimination can be realized by successive rank-1 modifications applied to smaller

and smaller lower-right blocks of the matrix. An implementation in this spirit is given in Code 2.3.1.13.



### C++ code 2.3.1.13: GE by rank-1 modification → GITLAB

```

2  /// in-situ Gaussian elimination, no pivoting
3  /// right hand side in rightmost column of A
4  /// back substitution is not done in this code!
5  void blockgs(Eigen::MatrixXd &A){
6      int n = A.rows();
7      for(int i = 1; i < n; ++i){
8          // rank-1 modification of C
9          A.bottomRightCorner(n-i, n-i+1) -=
10             A.col(i-1).tail(n-i) * A.row(i-1).tail(n-i+1) /
11             A(i-1, i-1);
12         A.col(i-1).tail(n-i).setZero(); // set d = 0
    }
}

```

In this code the Gaussian elimination is carried out *in situ*: the matrix  $A$  is replaced with the transformed matrices during elimination. If the matrix is not needed later this offers maximum efficiency. An in-situ LU-decomposition as described in Rem. 2.3.2.11 could also be performed by Code 2.3.1.13 after a modification of Line 10.  $\square$

**Remark 2.3.1.14 (Block Gaussian elimination)** Recall the “principle” from § 1.3.1.13: deal with block matrices (“matrices of matrices”) like regular matrices (except for commutativity of multiplication!). This suggests a block view of Gaussian elimination:

Given: regular matrix  $A \in \mathbb{K}^{n,n}$  with sub-matrices/blocks

$$\begin{aligned} A_{11} &:= (A)_{1:k, 1:k}, & A_{12} &:= (A)_{1:k, k+1:n}, \\ A_{21} &:= (A)_{k+1:n, 1:k}, & A_{22} &:= (A)_{k+1:n, k+1:n}, \end{aligned} \quad k < n,$$

and a right-hand-side vector  $b \in \mathbb{K}^n$  split into  $b_1 = (b)_{1:k}$ ,  $b_2 = (b)_{k+1:n}$

We apply the usual row transformations from (2.3.1.2) on the level of matrix block to this block-partitioned linear system using  $A_{11}$  as *pivot block*. Of course, we have to assume that  $A_{11}$  is invertible, generalizing the assumption that eligible pivot elements must not be zero. Again, the manipulations can be broken down into an elimination step ① and a backsubstitution step ②.

$$\begin{array}{c} \left[ \begin{array}{cc|c} A_{11} & A_{12} & b_1 \\ A_{21} & A_{22} & b_2 \end{array} \right] \xrightarrow{\textcircled{1}} \left[ \begin{array}{cc|c} A_{11} & A_{12} & b_1 \\ 0 & A_{22} - A_{21}A_{11}^{-1}A_{12} & b_2 - A_{21}A_{11}^{-1}b_1 \end{array} \right] \\ \xrightarrow{\textcircled{2}} \left[ \begin{array}{cc|c} I & 0 & A_{11}^{-1}(b_1 - A_{12}S^{-1}b_S) \\ 0 & I & S^{-1}b_S \end{array} \right]; \end{array}$$

where we abbreviated  $S := A_{22} - A_{21}A_{11}^{-1}A_{12}$ , a matrix known as *Schur complement*, see Rem. 2.3.2.19, and  $b_S := b_2 - A_{21}A_{11}^{-1}b_1$ .

We can read off the solution of the block-partitioned linear system from the above Gaussian elimination:

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \Rightarrow \begin{aligned} x_2 &= S^{-1}b_S, \\ x_1 &= A_{11}^{-1}(b_1 - A_{12}S^{-1}b_S). \end{aligned} \quad (2.3.1.15) \quad \square$$

## 2.3.2 LU-Decomposition

A **matrix factorization** (*ger.* Matrixzerlegung) expresses a general matrix  $A$  as product of two *special* (factor) matrices. Requirements for these special matrices define the matrix factorization. Matrix factorizations

come with the mathematical issue of existence & uniqueness, and pose the numerical challenge of finding algorithms for computing the factor matrices (efficiently and stably).

### Matrix factorizations

- ☞ often capture the essence of algorithms in compact form (here: Gaussian elimination),
- ☞ are important building blocks for complex algorithms,
- ☞ are key theoretical tools for algorithm analysis.

In this section the forward elimination step of Gaussian elimination will be related to a special matrix factorization, the so-called LU-decomposition or LU-factorization.

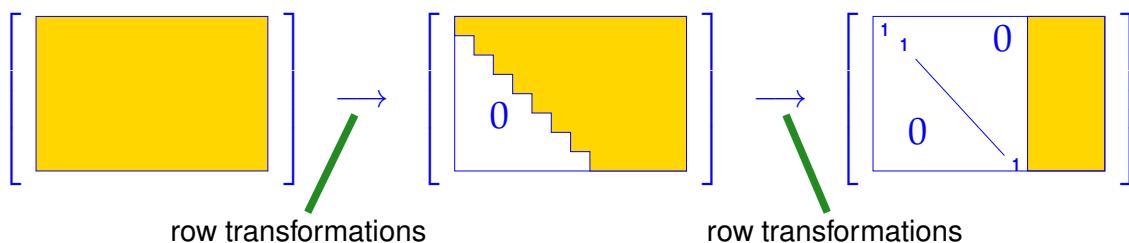


*Supplementary literature.* The LU-factorization should be well known from the introductory linear algebra course. In case you need to refresh your knowledge, please consult one of the following:

- textbook by Nipp & Stoffer [[NS02](#), Sect. 2.4],
- book by M. Hanke-Bourgeois [[Han02](#), p. II.4],
- linear algebra lecture notes by M. Gutknecht [[Gut09](#), Sect. 3.1],
- textbook by Quarteroni et al. [[QSS00](#), Sect.3.3.1],
- Sect. 3.5 of the book by Dahmen & Reusken,
- Sect. 5.1 of the textbook by Ascher & Greif [[AG11](#)].

See also (2.3.2.1) below.

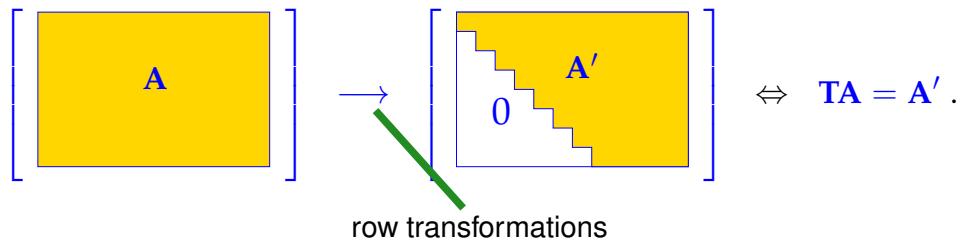
Recall the gist of Gaussian elimination split into the two steps of forward elimination and backsubstitution with (multiple) right-hand-side vectors appended to the coefficient matrix as rightmost columns:



Here: **row transformation** = adding a multiple of a matrix row to another row, or multiplying a row with a non-zero scalar (number) swapping two rows (more special row transformations are discussed in Rem. 1.3.1.12)

Note: Row transformations preserve regularity of a matrix and, thus, are suitable for transforming linear systems of equations: they will not affect the solution when applied to both the coefficient matrix and right-hand-side vector.

Rem. 1.3.1.12: row transformations can be realized by multiplication *from left* with suitable transformation matrices. When multiplying these transformation matrices we can emulate the effect to successive row transformations through left multiplication with a matrix  $\mathbf{T}$ :



Now we want to determine the  $T$  for the forward elimination step of Gaussian elimination.

**EXAMPLE 2.3.2.1 (Gaussian elimination and LU-factorization → [NS02, Sect. 2.4], [Han02, p. II.4], [Gut09, Sect. 3.1])** We revisit the LSE from Ex. 2.3.1.1 and carry out (forward) Gaussian elimination:

$$\begin{array}{c} \left[ \begin{array}{ccc} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 3 & -1 & -1 \end{array} \right] \left[ \begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] = \left[ \begin{array}{c} 4 \\ 1 \\ -3 \end{array} \right] \longleftrightarrow \begin{array}{l} x_1 + x_2 = 4 \\ 2x_1 + x_2 - x_3 = 1 \\ 3x_1 - x_2 - x_3 = -3 \end{array} \\ \left[ \begin{array}{ccc} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 3 & -1 & -1 \end{array} \right] \left[ \begin{array}{c} 4 \\ 1 \\ -3 \end{array} \right] \xrightarrow{} \left[ \begin{array}{ccc} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 0 & 1 & -1 \end{array} \right] \left[ \begin{array}{c} 4 \\ -7 \\ -3 \end{array} \right] \xrightarrow{} \\ \left[ \begin{array}{ccc} 1 & 1 & 0 \\ 0 & -1 & -1 \\ 3 & 4 & 1 \end{array} \right] \left[ \begin{array}{c} 4 \\ -7 \\ 13 \end{array} \right] \xrightarrow{} \left[ \begin{array}{ccc} 1 & 1 & 0 \\ 0 & -1 & -1 \\ 0 & 0 & 3 \end{array} \right] \left[ \begin{array}{c} 4 \\ -7 \\ 13 \end{array} \right] \end{array}$$

$=: L$        $=: U$

As before we highlight the pivot rows with █ and write the **pivot element** in **bold**. In addition, we let the *negative* multipliers take the places of matrix entries made to vanish; we color these entries red.

After this replacement we make the “surprising” observation that  $A = LU!$  □

The link between Gaussian elimination and **matrix factorization**, an explanations for the observation made in Ex. 2.3.2.1, becomes clear by recalling that row transformations result from multiplications with elimination matrices:

$$a_1 \neq 0 \quad \xrightarrow{} \quad \left[ \begin{array}{cccccc|c} 1 & 0 & \cdots & \cdots & 0 & a_1 \\ -\frac{a_2}{a_1} & 1 & & & 0 & a_2 \\ -\frac{a_3}{a_1} & & \ddots & & & a_3 \\ \vdots & & & \ddots & & \vdots \\ -\frac{a_n}{a_1} & 0 & & & 1 & a_n \end{array} \right] = \left[ \begin{array}{c} a_1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{array} \right]. \quad (2.3.2.2)$$

►  $n-1$  steps of Gaussian forward elimination immediately give rise to a matrix factorization (non-zero pivot elements assumed)

$$A = L_1 \cdots L_{n-1} U, \quad \text{with elimination matrices } L_i, i = 1, \dots, n-1, \text{ upper triangular matrix } U \in \mathbb{R}^{n,n}.$$

$$\left[ \begin{array}{ccccc|c} 1 & 0 & \cdots & \cdots & 0 & 1 \\ l_2 & 1 & & & 0 & 0 \\ l_3 & & \ddots & & & 0 \\ \vdots & & & \ddots & & \\ l_n & 0 & & & 1 & 1 \end{array} \right] \left[ \begin{array}{ccccc|c} 1 & 0 & \cdots & \cdots & 0 & 1 \\ 0 & 1 & & & 0 & h_2 \\ 0 & h_3 & 1 & & & 1 \\ \vdots & \vdots & & \ddots & & \\ 0 & h_n & 0 & & 1 & 1 \end{array} \right] = \left[ \begin{array}{ccccc|c} 1 & 0 & \cdots & \cdots & 0 & 1 \\ l_2 & 1 & & & & 0 \\ l_3 & h_3 & 1 & & & 0 \\ \vdots & \vdots & & \ddots & & \\ l_n & h_n & 0 & & 1 & 1 \end{array} \right]$$



The matrix products  $\mathbf{L}_1 \cdots \cdots \mathbf{L}_{n-1}$  yield **normalized lower triangular matrices**, whose entries are the multipliers  $-\frac{a_{ik}}{a_{kk}}$  from (2.3.1.2)  $\rightarrow$  Ex. 2.3.1.1.

The matrix factorization that “automatically” emerges during Gaussian forward elimination has a special name:

### Definition 2.3.2.3. LU-decomposition/LU-factorization

Given a square matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$ , an *upper triangular matrix*  $\mathbf{U} \in \mathbb{K}^{n,n}$  and a *normalized lower triangular matrix* ( $\rightarrow$  Def. 1.1.2.3) form an **LU-decomposition/LU-factorization** of  $\mathbf{A}$ , if  $\mathbf{A} = \mathbf{LU}$ .

$$\left[ \begin{array}{|c|} \hline \mathbf{A} \\ \hline \end{array} \right] = \left[ \begin{array}{|c|} \hline \mathbf{L} \\ \hline \end{array} \right] \cdot \left[ \begin{array}{|c|} \hline \mathbf{U} \\ \hline \end{array} \right],$$

Using this notion we can summarize what we have learned from studying elimination matrices:

The (forward) Gaussian elimination (without pivoting), for  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{n,n}$ , if possible, is *algebraically equivalent* to an **LU-factorization/LU-decomposition**  $\mathbf{A} = \mathbf{LU}$  of  $\mathbf{A}$  into a normalized lower triangular matrix  $\mathbf{L}$  and an upper triangular matrix  $\mathbf{U}$ , [DR08, Thm. 3.2.1], [NS02, Thm. 2.10], [Gut09, Sect. 3.1].

Algebraically equivalent  $\hat{=}$  when carrying out the forward elimination in situ as in Code 2.3.1.4 and storing the multipliers in a lower triangular matrix as in Ex. 2.3.2.1, then the latter will contain the **L-factor** and the original matrix will be replaced with the **U-factor**.

### Lemma 2.3.2.4. Existence of LU-decomposition

*The LU-decomposition of  $\mathbf{A} \in \mathbb{K}^{n,n}$  exists, if all submatrices  $(\mathbf{A})_{1:k, 1:k}$ ,  $1 \leq k \leq n$ , are regular.*

*Proof.* We adopt a block matrix perspective ( $\rightarrow$  Rem. 1.3.1.13) and employ induction w.r.t.  $n$ :

$n = 1$ : assertion trivial

$n - 1 \rightarrow n$ : Induction hypothesis ensures existence of normalized lower triangular matrix  $\tilde{\mathbf{L}}$  and regular upper triangular matrix  $\tilde{\mathbf{U}}$  such that  $\tilde{\mathbf{A}} = \tilde{\mathbf{L}}\tilde{\mathbf{U}}$ , where  $\tilde{\mathbf{A}}$  is the upper left  $(n - 1) \times (n - 1)$  block of  $\mathbf{A}$ :

$$\left[ \begin{array}{c|c} \tilde{\mathbf{A}} & \mathbf{b} \\ \hline \mathbf{a}^\top & \alpha \end{array} \right] = \left[ \begin{array}{c|c} \tilde{\mathbf{L}} & \mathbf{0} \\ \hline \mathbf{x}^\top & 1 \end{array} \right] \left[ \begin{array}{c|c} \tilde{\mathbf{U}} & \mathbf{y} \\ \hline \mathbf{0} & \xi \end{array} \right] =: \mathbf{LU}.$$

Then solve

- |   |  |  |
|---|--|--|
| ① | $\tilde{\mathbf{L}}\mathbf{y} = \mathbf{b}$            | $\rightarrow$ provides $\mathbf{y} \in \mathbb{K}^n$ , |
| ② | $\mathbf{x}^\top \tilde{\mathbf{U}} = \mathbf{a}^\top$ | $\rightarrow$ provides $\mathbf{x} \in \mathbb{K}^n$ , |
| ③ | $\mathbf{x}^\top \mathbf{y} + \xi = \alpha$            | $\rightarrow$ provides $\xi \in \mathbb{K}$ .          |

Regularity of  $\mathbf{A}$  involves  $\xi \neq 0$  (why?) so that  $\mathbf{U}$  will be regular, too.

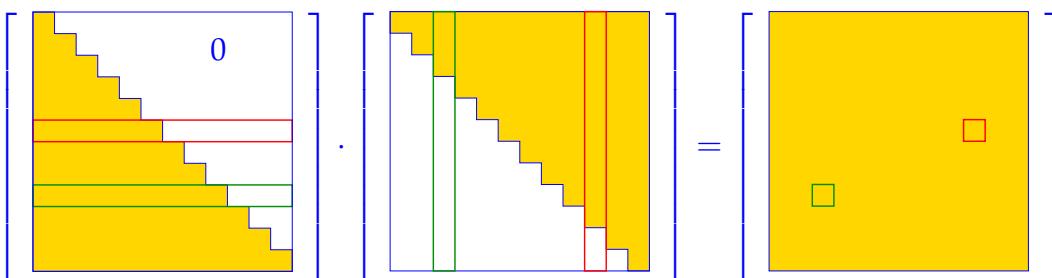
□

**§2.3.2.5 (Uniqueness of LU-decomposition)** Regular upper triangular matrices and normalized lower triangular matrices form *matrix groups* ( $\rightarrow$  Lemma 1.3.1.9). Their only common element is the identity matrix.

$$\mathbf{L}_1 \mathbf{U}_1 = \mathbf{L}_2 \mathbf{U}_2 \Rightarrow \mathbf{L}_2^{-1} \mathbf{L}_1 = \mathbf{U}_2 \mathbf{U}_1^{-1} = \mathbf{I}.$$

Since inverses of matrices are unique, so are the LU-factors:  $\mathbf{U}_1 = \mathbf{U}_2$  and  $\mathbf{L}_1 = \mathbf{L}_2$ . □

**§2.3.2.6 (Basic algorithm for computing LU-decomposition)** There are direct ways to determine the factor matrices of the LU-decomposition [Gut09, Sect. 3.1], [QSS00, Sect. 3.3.3] and, of course, they are closely related to forward Gaussian elimination. To derive the algorithm we study the entries of the product of a normalized lower triangular and an upper triangular matrix, see Def. 1.1.2.3:



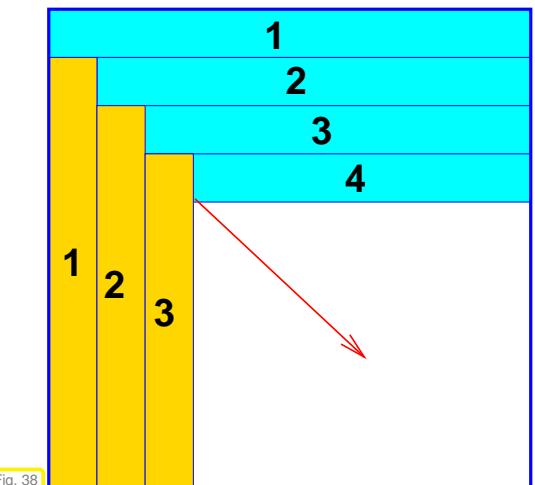
Taking into account the zero entries known a priori, we arrive at

$$\mathbf{LU} = \mathbf{A} \Rightarrow a_{ik} = \sum_{j=1}^{\min\{i,k\}} l_{ij} u_{jk} = \begin{cases} \sum_{j=1}^{i-1} l_{ij} u_{jk} + 1 \cdot u_{ik} & , \text{ if } i \leq k , \\ \sum_{j=1}^{k-1} l_{ij} u_{jk} + l_{ik} u_{kk} & , \text{ if } i > k . \end{cases} \quad (2.3.2.7)$$

This reveals how to compute the entries of  $\mathbf{L}$  and  $\mathbf{U}$  sequentially. We start with the top row of  $\mathbf{U}$ , which agrees with that of  $\mathbf{A}$ , and then work our way towards to bottom right corner:

- • row by row computation of  $\mathbf{U}$
  - column by column computation of  $\mathbf{L}$
- Entries of  $\mathbf{A}$  can be replaced with those of  $\mathbf{L}$ ,  $\mathbf{U}$  !  
(so-called **in situ/in place computation**)
- (Crout's algorithm, [Gut09, Alg. 3.1])

$\hat{=}$  rows of  $\mathbf{U}$   
  $\hat{=}$  columns of  $\mathbf{L}$



The following code follows this sequential computation scheme:

#### C++ code 2.3.2.8: LU-factorization → GITLAB

```

2 //! Algorithm of Crout: LU-factorization of  $\mathbf{A} \in \mathbb{K}^{n,n}$ 
3 std::pair<MatrixXd, MatrixXd> lufak(const MatrixXd &A) {
4     int n = A.rows();
5     assert(n == A.cols()); // Ensure matrix is square
6     MatrixXd L{MatrixXd::Identity(n, n)};
7     MatrixXd U{MatrixXd::Zero(n, n)};
8     for (int k = 0; k < n; ++k) {

```

```

9  // Compute row of U
10 for (int j = k; j < n; ++j)
11     U(k, j) = A(k, j) - (L.block(k, 0, 1, k) * U.block(0, j, k, 1))(0, 0);
12 // Compute column of L
13 for (int i = k + 1; i < n; ++i)
14     L(i, k) = (A(i, k) - (L.block(i, 0, 1, k) * U.block(0, k, k, 1))(0, 0)) /
15         U(k, k);
16 }
17 return { L, U };
18 }
```

It is instructive to compare this code with a simple implementation of the matrix product of a normalized lower triangular and an upper triangular matrix. From this perspective the LU-factorization looks like the “inversion” of matrix multiplication:

#### C++ code 2.3.2.9: matrix multiplication $L \cdot U \rightarrow$ GITLAB

```

2 //! Multiplication of normalized lower/upper triangular matrices
3 Eigen::MatrixXd lumult(const MatrixXd &L, const MatrixXd &U) {
4     int n = L.rows();
5     assert(n == L.cols() && n == U.cols() && n == U.rows());
6     Eigen::MatrixXd A{MatrixXd::Zero(n, n)};
7     for (int k = 0; k < n; ++k) {
8         for (int j = k; j < n; ++j)
9             A(k, j) = U(k, j) + (L.block(k, 0, 1, k) * U.block(0, j, k, 1))(0, 0);
10        for (int i = k + 1; i < n; ++i)
11            A(i, k) =
12                (L.block(i, 0, 1, k) * U.block(0, k, k, 1))(0, 0) + L(i, k) * U(k, k);
13    }
14    return A;
15 }
```

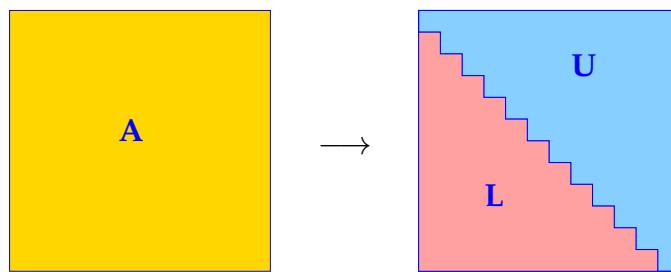
Observe: Solving for entries  $L(i, k)$  of  $L$  and  $U(k, j)$  of  $U$  in the multiplication of an upper triangular and normalized lower triangular matrix (Code 2.3.2.9) yields the algorithm for LU-factorization (Code 2.3.2.8).  $\dashv$

The computational cost of LU-factorization is immediate from Code 2.3.2.8 and the same as for Gaussian elimination, cf. § 2.3.1.5:

Asymptotic complexity of LU-factorization of  $A \in \mathbb{R}^{n,n}$  (2.3.2.10)

$$= \frac{2}{3}n^3 + O(n^2) = O(n^3) \text{ for } n \rightarrow \infty$$

**Remark 2.3.2.11 (In-situ LU-decomposition)** “In situ” is Latin and means “in place”. Many library routines provide routines that *overwrite* the matrix  $A$  with its LU-factors in order to save memory when the original matrix is no longer needed. This is possible because the number of unknown entries of the LU-factors combined exactly agrees with the number of entries of  $A$ . The convention is to replace the strict lower-triangular part of  $A$  with  $L$ , and the upper triangular part with  $VU$ :



]

**Remark 2.3.2.12 (Recursive LU-factorization)** Recall Rem. 2.3.1.11 and the recursive view of Gaussian elimination it suggests, because in (2.3.1.12) an analogous row transformation can be applied to the remaining right-lower block  $\mathbf{C}'$ .

In light of the close relationship between Gaussian elimination and LU-factorization there will also be a recursive version of LU-factorization.

The following code implements the recursive *in situ* (in place) LU-decomposition of  $\mathbf{A} \in \mathbb{R}^{n,n}$  (without pivoting). It is closely related to Code 2.3.1.13, but now both  $\mathbf{L}$  and  $\mathbf{U}$  are stored in place of  $\mathbf{A}$ :

#### C++ code 2.3.2.13: Recursive LU-factorization → [GITLAB](#)

```

2  ///!  in situ recursive LU-factorization
3  MatrixXd lurec(const MatrixXd &A){
4      int n = A.rows();
5      MatrixXd result(n,n);
6      if(n > 1){
7          VectorXd fac = A.col(0).tail(n-1) / A(0,0); ///
8          result.bottomRightCorner(n-1,n-1) = lurec(A.bottomRightCorner(n-1,n-1)
9              - fac * A.row(0).tail(n-1)); ///
10         result.row(0) = A.row(0); result.col(0).tail(n-1) = fac;
11     }
12     return result;
13 }
```

Refer to (2.3.1.12) to understand `lurec`: the rank-1 modification of the lower  $(n-1) \times (n-1)$ -block of the matrix is done in lines Line 7-Line 8 of the code.

#### C++ code 2.3.2.14: Driver for recursive LU-factorization of Code 2.3.2.13 → [GITLAB](#)

```

2  ///!  post-processing: extract L and U
3  void lurecdriver(const MatrixXd &A, MatrixXd &L, MatrixXd &U){
4      MatrixXd A_dec = lurec(A);
5      // post-processing:
6      //extract L and U
7      U = A_dec.triangularView<Upper>();
8      L.setIdentity();
9      L += A_dec.triangularView<StrictlyLower>();
10 }
```

]

**§2.3.2.15 (Using LU-factorization to solve a linear system of equations)** An intermediate LU-factorization paves the way for a *three-stage* procedure for solving an  $n \times n$  linear system of equations.

① LU-decomposition  $\mathbf{A} = \mathbf{L}\mathbf{U}$ , #elementary operations =  $\frac{1}{3}n(n-1)(n+1)$

$\mathbf{Ax} = \mathbf{b}$  : ② forward substitution, solve  $\mathbf{Lz} = \mathbf{b}$ , #elementary operations =  $\frac{1}{2}n(n-1)$

③ backward substitution, solve  $\mathbf{Ux} = \mathbf{z}$ , #elementary operations =  $\frac{1}{2}n(n+1)$

- The asymptotic complexity of the complete three-stage algorithm is (in leading order) the same as for Gaussian elimination (The bulk of computational cost is incurred in the factorization step ①).

However, the perspective of LU-factorization reveals that the solution of linear systems of equations can be split into two separate phases with different asymptotic complexity in terms of the number  $n$  of unknowns:

**setup phase**

(factorization)

Cost:  $O(n^3)$

+

**elimination phase**

(forward/backward substition)

Cost:  $O(n^2)$

**Remark 2.3.2.16 (Rationale for using LU-decomposition in algorithms)** Gauss elimination and LU-factorization for the solution of a linear system of equations (→ § 2.3.2.15) are equivalent and only differ in the ordering of the steps.

Then, why is it important to know about LU-factorization?

Because in the case of LU-factorization the expensive forward elimination and the less expensive (forward/backward) substitutions are separated, which sometimes can be exploited to reduce computational cost, as highlighted in Rem. 2.5.0.10 below.

**Remark 2.3.2.17 ("Partial LU-decompositions" of principal minors)** The algorithm from § 2.3.2.6 reveals that the computation of the LU-decomposition of a matrix proceeds from top-left to bottom-right. This implies the locality property discussed in this remark. To understand its heading we remind that a **principal minor** refers to the left upper block of a matrix

The following “visual rule” help identify the structure of the LU-factors of a matrix.

(2.3.2.18)

The left-upper blocks of both  $\mathbf{L}$  and  $\mathbf{U}$  in the LU-factorization of  $\mathbf{A}$  depend only on the corresponding left-upper block of  $\mathbf{A}$ !

**Remark 2.3.2.19 (Block LU-factorization)** In the spirit of § 1.3.1.13 we can also adopt a matrix-block perspective of LU-factorization. This is a natural idea in light of the close connection between matrix multiplication and matrix factorization, cf. the relationship between matrix factorization and matrix multiplication found in § 2.3.2.6:

Block matrix multiplication (1.3.1.14)  $\cong$  block LU-decomposition:

We consider a block-partitioned matrix

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}, \quad \mathbf{A}_{11} \in \mathbb{K}^{n,n} \text{ regular}, \quad \mathbf{A}_{12} \in \mathbb{K}^{n,m} \\ \text{and } \mathbf{A}_{21} \in \mathbb{K}^{m,n}, \quad \mathbf{A}_{22} \in \mathbb{K}^{m,m}.$$

The block LU-decomposition arises from the block Gaussian forward elimination of Rem. 2.3.1.14 in the same way as the standard LU-decomposition is spawned by the entry-wise Gaussian elimination:

With

$$\underbrace{\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}}_{\text{block LU-factorization}} = \underbrace{\begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{I} \end{bmatrix}}_{\text{with Schur complement}} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ 0 & \mathbf{S} \end{bmatrix}, \quad \boxed{\mathbf{S} := \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}} \quad (2.3.2.20)$$

Under the assumption that  $\mathbf{A}_{11}$  is invertible, the Schur complement matrix  $\mathbf{S}$  is invertible, if and only if this holds for  $\mathbf{A}$ .  $\square$

### 2.3.3 Pivoting

Known from linear algebra [NS02, Sect. 1.1]:

$$\begin{array}{ccc} \begin{array}{c} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \\ \text{breakdown of Gaussian elimination} \\ \text{pivot element} = 0 \end{array} & \quad & \begin{array}{c} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_2 \\ b_1 \end{bmatrix} \\ \text{Gaussian elimination feasible} \end{array} \end{array}$$

Idea (in linear algebra): Avoid zero pivot elements by **swapping rows**

#### EXAMPLE 2.3.3.1 (Pivoting and numerical stability) → [DR08, Example 3.2.3])

```

2 MatrixXd A(2, 2);
3 A << 5.0e-17, 1.0, 1.0, 1.0;
4 VectorXd b(2), x2(2);
5 b << 1.0, 2.0;
6 VectorXd x1 = A.fullPivLu().solve(b);
7 gausselimsolve(A, b, x2); // see Code 2.3.1.10
8 auto [L,U] = lufak(A); // see Code 2.3.2.8
9 VectorXd z = L.lu().solve(b);
10 VectorXd x3 = U.lu().solve(z);
11 cout << "x1 = \n"
12     << x1 << "\nx2 = \n"
13     << x2 << "\nx3 = \n"
14     << x3 << std::endl;

```

Output:

```

1 x1 =
2 1
3 1
4 x2 =
5 0
6 1
7 x3 =
8 0
9 1

```

$$\mathbf{A} = \begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \Rightarrow \mathbf{x} = \begin{bmatrix} \frac{1}{1-\epsilon} \\ \frac{1-2\epsilon}{1-\epsilon} \end{bmatrix} \approx \begin{bmatrix} 1 \\ 1 \end{bmatrix} \text{ for } |\epsilon| \ll 1.$$

**What is wrong with EIGEN?** Needed: insight into **roundoff errors**, which we already have → Section 1.5.3

Armed with knowledge about the behavior of machine numbers and roundoff errors we can now understand what is going on in Ex. 2.3.3.1

Straightforward LU-factorization: if  $\epsilon \leq \frac{1}{2}\text{EPS}$ , EPS  $\hat{=}$  machine precision,

$$\blacktriangleright \quad \mathbf{L} = \begin{bmatrix} 1 & 0 \\ \epsilon^{-1} & 1 \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} \epsilon & 1 \\ 0 & 1 - \epsilon^{-1} \end{bmatrix} \stackrel{(*)}{=} \tilde{\mathbf{U}} := \begin{bmatrix} \epsilon & 1 \\ 0 & -\epsilon^{-1} \end{bmatrix} \text{ in } \mathbb{M}! \quad (2.3.3.2)$$

(\*) because  $1/\tilde{\epsilon}^2/\text{EPS} = 2/\text{EPS}$ , see Exp. 1.5.3.14.

$$\blacktriangleright \quad \text{Solution of } \mathbf{L}\tilde{\mathbf{U}}\mathbf{x} = \mathbf{b}: \quad \mathbf{x} = \begin{bmatrix} 2\epsilon \\ 1-2\epsilon \end{bmatrix} \text{ (meaningless result !)}$$

LU-factorization after swapping rows:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ \epsilon & 1 \end{bmatrix} \Rightarrow \mathbf{L} = \begin{bmatrix} 1 & 0 \\ \epsilon & 1 \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} 1 & 1 \\ 0 & 1-\epsilon \end{bmatrix} = \tilde{\mathbf{U}} := \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \text{ in } \mathbb{M}. \quad (2.3.3.3)$$

$$\blacktriangleright \quad \text{Solution of } \mathbf{L}\tilde{\mathbf{U}}\mathbf{x} = \mathbf{b}: \quad \mathbf{x} = \begin{bmatrix} 1+2\epsilon \\ 1-2\epsilon \end{bmatrix} \text{ (sufficiently accurate result !)}$$

no row swapping,  $\rightarrow$  (2.3.3.2):  $\mathbf{L}\tilde{\mathbf{U}} = \mathbf{A} + \mathbf{E}$  with  $\mathbf{E} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$   $\blacktriangleright$  unstable!

after row swapping,  $\rightarrow$  (2.3.3.3):  $\mathbf{L}\tilde{\mathbf{U}} = \tilde{\mathbf{A}} + \mathbf{E}$  with  $\mathbf{E} = \begin{bmatrix} 0 & 0 \\ 0 & \epsilon \end{bmatrix}$   $\blacktriangleright$  stable!

Introduction to the notion of **stability**  $\rightarrow$  Section 1.5.5, Def. 1.5.5.19, see also [DR08, Sect. 2.3].

Suitable pivoting essential for controlling impact of roundoff errors  
on Gaussian elimination ( $\rightarrow$  Section 1.5.5, [NS02, Sect. 2.5])

#### EXAMPLE 2.3.3.4 (Gaussian elimination with pivoting for $3 \times 3$ -matrix)

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 2 \\ 2 & -3 & 2 \\ 1 & 24 & 0 \end{bmatrix} \xrightarrow{\textcircled{1}} \begin{bmatrix} 2 & -3 & 2 \\ 1 & 2 & 2 \\ 1 & 24 & 0 \end{bmatrix} \xrightarrow{\textcircled{2}} \begin{bmatrix} 2 & -3 & 2 \\ 0 & 3.5 & 1 \\ 0 & 25.5 & -1 \end{bmatrix} \xrightarrow{\textcircled{3}} \begin{bmatrix} 2 & -7 & 2 \\ 0 & 25.5 & -1 \\ 0 & 3.5 & 1 \end{bmatrix} \xrightarrow{\textcircled{4}} \begin{bmatrix} 2 & -7 & 2 \\ 0 & 25.5 & -1 \\ 0 & 0 & 1.373 \end{bmatrix}$$

- ①: swap rows 1 & 2.
- ②: elimination with top row as pivot row
- ③: swap rows 2 & 3
- ④: elimination with 2nd row as pivot row

#### §2.3.3.5 (Algorithm: Gaussian elimination with partial pivoting)

##### C++ code 2.3.3.6: Gaussian elimination with pivoting: extension of Code 2.3.1.4 $\rightarrow$ GITLAB

```

2 //! Solving an LSE  $\mathbf{Ax} = \mathbf{b}$  by Gaussian elimination with partial pivoting
3 //! A must be an  $n \times n$ -matrix, b an  $n$ -vector
4 void gepiv(const MatrixXd &A, const VectorXd& b, VectorXd& x){
5     int n = A.rows();
6     MatrixXd Ab(n, n+1);
7     Ab << A, b;      //
8     // Forward elimination by rank-1 modification, see Rem. 2.3.1.11
9     for(int k = 0; k < n-1; ++k){

```

```

10 |         int j; double p; // p = relatively largest pivot, j = pivot row
11 |         index
12 |         p = (Ab.col(k).tail(n-k).cwiseAbs().cwiseQuotient(
13 |             Ab.block(k,k,n-k,n-k).cwiseAbs().rowwise().maxCoeff() )
14 |             ).maxCoeff(&j);//
15 |         if( p < std::numeric_limits<double>::epsilon() *
16 |             Ab.block(k,k,n-k,n-k).norm() )
17 |             throw std::logic_error("nearly singular");// 
18 |         Ab.row(k).tail(n-k+1).swap(Ab.row(k+j).tail(n-k+1));//
19 |         Ab.bottomRightCorner(n-k-1,n-k) -= Ab.col(k).tail(n-k-1) *
20 |             Ab.row(k).tail(n-k) / Ab(k,k);//
21 |
22 |     }
23 |     // Back substitution (same as in Code 2.3.1.4)
24 |     Ab(n-1,n) = Ab(n-1,n) / Ab(n-1,n-1);
25 |     for(int i = n-2; i >= 0; --i){
26 |         for(int l = i+1; l < n; ++l){
27 |             Ab(i,n) -= Ab(l,n)*Ab(i,l);
28 |         }
29 |         Ab(i,n) /= Ab(i,i);
30 |     }
31 |     x = Ab.rightCols(1); //
32 |
33 }

```

choice of pivot row index  $j$  (Line 11 of code): relatively largest pivot [NS02, Sect. 2.5],

$$j \in \{k, \dots, n\} \text{ such that } \frac{|a_{ji}|}{\max\{|a_{jl}|, l = k, \dots, n\}} \rightarrow \max \quad (2.3.3.7)$$

for  $k = j, k \in \{i, \dots, n\}$ : partial pivoting

Explanations to Code 2.3.3.6:

Line 7: Augment matrix  $\mathbf{A}$  by right hand side vector  $\mathbf{b}$ , see comments on Code 2.3.1.4 for explanations.

Line 11: Select index  $j$  for pivot row according to the recipe of partial pivoting, see (2.3.3.7).

Note: Inefficient implementation above (too many comparisons)! Try to do better!

Line 13: If the pivot element is still very small relative to the norm of the matrix, then we have encountered an entire column that is close to zero. Gaussian elimination may not be possible in a stable fashion for this matrix; warn user and terminate.

Line 14: A way to swap rows of a matrix in EIGEN.

Line 15: Forward elimination by means of rank-1-update, see (2.3.1.12).

Line 25: As in Code 2.3.1.4: after back substitution last column of augmented matrix supplies solution  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ .

§2.3.3.8 (Algorithm: LU-factorization with pivoting) Recall: close relationship between Gaussian elimination and LU-factorization

➤ LU-factorization with pivoting? Of course, just by rearranging the operations of Gaussian forward elimination with pivoting.

EIGEN-code for in place LU-factorization of  $\mathbf{A} \in \mathbb{R}^{n,n}$  with **partial pivoting**:

#### C++ code 2.3.3.9: LU-factorization with partial pivoting → [GITLAB](#)

```

2 void lupiv(MatrixXd &A) { //insitu
3     int n = A.rows();
4     for(int k = 0; k < n-1; ++k){
5         int j; double p; // p = relatively largest pivot, j = pivot row
6         index
7         p = ( A.col(k).tail(n-k).cwiseAbs().cwiseQuotient(
8             A.block(k,k,n-k,n-k).cwiseAbs().rowwise().maxCoeff() ) ).maxCoeff(&j);
9         //
10        if( p < std::numeric_limits<double>::epsilon() *
11            A.block(k,k,n-k,n-k).norm() ) //
12            throw std::logic_error("nearly singular");
13        A.row(k).tail(n-k-1).swap(A.row(k+j).tail(n-k-1));//
14        VectorXd fac = A.col(k).tail(n-k-1) / A(k,k);//
15        A.bottomRightCorner(n-k-1,n-k-1) -= fac * A.row(k).tail(n-k-1);//
16        A.col(k).tail(n-k-1) = fac;//
17    }
18}

```

Notice that the recursive call is omitted as in 2.3.1.11

Explanations to Code 2.3.3.9:

- Line 6: Find the *relatively largest* pivot element  $p$  and the index  $j$  of the corresponding row of the matrix, see (2.3.3.7)
- Line 7: If the pivot element is still very small relative to the norm of the matrix, then we have encountered an entire column that is close to zero. The matrix is (close to) singular and LU-factorization does not exist.
- Line 9: Swap the first and the  $j$ -th row of the matrix.
- Line 10: Initialize the vector of multiplier.
- Line 11: Call the routine for the lower right  $(n-1) \times (n-1)$ -block of the matrix after subtracting suitable multiples of the first row from the other rows, cf. Rem. 2.3.1.11 and Rem. 2.3.2.12.
- Line 12: Reassemble the parts of the LU-factors. The vector of multipliers yields a column of  $\mathbf{L}$ , see Ex. 2.3.2.1.

**Remark 2.3.3.10 (Rationale for partial pivoting policy (2.3.3.7) → [NS02, Page 47])** Why *relatively largest* pivot element in (2.3.3.7)? **scaling invariance** desirable

Scale linear system of equations from Ex. 2.3.3.1:

$$\begin{bmatrix} 2/\epsilon & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 & 2/\epsilon \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2/\epsilon \\ 2 \end{bmatrix} := \tilde{\mathbf{b}}$$

No row swapping, if absolutely largest pivot element is used:

$$\begin{bmatrix} 2 & 2/\epsilon \\ 1 & 1 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}}_{\tilde{\mathbf{L}}} \underbrace{\begin{bmatrix} 2 & 2/\epsilon \\ 0 & 1-2/\epsilon \end{bmatrix}}_{\tilde{\mathbf{U}}} = \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}}_{\tilde{\mathbf{L}}} \underbrace{\begin{bmatrix} 2 & 2/\epsilon \\ 0 & -2/\epsilon \end{bmatrix}}_{\tilde{\mathbf{U}}} \quad \text{in } \mathbf{M}.$$

Using the rules of arithmetic in  $\mathbf{M}$  (→ Exp. 1.5.3.14), we find

$$\tilde{\mathbf{U}}^{-1} \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} = -\frac{1}{2} \begin{bmatrix} -1 & -1 \\ 0 & \epsilon \end{bmatrix} \frac{2}{\epsilon} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix},$$

which is not an acceptable result.

### Pivoting: Theoretical perspective

#### Definition 2.3.3.11. Permutation matrix

An  $n$ -permutation,  $n \in \mathbb{N}$ , is a bijective mapping  $\pi : \{1, \dots, n\} \mapsto \{1, \dots, n\}$ . The corresponding permutation matrix  $\mathbf{P}_\pi \in \mathbb{K}^{n,n}$  is defined by

$$(\mathbf{P}_\pi)_{ij} = \begin{cases} 1 & , \text{ if } j = \pi(i) , \\ 0 & \text{else.} \end{cases}$$

$$\text{permutation } (1, 2, 3, 4) \mapsto (1, 3, 2, 4) \hat{=} \mathbf{P} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} .$$

Note:

- ◆  $\mathbf{P}^\top = \mathbf{P}^{-1}$  for any permutation matrix  $\mathbf{P}$  ( $\rightarrow$  permutation matrices orthogonal/unitary)
- ◆  $\mathbf{P}_\pi \mathbf{A}$  effects  $\pi$ -permutation of rows of  $\mathbf{A} \in \mathbb{K}^{n,m}$
- ◆  $\mathbf{A}\mathbf{P}_\pi$  effects  $\pi$ -permutation of columns of  $\mathbf{A} \in \mathbb{K}^{m,n}$

#### Lemma 2.3.3.12. Existence of LU-factorization with pivoting $\rightarrow$ [DR08, Thm. 3.25], [Han02, Thm. 4.4]

For any regular  $\mathbf{A} \in \mathbb{K}^{n,n}$  there is a permutation matrix ( $\rightarrow$  Def. 2.3.3.11)  $\mathbf{P} \in \mathbb{K}^{n,n}$ , a normalized lower triangular matrix  $\mathbf{L} \in \mathbb{K}^{n,n}$ , and a regular upper triangular matrix  $\mathbf{U} \in \mathbb{K}^{n,n}$  ( $\rightarrow$  Def. 1.1.2.3), such that  $\mathbf{PA} = \mathbf{LU}$ .

*Proof.* (by induction)

Every regular matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$  admits a row permutation encoded by the permutation matrix  $\mathbf{P} \in \mathbb{K}^{n,n}$ , such that  $\mathbf{A}' := (\mathbf{A})_{1:n-1, 1:n-1}$  is regular (why?).

By induction assumption there is a permutation matrix  $\mathbf{P}' \in \mathbb{K}^{n-1,n-1}$  such that  $\mathbf{P}'\mathbf{A}'$  possesses a LU-factorization  $\mathbf{A}' = \mathbf{L}'\mathbf{U}'$ . There are  $\mathbf{x}, \mathbf{y} \in \mathbb{K}^{n-1}$ ,  $\gamma \in \mathbb{K}$  such that

$$\begin{bmatrix} \mathbf{P}' & 0 \\ 0 & 1 \end{bmatrix} \mathbf{PA} = \begin{bmatrix} \mathbf{P}' & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{A}' & \mathbf{x} \\ \mathbf{y}^\top & \gamma \end{bmatrix} = \begin{bmatrix} \mathbf{L}'\mathbf{U}' & \mathbf{x} \\ \mathbf{y}^\top & \gamma \end{bmatrix} = \begin{bmatrix} \mathbf{L}' & 0 \\ \mathbf{c}^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{U} & \mathbf{d} \\ 0 & \alpha \end{bmatrix} ,$$

if we choose

$$\mathbf{d} = (\mathbf{L}')^{-1}\mathbf{x} , \quad \mathbf{c} = (\mathbf{u}')^{-T}\mathbf{y} , \quad \alpha = \gamma - \mathbf{c}^\top \mathbf{d} ,$$

which is always possible.  $\square$

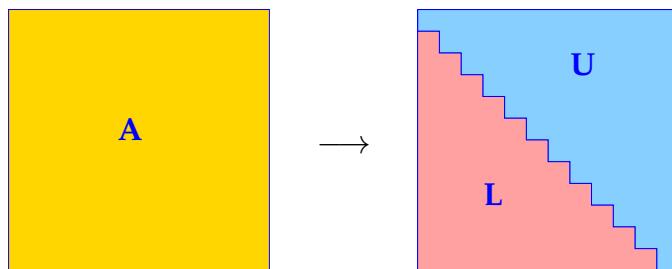
#### EXAMPLE 2.3.3.13 (Ex. 2.3.3.4 cont'd)

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 2 \\ 2 & -3 & 2 \\ 1 & 24 & 0 \end{bmatrix} \xrightarrow{\textcircled{1}} \begin{bmatrix} 2 & -3 & 2 \\ 1 & 2 & 2 \\ 1 & 24 & 0 \end{bmatrix} \xrightarrow{\textcircled{2}} \begin{bmatrix} 2 & -3 & 2 \\ 0 & 3.5 & 1 \\ 0 & 25.5 & -1 \end{bmatrix} \xrightarrow{\textcircled{3}} \begin{bmatrix} 2 & -7 & 2 \\ 0 & 25.5 & -1 \\ 0 & 3.5 & 1 \end{bmatrix} \xrightarrow{\textcircled{4}} \begin{bmatrix} 2 & -7 & 2 \\ 0 & 25.5 & -1 \\ 0 & 0 & 1.373 \end{bmatrix}$$

$$\mathbf{U} = \begin{bmatrix} 2 & -3 & 2 \\ 0 & 25.5 & -1 \\ 0 & 0 & 1.1373 \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0.5 & 0.1373 & 1 \end{bmatrix}, \quad \mathbf{P} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}.$$

Two permutations: in step ❶ swap rows #1 and #2, in step ❷ swap rows #2 and #3. Apply these swaps to the identity matrix and you will recover  $\mathbf{P}$ . See also [DR08, Ex. 3.30]. □

**§2.3.3.14 (LU-decomposition in EIGEN)** EIGEN provides various functions for computing the LU-decomposition of a given matrix. They all perform the factorization **in-situ** → Rem. 2.3.2.11:



The resulting matrix can be retrieved and used to recover the LU-factors, as demonstrated in the next code snippet. Note that the method `matrixLU` returns just a single matrix, from which the LU-factors have to be extracted using special view methods.

**C++ code 2.3.3.15: Performing explicit LU-factorization in EIGEN → GITLAB**

```
const Eigen::MatrixXd::Index n = A.cols();
assert(n == A.rows()); // ensure square matrix
Eigen::PartialPivLU<MatrixXd> lu(A);
// Normalized lower-triangular factor
MatrixXd L = MatrixXd::Identity(n, n);
L.triangularView<StrictlyLower>() += lu.matrixLU();
// Upper triangular factor
MatrixXd U = lu.matrixLU().triangularView<Upper>();
// Permutation matrix, see Def. 2.3.3.11
MatrixXd P = lu.permutationP();
```

Note that for solving a linear system of equations by means of LU-decomposition (the standard algorithm) we never have to extract the LU-factors. □

**Remark 2.3.3.16 (Row swapping commutes with forward elimination)** Any kind of pivoting only involves comparisons and row/column permutations, but no arithmetic operations on the matrix entries. This makes the following observation plausible:

The LU-factorization of  $\mathbf{A} \in \mathbb{K}^{n,n}$  with partial pivoting by § 2.3.3.8 is *numerically equivalent* to the LU-factorization of  $\mathbf{PA}$  without pivoting (→ Code in § 2.3.2.6), when  $\mathbf{P}$  is a permutation matrix gathering the row swaps entailed by partial pivoting.

*numerically equivalent*  $\hat{=}$  same result when executed with the same machine arithmetic

- The above statement means that whenever we study the impact of roundoff errors on LU-factorization it is safe to consider only the basic version without pivoting, because we can always assume that row swaps have been conducted beforehand.

## 2.4 Stability of Gaussian Elimination

It will turn out that when investigating the stability of algorithms meant to solve linear systems of equations, a key quantity is the residual.

### Definition 2.4.0.1. Residual

Given an approximate solution  $\tilde{\mathbf{x}} \in \mathbb{K}^n$  of the LSE  $\mathbf{A}\mathbf{x} = \mathbf{b}$  ( $\mathbf{A} \in \mathbb{K}^{n,n}$ ,  $\mathbf{b} \in \mathbb{K}^n$ ), its **residual** is the vector

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}.$$

**§2.4.0.2 (Probing stability of a direct solver for LSE)** Assume that you have downloaded a direct solver for a general (dense) linear system of equations  $\mathbf{A}\mathbf{x} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{K}^{n,n}$  regular,  $\mathbf{b} \in \mathbb{K}^n$ . When given the data  $\mathbf{A}$  and  $\mathbf{b}$  it returns the perturbed solution  $\tilde{\mathbf{x}}$ . How can we tell that  $\tilde{\mathbf{x}}$  is the exact solution of a linear system with slightly perturbed data (in the sense of a tiny relative error of size  $\approx \text{EPS}$ ,  $\text{EPS}$  the machine precision, see § 1.5.3.8). That is, how can we tell that  $\tilde{\mathbf{x}}$  is an **acceptable** solution in the sense of backward error analysis, cf. Def. 1.5.5.19. A similar question was explored in Ex. 1.5.5.20 for matrix  $\times$  vector multiplication.

- ①  $\mathbf{x} - \tilde{\mathbf{x}}$  accounted for by perturbation of right hand side:

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{b} \\ \mathbf{A}\tilde{\mathbf{x}} &= \mathbf{b} + \Delta\mathbf{b} \end{aligned} \quad \Rightarrow \quad \Delta\mathbf{b} = \mathbf{A}\tilde{\mathbf{x}} - \mathbf{b} =: -\mathbf{r} \quad (\text{residual, Def. 2.4.0.1}).$$

Hence,  $\tilde{\mathbf{x}}$  can be accepted as a solution, if  $\frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq Cn^3 \cdot \text{EPS}$ , for some small constant  $C \approx 1$ , see Def. 1.5.5.19. Here,  $\|\cdot\|$  can be *any* vector norm on  $\mathbb{K}^n$ .

- ②  $\mathbf{x} - \tilde{\mathbf{x}}$  accounted for by perturbation of system matrix:

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{b}, \quad (\mathbf{A} + \Delta\mathbf{A})\tilde{\mathbf{x}} = \mathbf{b} \\ [\quad \text{try perturbation} \quad \Delta\mathbf{A} = \mathbf{u}\tilde{\mathbf{x}}^H, \quad \mathbf{u} \in \mathbb{K}^n \quad] \\ \blacktriangleright \quad \mathbf{u} &= \frac{\mathbf{r}}{\|\mathbf{x}\|_2^2} \quad \Rightarrow \quad \Delta\mathbf{A} = \frac{\mathbf{r}\tilde{\mathbf{x}}^H}{\|\mathbf{x}\|_2^2}. \end{aligned}$$

As in Ex. 1.5.5.20 we find

$$\frac{\|\Delta\mathbf{A}\|_2}{\|\mathbf{A}\|_2} = \frac{\|\mathbf{r}\|}{\|\mathbf{A}\|_2 \|\tilde{\mathbf{x}}\|_2} \leq \frac{\|\mathbf{r}\|_2}{\|\mathbf{A}\tilde{\mathbf{x}}\|_2}.$$

Thus,  $\tilde{\mathbf{x}}$  is ok in the sense of backward error analysis, if  $\frac{\|\mathbf{r}\|}{\|\mathbf{A}\tilde{\mathbf{x}}\|} \leq Cn^3 \cdot \text{EPS}$ .

A stable algorithm for solving an LSE yields a **residual**  $\mathbf{r} := \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$  small (in norm) **relative to  $\mathbf{b}$** .

The roundoff error analysis of Gaussian elimination based Ass. 1.5.3.11 is rather involved. Here we merely summarise the results:

Simplification: equivalence of Gaussian elimination and LU-factorization extends to machine arithmetic, cf. Sect. 2.3.2

#### Lemma 2.4.0.3. Equivalence of Gaussian elimination and LU-factorization

The following algorithms for solving the LSE  $\mathbf{Ax} = \mathbf{b}$  ( $\mathbf{A} \in \mathbb{K}^{n,n}$ ,  $\mathbf{b} \in \mathbb{K}^n$ ) are numerically equivalent:

- ① Gaussian elimination (forward elimination and back substitution) without pivoting, see Algorithm 2.3.1.3.
- ② LU-factorization of  $\mathbf{A}$  ( $\rightarrow$  Code in § 2.3.2.6) followed by forward and backward substitution, see Algorithm 2.3.2.15.

Rem. 2.3.3.16  $\Rightarrow$  sufficient to consider LU-factorization without pivoting

A profound roundoff analysis of Gaussian elimination/LU-factorization can be found in [GV89, Sect. 3.3 & 3.5] and [Hig02, Sect. 9.3]. A less rigorous, but more lucid discussion is given in [TB97, Lecture 22].

Here we only quote a result due to Wilkinson, [Hig02, Thm. 9.5]:

#### Theorem 2.4.0.4. Stability of Gaussian elimination with partial pivoting

Let  $\mathbf{A} \in \mathbb{R}^{n,n}$  be regular and  $\mathbf{A}^{(k)} \in \mathbb{R}^{n,n}$ ,  $k = 1, \dots, n-1$ , denote the intermediate matrix arising in the  $k$ -th step of § 2.3.3.8 (Gaussian elimination with partial pivoting) when carried out with exact arithmetic.

For the approximate solution  $\tilde{\mathbf{x}} \in \mathbb{R}^n$  of the LSE  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{b} \in \mathbb{R}^n$ , computed as in § 2.3.3.8 (based on machine arithmetic with machine precision  $\text{EPS}$ ,  $\rightarrow$  Ass. 1.5.3.11) there is  $\Delta\mathbf{A} \in \mathbb{R}^{n,n}$  with

$$\|\Delta\mathbf{A}\|_\infty \leq n^3 \frac{3\text{EPS}}{1 - 3n\text{EPS}} \rho \|\mathbf{A}\|_\infty, \quad \rho := \frac{\max_{i,j,k} |(\mathbf{A}^{(k)})_{ij}|}{\max_{i,j} |(\mathbf{A})_{ij}|},$$

such that  $(\mathbf{A} + \Delta\mathbf{A})\tilde{\mathbf{x}} = \mathbf{b}$ .

$\rho$  “small”  $\rightarrow$  Gaussian elimination with partial pivoting is stable ( $\rightarrow$  Def. 1.5.5.19)

If  $\rho$  is “small”, the computed solution of a LSE can be regarded as the exact solution of a LSE with “slightly perturbed” system matrix (perturbations of size  $O(n^3\text{EPS})$ ).



Bad news:

exponential growth  $\rho \sim 2^n$  is possible !

**EXAMPLE 2.4.0.5 (Wilkinson's counterexample)** We confirm the bad news by means of a famous ex-

ample, known as the so-called Wilkinson matrix.

$$n=10: \quad a_{ij} = \begin{cases} 1 & , \text{ if } i = j \vee j = n , \\ -1 & , \text{ if } i > j , \\ 0 & \text{ else.} \end{cases}, \quad \mathbf{A} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & 1 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & 1 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & 1 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 \end{bmatrix}$$

Partial pivoting does not trigger row permutations !

$$\Rightarrow \mathbf{A} = \mathbf{LU}, \quad l_{ij} = \begin{cases} 1 & , \text{ if } i = j , \\ -1 & , \text{ if } i > j , \\ 0 & \text{ else} \end{cases}, \quad u_{ij} = \begin{cases} 1 & , \text{ if } i = j , \\ 2^{i-1} & , \text{ if } j = n , \\ 0 & \text{ else.} \end{cases}$$

► Exponential blow-up of entries of  $\mathbf{U}$  !

Blow-up of entries of  $\mathbf{U}$  !

$\Updownarrow (*)$   
However,  $\text{cond}_2(\mathbf{A})$  is small!

► Evidence of **Instability of Gaussian elimination!**

#### C++ code 2.4.0.6: Gaussian elimination for “Wilkinson system” in EIGEN → GITLAB

```

2 MatrixXd res(100,2);
3 for(int n = 10; n <= 100*10; n += 10){
4     MatrixXd A(n,n); A.setIdentity();
5     A.triangularView<StrictlyLower>().setConstant(-1);
6     A.rightCols<1>().setOnes();
7     VectorXd x = VectorXd::Constant(n,-1).binaryExpr( VectorXd::LinSpaced(n,1,n),
8         [] (double x, double y){return pow(x,y);} );
9     double relerr = (A.lu().solve(A*x)-x).norm()/x.norm();
10    res(n/10-1,0) = n; res(n/10-1,1) = relerr;
11 }
// ... different solver(e.g. colPivHouseholderQr()), plotting

```

- (\*) If  $\text{cond}_2(\mathbf{A})$  was huge, then big errors in the solution of a linear system can be caused by small perturbations of either the system matrix or the right hand side vector, see (2.4.0.12) and the message of Thm. 2.2.2.4, (2.2.2.8). In this case, a stable algorithm can obviously produce a grossly “wrong” solution, as was already explained after (2.2.2.8).

Hence, lack of stability of Gaussian elimination will only become apparent for linear systems with well-conditioned system matrices.

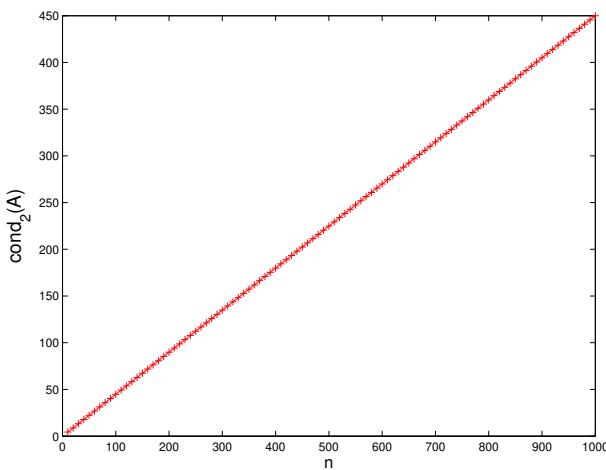
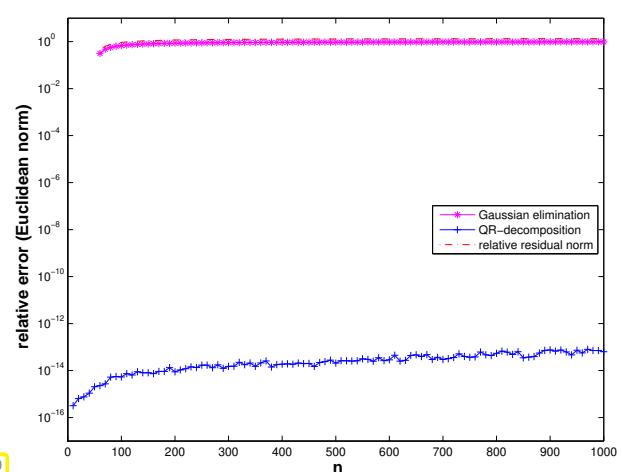


Fig. 39



These observations match Thm. 2.4.0.4, because in this case we encounter an exponential growth of  $\rho = \rho(n)$ , see Ex. 2.4.0.5.  $\square$



Observation: In practice  $\rho$  (almost) always grows only mildly (like  $O(\sqrt{n})$ ) with  $n$

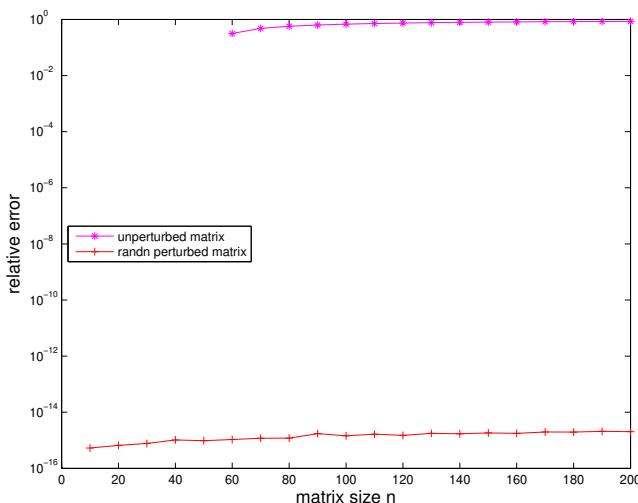
Discussion in [TB97, Lecture 22]: growth factors larger than the order  $O(\sqrt{n})$  are exponentially rare in certain relevant classes of *random matrices*.

**EXAMPLE 2.4.0.7 (Stability by small random perturbations)** Spielman and Teng [ST96] discovered that a tiny relative random perturbation of the Wilkinson matrix on the scale of the machine precision EPS ( $\rightarrow$  § 1.5.3.8) already remedies the instability of Gaussian elimination.

#### C++ code 2.4.0.8: Stabilization of Gaussian elimination with partial pivoting by **small** random perturbations $\Rightarrow$ [GITLAB](#)

```

2  //! Curing Wilkinson's counterexample by random perturbation
3  MatrixXd res(20, 3);
4  mt19937 gen(42); // seed
5  // normal distribution, mean = 0.0, stddev = 1.0
6  std::normal_distribution<> bellcurve;
7  for (int n = 10; n <= 10 * 20; n += 10) {
8      // Build Wilkinson matrix
9      MatrixXd A(n, n); A.setIdentity();
10     A.triangularView<StrictlyLower>().setConstant(-1);
11     A.rightCols<1>().setOnes();
12     // imposed solution
13     VectorXd x = VectorXd::Constant(n, -1).binaryExpr(
14         VectorXd::LinSpaced(n, 1, n),
15         [] (double x, double y) { return pow(x, y); });
16     double relerr = (A.lu().solve(A * x) - x).norm() / x.norm();
17     // Randomly perturbed Wilkinson matrix by matrix with iid
18     // N(0,eps) distributed entries
19     MatrixXd Ap = A.unaryExpr([&](double x) {
20         return x + numeric_limits<double>::epsilon() * bellcurve(gen);
21     });
22     double relerrp = (Ap.lu().solve(Ap * x) - x).norm() / x.norm();
23     res(n / 10 - 1, 0) = n;
24     res(n / 10 - 1, 1) = relerr;
25     res(n / 10 - 1, 2) = relerrp;
26 }
```



Recall the statement made above about “improbability” of matrices for which Gaussian elimination with partial pivoting is unstable. This is now matched by the observation that a *tiny random* perturbation of the matrix (almost certainly) cures the problem. This is investigated by the brand-new field of **smoothed analysis** of numerical algorithms, see [SST06].

*Gaussian elimination/LU-factorization with partial pivoting is stable (\*)  
(for all practical purposes) !*

(\*): stability refers to maximum norm  $\|\cdot\|_\infty$ .

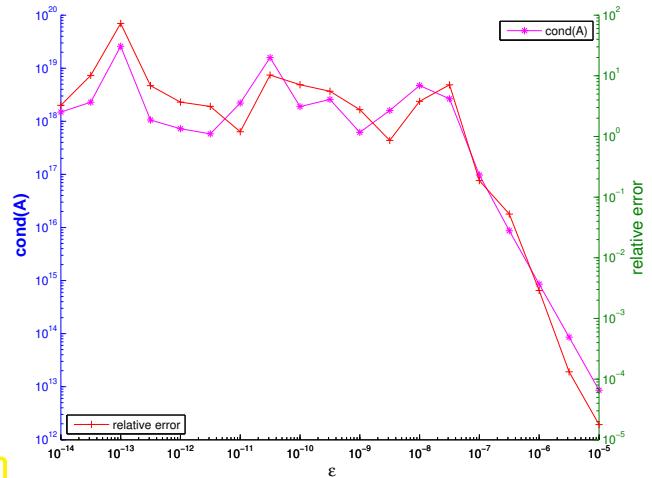
#### EXPERIMENT 2.4.0.9 (Conditioning and relative error → Ex. 2.4.0.10 cnt'd)

In the discussion of numerical stability (→ Def. 1.5.5.19, Rem. 1.5.5.22) we have seen that a stable algorithm may produce results with large errors for ill-conditioned problems. The conditioning of the problem of solving a linear system of equations is determined by the condition number (→ Def. 2.2.2.7) of the system matrix, see Thm. 2.2.2.4.

Hence, for an **ill-conditioned** linear system, whose system matrix is beset with a huge condition number, (stable) Gaussian elimination may return “solutions” with large errors. This will be demonstrated in this experiment.

Numerical experiment with *nearly singular matrix* from Ex. 2.4.0.10

$$\begin{aligned} \mathbf{A} &= \mathbf{u}\mathbf{v}^T + \epsilon \mathbf{I}, \\ \mathbf{u} &= \frac{1}{3}(1, 2, 3, \dots, 10)^T, \\ \mathbf{v} &= (-1, \frac{1}{2}, -\frac{1}{3}, \frac{1}{4}, \dots, \frac{1}{10})^T \end{aligned}$$



The practical stability of Gaussian elimination for  $\mathbf{Ax} = \mathbf{b}$  is reflected by the size of a particular vector, the residual  $\mathbf{r} := \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$ ,  $\tilde{\mathbf{x}}$  the computed solution, that can easily be computed after the elimination solver has finished:

In practice *Gaussian elimination/LU-factorization with partial pivoting*  
produces “relatively small residual vectors”

Simple consideration as in § 2.4.0.2:

$$(\mathbf{A} + \Delta\mathbf{A})\tilde{\mathbf{x}} = \mathbf{b} \Rightarrow \mathbf{r} = \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}} = \Delta\mathbf{A}\tilde{\mathbf{x}} \Rightarrow \|\mathbf{r}\| \leq \|\Delta\mathbf{A}\|\|\tilde{\mathbf{x}}\|,$$

for any vector norm  $\|\cdot\|$ . This means that, if a direct solver for an LSE is stable in the sense of backward error analysis, that is, the perturbed solution could be obtained as the exact solution for a slightly relatively perturbed system matrix, then the *residual will be (relatively) small*.

**EXPERIMENT 2.4.0.10 (Small residuals by Gaussian elimination)** Gaussian elimination works miraculously about delivering small residuals!

$$\mathbf{A} = \mathbf{u}\mathbf{v}^T + \epsilon\mathbf{I},$$

↑  
singular rank-1 matrix

with

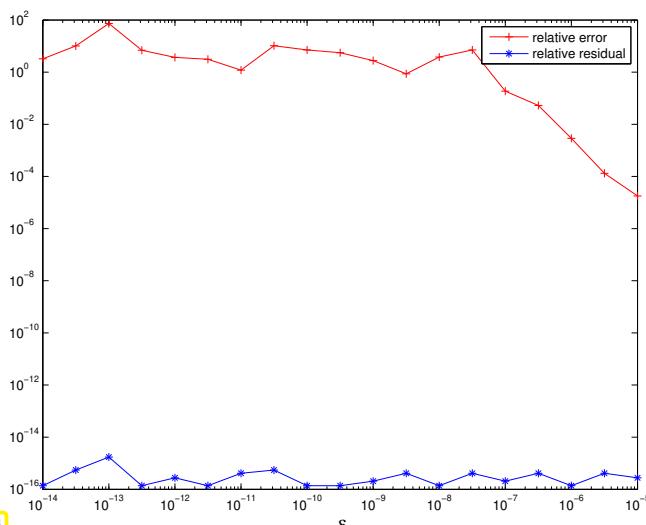
$$\mathbf{u} = \frac{1}{3}(1, 2, 3, \dots, 10)^T,$$

$$\mathbf{v} = (-1, \frac{1}{2}, -\frac{1}{3}, \frac{1}{4}, \dots, \frac{1}{10})^T$$

#### C++ code 2.4.0.11: Small residuals for Gauss elimination → GITLAB

```

2 int n = 10;
3 // Initialize vectors u and v
4 VectorXd u = VectorXd::LinSpaced(n, 1, n) / 3.0;
5 VectorXd v = u.cwiseInverse().array() *
6     VectorXd::LinSpaced(n, 1, n)
7     .unaryExpr([](double x) { return pow(-1, x); })
8     .array();
9 VectorXd x = VectorXd::Ones(n);
10 double nx = x.lpNorm<Infinity>();
11 VectorXd expo = VectorXd::LinSpaced(19, -5, -14);
12 Eigen::MatrixXd res(expo.size(), 4);
13 for (int i = 0; i <= expo.size(); ++i) {
14     // Build coefficient matrix A
15     double epsilon = std::pow(10, expo(i));
16     MatrixXd A = u * v.transpose() + epsilon * MatrixXd::Identity(n, n);
17     VectorXd b = A * x; // right-hand-side vector
18     double nb = b.lpNorm<Infinity>(); // maximum norm
19     VectorXd xt = A.lu().solve(b); // Gaussian elimination
20     VectorXd r = b - A * xt; // residual vector
21     res(i, 0) = epsilon;
22     res(i, 1) = (x - xt).lpNorm<Infinity>() / nx;
23     res(i, 2) = r.lpNorm<Infinity>() / nb;
24     // L-infinity condition number
25     res(i, 3) = A.inverse().cwiseAbs().rowwise().sum().maxCoeff() *
26                 A.cwiseAbs().rowwise().sum().maxCoeff();
27 }
```



Observations (w.r.t  $\|\cdot\|_\infty$ -norm)

- ◆ for  $\epsilon \ll 1$  large relative error in computed solution  $\tilde{\mathbf{x}}$
- ◆ small residuals for any  $\epsilon$

How can a *large* relative error be reconciled with a *small* relative residual ?

$$\begin{aligned} \mathbf{Ax} = \mathbf{b} &\leftrightarrow \mathbf{A}\tilde{\mathbf{x}} \approx \mathbf{b} \\ \left\{ \begin{array}{l} \mathbf{A}(\mathbf{x} - \tilde{\mathbf{x}}) = \mathbf{r} \\ \mathbf{Ax} = \mathbf{b} \end{array} \right. \Rightarrow \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}^{-1}\| \|\mathbf{r}\| \quad \Rightarrow \quad \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}. \end{aligned} \quad (2.4.0.12)$$

➤ If  $\text{cond}(\mathbf{A}) := \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \gg 1$ , then a small relative residual may not imply a small relative error. Also recall the discussion in Exp. 2.4.0.9. □

**EXPERIMENT 2.4.0.13 (Instability of multiplication with inverse)** An important justification for Rem. 2.2.1.6 that advised us not to compute the inverse of a matrix in order to solve a linear system of equations is conveyed by this experiment. We again consider the nearly singular matrix from Exp. 2.4.0.10.

$$\mathbf{A} = \mathbf{uv}^T + \epsilon \mathbf{I}, \quad \begin{matrix} \uparrow \\ \text{singular rank-1 matrix} \end{matrix} \quad \text{with} \quad \mathbf{u} = \frac{1}{3}(1, 2, 3, \dots, 10)^T, \\ \mathbf{v} = (-1, \frac{1}{2}, -\frac{1}{3}, \frac{1}{4}, \dots, \frac{1}{10})^T$$

#### C++ code 2.4.0.14: Instability of multiplication with inverse → GITLAB

```

2 int n = 10;
3 VectorXd u = VectorXd::LinSpaced(n, 1, n) / 3.0;
4 VectorXd v = u.cwiseInverse().array() * VectorXd::LinSpaced(n, 1, n).unaryExpr(
5     [] (double x){return pow(-1,x);} ).array();
6 VectorXd x = VectorXd::Ones(n);
7 VectorXd expo = VectorXd::LinSpaced(19, -5, -14);
8 MatrixXd res(expo.size(), 4);
9 for(int i = 0; i < expo.size(); ++i){
10    double epsilon = std::pow(10, expo(i));
11    MatrixXd A = u*v.transpose() + epsilon * (MatrixXd::Random(n,n) +
12        MatrixXd::Ones(n,n))/2;
13    VectorXd b = A * x;
14    double nb = b.lpNorm<Infinity>();
15    VectorXd xt = A.lu().solve(b); // stable solving
16    VectorXd r = b - A*xt; // residual
17    MatrixXd B = A.inverse();
18    VectorXd xi = B*b; // solving via inverse
19    VectorXd ri = b - A*xi; // residual
20    MatrixXd R = MatrixXd::Identity(n,n) - A*B; // residual
21    res(i,0) = epsilon; res(i,1) = (r).lpNorm<Infinity>()/nb;
22    res(i,2) = ri.lpNorm<Infinity>()/nb;
23    // L-infinity condition number

```

```

22     res(i,3) = R.IpNorm<Infinity>() / B.IpNorm<Infinity>();
23 }
  
```

The computation of the inverse is affected by round-off errors, but does not benefit from the same favorable cancellation of roundoff errors as Gaussian elimination.

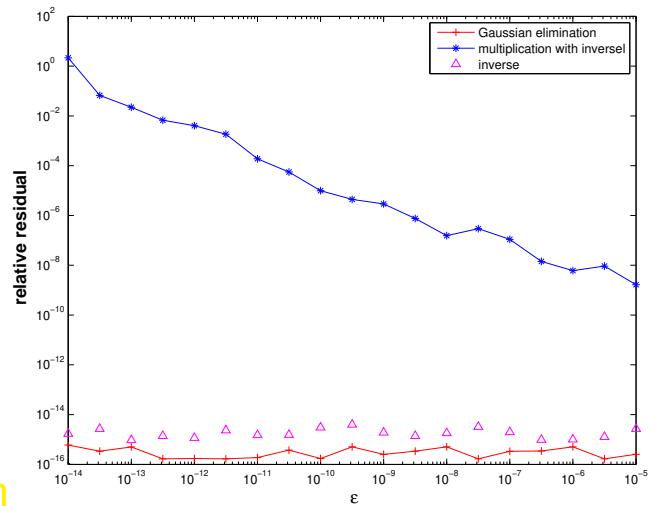


Fig. 44

## 2.5 Survey: Elimination solvers for linear systems of equations

All direct (\*) solver algorithms for square linear systems of equations  $\mathbf{Ax} = \mathbf{b}$  with given matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$ , right hand side vector  $\mathbf{b} \in \mathbb{K}^n$  and unknown  $\mathbf{x} \in \mathbb{K}^n$  rely on variants of Gaussian elimination with **pivoting**, see Section 2.3.3. Sophisticated, optimised and verified implementations are available in numerical libraries like LAPACK/MKL.

(\*): a direct solver terminates after a predictable finite number of elementary operations for every admissible input.

Never contemplate implementing a general solver for linear systems of equations!

If possible, use algorithms from numerical libraries! ( $\rightarrow$  Exp. 2.3.1.7)

Therefore, familiarity with details of Gaussian elimination is not required, but one must know when and how to use the library functions and one must be able to assess the computational effort they involve.

**§2.5.0.1 (Computational effort for direct elimination)** We repeat the reasoning of § 2.3.1.5: Gaussian elimination for a general (dense) matrix invariably involves *three nested loops* of length  $n$ , see Code 2.3.1.4, Code 2.3.3.6.

**Theorem 2.5.0.2. Cost of Gaussian elimination**  $\rightarrow$  § 2.3.1.5

Given a linear system of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{K}^{n,n}$  regular,  $\mathbf{b} \in \mathbb{K}^n$ ,  $n \in \mathbb{N}$ , the asymptotic computational effort ( $\rightarrow$  Def. 1.4.0.1) for its direct solution by means of Gaussian elimination in terms of the problem size parameter  $n$  is  $O(n^3)$  for  $n \rightarrow \infty$ .

The constant hidden in the Landau symbol can be expected to be rather small ( $\approx 1$ ) as is clear from

(2.3.1.6).

The cost for solving are substantially lower, if certain properties of the matrix  $\mathbf{A}$  are known. This is clear, if  $\mathbf{A}$  is diagonal or orthogonal/unitary. It is also true for triangular matrices ( $\rightarrow$  Def. 1.1.2.3), because they can be solved by simple back substitution or forward elimination. We recall the observation made in see § 2.3.2.15.

### Theorem 2.5.0.3. Cost for solving triangular systems $\rightarrow$ § 2.3.1.5

In the setting of Thm. 2.5.0.2, the asymptotic computational effort for solving a *triangular* linear system of equations is  $O(n^2)$  for  $n \rightarrow \infty$ .

**§2.5.0.4 (Direct solution of linear systems of equations in EIGEN)** EIGEN supplies a rich suite of functions for matrix decompositions and solving LSEs, see  [EIGEN documentation](#). The default solver is Gaussian elimination with partial pivoting, accessible through the methods `lu()` and `solve()` of dense matrix types:

Given: system/coefficient matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$  regular  $\leftrightarrow$  A ( $n \times n$  EIGEN matrix)  
right hand side vectors  $\mathbf{B} \in \mathbb{K}^{n,\ell}$   $\leftrightarrow$  B ( $n \times \ell$  EIGEN matrix)  
(corresponds to multiple right hand sides, cf. Code 2.3.1.10)

linear algebra	EIGEN
$\mathbf{X} = \mathbf{A}^{-1}\mathbf{B} = [\mathbf{A}^{-1}(\mathbf{B})_{:,1}, \dots, \mathbf{A}^{-1}(\mathbf{B})_{:,\ell}]$	<code>X = A.lu().solve(B)</code>

Summarizing the detailed information given in § 2.3.2.15:

$$\text{cost}(X = A.lu().solve(B)) = O(n^3 + ln^2) \quad \text{for } n, l \rightarrow \infty$$

**Remark 2.5.0.5 (Communicating special properties of system matrices in EIGEN)** Sometimes, the coefficient matrix of a linear system of equations is known to have certain analytic properties that a direct solver can exploit to perform elimination more efficiently. These properties may even be impossible to detect by an algorithm, because matrix entries that should vanish exactly might have been perturbed due to roundoff.

Thus one needs to pass EIGEN these informations as follows:

```

2 // A is lower triangular
3 x = A.triangularView<Eigen::Lower>().solve(b);
4 // A is upper triangular
5 x = A.triangularView<Eigen::Upper>().solve(b);
6 // A is hermitian / self adjoint and positive definite
7 x = A.selfadjointView<Eigen::Upper>().lLt() .solve(b);
8 // A is hermitian / self adjoint (positive or negative semidefinite)
9 x = A.selfadjointView<Eigen::Upper>().ldLt() .solve(b);
```

The methods `lLt()` and `ldLt()` rely on special factorizations for *symmetric matrices*, see § 2.8.0.13 below.

**EXPERIMENT 2.5.0.6 (Standard EIGEN lu() operator versus triangularView() )** In this numerical experiment we study the gain in efficiency achievable by make the direct solver aware of important matrix properties.

### C++ code 2.5.0.7: Direct solver applied to a upper triangular matrix → GITLAB

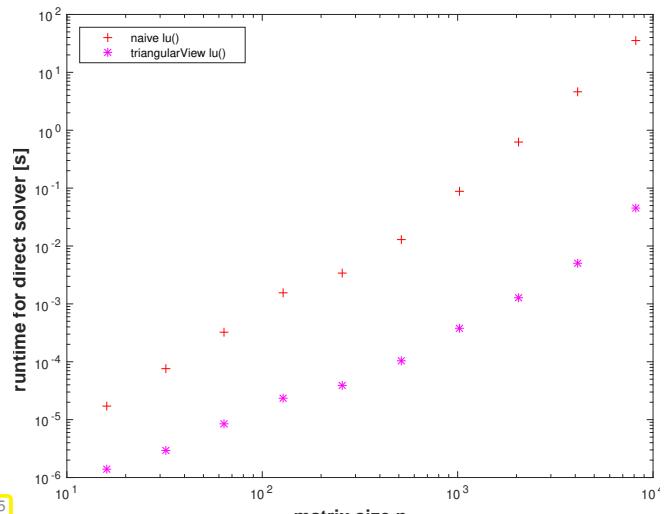
```

2  ///! Eigen code: assessing the gain from using special properties
3  ///! of system matrices in Eigen
4  MatrixXd timing(){
5      std::vector<int> n = {16,32,64,128,256,512,1024,2048,4096,8192};
6      int nrungs = 3;
7      MatrixXd times(n.size(),3);
8      for(int i = 0; i < n.size(); ++i){
9          Timer t1, t2; // timer class
10         MatrixXd A = VectorXd::LinSpaced(n[i],1,n[i]).asDiagonal();
11         A += MatrixXd::Ones(n[i],n[i]).triangularView<Upper>();
12         VectorXd b = VectorXd::Random(n[i]);
13         VectorXd x1(n[i]), x2(n[i]);
14         for(int j = 0; j < nrungs; ++j){
15             t1.start(); x1 = A.lu().solve(b); t1.stop();
16             t2.start(); x2 = A.triangularView<Upper>().solve(b); t2.stop();
17         }
18         times(i,0) = n[i]; times(i,1) = t1.min(); times(i,2) = t2.min();
19     }
20     return times;
21 }
```

Observation:



Being told that only the upper triangular part of the matrix needs to be taken into account, Gaussian elimination reduces to cheap backward elimination, which is much faster than full elimination, cf. Thm. 2.5.0.2 vs. Thm. 2.5.0.3.



**§2.5.0.8 (Direct solvers for LSE in EIGEN)** Invocation of direct solvers in EIGEN is a two stage process:

- ① Request a decomposition (LU,QR,LDLT) of the matrix and store it in a temporary “decomposition object”.
- ② Perform backward & forward substitutions by calling the `solve()` method of the decomposition object.

The general format for invoking linear solvers in EIGEN is as follows:

```
Eigen::SolverType<Eigen::MatrixXd> solver(A);
Eigen::VectorXd x = solver.solve(b);
```

This can be reduced to one line, as the solvers can also be used as methods acting on matrices:

```
Eigen::VectorXd x = A.solverType().solve(b);
```

A full list of solvers can be found in the [EIGEN documentation](#). The next code demonstrates a few of the available decompositions that can serve as the basis for a linear solver:

**C++-code 2.5.0.9: EIGEN based function solving a LSE → GITLAB**

```

2 // Gaussian elimination with partial pivoting, Code 2.3.3.6
3 void lu_solve(const MatrixXd &A, const VectorXd &b, VectorXd &x) {
4     x = A.lu().solve(b); // 'lu()' is short for 'partialPivLu()'
5 }
6
7 // Gaussian elimination with total pivoting
8 void fullpivlu_solve(const MatrixXd &A, const VectorXd &b, VectorXd &x) {
9     x = A.fullPivLu().solve(b); // total pivoting
10 }
11
12 // An elimination solver based on Householder transformations
13 void qr_solve(const MatrixXd &A, const VectorXd &b, VectorXd &x) {
14     Eigen::HouseholderQR<MatrixXd> solver(A); // see Section 3.3.3
15     x = solver.solve(b);
16 }
17
18 // Use singular value decomposition (SVD)
19 void svd_solve(const MatrixXd &A, const VectorXd &b, VectorXd &x) {
20     // SVD based solvers, see Section 3.4
21     x = A.jacobiSvd(Eigen::ComputeThinU | Eigen::ComputeThinV).solve(b);
22 }
```

The [different decompositions](#) trade speed for stability and accuracy: fully pivoted and QR-based decompositions also work for nearly singular matrices, for which the standard LU-factorization may no longer be reliable.

**Remark 2.5.0.10 (Many sequential solutions of LSE)** As we have seen in Code 2.5.0.9, EIGEN provides functions that return decompositions of matrices. For instance, we can get an object “containing” the LU-decomposition (→ Section 2.3.2) of a matrix by the following commands:

```
Eigen::MatrixXd A(n,n); // A dense square matrix object
.....
auto ludec = A.lu(); // Perform LU-decomposition and store the
factors.
```

Based on the *precomputed decompositions*, a linear system of equations with coefficient matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$  can be solved with asymptotic computational effort  $O(n^2)$ , cf. § 2.3.2.15.

The following example illustrates a special situation, in which matrix decompositions can curb computational cost:

**C++ code 2.5.0.11: Wasteful approach!**→ [GITLAB](#)

```

2 // Setting: N ≫ 1,
3 // large matrix A ∈ K^{n,n}
4 for(int j = 0; j < N; ++j){
5     x = A.lu().solve(b);
6     b = some_function(x);
7 }
```

computational effort  $O(Nn^3)$

**C++ code 2.5.0.12: Smart approach!**→ [GITLAB](#)

```

2 // Setting: N ≫ 1,
3 // large matrix A ∈ K^{n,n}
4 auto A_lu_dec = A.lu();
5 for(int j = 0; j < N; ++j){
6     x = A_lu_dec.solve(b);
7     b = some_function(x);
8 }
```

computational effort  $O(n^3 + Nn^2)$

A concrete example is the so-called [inverse power iteration](#), see Chapter 9, for which a skeleton code is

given next. It computes the iterates

$$\mathbf{x}^* := \mathbf{A}^{-1} \mathbf{x}^{(k)}, \quad \mathbf{x}^{(k+1)} := \frac{\mathbf{x}^*}{\|\mathbf{x}^*\|_2}, \quad k = 0, 1, 2, \dots, \quad (2.5.0.13)$$

#### C++-code 2.5.0.14: Efficient implementation of inverse power method in EIGEN → GITLAB

```

2 template<class VecType, class MatType>
3 VecType invpowit(const Eigen::MatrixBase<MatType> &A, double tol)
4 {
5     using index_t = typename MatType::Index;
6     using scalar_t = typename VecType::Scalar;
7     // Make sure that the function is called with a square matrix
8     const index_t n = A.cols();
9     const index_t m = A.rows();
10    eigen_assert(n == m);
11    // Request LU-decomposition
12    auto A_lu_dec = A.lu();
13    // Initial guess for inverse power iteration
14    VecType xo = VecType::Zero(n);
15    VecType xn = VecType::Random(n);
16    // Normalize vector
17    xn /= xn.norm();
18    // Terminate if relative (normwise) change below threshold
19    while ((xo-xn).norm() > xn.norm()*tol) {
20        xo = xn;
21        xn = A_lu_dec.solve(xo);
22        xn /= xn.norm();
23    }
24    return (xn);
25 }
```

The use of **Eigen**::MatrixBase<MatType> makes it possible to call `invpowit` with an expression argument:

```

2 MatrixXd A = MatrixXd::Random(n,n);
3 MatrixXd B = MatrixXd::Random(n,n);
4 VectorXd ev = invpowit<VectorXd>(A+B, tol);
```

This is necessary, because `A+B` will spawn an auxiliary object of a “strange” type determined by the **expression template** mechanism. ↴

## 2.6 Exploiting Structure when Solving Linear Systems

By “structure” of a linear system we mean prior knowledge that

- ◆ either certain entries of the system matrix vanish,
- ◆ or the system matrix is generated by a particular formula.

**§2.6.0.1 (Triangular linear systems)** Triangular linear systems are linear systems of equations whose system matrix is a *triangular matrix* (→ Def. 1.1.2.3).

Thm. 2.5.0.3 tells us that (dense) triangular linear systems can be solved by backward/forward elimination with  $O(n^2)$  asymptotic computational effort ( $n \doteq$  number of unknowns) compared to an asymptotic complexity of  $O(n^3)$  for solving a generic (dense) linear system of equations (→ Thm. 2.5.0.2, Exp. 2.5.0.6).

This is the simplest case where exploiting special structure of the system matrix leads to faster algorithms for the solution of a special class of linear systems.  $\square$

**§2.6.0.2 (Block elimination)** Remember that thanks to the possibility to compute the matrix product in a block-wise fashion ( $\rightarrow$  § 1.3.1.13), Gaussian elimination can be conducted on the level of matrix blocks. We recall Rem. 2.3.1.14 and Rem. 2.3.2.19.

For  $k, \ell \in \mathbb{N}$  consider the block partitioned square  $n \times n$ ,  $n := k + \ell$ , linear system

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}, \quad \mathbf{A}_{11} \in \mathbb{K}^{k,k}, \mathbf{A}_{12} \in \mathbb{K}^{k,\ell}, \mathbf{A}_{21} \in \mathbb{K}^{\ell,k}, \mathbf{A}_{22} \in \mathbb{K}^{\ell,\ell}, \quad (2.6.0.3)$$

Using block matrix multiplication (applied to the matrix  $\times$  vector product in (2.6.0.3)) we find an equivalent way to write the block partitioned linear system of equations:

$$\begin{aligned} \mathbf{A}_{11}\mathbf{x}_1 + \mathbf{A}_{12}\mathbf{x}_2 &= \mathbf{b}_1, \\ \mathbf{A}_{21}\mathbf{x}_1 + \mathbf{A}_{22}\mathbf{x}_2 &= \mathbf{b}_2. \end{aligned} \quad (2.6.0.4)$$

We assume that  $\mathbf{A}_{11}$  is *regular* (invertible) so that we can solve for  $\mathbf{x}_1$  from the first equation.

➤ By elementary algebraic manipulations (“block Gaussian elimination”) we find

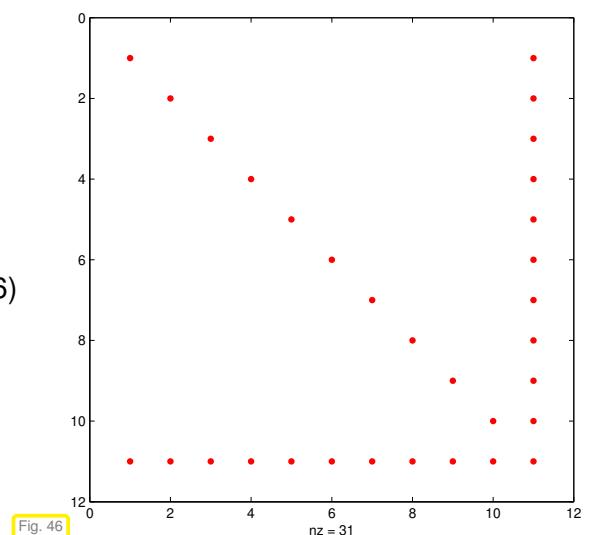
$$\begin{aligned} \mathbf{x}_1 &= \mathbf{A}_{11}^{-1}(\mathbf{b}_1 - \mathbf{A}_{12}\mathbf{x}_2), \\ \mathbf{x}_2 &= \underbrace{(\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12})}_{\text{Schur complement}} \mathbf{x}_2 = \mathbf{b}_2 - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{b}_1, \end{aligned}$$

The resulting  $\ell \times \ell$  linear system of equations for the unknown vector  $\mathbf{x}_2$  is called the **Schur complement system** for (2.6.0.3).

Unless  $\mathbf{A}$  has a special structure that allows the efficient solution of linear systems with system matrix  $\mathbf{A}_{11}$ , the Schur complement system is mainly of theoretical interest.  $\square$

**EXAMPLE 2.6.0.5 (Linear systems with arrow matrices)** From  $n \in \mathbb{N}$ , a diagonal matrix  $\mathbf{D} \in \mathbb{K}^{n,n}$ ,  $\mathbf{c} \in \mathbb{K}^n$ ,  $\mathbf{b} \in \mathbb{K}^n$ , and  $\alpha \in \mathbb{K}$ , we can build an  $(n+1) \times (n+1)$  **arrow matrix**.

$$\mathbf{A} = \left[ \begin{array}{c|c} \mathbf{D} & \mathbf{c} \\ \hline \mathbf{b}^\top & \alpha \end{array} \right] \quad (2.6.0.6)$$



We can apply the block partitioning (2.6.0.3) with  $k = n$  and  $\ell = 1$  to a linear system  $\mathbf{Ax} = \mathbf{y}$  with system matrix  $\mathbf{A}$  and obtain  $\mathbf{A}_{11} = \mathbf{D}$ , which can be inverted easily, provided that all diagonal entries of  $\mathbf{D}$  are different from zero. In this case

$$\mathbf{Ax} = \begin{bmatrix} \mathbf{D} & \mathbf{c} \\ \mathbf{b}^\top & \alpha \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \xi \end{bmatrix} = \mathbf{y} := \begin{bmatrix} \mathbf{y}_1 \\ \eta \end{bmatrix}, \quad (2.6.0.7)$$

$$\Rightarrow \begin{aligned} \xi &= \frac{\eta - \mathbf{b}^T \mathbf{D}^{-1} \mathbf{y}_1}{\alpha - \mathbf{b}^T \mathbf{D}^{-1} \mathbf{c}}, \\ \mathbf{x}_1 &= \mathbf{D}^{-1}(\mathbf{y}_1 - \xi \mathbf{c}). \end{aligned} \quad (2.6.0.8)$$

These formulas make sense, if  $\mathbf{D}$  is regular and  $\alpha - \mathbf{b}^T \mathbf{D}^{-1} \mathbf{c} \neq 0$ , which is another condition for the invertibility of  $\mathbf{A}$ .

Using the formula (2.6.0.8) we can solve the linear system (2.6.0.7) with an asymptotic complexity  $O(n)!$  This superior speed compared to Gaussian elimination applied to the (dense) linear system is evident in runtime measurements.

#### C++ code 2.6.0.9: Dense Gaussian elimination applied to arrow system → [GITLAB](#)

```

2 VectorXd arrowsys_slow(const VectorXd &d, const VectorXd &c, const VectorXd &b,
3                               const double alpha, const VectorXd &y) {
4     int n = d.size();
5     MatrixXd A(n + 1, n + 1); // Empty dense matrix
6     A.setZero(); // Initialize with all zeros.
7     A.diagonal().head(n) = d; // Initialize matrix diagonal from a vector.
8     A.col(n).head(n) = c; // Set rightmost column c.
9     A.row(n).head(n) = b; // Set bottom row b^T.
10    A(n, n) = alpha; // Set bottom-right entry alpha.
11    return A.lu().solve(y); // Gaussian elimination
12 }
```

► Asymptotic complexity  $O(n^3)!$

(Due to the serious blunder of accidentally creating a matrix full of zeros, cf. Exp. 1.3.1.10.)

#### C++ code 2.6.0.10: Solving an arrow system according to (2.6.0.8) → [GITLAB](#)

```

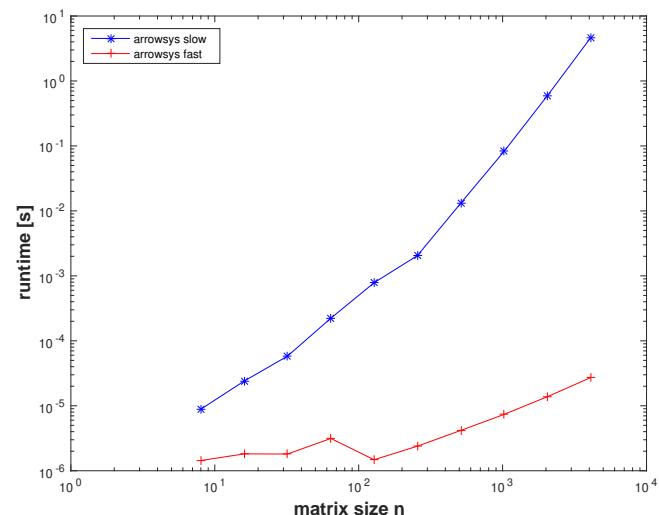
2 VectorXd arrowsys_fast(const VectorXd &d, const VectorXd &c, const VectorXd &b,
3                               const double alpha, const VectorXd &y) {
4     int n = d.size();
5     VectorXd z = c.array() / d.array(); // z = D^-1 c
6     VectorXd w = y.head(n).array() / d.array(); // w = D^-1 y_1
7     const double den = alpha - b.dot(z); // denominator in (2.6.0.8)
8     // Check for (relatively!) small denominator
9     if (std::abs(den) <
10         std::numeric_limits<double>::epsilon() * (b.norm() + std::abs(alpha))) {
11         throw std::runtime_error("Nearly singular system");
12     }
13     const double xi = (y(n) - b.dot(w)) / den;
14     return (VectorXd(n + 1) << w - xi * z, xi).finished();
15 }
```

► Asymptotic complexity  $O(n)$  for  $n \rightarrow \infty$ !

Code for Runtime measurements can be obtained from → [GITLAB](#).

(Intel i7-3517U CPU @ 1.90GHz—4, 64-bit, ubuntu 14.04 LTS, gcc 4.8.4, -O3)

No comment!



**Remark 2.6.0.11 (Sacrificing numerical stability for efficiency)** The vector based implementation of the solver of Code 2.6.0.10 can be vulnerable to roundoff errors, because, upon closer inspection, the algorithm turns out to be equivalent to Gaussian elimination **without pivoting**, cf. Section 2.3.3, Ex. 2.3.3.1.



Caution:

stability at risk in Code 2.6.0.10

Yet, there are classes of matrices for which Gaussian elimination without pivoting is guaranteed to be stable. For such matrices algorithms like that of Code 2.6.0.10 are safe.

**§2.6.0.12 (Solving LSE with low-rank modification of system matrix)** Given a regular matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$ , let us assume that at some point in a code we are in a position to solve any linear system  $\mathbf{Ax} = \mathbf{b}$  “fast”, because

- ◆ either  $\mathbf{A}$  has a favorable structure, eg. triangular, see § 2.6.0.1,
- ◆ or an LU-decomposition of  $\mathbf{A}$  is already available, see § 2.3.2.15.

Now, a  $\tilde{\mathbf{A}}$  is obtained by changing a *single entry* of  $\mathbf{A}$ :

$$\mathbf{A}, \tilde{\mathbf{A}} \in \mathbb{K}^{n,n}: \quad \tilde{a}_{ij} = \begin{cases} a_{ij} & , \text{ if } (i,j) \neq (i^*, j^*) , \\ z + a_{ij} & , \text{ if } (i,j) = (i^*, j^*) , \end{cases} \quad i^*, j^* \in \{1, \dots, n\} . \quad (2.6.0.13)$$

$$\Rightarrow \quad \tilde{\mathbf{A}} = \mathbf{A} + z \cdot \mathbf{e}_{i^*} \mathbf{e}_{j^*}^T . \quad (2.6.0.14)$$

(Recall:  $\mathbf{e}_i \doteq i$ -th unit vector.) The question is whether we can reuse some of the computations spent on solving  $\mathbf{Ax} = \mathbf{b}$  in order to solve  $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{b}$  with less effort than entailed by a direct Gaussian elimination from scratch.

We may also consider a matrix modification affecting a single row: Changing a single row: given  $\mathbf{z} \in \mathbb{K}^n$

$$\mathbf{A}, \tilde{\mathbf{A}} \in \mathbb{K}^{n,n}: \quad \tilde{a}_{ij} = \begin{cases} a_{ij} & , \text{ if } i \neq i^* , \\ (\mathbf{z})_j + a_{ij} & , \text{ if } i = i^* , \end{cases} \quad i^*, j^* \in \{1, \dots, n\} .$$

$$\Rightarrow \quad \tilde{\mathbf{A}} = \mathbf{A} + \mathbf{e}_{i^*} \mathbf{z}^T . \quad (2.6.0.15)$$

Both matrix modifications (2.6.0.13) and (2.6.0.15) represent **rank-1-modifications** of  $\mathbf{A}$ . A generic rank-1-modification reads

$$\mathbf{A} \in \mathbb{K}^{n,n} \mapsto \tilde{\mathbf{A}} := \mathbf{A} + \boxed{\mathbf{u}\mathbf{v}^H}, \quad \mathbf{u}, \mathbf{v} \in \mathbb{K}^n. \quad (2.6.0.16)$$

general rank-1-matrix



Trick:

**Block elimination** of an extended linear system, see § 2.6.0.2

We consider the block partitioned linear system

$$\begin{bmatrix} \mathbf{A} & \mathbf{u} \\ \mathbf{v}^H & -1 \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{x}} \\ \xi \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ 0 \end{bmatrix}. \quad (2.6.0.17)$$

The **Schur complement** system after elimination of  $\xi$  reads

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^H)\tilde{\mathbf{x}} = \mathbf{b} \Leftrightarrow \tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{b}. ! \quad (2.6.0.18)$$

Hence, we have solved the modified LSE, once we have found the component  $\tilde{\mathbf{x}}$  of the solution of the extended linear system (2.6.0.17). We do block elimination again, now getting rid of  $\tilde{\mathbf{x}}$  first, which yields the other **Schur complement** system

$$(1 + \mathbf{v}^H \mathbf{A}^{-1} \mathbf{u})\xi = \mathbf{v}^H \mathbf{A}^{-1} \mathbf{b}. \quad (2.6.0.19)$$

$$\Rightarrow \tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{b} - \frac{\mathbf{u}\mathbf{v}^H \mathbf{A}^{-1}}{1 + \mathbf{v}^H \mathbf{A}^{-1} \mathbf{u}} \mathbf{b}. \quad (2.6.0.20)$$

The generalization of this formula to rank- $k$ -perturbations is given in the following lemma:

### Lemma 2.6.0.21. Sherman-Morrison-Woodbury formula

For regular  $\mathbf{A} \in \mathbb{K}^{n,n}$ , and  $\mathbf{U}, \mathbf{V} \in \mathbb{K}^{n,k}$ ,  $n, k \in \mathbb{N}$ ,  $k \leq n$ , holds

$$(\mathbf{A} + \mathbf{U}\mathbf{V}^H)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{V}^H \mathbf{A}^{-1} \mathbf{U})^{-1} \mathbf{V}^H \mathbf{A}^{-1},$$

if  $\mathbf{I} + \mathbf{V}^H \mathbf{A}^{-1} \mathbf{U}$  is regular.

*Proof.* Straightforward algebra:

$$\begin{aligned} & (\mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{V}^H \mathbf{A}^{-1} \mathbf{U})^{-1} \mathbf{V}^H \mathbf{A}^{-1})(\mathbf{A} + \mathbf{U}\mathbf{V}^H) = \\ & \mathbf{I} - \mathbf{A}^{-1}\mathbf{U} \underbrace{(\mathbf{I} + \mathbf{V}^H \mathbf{A}^{-1} \mathbf{U})^{-1}(\mathbf{I} + \mathbf{V}^H \mathbf{A}^{-1} \mathbf{U})}_{=\mathbf{I}} \mathbf{V}^H + \mathbf{A}^{-1}\mathbf{U}\mathbf{V}^H = \mathbf{I}. \end{aligned}$$

Uniqueness of the inverse settles the case.  $\square$

We use this result to solve  $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{b}$  with  $\tilde{\mathbf{A}}$  from (2.6.0.16) more efficiently than straightforward elimination could deliver, provided that the LU-factorisation  $\mathbf{a} = \mathbf{LU}$  is already known.

Apply Lemma 2.6.0.21 for  $k = 1$ :

$$\tilde{\mathbf{x}} = \mathbf{A}^{-1}\mathbf{b} - \frac{\mathbf{A}^{-1}\mathbf{u}(\mathbf{v}^H(\mathbf{A}^{-1}\mathbf{b}))}{1 + \mathbf{v}^H(\mathbf{A}^{-1}\mathbf{u})}. \quad (2.6.0.22)$$

LU-decomposition of  $\mathbf{A}$  is available.

Efficient implementation ! Asymptotic complexity  $\mathcal{O}(n^2)$  (back substitutions on Line 6-7)

### C++ code 2.6.0.23: Solving a rank-1 modified LSE → GITLAB

```

2 // Solving rank-1 updated LSE based on (2.6.0.22)
3 template <class LUDec>
4 Eigen::VectorXd smw(const LUDec &lu, const Eigen::VectorXd &u,
5                      const Eigen::VectorXd &v, const Eigen::VectorXd &b) {
6     const Eigen::VectorXd z = lu.solve(b); // z = A-1b
7     const Eigen::VectorXd w = lu.solve(u); // w = A-1u
8     double alpha = 1.0 + v.dot(w); // Compute denominator of (2.6.0.22)
9     double beta = v.dot(z); // Factor for numerator of (2.6.0.22)
10    if (std::abs(alpha) < std::numeric_limits<double>::epsilon() * std::abs(beta))
11        throw std::runtime_error("A nearly singular");
12    else
13        return (z - w * beta / alpha); // see (2.6.0.22)
14 }
```

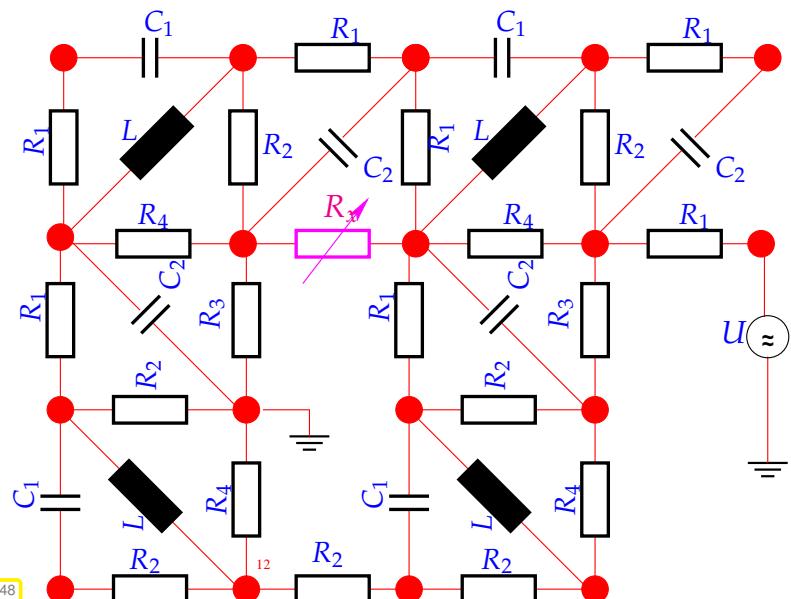
**EXAMPLE 2.6.0.24 (Resistance to currents map)** Many lineare systems with system matrices that differ in a single entry only have to be solved when we want to determine the dependence of the total impedance of a (linear) circuit from the parameters of a single component.

Large (linear) electric circuit (modelling → Ex. 2.1.0.3)

Sought:

Dependence of (certain) branch currents on “continuously varying” resistance  $R_x$

(> currents for many different values of  $R_x$ )



► Only a few entries of the nodal analysis matrix  $\mathbf{A}$  (→ Ex. 2.1.0.3) are affected by variation of  $R_x$ !  
(If  $R_x$  connects nodes  $i$  &  $j$  ⇒ only entries  $a_{ii}, a_{jj}, a_{ij}, a_{ji}$  of  $\mathbf{A}$  depend on  $R_x$ )

## 2.7 Sparse Linear Systems

We start with a (rather fuzzy) classification of matrices according to their numbers of zeros:

Dense(ly populated) matrices



sparse(ly populated) matrices

### Notion 2.7.0.1. Sparse matrix

$\mathbf{A} \in \mathbb{K}^{m,n}$ ,  $m, n \in \mathbb{N}$ , is **sparse**, if

$$\text{nnz}(\mathbf{A}) := \#\{(i,j) \in \{1, \dots, m\} \times \{1, \dots, n\} : a_{ij} \neq 0\} \ll mn .$$

Sloppy parlance: matrix **sparse** : $\Leftrightarrow$  “almost all” entries  $= 0$  /“only a few percent of” entries  $\neq 0$

J.H. Wilkinson's informal working definition for a developer of simulation codes:

### Notion 2.7.0.2. Sparse matrix

A matrix with enough zeros that it pays to take advantage of them should be *treated as sparse*.

A more rigorous “mathematical” definition:

### Definition 2.7.0.3. Families of sparse matrices

Given a strictly increasing sequences  $m : \mathbb{N} \mapsto \mathbb{N}$ ,  $n : \mathbb{N} \mapsto \mathbb{N}$ , a family  $(\mathbf{A}^{(l)})_{l \in \mathbb{N}}$  of matrices with  $\mathbf{A}^{(l)} \in \mathbb{K}^{m_l, n_l}$  is **sparse** (opposite: dense), if

$$\lim_{l \rightarrow \infty} \frac{\text{nnz}(\mathbf{A}^{(l)})}{n_l m_l} = 0 .$$

Simple example: families of diagonal matrices ( $\rightarrow$  Def. 1.1.2.3) of increasing size.

**EXAMPLE 2.7.0.4 (Sparse LSE in circuit modelling)** See Ex. 2.1.0.3 for the description of a linear electric circuit by means of a linear system of equations for nodal voltages. For large circuits the system matrices will invariably be huge and sparse.

Modern electric circuits (VLSI chips):

$10^5 - 10^7$  circuit elements

- Each element is connected to only a few nodes
- Each node is connected to only a few elements  
[In the case of a linear circuit]

nodal analysis  $\rightarrow$  **sparse** circuit matrix  
(Impossible to even store as dense matrices)

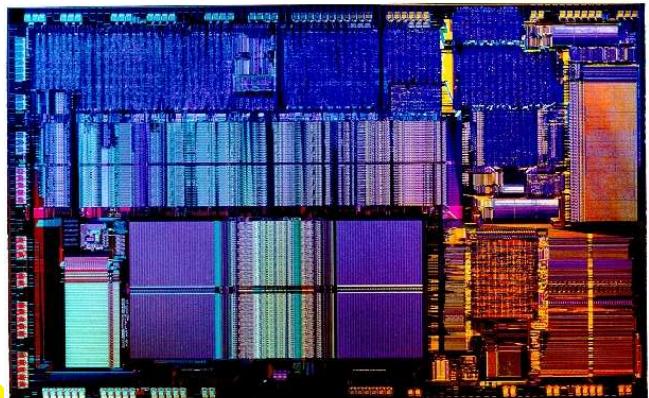


Fig. 49

**Remark 2.7.0.5 (Sparse matrices from the discretization of linear partial differential equations)**  
Another important context in which sparse matrices usually arise:

- spatial discretization of linear boundary value problems for partial differential equations by means of finite element (FE), finite volume (FV), or finite difference (FD) methods ( $\rightarrow$  4th semester course “Numerical methods for PDEs”).

## 2.7.1 Sparse matrix storage formats

Sparse matrix storage formats for storing a “sparse matrix”  $\mathbf{A} \in \mathbb{K}^{m,n}$  are designed to achieve two objectives:

- ① Amount of memory required is only slightly more than  $\text{nnz}(\mathbf{A})$  scalars.
- ② Computational effort for matrix  $\times$  vector multiplication is proportional to  $\text{nnz}(\mathbf{A})$ .

In this section we see a few schemes used by numerical libraries.

**§2.7.1.1 (Triplet/coordinate list (COO) format)** In the case of a sparse matrix  $\mathbf{A} \in \mathbb{K}^{m,n}$ , this format stores triplets  $(i,j,\alpha_{i,j}), 1 \leq i \leq m, 1 \leq j \leq n$ :

```
struct TripletMatrix {
    std::size_t m,n;           // Number of rows and columns
    std::vector<std::size_t> I; // row indices
    std::vector<std::size_t> J; // column indices
    std::vector<scalar_t> a;   // values associated with index pairs
};
```

All vectors have the same size  $\geq \text{nnz}(\mathbf{A})$ .

We write “ $\geq$ ”, because *repetitions* of index pairs  $(i,j)$  are *allowed*. The matrix entry  $(\mathbf{A})_{i,j}$  is defined to be the *sum* of all values  $\alpha_{i,j}$  associated with the index pair  $(i,j)$ . The next code clearly demonstrates this summation.

### C++-code 2.7.1.2: Matrix $\times$ vector product $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{y}$ in triplet format

```
1 void multTripIMatvec (const TripletMatrix &A,
2                         const vector<scalar_t> &x,
3                         vector<scalar_t> &y)
4     for (size_t l=0; l<A.a.size(); l++) {
5         y[A.I[l]] += A.a[l]*x[A.J[l]];
6     }
```

Note that this code assumes that the result vector  $\mathbf{y}$  has the appropriate length; no index checks are performed.

Code 2.7.1.2: computational effort is proportional to the number of triplets. (This might be much larger than  $\text{nnz}(\mathbf{A})$  in case of many repetitions of triplets.)  $\dashv$

**Remark 2.7.1.3 (The zoo of sparse matrix formats)** Special *sparse matrix storage formats* store *only* non-zero entries:

- Compressed Row Storage (CRS)
- Compressed Column Storage (CCS)  $\rightarrow$  used by MATLAB
- Block Compressed Row Storage (BCRS)
- Compressed Diagonal Storage (CDS)
- Jagged Diagonal Storage (JDS)
- Skyline Storage (SKS)

All of these formats achieve the two objectives stated above. Some have been designed for sparse matrices with additional structure or for seamless cooperation with direct elimination algorithms (JDS,SKS). ↴

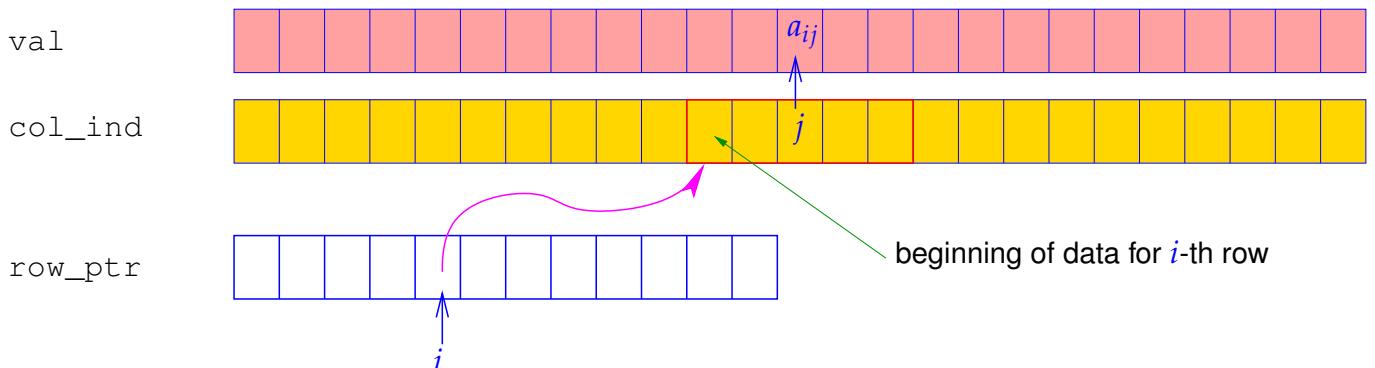
**§2.7.1.4 (Compressed row-storage (CRS) format)** The CRS format for a sparse matrix  $\mathbf{A} = [a_{ij}] \in \mathbb{K}^{n,n}$  keeps the data in three *contiguous* arrays:

<code>std::vector&lt;scalar_t&gt; val</code>	size $\text{nnz}(\mathbf{A}) := \#\{(i,j) \in \{1, \dots, n\}^2, a_{ij} \neq 0\}$
<code>std::vector&lt;size_t&gt; col_ind</code>	size $\text{nnz}(\mathbf{A})$
<code>std::vector&lt;size_t&gt; row_ptr</code>	size $n+1$ & $\text{row\_ptr}[n+1] = \text{nnz}(\mathbf{A}) + 1$ (sentinel value)

As above we write  $\text{nnz}(\mathbf{A}) \triangleq (\text{number of nonzeros})$  of  $\mathbf{A}$

Access to matrix entry  $a_{ij} \neq 0, 1 \leq i, j \leq n$  (“mathematical indexing”)

$$\text{val}[k] = a_{ij} \Leftrightarrow \begin{cases} \text{col\_ind}[k] = j, \\ \text{row\_ptr}[i] \leq k < \text{row\_ptr}[i+1], \end{cases} \quad 1 \leq k \leq \text{nnz}(\mathbf{A}).$$



$$\mathbf{A} = \begin{bmatrix} 10 & 0 & 0 & 0 & -2 & 0 \\ 3 & 9 & 0 & 0 & 0 & 3 \\ 0 & 7 & 8 & 7 & 0 & 0 \\ 3 & 0 & 8 & 7 & 5 & 0 \\ 0 & 8 & 0 & 9 & 9 & 13 \\ 0 & 4 & 0 & 0 & 2 & -1 \end{bmatrix} \quad \begin{array}{l} \text{val-vector:} \\ \begin{bmatrix} 10 & -2 & 3 & 9 & 3 & 7 & 8 & 7 & 3 \dots 9 & 13 & 4 & 2 & -1 \end{bmatrix} \\ \text{col\_ind-array:} \\ \begin{bmatrix} 1 & 5 & 1 & 2 & 6 & 2 & 3 & 4 & 1 \dots 5 & 6 & 2 & 5 & 6 \end{bmatrix} \\ \text{row\_ptr-array:} \\ \begin{bmatrix} 1 & 3 & 6 & 9 & 13 & 17 & 20 \end{bmatrix} \end{array}$$

Variant: diagonal CRS format (matrix diagonal stored in separate array)

The CCS format is equivalent to CRS format for the transposed matrix. ↴

## 2.7.2 Sparse matrices in EIGEN

Eigen can handle **sparse matrices** in the standard **Compressed Row Storage (CRS)** and **Compressed Column Storage (CCS)** format, see § 2.7.1.4 and the  **EIGEN** documentation:

```
#include <Eigen/Sparse>
Eigen::SparseMatrix<int, Eigen::ColMajor> Asp(rows,cols); // CCS
format
Eigen::SparseMatrix<double, Eigen::RowMajor> Bsp(rows,cols); // CRS
format
```

Usually sparse matrices must not be filled by setting entries through index-pair access. The matrix should first be assembled in [triplet format](#)  [EIGEN documentation](#), from which a sparse matrix is built. EIGEN offers special facilities for handling triplets.

```
std::vector <Eigen::Triplet <double> > triplets;
// .. fill the std::vector triplets ..
Eigen::SparseMatrix<double, Eigen::RowMajor> spMat(rows, cols);
spMat.setFromTriplets(triplets.begin(), triplets.end());
```

A triplet object can be initialized as demonstrated in the following example:

```
unsigned int row_idx = 2;
unsigned int col_idx = 4;
double value = 2.5;
Eigen::Triplet<double> triplet(row_idx, col_idx, value);
std::cout << '(' << triplet.row() << ',' << triplet.col()
     << ',' << triplet.value() << ')' << std::endl;
```

As shown, a `Triplet` object offers the access member functions `row()`, `col()`, and `value()` to fetch the row index, column index, and scalar value stored in a `Triplet`.

The statement that entry-wise initialization of sparse matrices is not efficient has to be qualified in Eigen. Entries can be set, provided that *enough space* for each row (in `RowMajor` format) is reserved in advance. This done by the `reserve()` method that takes an integer vector of maximal expected numbers of non-zero entries per row:

#### C++-code 2.7.2.1: Accessing entries of a sparse matrix: potentially inefficient!

```
1 unsigned int rows, cols, max_no_nnz_per_row;
2 ....
3 SparseMatrix<double, RowMajor> mat(rows, cols);
4 mat.reserve(RowVectorXi::Constant(cols, max_no_nnz_per_row));
5 // do many (incremental) initializations
6 for ( ) {
7     mat.insert(i, j) = value_ij;
8     mat.coeffRef(i, j) += increment_ij;
9 }
10 mat.makeCompressed();
```

`insert(i, j)` sets an entry of the sparse matrix, which is rather efficient, provided that enough space has been reserved. `coeffRef(i, j)` gives l-value and r-value access to any matrix entry, creating a non-zero entry, if needed: costly!

The usual matrix operations are supported for sparse matrices; addition and subtraction may involve only sparse matrices stored in the *same format*. These operations may incur large hidden costs and have to be used with care!

**EXPERIMENT 2.7.2.2 (Initialization of sparse matrices in Eigen)** We study the runtime behavior of the initialization of a sparse matrix in Eigen. We use the methods described above. The code is available from → [GITLAB](#).

Runtimes (in ms) for the initialization of a banded matrix (with 5 non-zero diagonals, that is, a maximum of 5 non-zero entries per row) using different techniques in Eigen.

**Green** line: timing for entry-wise initialization with only 4 non-zero entries per row reserved in advance.

(OS: Ubuntu Linux 14.04, CPU: Intel i5@1.80 Ghz, Compiler: g++-4.8.2, -O2)

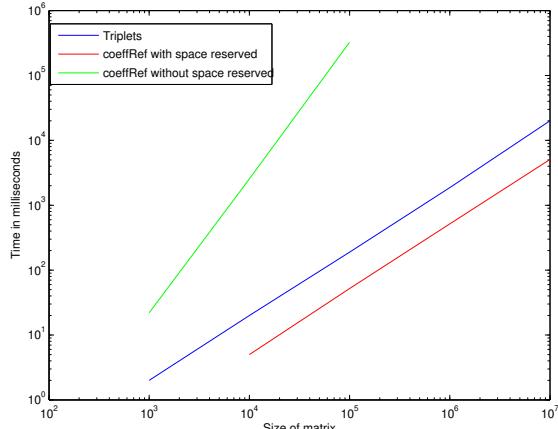


Fig. 50

Observation: insufficient advance allocation of memory massively slows down the set-up of a sparse matrix in the case of direct entry-wise initialization.

Reason: Massive internal copying of data required to created space for “unexpected” entries.  $\square$

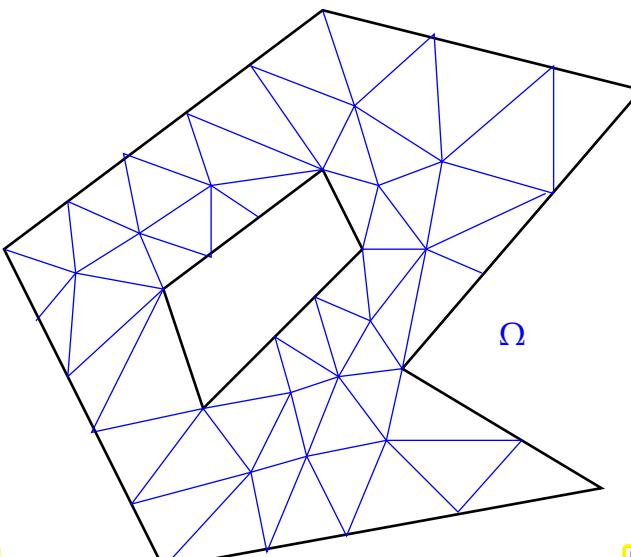
**EXAMPLE 2.7.2.3 (Smoothing of a triangulation)** This example demonstrates that sparse linear systems of equations naturally arise in the handling of triangulations.

#### Definition 2.7.2.4. Planar triangulation

A **planar triangulation** (mesh)  $\mathcal{M}$  consists of a set  $\mathcal{N}$  of  $N \in \mathbb{N}$  distinct points  $\in \mathbb{R}^2$  and a set  $\mathcal{T}$  of triangles with vertices in  $\mathcal{N}$ , such that the following two conditions are satisfied:

1. the interiors of the triangles are mutually disjoint (“no overlap”),
2. for every two *closed* distinct triangles  $\in \mathcal{T}$  their intersection satisfies exactly one of the following conditions:
  - (a) it is empty
  - (b) it is exactly one vertex from  $\mathcal{N}$ ,
  - (c) it is a **common edge** of both triangles

The points in  $\mathcal{N}$  are also called the **nodes** of the mesh, the triangles the **cells**, and all line segments connecting two nodes and occurring as a side of a triangle form the set of **edges**. We always assume a consecutive numbering of the nodes and cells of the triangulation (starting from 1, MATLAB's convention).



Valid planar triangulation

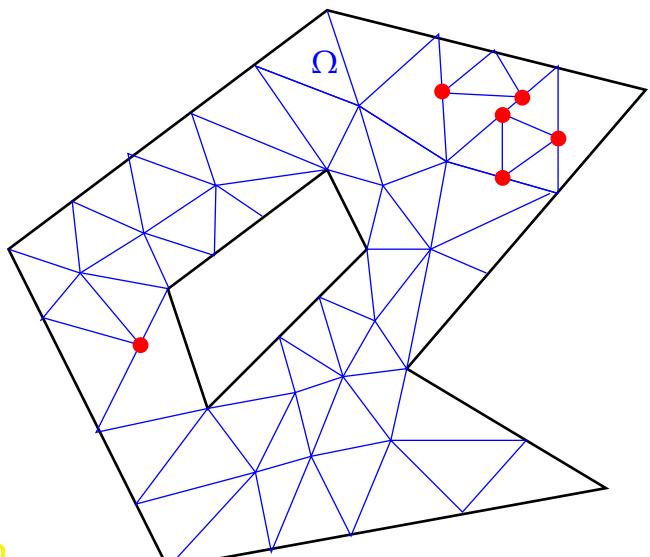
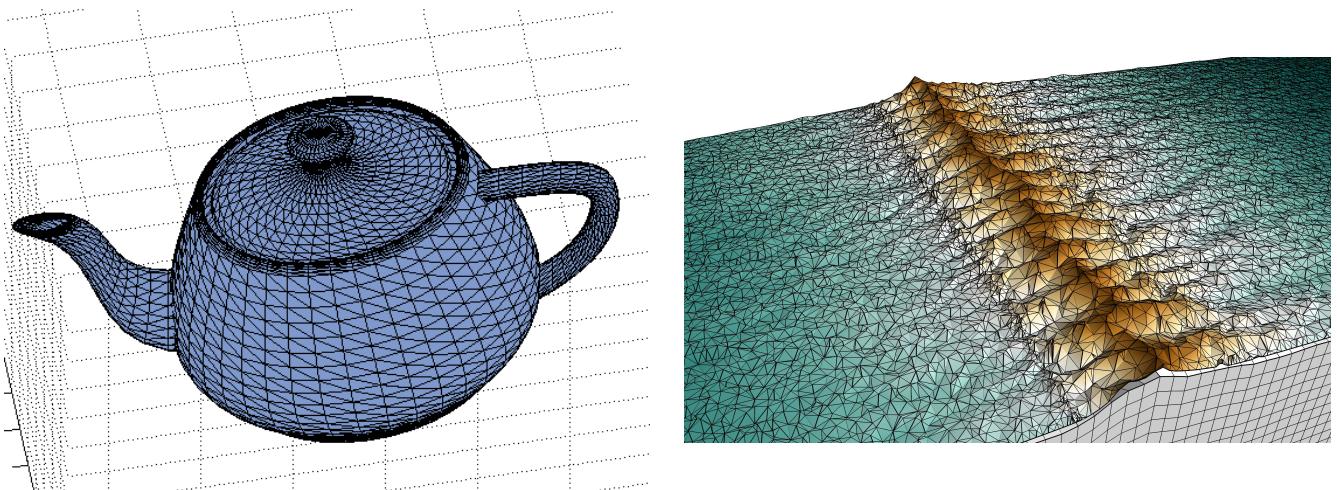
Mesh with “illegal” **hanging nodes**

Fig. 51

Fig. 52

Triangulations are of fundamental importance for computer graphics, landscape models, geodesy, and numerical methods. They need not be planar, but the algorithmic issues remain the same.



**MATLAB data structure** for describing a triangulation with  $N$  nodes and  $M$  cells:

- column vector  $\mathbf{x} \in \mathbb{R}^N$ :  $x$ -coordinates of nodes
- column vector  $\mathbf{y} \in \mathbb{R}^N$ :  $y$ -coordinates of nodes
- $M \times 3$ -matrix  $\mathbf{T}$  whose rows contain the index numbers of the vertices of the cells.  
(This matrix is a so-called triangle-node **incidence matrix**.)

The **Figure** library provides the function `triplot` for drawing planar triangulations:

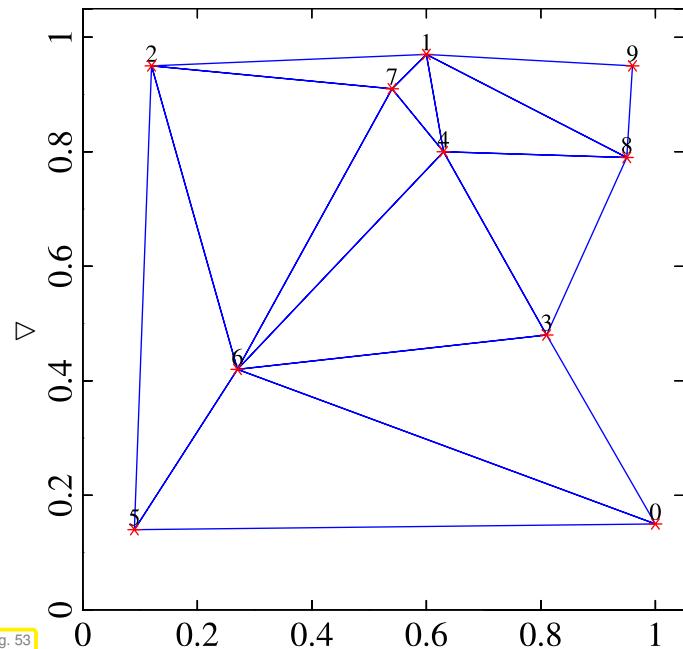
#### C++-code 2.7.2.5: Initializing and drawing a simple planar triangulations → [GITLAB](#)

```

2 // Demonstration for visualizing a plane triangular mesh
3 // Initialize node coordinates
4 Eigen::VectorXd x(10), y(10);
5 // x and y coordinate of mesh
6 x << 1.0,0.60,0.12,0.81,0.63,0.09,0.27,0.54,0.95,0.96;
7 y << 0.15,0.97,0.95,0.48,0.80,0.14,0.42,0.91,0.79,0.95;
8 // Then specify triangles through the indices of their vertices.
9 // These
10 // indices refer to the ordering of the coordinates as given in the
11 // vectors x and y.
12 Eigen::MatrixXi T(11,3);
13 T << 7, 1, 2,      5, 6, 2,      4, 1, 7,      6, 7, 2,
14     6, 4, 7,      6, 5, 0,      3, 6, 0,      8, 4, 3,
15     3, 4, 6,      8, 1, 4,      9, 1, 8;
16 // Call the Figure plotting routine, draw mesh with blue edges
17 // red vertices and a numbering/ordering
18 mgl::Figure fig1; fig1.setFontSize(8);
19 fig1.ranges(0.0, 1.05, 0.0, 1.05);
20 fig1.triplot(T, x, y, "b?"); // drawing triangulation with numbers
21 fig1.plot(x, y, "*r"); // mark vertices
22 fig1.save("meshplot_cpp");

```

Output of Code 2.7.2.5



The cells of a mesh may be rather distorted triangles (with very large and/or small angles), which is usually not desirable. We study an algorithm for **smoothing** a mesh without changing the planar domain covered by it.

#### Definition 2.7.2.6. Boundary edge

Every edge that is adjacent to only one cell is a **boundary edge** of the triangulation. Nodes that are endpoints of boundary edges are **boundary nodes**.

- ☞ Notation:  $\Gamma \subset \{1, \dots, N\} \hat{=} \text{set of indices of boundary nodes.}$
- ☞ Notation:  $\mathbf{p}^i = (p_1^i, p_2^i) \in \mathbb{R}^2 \hat{=} \text{coordinate vector of node } \#i, i = 1, \dots, N$   
 $(p_1^i \leftrightarrow x(i), p_2^i \leftrightarrow y(i) \text{ in MATLAB})$

We define

$$S(i) := \{j \in \{1, \dots, N\} : \text{nodes } i \text{ and } j \text{ are connected by an edge}\}, \quad (2.7.2.7)$$

as the set of node indices of the “neighbours” of the node with index number  $i$ .

#### Definition 2.7.2.8. Smoothed triangulation

A triangulation is called **smoothed**, if

$$\mathbf{p}^i = \frac{1}{\#S(i)} \sum_{j \in S(i)} \mathbf{p}^j \Leftrightarrow \#S(i)p_d^i = \sum_{j \in S(i)} p_d^j, d = 1, 2, \text{ for all } i \in \{1, \dots, N\} \setminus \Gamma, \quad (2.7.2.9)$$

that is, every interior node is located in the center of gravity of its neighbours.

The relations (2.7.2.9) correspond to the lines of a sparse linear system of equations! In order to state it, we insert the coordinates of all nodes into a column vector  $\mathbf{z} \in \mathbb{K}^{2N}$ , according to

$$z_i = \begin{cases} p_1^i & , \text{ if } 1 \leq i \leq N, \\ p_2^i & , \text{ if } N+1 \leq i \leq 2N. \end{cases} \quad (2.7.2.10)$$

For the sake of ease of presentation, in the sequel we *assume* (which is not the case in usual triangulation data) that interior nodes have index numbers smaller than that of boundary nodes.

From (2.7.2.7) we infer that the system matrix  $\mathbf{C} \in \mathbb{R}^{2n, 2N}$ ,  $n := N - |\Gamma|$ , of that linear system has the following structure:

$$\mathbf{C} = \begin{bmatrix} \mathbf{A} & \mathbf{O} \\ \mathbf{O} & \mathbf{A} \end{bmatrix}, \quad (\mathbf{A})_{i,j} = \begin{cases} \#S(i) & , \text{if } i = j, \\ -1 & , \text{if } j \in S(i), \\ 0 & \text{else.} \end{cases} \quad i \in \{1, \dots, n\}, \quad j \in \{1, \dots, N\}. \quad (2.7.2.11)$$

$$\blacktriangleright \quad (2.7.2.9) \Leftrightarrow \mathbf{Cz} = \mathbf{0}. \quad (2.7.2.12)$$

- $\text{nnz}(\mathbf{A}) \leq$  number of edges of  $\mathcal{M}$  + number of interior nodes of  $\mathcal{M}$ .
- The matrix  $\mathbf{C}$  associated with  $\mathcal{M}$  according to (2.7.2.11) is clearly sparse.
- The sum of the entries in every row of  $\mathbf{C}$  vanishes.

We partition the vector  $\mathbf{z}$  into coordinates of nodes in the interior and of nodes on the boundary

$$\mathbf{z}^T = \begin{bmatrix} \mathbf{z}_1^{\text{int}} \\ \mathbf{z}_1^{\text{bd}} \\ \mathbf{z}_2^{\text{int}} \\ \mathbf{z}_2^{\text{bd}} \end{bmatrix} := [z_1, \dots, z_n, z_{n+1}, \dots, z_N, z_{N+1}, \dots, z_{N+n}, z_{N+n+1}, \dots, z_{2N}]^\top.$$

This induces the following block partitioning of the linear system (2.7.2.12):

$$\begin{bmatrix} \mathbf{A}_{\text{int}} & \mathbf{A}_{\text{bd}} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{A}_{\text{int}} & \mathbf{A}_{\text{bd}} \end{bmatrix} \begin{bmatrix} \mathbf{z}_1^{\text{int}} \\ \mathbf{z}_1^{\text{bd}} \\ \mathbf{z}_2^{\text{int}} \\ \mathbf{z}_2^{\text{bd}} \end{bmatrix} = \mathbf{0}, \quad \mathbf{A}_{\text{int}} \in \mathbb{R}^{n,n}, \quad \mathbf{A}_{\text{bd}} \in \mathbb{R}^{n,N-n}.$$

$\Updownarrow$

$= \mathbf{0}.$ 
(2.7.2.13)

The linear system (2.7.2.13) holds the key to the algorithmic realization of mesh smoothing; when smoothing the mesh

- (i) the node coordinates belonging to interior nodes have to be adjusted to satisfy the equilibrium condition (2.7.2.9), they are **unknowns**,
- (ii) the coordinates of nodes located on the boundary are fixed, that is, their values are known.

$$\blacktriangleright \quad \text{unknown } \mathbf{z}_1^{\text{int}}, \mathbf{z}_2^{\text{int}}, \quad \text{known } \mathbf{z}_1^{\text{bd}}, \mathbf{z}_2^{\text{bd}} \\ (\text{yellow in (2.7.2.13)}) \quad (\text{pink in (2.7.2.13)})$$

$$(2.7.2.12)/(2.7.2.13) \Leftrightarrow \mathbf{A}_{\text{int}} [\mathbf{z}_1^{\text{int}} \mathbf{z}_2^{\text{int}}] = -[\mathbf{A}_{\text{bd}} \mathbf{z}_1^{\text{bd}} \mathbf{A}_{\text{bd}} \mathbf{z}_2^{\text{bd}}]. \quad (2.7.2.14)$$

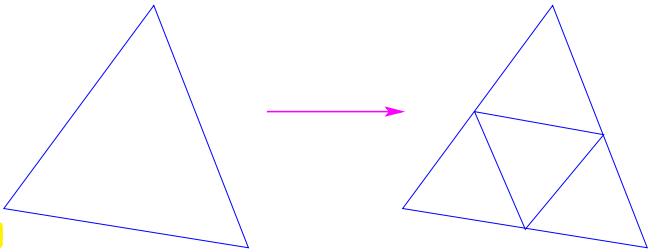
This is a square linear system with an  $n \times n$  system matrix, to be solved for two different right hand side vectors. The matrix  $\mathbf{A}_{\text{int}}$  is also known as the matrix of the **combinatorial graph Laplacian**.

We examine the sparsity pattern of the system matrices  $\mathbf{A}_{\text{int}}$  for a sequence of triangulations created by regular refinement.

### Definition 2.7.2.15. Regular refinement of a planar triangulation

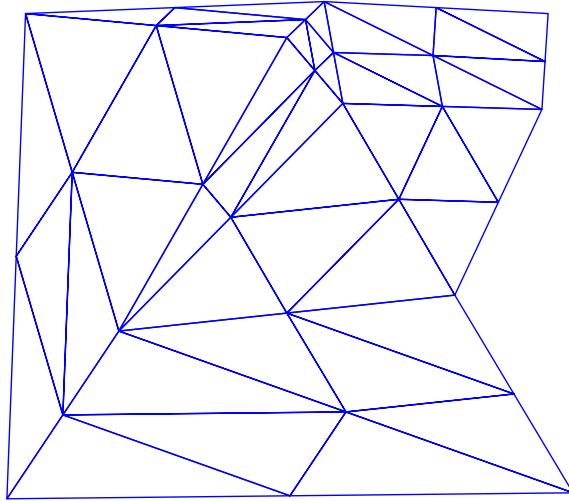
The planar triangulation with cells obtained by splitting all cells of a planar triangulation  $\mathcal{M}$  into four congruent triangles is called the **regular refinement** of  $\mathcal{M}$ .

Fig 54

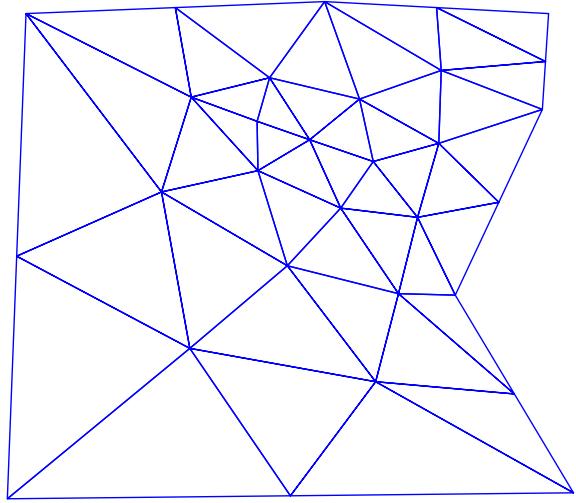


We start from the triangulation of Fig. 53 and in turns perform regular refinement and smoothing (left  $\leftrightarrow$  after refinement, right  $\leftrightarrow$  after smoothing)

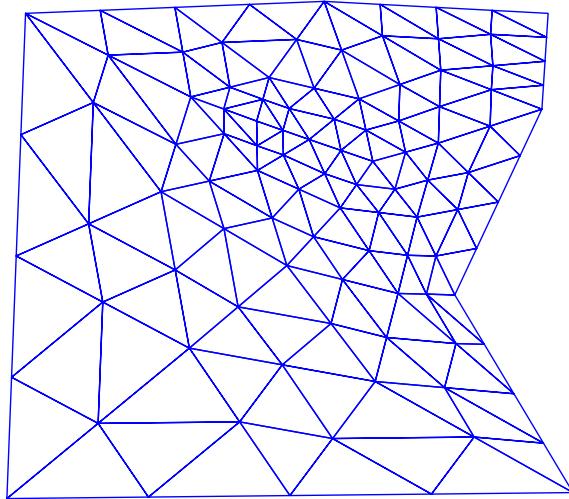
Refined mesh level 1



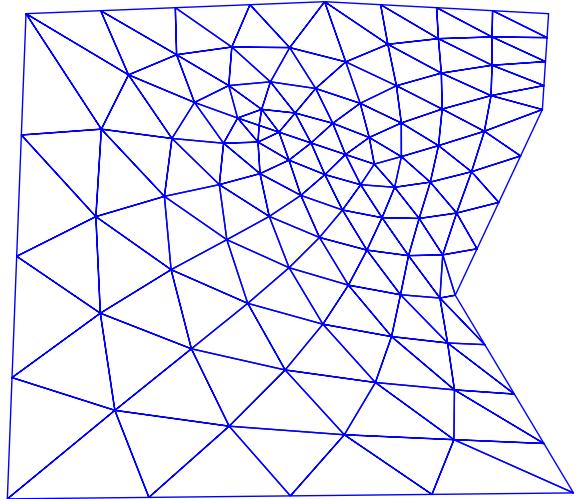
Smoothed mesh level 1

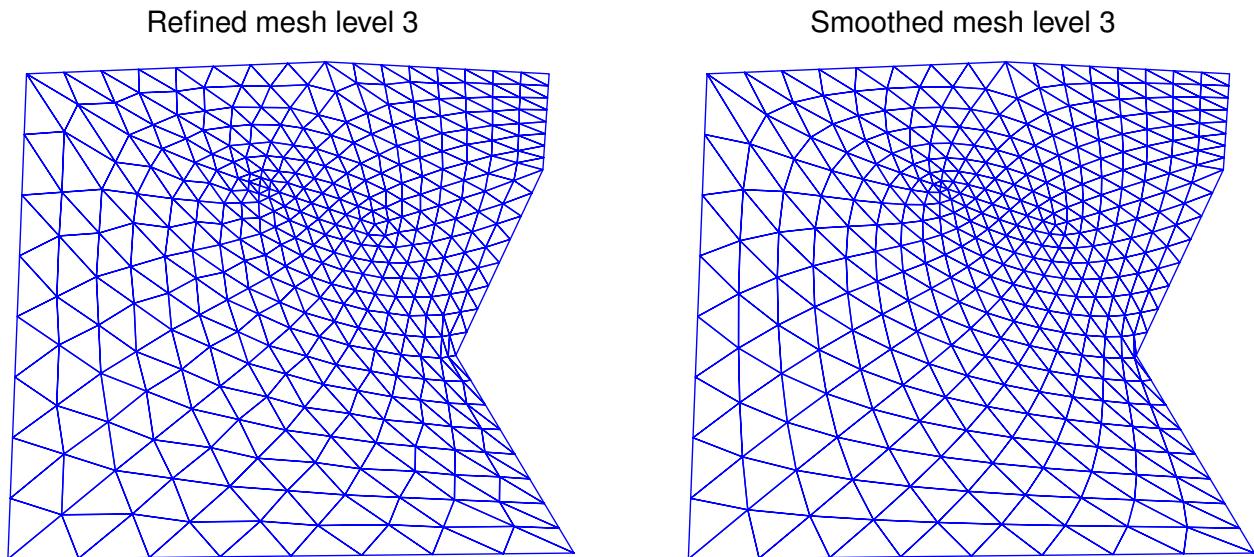


Refined mesh level 2



Smoothed mesh level 2





Below we give spy plots of the system matrices  $\mathbf{A}_{\text{int}}$  for the first three triangulations of the sequence:

level 1: 30 points, 73 edges, 44 cells

level 2: 103 points, 278 edges, 176 cells

level 3: 381 points, 1084 edges, 704 cells

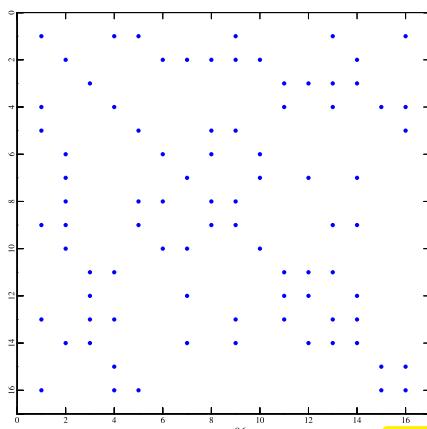


Fig. 55

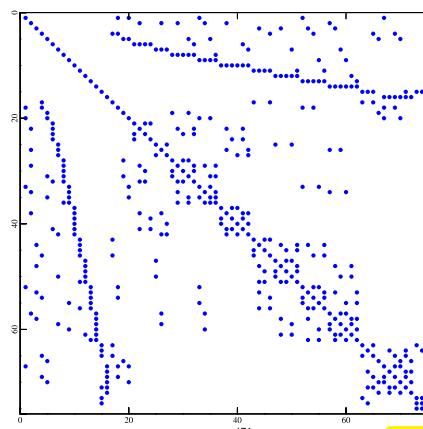


Fig. 56

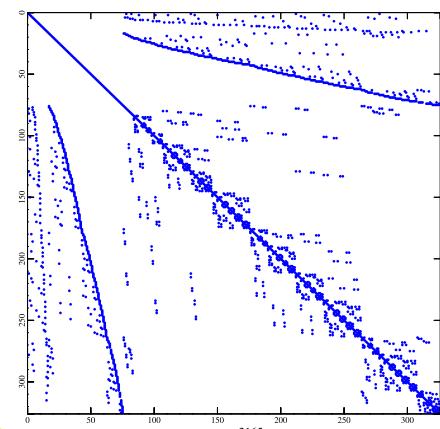


Fig. 57

## 2.7.3 Direct Solution of Sparse Linear Systems of Equations

Efficient Gaussian elimination for sparse matrices requires sophisticated algorithms that are encapsulated in special types of solvers in EIGEN. Their calling syntax remains unchanged, however:

```
Eigen::SolverType<Eigen::SparseMatrix<double>> solver(A);
Eigen::VectorXd x = solver.solve(b);
```

The standard sparse solver is SparseLU.

### C++-code 2.7.3.1: Function for solving a sparse LSE with EIGEN → GITLAB

```
2  using SparseMatrix = Eigen::SparseMatrix<double>;
3  // Perform sparse elimination
4  void sparse_solve(const SparseMatrix &A, const VectorXd &b, VectorXd &x) {
5      Eigen::SparseLU<SparseMatrix> solver(A);
6      if (solver.info() != Eigen::Success) {
7          throw "Matrix factorization failed";
8      }
9      x = solver.solve(b);
10 }
```

The constructor of the `solver` object builds the actual sparse LU-decomposition. The `solve` method then does forward and backward elimination, cf. § 2.3.2.15. It can be called multiple times, see Rem. 2.5.0.10. For more sample codes see → [GITLAB](#).

**EXPERIMENT 2.7.3.2 (Sparse elimination for arrow matrix)** In Ex. 2.6.0.5 we saw that applying the standard `lu()` solver to a sparse arrow matrix results in an extreme waste of computational resources.

Yet, EIGEN can do much better! The main mistake was the creation of a dense matrix instead of storing the arrow matrix in sparse format. There are EIGEN solvers which rely on particular **sparse elimination techniques**. They still rely of Gaussian elimination with (partial) pivoting (→ Code 2.3.3.6), but take pains to operate on non-zero entries only. This can greatly boost the speed of the elimination.

### C++ code 2.7.3.3: Invoking sparse elimination solver for arrow matrix → [GITLAB](#)

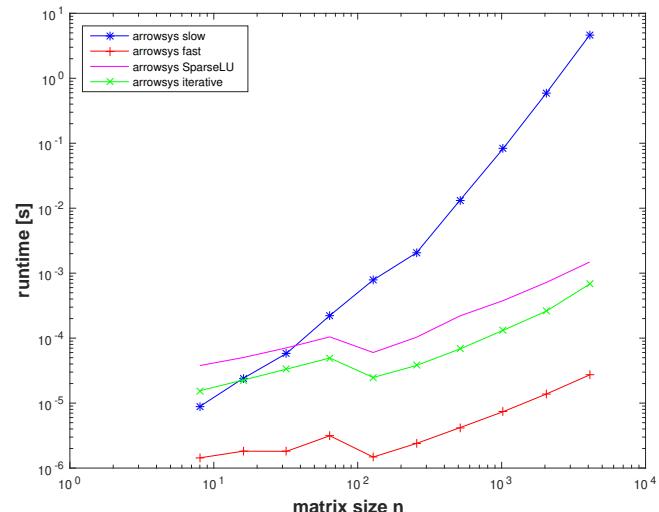
```

2  template <class solver_t>
3  VectorXd arrowsys_sparse(const VectorXd &d, const VectorXd &c, const VectorXd &b,
4      const double alpha, const VectorXd &y){
5      int n = d.size();
6      SparseMatrix<double> A(n+1, n+1); // default: column-major
7      VectorXi reserveVec = VectorXi::Constant(n+1, 2); // nnz per col
8      reserveVec(n) = n+1; // last full col
9      A.reserve(reserveVec);
10     for(int j = 0; j < n; ++j){ // initialize along cols for efficiency
11         A.insert(j, j) = d(j); // diagonal entries
12         A.insert(n, j) = b(j); // bottom row entries
13     }
14     for(int i = 0; i < n; ++i){
15         A.insert(i, n) = c(i); // last col
16     }
17     A.insert(n, n) = alpha; // bottomRight entry
18     A.makeCompressed();
19     return solver_t(A).solve(y);
}
```

Observation:

The sparse elimination solver is several orders of magnitude faster than `lu()` operating on a dense matrix.

The sparse solver is still slower than Code 2.6.0.10. The reason is that it is a general algorithm that has to keep track of non-zero entries and has to be prepared to do pivoting.



**EXPERIMENT 2.7.3.4 (Timing sparse elimination for the combinatorial graph Laplacian)** We consider a sequence of planar triangulations created by successive regular refinement (→ Def. 2.7.2.15) of the planar triangulation of Fig. 53, see Ex. 2.7.2.3. We use different EIGEN and MKL sparse solver for the linear system of equations (2.7.2.14) associated with each mesh.

## Timing results

Platform:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB, Mem 8 GB
- ◆ gcc 4.8.4, -O3

We observe an *empirical* asymptotic complexity ( $\rightarrow$  Def. 1.4.1.1) of  $O(n^{1.5})$ , way better than the asymptotic complexity of  $O(n^3)$  expected for Gaussian elimination in the case of dense matrices.

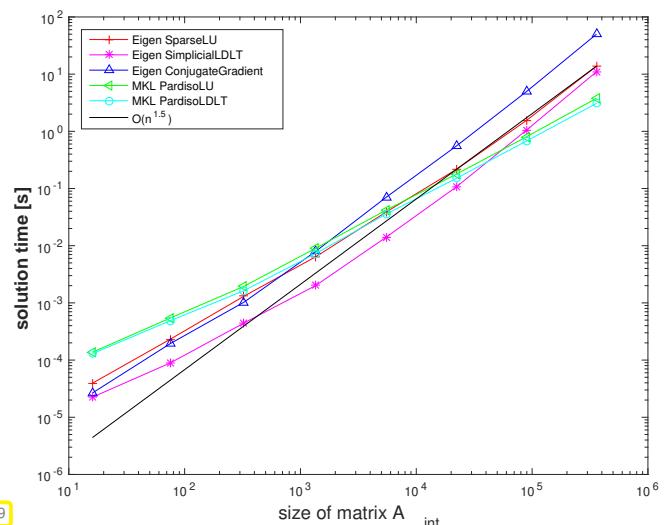


Fig. 59

↓

When solving linear systems of equations directly **dedicated sparse elimination solvers** from *numerical libraries* have to be used!

System matrices are passed to these algorithms in sparse storage formats ( $\rightarrow$  2.7.1) to convey information about zero entries.

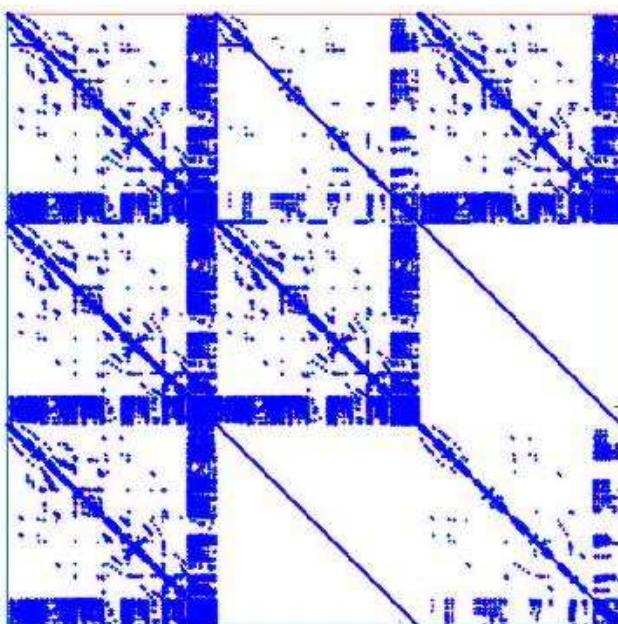


Never ever even think about implementing a general sparse elimination solver by yourself!

For an survey of sparse solvers available in EIGEN see  [EIGEN documentation](#).

**§2.7.3.5 (Implementations of sparse solvers)** Widely used implementations of sparse solvers are:

- SuperLU (<http://www.cs.berkeley.edu/~demmel/SuperLU.html>),
- UMFPACK (<http://www.cise.ufl.edu/research/sparse/umfpack/>), used by MATLAB's \,
- **PARDISO [SG04]** (<http://www.pardiso-project.org/>), incorporated into MKL



▷ fill-in ( $\rightarrow$  Def. 2.7.4.3) during sparse elimination with PARDISO

**PARDISO** has been developed by Prof. O. Schenk and his group (formerly University of Basel, now USI Lugano).



Fig. 60

**C++-code 2.7.3.6: Example code demonstrating the use of PARDISO with EIGEN → GITLAB**

```

2 void solveSparsePardiso(size_t n) {
3     using SpMat = Eigen::SparseMatrix<double>;
4     // Initialize a sparse matrix
5     const SpMat M = initSparseMatrix<SpMat>(n);
6     const Eigen::VectorXd b = Eigen::VectorXd::Random(n);
7     Eigen::VectorXd x(n);
8     // Initialization of the sparse direct solver based on the Pardiso
9     // library with
10    // directly passing the matrix M to the solver Pardiso is part of the
11    // Intel
12    // MKL library, see also Ex. 1.3.2.6
13    Eigen::PardisoLU<SpMat> solver(M);
14    // The checks of ?? are omitted
15    // solve the LSE
16    x = solver.solve(b);
17 }
```

Required is `#include <Eigen/PardisoSupport>`, the compilation flag `-DEIGEN_USE_MKL_ALL`, and the inclusion of MKL libraries during the linking phase.

```

1 COMPILER = clang++
2 FLAGS = -std=c++11 -m64 -I/usr/include/eigen3 -I${MKLROOT}/include -O3 -DNDEBUG
3
4 # Intel(R) MKL 11.3.2, Linux, None, GNU C/C++, Intel(R) 64, Static, LP64,
5 # Sequential
6 FLAGS_LINK = -WI,--start-group ${MKLROOT}/lib/intel64/libmkl_intel_lp64.a \
7             ${MKLROOT}/lib/intel64/libmkl_core.a \
8                 ${MKLROOT}/lib/intel64/libmkl_sequential.a \
9 -WI,--end-group -lpthread -lm -ldl
10
11 all: main.cpp
12     $(COMPILER) $(FLAGS) -DEIGEN_USE_MKL_ALL $< -o main $(FLAGS_LINK)
```

## 2.7.4 LU-factorization of sparse matrices

In Sect. 2.7.1 we have seen, how sparse matrices can be stored requiring  $O(\text{nnz}(\mathbf{A}))$  memory.

However, simple examples show that the product of sparse matrices need not be sparse, which means that the multiplication of large sparse matrices will usually require an effort way bigger than the sum of the numbers of their non-zero entries.

What is the situation concerning the solution of square linear systems of equations with sparse system matrices? Generically, we have to brace for a computational effort  $O(n^3)$  for matrix size  $n \rightarrow \infty$ . Yet Section 2.7.3 sends the message that a better asymptotic complexity can often be achieved, if the sparse matrix has a particular structure and sophisticated library routines are used. In this section, we examine some aspects of Gaussian elimination  $\leftrightarrow$  LU-factorisation when applied in a sparse-matrix context.

**EXAMPLE 2.7.4.1 (LU-factorization of sparse matrices)** We examine the following “sparse” matrix with

a typical structure and inspect the **pattern** of the LU-factors returned by EIGEN, see Code 2.7.4.2.

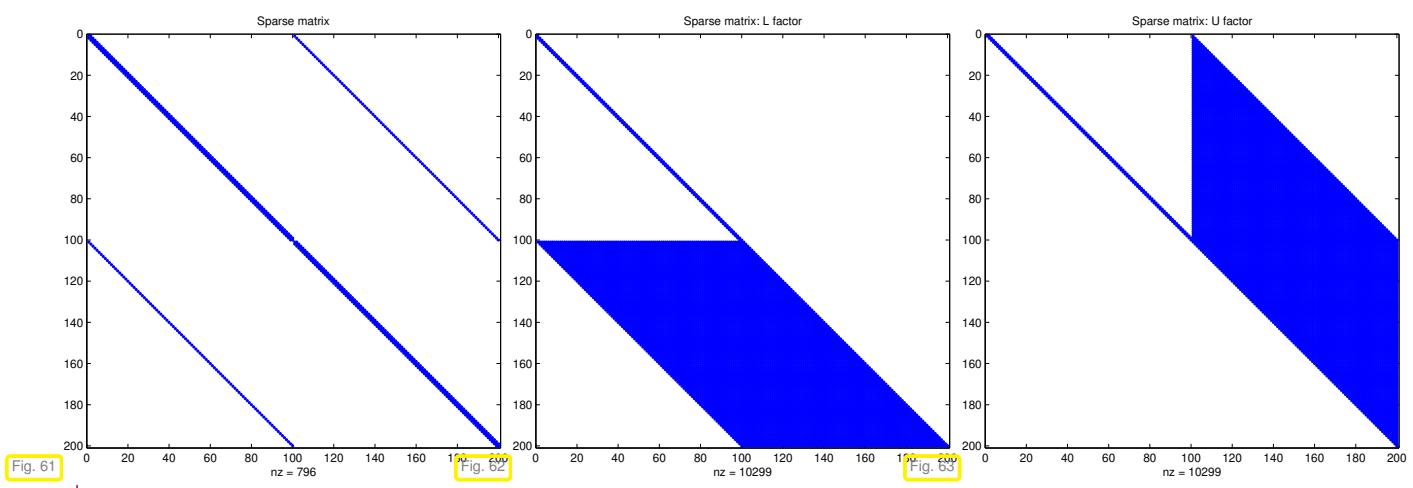
$$\mathbf{A} = \left( \begin{array}{cc|cc} 3 & -1 & -1 & \\ -1 & \ddots & \ddots & \\ \ddots & \ddots & -1 & \\ \hline -1 & -1 & 3 & -1 \\ \ddots & \ddots & -1 & \ddots \\ \ddots & \ddots & -1 & -1 \end{array} \right) \in \mathbb{R}^{n,n}, n \in \mathbb{N}$$

### C++ code 2.7.4.2: Visualizing LU-factors of a sparse matrix → GITLAB

```

2 // Build matrix
3 int n = 100;
4 RowVectorXd diag_el(5); diag_el << -1,-1, 3, -1,-1;
5 VectorXi diag_no(5); diag_no << -n, -1, 0, 1, n;
6 MatrixXd B = diag_el.replicate(2*n,1);
7 B(n-1,1) = 0; B(n,3) = 0; // delete elements
8 // A custom function from the Utils folder
9 SparseMatrix<double> A = spdiags(B, diag_no, 2*n, 2*n);
10 // It is not possible to access the LU-factors in the case of
11 // EIGEN's LU-decomposition for sparse matrices.
12 // Therefore we have to resort to the dense version.
13 auto solver = MatrixXd(A).lu();
14 MatrixXd L = MatrixXd::Identity(2*n,2*n);
15 L += solver.matrixLU().triangularView<StrictlyLower>();
16 MatrixXd U = solver.matrixLU().triangularView<Upper>();
17 // Plotting
18 mgl::Figure fig1, fig2, fig3;
19 fig1.spy(A); fig1.setFontSize(4);
20 fig1.title("Sparse matrix"); fig1.save("sparseA_cpp");
21 fig2.spy(L); fig2.setFontSize(4);
22 fig2.title("Sparse matrix: L factor"); fig2.save("sparseL_cpp");
23 fig3.spy(U); fig3.setFontSize(4);
24 fig3.title("Sparse matrix: U factor"); fig3.save("sparseU_cpp");

```



Observation:

$\mathbf{A}$  sparse  $\not\Rightarrow$  LU-factors sparse

Of course, in case the LU-factors of a sparse matrix possess many more non-zero entries than the matrix itself, the effort for solving a linear system with direct elimination will increase significantly. This can be

quantified by means of the following concept:

#### Definition 2.7.4.3. Fill-in

Let  $\mathbf{A} = \mathbf{LU}$  be an  $LULU$ -factorization ( $\rightarrow$  Sect. 2.3.2) of  $\mathbf{A} \in \mathbb{K}^{n,n}$ . If  $l_{ij} \neq 0$  or  $u_{ij} \neq 0$  though  $a_{ij} = 0$ , then we encounter **fill-in** at position  $(i,j)$ .

#### EXAMPLE 2.7.4.4 (Sparse LU-factors) Ex. 2.7.4.1 > massive fill-in can occur for sparse matrices

This example demonstrates that fill-in can largely be avoided, if the matrix has favorable structure. In this case a LSE with this particular system matrix  $\mathbf{A}$  can be solved efficiently, that is, with a computational effort  $\mathcal{O}(\text{nnz}(\mathbf{A}))$  by Gaussian elimination.

#### C++ code 2.7.4.5: LU-factorization of sparse matrix → GITLAB

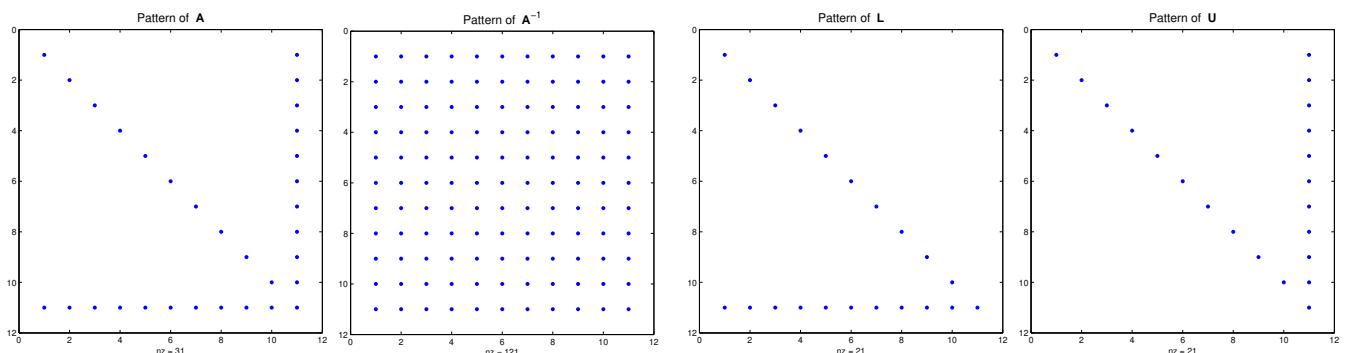
```

2 // Build matrix
3 MatrixXd A(11,11); A.setIdentity();
4 A.col(10).setOnes(); A.row(10).setOnes();
5 // A.reverseInPlace(); // used inEx. 2.7.4.6
6 auto solver = A.lu();
7 MatrixXd L = MatrixXd::Identity(11,11);
8 L += solver.matrixLU().triangularView<StrictlyLower>();
9 MatrixXd U = solver.matrixLU().triangularView<Upper>();
10 MatrixXd Ainv = A.inverse();
11 // Plotting
12 mgl::Figure fig1, fig2, fig3, fig4;
13 fig1.spy(A); fig1.setFontSize(4);
14 fig1.title("Pattern of A"); fig1.save("Apat_cpp");
15 fig2.spy(L); fig2.setFontSize(4);
16 fig2.title("Pattern of L"); fig2.save("Lpat_cpp");
17 fig3.spy(U); fig3.setFontSize(4);
18 fig3.title("Pattern of U"); fig3.save("Upat_cpp");
19 fig4.spy(Ainv); fig4.setFontSize(4);
20 fig4.title("Pattern of A^{-1}"); fig4.save("Ainvpat_cpp");

```

$\mathbf{A}$  is called an “arrow matrix”, see the pattern of non-zero entries below and Ex. 2.6.0.5.

Recalling Rem. 2.3.2.17 it is easy to see that the LU-factors of  $\mathbf{A}$  will be sparse and that their sparsity patterns will be as depicted below. Observe that despite sparse LU-factors,  $\mathbf{A}^{-1}$  will be densely populated.



$\mathbf{L}, \mathbf{U}$  sparse  $\not\Rightarrow \mathbf{A}^{-1}$  sparse !

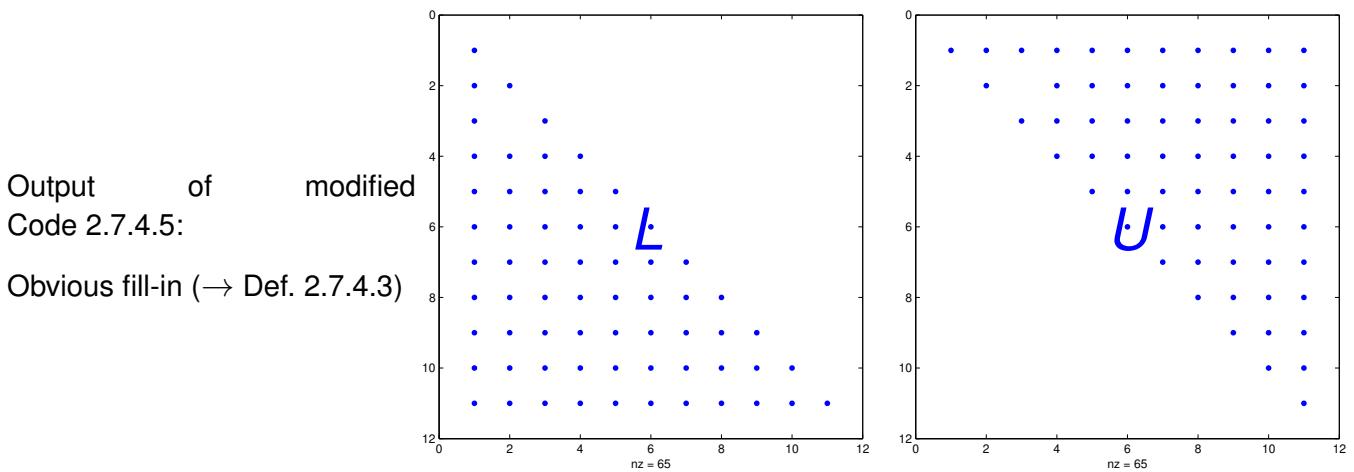
Besides stability and efficiency issues, see Exp. 2.4.0.10, this is another reason why using `x = A.inverse() * y` instead of `y = A.lu().solve(b)` is usually a major blunder. □

**EXAMPLE 2.7.4.6 (LU-decomposition of flipped “arrow matrix”)** Recall the discussion in Ex. 2.6.0.5. Here we look at an arrow matrix in a slightly different form:

$$\mathbf{M} = \begin{bmatrix} \alpha & \mathbf{b}^\top \\ \mathbf{c} & \mathbf{D} \end{bmatrix}, \quad \begin{array}{l} \alpha \in \mathbb{R}, \\ \mathbf{b}, \mathbf{c} \in \mathbb{R}^{n-1}, \\ \mathbf{D} \in \mathbb{R}^{n-1, n-1} \text{ regular diagonal matrix, } \rightarrow \text{Def. 1.1.2.3} \end{array}$$

Run the algorithm from § 2.3.2.6 (LU decomposition without pivoting):

- ◆ LU-decomposition **dense** factor matrices with  $O(n^2)$  non-zero entries.
  - ◆ asymptotic computational cost:  $O(n^3)$



Now it comes as a surprise that the arrow matrix **A** from Ex. 2.6.0.5, (2.6.0.6) has sparse LU-factors!

## Arrow matrix (2.6.0.6)

$$A \triangleright A = \begin{pmatrix} & & D \\ & & c \\ b^\top & & \alpha \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} & \mathbf{I} & & 0 \\ & & & | \\ & & & 0 \\ \hline & \mathbf{b}^\top \mathbf{D}^{-1} & & 1 \end{pmatrix} \cdot \begin{pmatrix} & \mathbf{D} & & \mathbf{c} \\ & & & | \\ & & & 0 \\ \hline & 0 & & \sigma \end{pmatrix}, \quad \sigma := \alpha - \mathbf{b}^\top \mathbf{D}^{-1} \mathbf{c}.$$

$=: \mathbf{L}$        $=: \mathbf{U}$

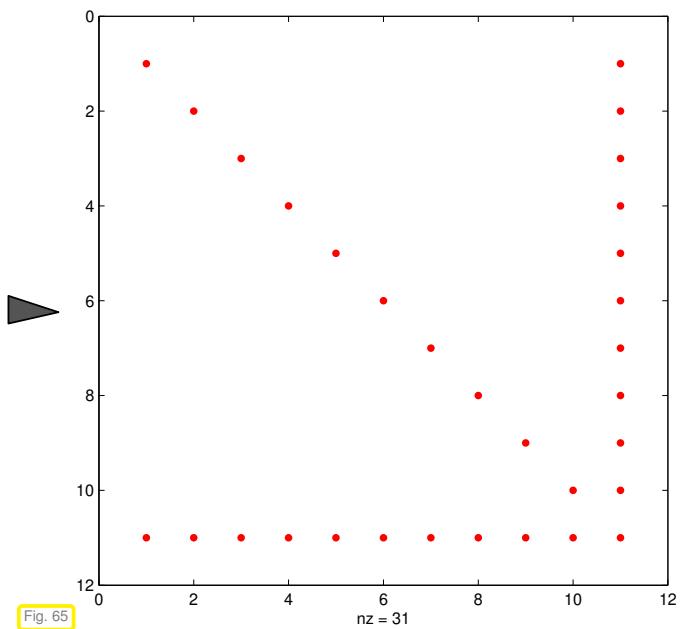
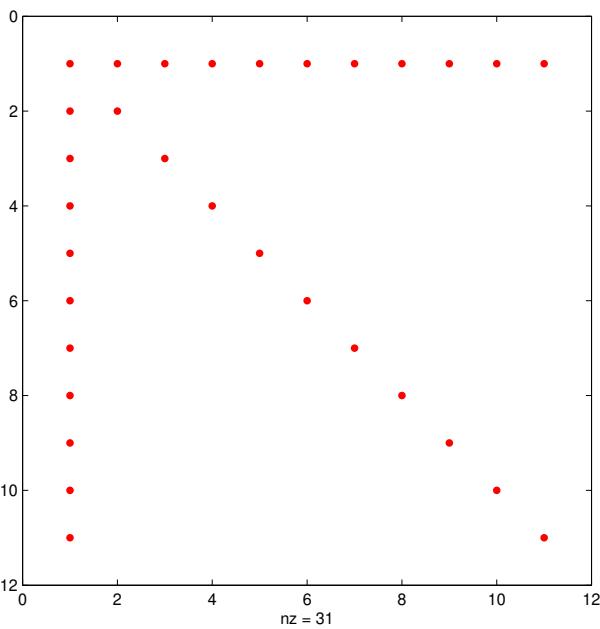
➤ In this case LU-factorisation is possible **without fill-in**, cost merely  $\mathcal{O}(n)$ !



Idea: Transform  $\mathbf{A}$  into  $\mathbf{A}$  by row and column **permutations** before performing LU-decomposition.

Details: Apply a **cyclic permutation** of rows/columns:

- 1st row/column  $\rightarrow n$ -th row/column
- $i$ -th row/column  $\rightarrow i-1$ -th row/column,  $i = 2, \dots, n$



➤ Then LU-factorization (*without pivoting*) of the resulting matrix requires  $\mathcal{O}(n)$  operations.

#### C++ code 2.7.4.8: Permuting arrow matrix, see Figs. 64, 65 → GITLAB

```

2 MatrixXd A(11,11); A.setIdentity();
3 A.col(0).setOnes(); A.row(0) = RowVectorXd::LinSpaced(11,11,1);
4 // Permutation matrix (→ Def. 2.3.3.11) encoding cyclic permutation
5 MatrixXd P(11,11); P.setZero();

```

```

6 P.topRightCorner(10,10).setIdentity(); P(10,0) = 1;
7 mgl::Figure fig1, fig2;
8 fig1.spy(A); fig1.setFontSize(4);
9 fig1.save("InvArrowSpy_cpp");
10 fig2.spy((P*A*P.transpose()).eval()); fig2.setFontSize(4);
11 fig2.save("ArrowSpy_cpp");

```

**EXAMPLE 2.7.4.9 (Pivoting destroys sparsity)** In Ex. 2.7.4.6 we found that permuting a matrix can make it amenable to Gaussian elimination/LU-decomposition with much less fill-in ( $\rightarrow$  Def. 2.7.4.3). However, recall from Section 2.3.3 that pivoting, which may be essential for achieving numerical stability, amounts to permuting the rows (or even columns) of the matrix. Thus, we may face the awkward situation that pivoting tries to reverse the very permutation we applied to minimize fill-in! The next example shows that this can happen for an arrow matrix.

#### C++ code 2.7.4.10: fill-in due to pivoting $\rightarrow$ GITLAB

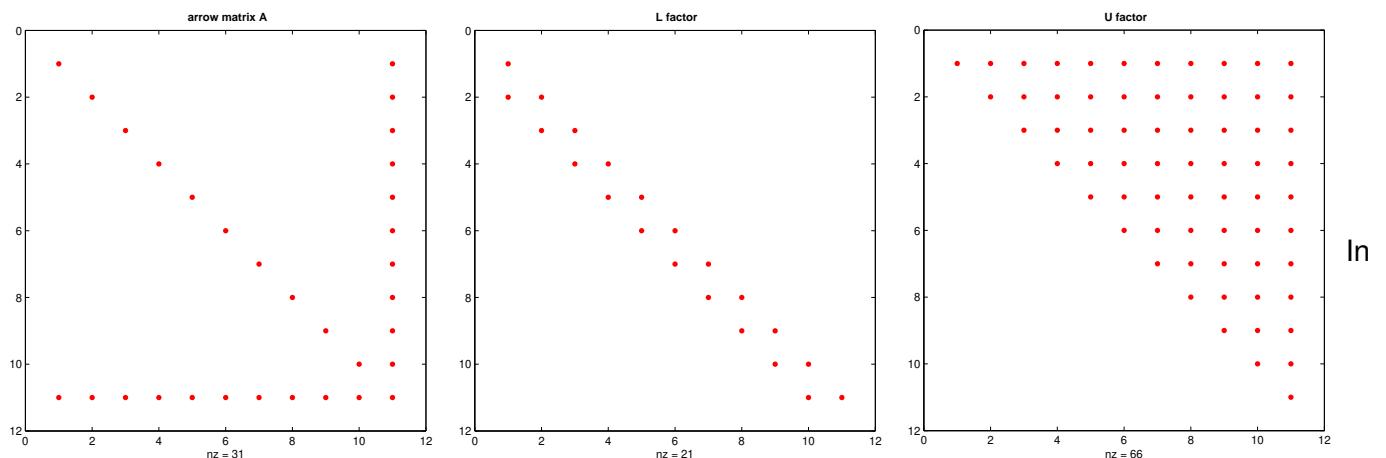
```

2 // Study of fill-in with LU-factorization due to pivoting
3 MatrixXd A(11,11); A.setZero();
4 A.diagonal() = VectorXd::LinSpaced(11,1,11).cwiseInverse();
5 A.col(10).setConstant(2); A.row(10).setConstant(2);
6 auto solver = A.lu();
7 MatrixXd L = MatrixXd::Identity(11,11);
8 L += solver.matrixLU().triangularView<StrictlyLower>();
9 MatrixXd U = solver.matrixLU().triangularView<Upper>();
10 // Plotting
11 mgl::Figure fig1, fig2, fig3, fig4;
12 fig1.spy(A); fig1.setFontSize(4);
13 fig1.title("arrow matrix A"); fig1.save("fillinpivotA");
14 fig2.spy(L); fig2.setFontSize(4);
15 fig2.title("L factor"); fig2.save("fillinpivotL");
16 fig3.spy(U); fig3.setFontSize(4);
17 fig3.title("U factor"); fig3.save("fillinpivotU");
18 std::cout << A << std::endl;

```

$$\mathbf{A} = \begin{bmatrix} 1 & & & 2 \\ & \frac{1}{2} & & 2 \\ & \ddots & & \vdots \\ 2 & \dots & \frac{1}{10} & 2 \end{bmatrix} \rightarrow \text{arrow matrix, Ex. 2.7.4.4}$$

The distributions of non-zero entries of the computed LU-factors (“spy-plots”) are as follows:



**A** **L** **U**  
this case the solution of a LSE with system matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  of the above type by means of Gaussian elimination with partial pivoting would incur costs of  $O(n^3)$ .  $\square$

## 2.7.5 Banded matrices [DR08, Sect. 3.7]

Banded matrices are a special class of sparse matrices ( $\rightarrow$  Notion 2.7.0.1 with extra structure):

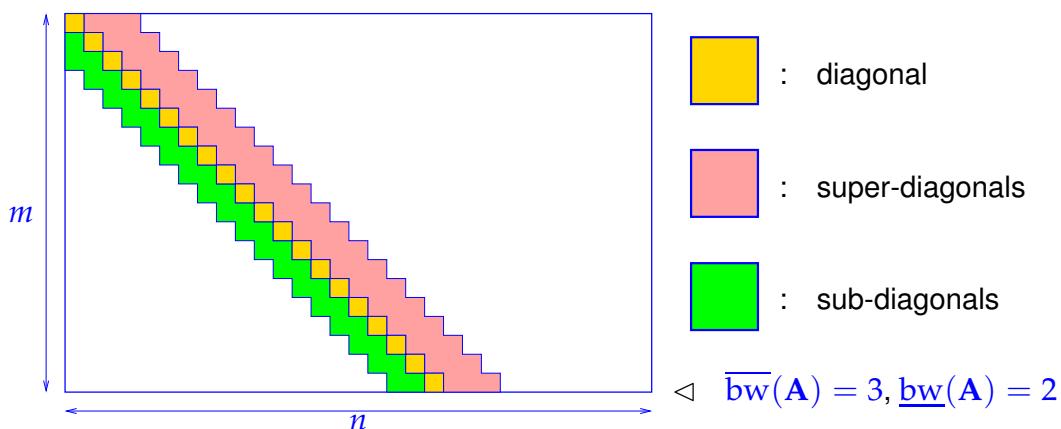
### Definition 2.7.5.1. Bandwidth

For  $\mathbf{A} = (a_{ij})_{i,j} \in \mathbb{K}^{m,n}$  we call

$$\begin{aligned}\overline{\text{bw}}(\mathbf{A}) &:= \min\{k \in \mathbb{N}: j - i > k \Rightarrow a_{ij} = 0\} \text{ upper bandwidth ,} \\ \underline{\text{bw}}(\mathbf{A}) &:= \min\{k \in \mathbb{N}: i - j > k \Rightarrow a_{ij} = 0\} \text{ lower bandwidth .}\end{aligned}$$

$$\text{bw}(\mathbf{A}) := \overline{\text{bw}}(\mathbf{A}) + \underline{\text{bw}}(\mathbf{A}) + 1 = \text{bandwidth von } \mathbf{A}.$$

- $\text{bw}(\mathbf{A}) = 1 \Rightarrow \mathbf{A}$  diagonal matrix,  $\rightarrow$  Def. 1.1.2.3
- $\overline{\text{bw}}(\mathbf{A}) = \underline{\text{bw}}(\mathbf{A}) = 1 \Rightarrow \mathbf{A}$  tridiagonal matrix
- More general:  $\mathbf{A} \in \mathbb{R}^{n,n}$  with  $\text{bw}(\mathbf{A}) \ll n \hat{=} \text{banded matrix}$



► for banded matrix  $\mathbf{A} \in \mathbb{K}^{m,n}$ :  $\text{nnz}(\mathbf{A}) \leq \min\{m, n\} \text{bw}(\mathbf{A})$

We now examine a generalization of the concept of a banded matrix that is particularly useful in the context of Gaussian elimination:

### Definition 2.7.5.2. Matrix envelope

For  $\mathbf{A} \in \mathbb{K}^{n,n}$  define

$$\text{row bandwidth} \quad \text{bw}_i^R(\mathbf{A}) := \max\{0, i - j : a_{ij} \neq 0, 1 \leq j \leq n\}, i \in \{1, \dots, n\}$$

$$\text{column bandwidth} \quad \text{bw}_j^C(\mathbf{A}) := \max\{0, j - i : a_{ij} \neq 0, 1 \leq i \leq n\}, j \in \{1, \dots, n\}$$

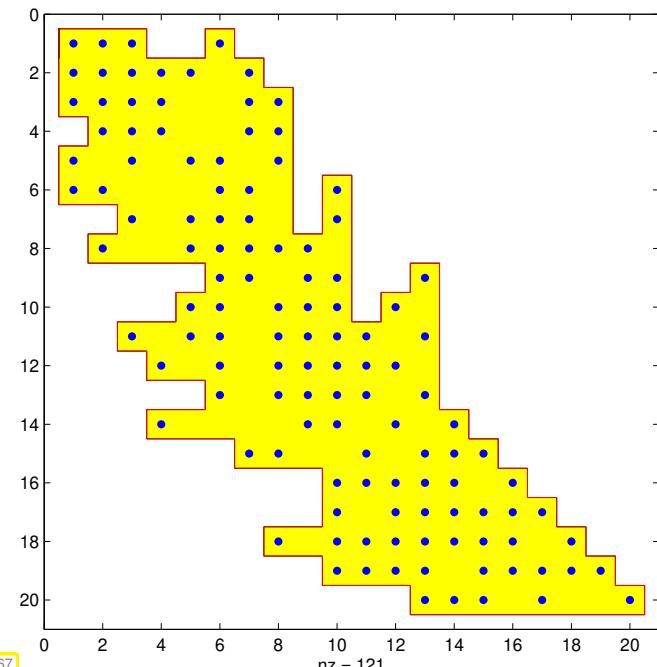
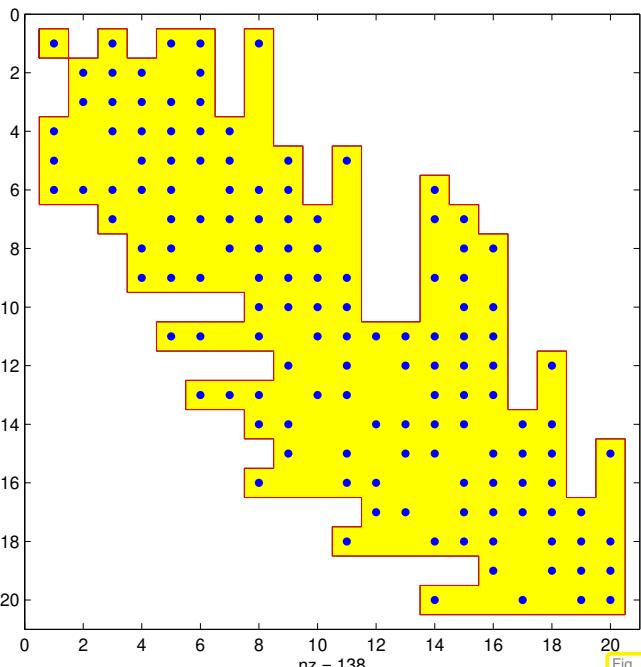
envelope

$$\text{env}(\mathbf{A}) := \left\{ (i, j) \in \{1, \dots, n\}^2 : \begin{array}{l} i - \text{bw}_i^R(\mathbf{A}) \leq j \leq i, \\ j - \text{bw}_j^C(\mathbf{A}) \leq i \leq j \end{array} \right\}$$

**EXAMPLE 2.7.5.3 (Envelope of a matrix)** We give an example illustrating Def. 2.7.5.2.

$$\mathbf{A} = \begin{pmatrix} * & 0 & * & 0 & 0 & 0 & 0 \\ 0 & * & 0 & 0 & * & 0 & 0 \\ * & 0 & * & 0 & 0 & 0 & * \\ 0 & 0 & 0 & * & * & 0 & * \\ 0 & * & 0 & * & * & * & 0 \\ 0 & 0 & 0 & 0 & * & * & 0 \\ 0 & 0 & * & * & 0 & 0 & * \end{pmatrix} \quad \begin{aligned} \text{bw}_1^R(\mathbf{A}) &= 0 \\ \text{bw}_2^R(\mathbf{A}) &= 0 \\ \text{bw}_3^R(\mathbf{A}) &= 2 \\ \text{bw}_4^R(\mathbf{A}) &= 0 \quad \text{env}(\mathbf{A}) = \text{red entries} \\ \text{bw}_5^R(\mathbf{A}) &= 3 \quad * \hat{=} \text{non-zero matrix entry } a_{ij} \neq 0 \\ \text{bw}_6^R(\mathbf{A}) &= 1 \\ \text{bw}_7^R(\mathbf{A}) &= 4 \end{aligned}$$

Starting from a “spy-plot”, it is easy to find the envelope:



Note: the envelope of the arrow matrix from Ex. 2.7.4.4 is just the set of index pairs of its non-zero entries. Hence, the following theorem provides another reason for the sparsity of the LU-factors in that example.

↓

### Theorem 2.7.5.4. Envelope and fill-in → [QSS00, Sect. 3.9]

If  $\mathbf{A} \in \mathbb{K}^{n,n}$  is regular with LU-factorization  $\mathbf{A} = \mathbf{LU}$ , then fill-in (→ Def. 2.7.4.3) is confined to  $\text{env}(\mathbf{A})$ .

Gaussian elimination without pivoting

*Proof.* (by induction, version I) Examine first step of Gaussian elimination without pivoting,  $a_{11} \neq 0$

$$\mathbf{A} = \begin{pmatrix} a_{11} & \mathbf{b}^\top \\ \mathbf{c} & \tilde{\mathbf{A}} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 \\ -\frac{\mathbf{c}}{a_{11}} & \mathbf{I} \end{pmatrix}}_{\mathbf{L}^{(1)}} \underbrace{\begin{pmatrix} a_{11} & \mathbf{b}^\top \\ 0 & \tilde{\mathbf{A}} - \frac{\mathbf{c}\mathbf{b}^\top}{a_{11}} \end{pmatrix}}_{\mathbf{U}^{(1)}}$$

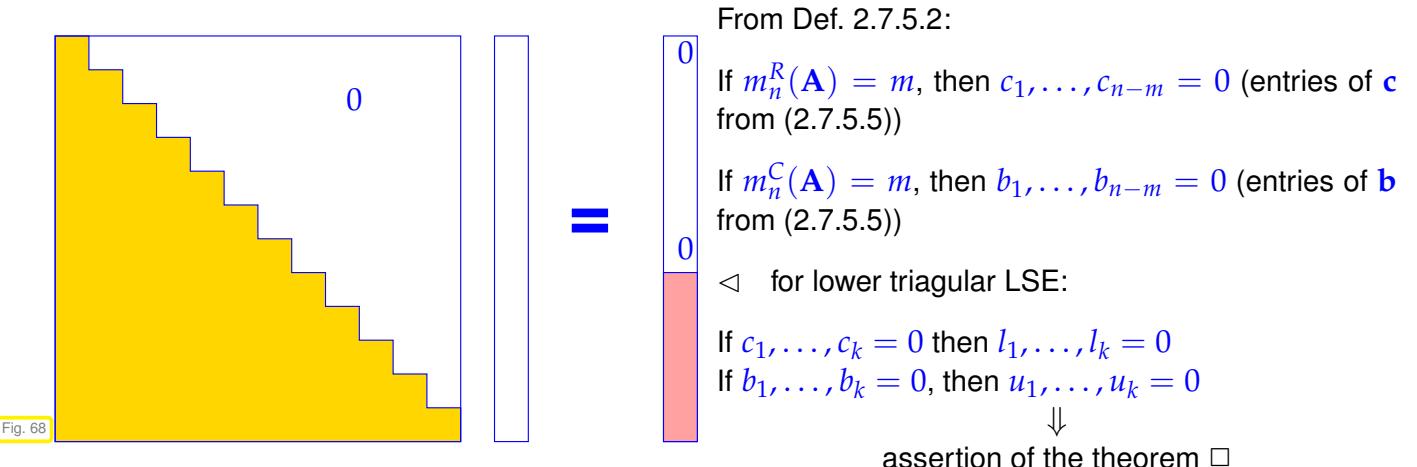
If  $(i, j) \notin \text{env}(\mathbf{A}) \Rightarrow \begin{cases} c_{i-1} = 0 & \text{, if } i > j, \\ b_{j-1} = 0 & \text{, if } i < j. \end{cases}$

$$\Rightarrow \text{env}(\mathbf{L}^{(1)}) \subset \text{env}(\mathbf{A}), \quad \text{env}(\mathbf{U}^{(1)}) \subset \text{env}(\mathbf{A}).$$

Moreover,  $\text{env}(\tilde{\mathbf{A}} - \frac{\mathbf{c}\mathbf{b}^\top}{a_{11}}) = \text{env}((\mathbf{A})_{2:n, 2:n})$  □

*Proof.* (by induction, version II) Use block-LU-factorization, cf. Rem. 2.3.2.19 and proof of Lemma 2.3.2.4:

$$\left( \begin{array}{c|c} \tilde{\mathbf{A}} & \mathbf{b} \\ \hline \mathbf{c}^\top & \alpha \end{array} \right) = \left( \begin{array}{c|c} \tilde{\mathbf{L}} & 0 \\ \hline \mathbf{1}^\top & 1 \end{array} \right) \left( \begin{array}{c|c} \tilde{\mathbf{U}} & \mathbf{u} \\ \hline 0 & \xi \end{array} \right) \Rightarrow \begin{array}{l} \tilde{\mathbf{U}}^\top \mathbf{1} = \mathbf{c}, \\ \tilde{\mathbf{L}} \mathbf{u} = \mathbf{b}. \end{array} \quad (2.7.5.5)$$



Thm. 2.7.5.4 immediately suggests a policy for saving computational effort when solving linear system whose system matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$  is sparse due to *small envelope*:

$$\#\text{env}(\mathbf{A}) \ll n^2 :$$

► Policy Confine elimination to envelope!

Details will be given now:

► **Envelope-aware LU-factorization:**

#### C++ code 2.7.5.6: Computing row bandwidths, → Def. 2.7.5.2 → GITLAB

```

2 //! computes rowbandwidth numbers  $m_i^R(\mathbf{A})$  of  $\mathbf{A}$  (sparse matrix) according
   to Def. 2.7.5.2
3 template <class numeric_t>
4 VectorXi rowbandwidth(const SparseMatrix<numeric_t> &A) {
5     VectorXi m = VectorXi::Zero(A.rows());
6     for(int k = 0; k < A.outerSize(); ++k)
7         for(typename SparseMatrix<numeric_t>::InnerIterator it(A,k); it; ++it)
8             m(it.row()) = std::max<VectorXi::Scalar>(m(it.row()),

```

```

9     it.row() - it.col() );
10    return m;
11 }
11 //! computes row bandwidth numbers  $m_i^R(\mathbf{A})$  of  $\mathbf{A}$  (dense matrix) according
11 //! to Def. 2.7.5.2
12 template <class Derived>
13 VectorXi rowbandwidth(const MatrixBase<Derived> &A) {
14     VectorXi m = VectorXi::Zero(A.rows());
15     for(int i = 1; i < A.rows(); ++i)
16         for(int j = 0; j < i; ++j)
17             if(A(i, j) != 0){
18                 m(i) = i - j;
19                 break;
20             }
21     return m;
22 }
```

### C++ code 2.7.5.7: Envelope aware forward substitution → GITLAB

```

2 //! envelope aware forward substitution for  $\mathbf{L}\mathbf{x} = \mathbf{y}$ 
3 //! ( $\mathbf{L}$  = lower triangular matrix)
4 //! argument mr: row bandwidth vector
5 VectorXd substenv(const MatrixXd &L, const VectorXd &y, const VectorXi &mr) {
6     int n = L.cols(); VectorXd x(n);
7     x(0) = y(0)/L(0,0);
8     for(int i = 1; i < n; ++i){
9         if( mr(i) > 0){
10             double zeta = L.row(i).segment(i-mr(i),mr(i))
11                     * x.segment(i-mr(i), mr(i));
12             x(i) = (y(i) - zeta)/L(i,i);
13         }
14         else x(i) = y(i) / L(i,i);
15     }
16     return x;
17 }
```

Asymptotic complexity of envelope aware forward substitution, cf. § 2.3.2.15, for  $\mathbf{Lx} = \mathbf{y}$ ,  $\mathbf{L} \in \mathbb{K}^{n,n}$  regular lower triangular matrix is

$$O(\#\text{env}(\mathbf{L})) !$$

By block LU-factorization (→ Rem. 2.3.2.19) we find

$$\left( \begin{array}{c|c} (\mathbf{A})_{1:n-1,1:n-1} & (\mathbf{A})_{1:n-1,n} \\ \hline (\mathbf{A})_{n,1:n-1} & (\mathbf{A})_{n,n} \end{array} \right) = \left( \begin{array}{c|c} \mathbf{L}_1 & 0 \\ \hline \mathbf{1}^\top & 1 \end{array} \right) \left( \begin{array}{c|c} \mathbf{U}_1 & \mathbf{u} \\ \hline 0 & \gamma \end{array} \right), \quad (2.7.5.8)$$

$$\Rightarrow (\mathbf{A})_{1:n-1,1:n-1} = \mathbf{L}_1 \mathbf{U}_1, \quad \mathbf{L}_1 \mathbf{u} = (\mathbf{A})_{1:n-1,n}, \quad \mathbf{U}_1^\top \mathbf{1} = (\mathbf{A})_{n,1:n-1}^\top, \quad \mathbf{1}^\top \mathbf{u} + \gamma = (\mathbf{A})_{n,n}. \quad (2.7.5.9)$$

### C++ code 2.7.5.10: Envelope aware recursive LU-factorization → GITLAB

```

2 //! envelope aware recursive LU-factorization
3 //! of structurally symmetric matrix
4 void luenv(const MatrixXd &A, MatrixXd &L, MatrixXd &U) {
5     int n = A.cols();
```

```

6   assert(n == A.rows() && "A must be square");
7   if(n == 1){ L.setIdentity(); U = A;}
8   else{
9     VectorXi mr = rowbandwidth(A); // = colbandwidth thanks to symmetry
10    double gamma;
11    MatrixXd L1(n-1,n-1), U1(n-1,n-1);
12    luenv(A.topLeftCorner(n-1, n-1), L1, U1);
13    VectorXd u = substenv(L1, A.col(n-1).head(n-1), mr);
14    VectorXd l = substenv(U1.transpose(), A.row(n-1).head(n-1).transpose(),
15                           mr);
16    if(mr(n-1) > 0)
17      gamma = A(n-1,n-1) - l.tail(mr(n-1)).dot(u.tail(mr(n-1)));
18    else
19      gamma = A(n-1,n-1);
20    L.topLeftCorner(n-1,n-1) = L1;      L.col(n-1).setZero();
21    L.row(n-1).head(n-1) = l.transpose();  L(n-1,n-1) = 1;
22    U.topLeftCorner(n-1,n-1) = U1;      U.col(n-1).head(n-1) = u;
23    U.row(n-1).setZero();      U(n-1,n-1) = gamma;
24  }
}

```

recursive implementation of envelope aware recursive LU-factorization (no pivoting !)

Assumption:  $\mathbf{A} \in \mathbb{K}^{n,n}$  is structurally symmetric

► Asymptotic complexity ( $\mathbf{A} \in \mathbb{K}^{n,n}$ )  $O(n \cdot \# \text{env}(\mathbf{A}))$  for  $n \rightarrow \infty$ .

### Definition 2.7.5.11. Structurally symmetric matrix

$\mathbf{A} \in \mathbb{K}^{n,n}$  is structurally symmetric, if

$$(\mathbf{A})_{i,j} \neq 0 \Leftrightarrow (\mathbf{A})_{j,i} \neq 0 \quad \forall i, j \in \{1, \dots, n\}.$$

Since by Thm. 2.7.5.4 fill-in is confined to the envelope, we need store only the matrix entries  $a_{ij}$ ,  $(i, j) \in \text{env}(\mathbf{A})$  when computing (in situ) LU-factorization of structurally symmetric  $\mathbf{A} \in \mathbb{K}^{n,n}$

- Storage required:  $n + 2 \sum_{i=1}^n m_i(\mathbf{A})$  floating point numbers
- terminology: envelope oriented matrix storage

**EXAMPLE 2.7.5.12 (Envelope oriented matrix storage)** Linear envelope oriented matrix storage of symmetric  $\mathbf{A} = \mathbf{A}^\top \in \mathbb{R}^{n,n}$ :

Two arrays:

scalar\_t \* val size  $P$ ,

size\_t \* dptr size  $n$

$$P := n + \sum_{i=1}^n m_i(\mathbf{A}). \quad (2.7.5.13)$$

$$\mathbf{A} = \begin{pmatrix} * & 0 & * & 0 & 0 & 0 & 0 \\ 0 & * & 0 & 0 & * & 0 & 0 \\ * & 0 & * & 0 & 0 & * & 0 \\ 0 & 0 & 0 & * & * & * & 0 \\ 0 & * & 0 & * & * & * & 0 \\ 0 & 0 & 0 & 0 & * & * & 0 \\ 0 & 0 & * & * & 0 & 0 & * \end{pmatrix}$$

Indexing rule:

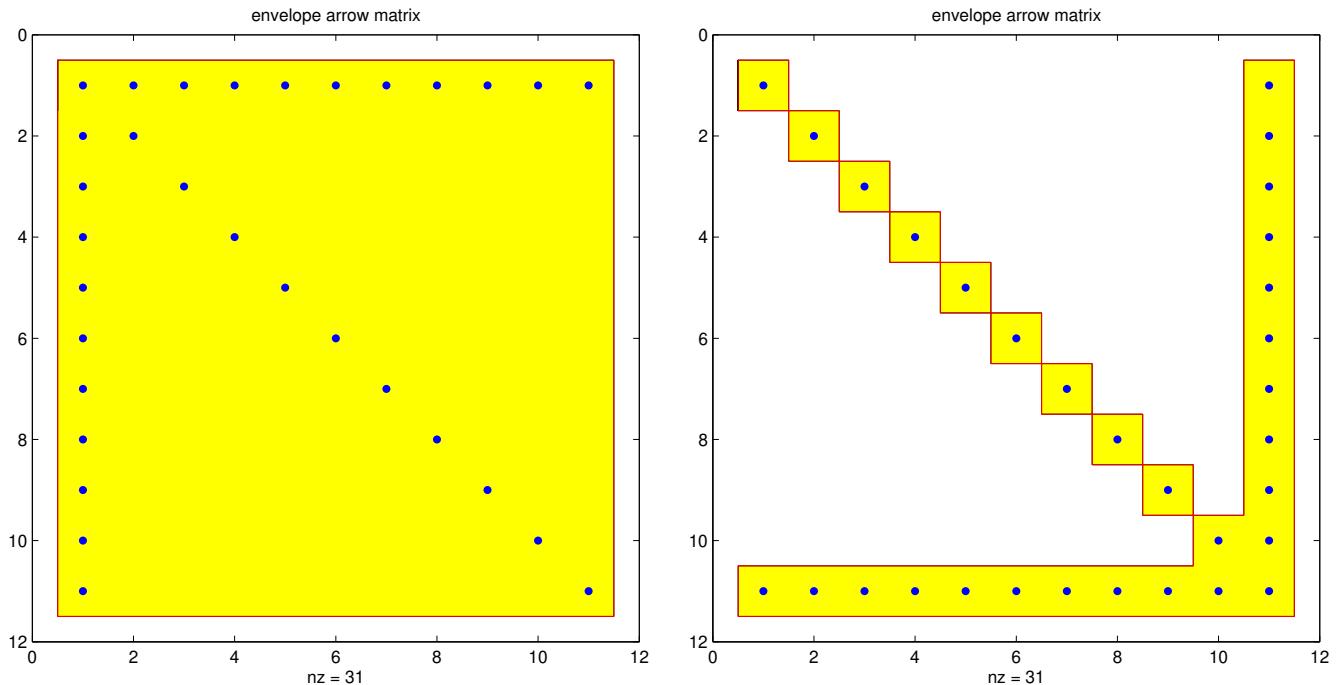
$$\begin{aligned} \text{dptr}[j] &= k \\ \Updownarrow \\ \text{val}[k] &= a_{jj} \end{aligned}$$

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
val	$a_{11}$	$a_{22}$	$a_{31}$	$a_{32}$	$a_{33}$	$a_{44}$	$a_{52}$	$a_{53}$	$a_{54}$	$a_{55}$	$a_{65}$	$a_{66}$	$a_{73}$	$a_{74}$	$a_{75}$	$a_{76}$	$a_{77}$	
dptr	0	1	2	5	6	10	12	17										

### Minimizing bandwidth/envelope:

Goal: Minimize  $m_i(\mathbf{A})$ ,  $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{N,N}$ , by permuting rows/columns of  $\mathbf{A}$

**EXAMPLE 2.7.5.14 (Reducing bandwidth by row/column permutations)** Recall: cyclic permutation of rows/columns of arrow matrix applied in Ex. 2.7.4.6. This can be viewed as a drastic shrinking of the envelope:



Another example: Reflection at cross diagonal  $\rightarrow$  reduction of  $\# \text{env}(\mathbf{A})$

$$\begin{bmatrix} * & 0 & 0 & * & * & * \\ 0 & * & 0 & 0 & 0 & 0 \\ 0 & 0 & * & 0 & 0 & 0 \\ * & 0 & 0 & 0 & * & * \\ * & 0 & 0 & * & * & * \\ * & 0 & 0 & * & * & * \end{bmatrix} \xrightarrow{i \leftarrow N+1-i} \begin{bmatrix} * & * & * & 0 & 0 & * \\ * & * & * & 0 & 0 & * \\ * & * & * & 0 & 0 & * \\ 0 & 0 & 0 & * & 0 & 0 \\ 0 & 0 & 0 & * & 0 & 0 \\ * & * & * & 0 & 0 & * \end{bmatrix}$$

$\# \text{env}(\mathbf{A}) = 30$        $\# \text{env}(\mathbf{A}) = 22$

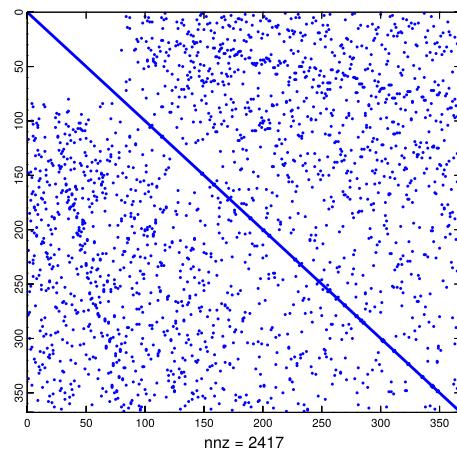
### EXAMPLE 2.7.5.15 (Reducing fill-in by reordering)

Envelope reducing permutations are at the heart of all modern sparse solvers ( $\rightarrow$  § 2.7.3.5). They employ elaborate algorithms for the analysis of **matrix graphs**, that is, the connections between components of the vector of unknowns defined by non-zero entries of the matrix. For further discussion see [AG11, Sect. 5.7].

EIGEN supplies a few ordering methods for sparse matrices. These methods use permutations to aim for minimal bandwidth/envelope of a given sparse matrix. We study an example with a  $347 \times 347$  matrix **M** originating in the numerical solution of partial differential equations, *cf.* Rem. 2.7.0.5.

Pattern of **M** 

(Here: no row swaps from pivoting !)



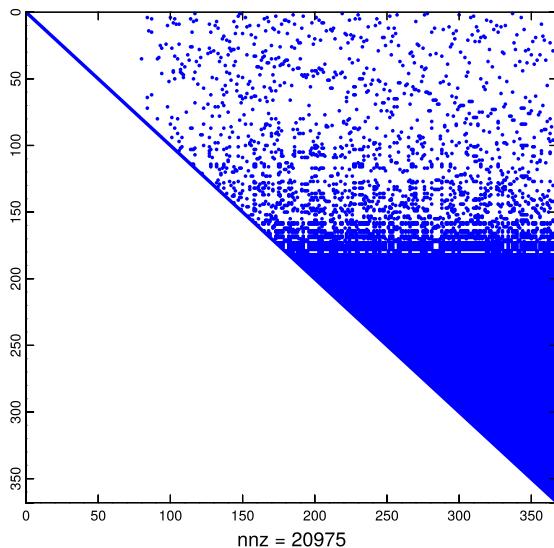
#### C++ code 2.7.5.16: preordering in EIGEN $\rightarrow$ GITLAB

```

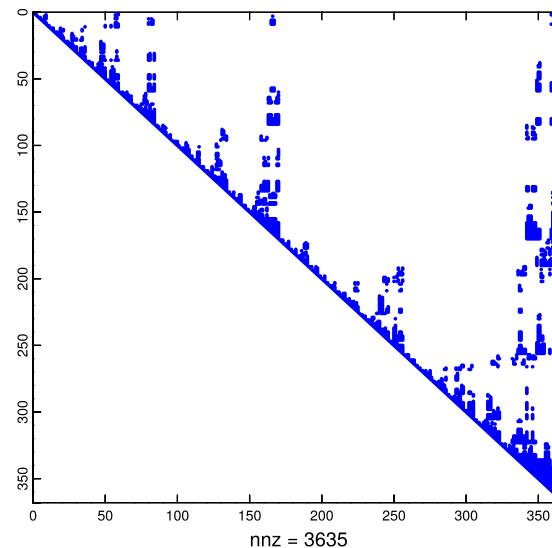
2 // L and U cannot be extracted from SparseLU -> LDLT
3 SimplicialLDLT<SpMat_t, Lower, AMDOrdering<int> > solver1 (M) ;
4 SimplicialLDLT<SpMat_t, Lower, NaturalOrdering <int> > solver2 (M) ;
5 MatrixXd U1 = MatrixXd(solver1.matrixU());
6 MatrixXd U2 = MatrixXd(solver2.matrixU());
7 // Plotting
8 mgl::Figure fig1, fig2, fig3;
9 fig1.spy(M); fig1.setFontSize(4); fig1.save("Mspy");
10 fig2.spy(U1); fig2.setFontSize(4); fig2.save("AMDUSpy");
11 fig3.spy(U2); fig3.setFontSize(4); fig3.save("NaturalUSpy");

```

Examine patterns of LU-factors ( $\rightarrow$  Sect. 2.3.2) after reordering:



no reordering



approximate minimum degree

## 2.8 Stable Gaussian elimination without pivoting

Recall insights, examples and experiments:

- Thm. 2.7.5.4 ➤ special structure of the matrix helps avoid fill-in in Gaussian elimination/LU-factorization *without* pivoting.
  - Ex. 2.7.4.9 ➤ pivoting can trigger huge fill-in that would not occur without it.
  - Ex. 2.7.5.15 ➤ fill-in reducing effect of reordering can be thwarted by later row swapping in the course of pivoting.
  - **BUT** pivoting is essential for stability of Gaussian elimination/LU-factorization → Ex. 2.3.3.1.
- Very desirable: a priori criteria, when Gaussian elimination/LU-factorization remains stable even without pivoting. This can help avoid the extra work for partial pivoting and makes it possible to exploit structure without worrying about stability.

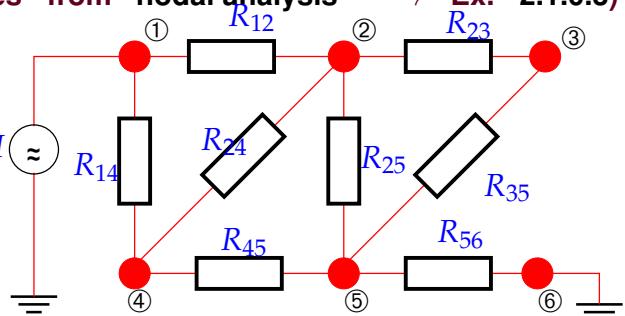
This section will introduce classes of matrices that allow Gaussian elimination without pivoting. Fortunately, linear systems of equations featuring system matrices from these classes are very common in applications.

### EXAMPLE 2.8.0.1 (Diagonally dominant matrices from nodal analysis → Ex. 2.1.0.3)

Consider:

electrical circuit entirely composed of Ohmic resistors.

Circuit equations from nodal analysis, see Ex. 2.1.0.3:



$$\textcircled{2}: R_{12}^{-1}(U_2 - U_1) + R_{23}^{-1}(U_2 - U_3) + R_{24}^{-1}(U_2 - U_4) + R_{25}^{-1}(U_2 - U_5) = 0,$$

$$\textcircled{3}: R_{23}^{-1}(U_3 - U_2) + R_{35}^{-1}(U_3 - U_5) = 0,$$

$$\textcircled{4}: R_{14}^{-1}(U_4 - U_1) + R_{24}^{-1}(U_4 - U_2) + R_{45}^{-1}(U_4 - U_5) = 0,$$

$$\textcircled{5}: R_{25}^{-1}(U_5 - U_2) + R_{35}^{-1}(U_5 - U_3) + R_{45}^{-1}(U_5 - U_4) + R_{56}^{-1}(U_5 - U_6) = 0,$$

$$U_1 = U, \quad U_6 = 0.$$



$$\begin{pmatrix} \frac{1}{R_{12}} + \frac{1}{R_{23}} + \frac{1}{R_{24}} + \frac{1}{R_{25}} & -\frac{1}{R_{23}} & -\frac{1}{R_{24}} & -\frac{1}{R_{25}} \\ -\frac{1}{R_{23}} & \frac{1}{R_{23}} + \frac{1}{R_{35}} & 0 & -\frac{1}{R_{35}} \\ -\frac{1}{R_{24}} & 0 & \frac{1}{R_{24}} + \frac{1}{R_{45}} & -\frac{1}{R_{45}} \\ -\frac{1}{R_{25}} & -\frac{1}{R_{35}} & -\frac{1}{R_{45}} & \frac{1}{R_{22}} + \frac{1}{R_{35}} + \frac{1}{R_{45}} + \frac{1}{R_{56}} \end{pmatrix} \begin{pmatrix} U_2 \\ U_3 \\ U_4 \\ U_5 \end{pmatrix} = \begin{pmatrix} \frac{1}{R_{12}} \\ 0 \\ \frac{1}{R_{14}} \\ 0 \end{pmatrix} U$$

➤ Matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  arising from nodal analysis satisfies

- $\mathbf{A} = \mathbf{A}^\top, \quad a_{kk} > 0, \quad a_{kj} \leq 0 \text{ for } k \neq j, \quad (2.8.0.2)$

$$\bullet \quad \sum_{j=1}^n a_{kj} \geq 0, \quad k = 1, \dots, n, \quad (2.8.0.3)$$

- $\mathbf{A}$  is regular.  $(2.8.0.4)$

All these properties are obvious except for the fact that  $\mathbf{A}$  is regular.

Proof of (2.8.0.4): By Thm. 2.2.1.4 it suffices to show that the nullspace of  $\mathbf{A}$  is trivial:  $\mathbf{Ax} = \mathbf{0} \Rightarrow \mathbf{x} = \mathbf{0}$ .

Pick  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{Ax} = \mathbf{0}$ , and  $i \in \{1, \dots, n\}$  so that

$$|x_i| = \max\{|x_j|, j = 1, \dots, n\}.$$

Intermediate goal: show that all entries of  $\mathbf{x}$  are the same

$$\mathbf{Ax} = \mathbf{0} \Rightarrow x_i = \sum_{j \neq i} \frac{a_{ij}}{a_{ii}} x_j \Rightarrow |x_i| \leq \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|} |x_j|. \quad (2.8.0.5)$$

By (2.8.0.3) and the sign condition from (2.8.0.2) we conclude

$$\sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|} \leq 1. \quad (2.8.0.6)$$

Hence, (2.8.0.6) combined with the above estimate (2.8.0.5) that tells us that the maximum is smaller equal than a mean implies  $|x_j| = |x_i|$  for all  $j = 1, \dots, n$ . Finally, the sign condition  $a_{kj} \leq 0$  for  $k \neq j$  enforces the same sign of all  $x_i$ . Thus, we conclude, w.l.o.g.,  $x_1 = x_2 = \dots = x_n$ . As

$$\exists i \in \{1, \dots, n\}: \sum_{j=1}^n a_{ij} > 0 \quad (\text{strict inequality}),$$

$\mathbf{Ax} = \mathbf{0}$  is only possible for  $\mathbf{x} = \mathbf{0}$ . □

### §2.8.0.7 (Diagonally dominant matrices)

**Definition 2.8.0.8. Diagonally dominant matrix** → [QSS00, Def. 1.24]

$\mathbf{A} \in \mathbb{K}^{n,n}$  is **diagonally dominant**, if

$$\forall k \in \{1, \dots, n\}: \sum_{j \neq k} |a_{kj}| \leq |a_{kk}|.$$

The matrix  $\mathbf{A}$  is called **strictly diagonally dominant**, if

$$\forall k \in \{1, \dots, n\}: \sum_{j \neq k} |a_{kj}| < |a_{kk}|.$$

### Lemma 2.8.0.9. LU-factorization of diagonally dominant matrices

$$\mathbf{A} \text{ regular, diagonally dominant with positive diagonal} \Leftrightarrow \left\{ \begin{array}{l} \mathbf{A} \text{ has LU-factorization} \\ \Updownarrow \\ \text{Gaussian elimination feasible without pivoting}^{(*)} \end{array} \right.$$

(\*) partial pivoting & diagonally dominant matrices  $\succ$  triggers no row permutations !

((2.3.3.7) will always be satisfied for  $j = k$ )

$\succ$  Pivoting can be dispensed with without compromising stability.

*Proof.*(of Lemma 2.8.0.9)

(2.3.1.12) → induction w.r.t.  $n$ :

Clear: partial pivoting in the first step selects  $a_{11}$  as pivot element, cf. (2.3.3.7).

After 1st step of elimination:

$$a_{ij}^{(1)} = a_{ij} - \frac{a_{i1}}{a_{11}} a_{1j}, \quad i, j = 2, \dots, n \Rightarrow a_{ii}^{(1)} > 0.$$

$$\begin{aligned} \blacktriangleright |a_{ii}^{(1)}| - \sum_{\substack{j=2 \\ j \neq i}}^n |a_{ij}^{(1)}| &= \left| a_{ii} - \frac{a_{i1}}{a_{11}} a_{1i} \right| - \sum_{\substack{j=2 \\ j \neq i}}^n \left| a_{ij} - \frac{a_{i1}}{a_{11}} a_{1j} \right| \\ &\geq a_{ii} - \frac{|a_{i1}| |a_{1i}|}{a_{11}} - \sum_{\substack{j=2 \\ j \neq i}}^n |a_{ij}| - \frac{|a_{i1}|}{a_{11}} \sum_{\substack{j=2 \\ j \neq i}}^n |a_{1j}| \\ &\geq a_{ii} - \frac{|a_{i1}| |a_{1i}|}{a_{11}} - \sum_{\substack{j=2 \\ j \neq i}}^n |a_{ij}| - |a_{i1}| \frac{a_{11} - |a_{1i}|}{a_{11}} \geq a_{ii} - \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \geq 0. \end{aligned}$$

A regular, diagonally dominant  $\Rightarrow$  partial pivoting according to (2.3.3.7) selects  $i$ -th row in  $i$ -th step.  $\square$

**§2.8.0.10 (Gaussian elimination for symmetric positive definite (s.p.d.) matrices)** The class of symmetric positive definite (s.p.d.) matrices has been defined in Def. 1.1.2.6. They permit stable Gaussian elimination without pivoting:

**Theorem 2.8.0.11. Gaussian elimination for s.p.d. matrices**

Every symmetric/Hermitian positive definite matrix ( $\rightarrow$  Def. 1.1.2.6) possesses an LU-decomposition ( $\rightarrow$  Sect. 2.3.2).

Equivalent to the assertion of the theorem: Gaussian elimination is feasible *without pivoting*

In fact, this theorem is a corollary of Lemma 2.3.2.4, because all principal minors of an s.p.d. matrix are s.p.d. themselves.

*Sketch of alternative self-contained proof.*

Proof by induction: consider first step of elimination

$$\mathbf{A} = \begin{pmatrix} a_{11} & \mathbf{b}^\top \\ \mathbf{b} & \tilde{\mathbf{A}} \end{pmatrix} \xrightarrow[\text{Gaussian elimination}]{\text{1. step}} \begin{pmatrix} a_{11} & \mathbf{b}^\top \\ 0 & \tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^\top}{a_{11}} \end{pmatrix}.$$

$\gg$  to show:  $\tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^\top}{a_{11}}$  s.p.d. ( $\rightarrow$  step of induction argument)



Evident: symmetry of  $\tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^\top}{a_{11}} \in \mathbb{R}^{n-1, n-1}$



As  $\mathbf{A}$  s.p.d. ( $\rightarrow$  Def. 1.1.2.6), for every  $\mathbf{y} \in \mathbb{R}^{n-1} \setminus \{0\}$

$$0 < \left( -\frac{\mathbf{b}^\top \mathbf{y}}{a_{11}} \right)^\top \begin{pmatrix} a_{11} & \mathbf{b}^\top \\ \mathbf{b} & \tilde{\mathbf{A}} \end{pmatrix} \begin{pmatrix} -\frac{\mathbf{b}^\top \mathbf{y}}{a_{11}} \\ \mathbf{y} \end{pmatrix} = \mathbf{y}^\top (\tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^\top}{a_{11}})\mathbf{y}.$$

►  $\tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^\top}{a_{11}}$  positive definite. □

The proof can also be based on the identities

$$\begin{aligned} \left( \begin{array}{c|c} (\mathbf{A})_{1:n-1,1:n-1} & (\mathbf{A})_{1:n-1,n} \\ \hline (\mathbf{A})_{n,1:n-1} & (\mathbf{A})_{n,n} \end{array} \right) &= \left( \begin{array}{c|c} \mathbf{L}_1 & 0 \\ \hline \mathbf{1}^\top & 1 \end{array} \right) \left( \begin{array}{c|c} \mathbf{U}_1 & \mathbf{u} \\ \hline 0 & \gamma \end{array} \right), \\ \Rightarrow (\mathbf{A})_{1:n-1,1:n-1} &= \mathbf{L}_1 \mathbf{U}_1, \quad \mathbf{L}_1 \mathbf{u} = (\mathbf{A})_{1:n-1,n}, \quad \mathbf{U}_1^\top \mathbf{1} = (\mathbf{A})_{n,1:n-1}^\top, \quad \mathbf{1}^\top \mathbf{u} + \gamma = (\mathbf{A})_{n,n}, \end{aligned} \quad (2.7.5.8)$$

noticing that the principal minor  $(\mathbf{A})_{1:n-1,1:n-1}$  is also s.p.d. This allows a simple induction argument.

Note: no pivoting required ( $\rightarrow$  Sect. 2.3.3)  
 (partial pivoting always picks current pivot row)

The next result gives a useful criterion for telling whether a given symmetric/Hermitian matrix is s.p.d.:

### Lemma 2.8.0.12. Diagonal dominance and definiteness

A diagonally dominant Hermitian/symmetric matrix with non-negative diagonal entries is positive semi-definite.

A strictly diagonally dominant Hermitian/symmetric matrix with positive diagonal entries is positive definite.

*Proof.* For  $\mathbf{A} = \mathbf{A}^H$  diagonally dominant, use inequality between arithmetic and geometric mean (AGM)  $ab \leq \frac{1}{2}(a^2 + b^2)$ :

$$\begin{aligned} \mathbf{x}^H \mathbf{A} \mathbf{x} &= \sum_{i=1}^n \left( a_{ii} |x_i|^2 + \sum_{i \neq j} a_{ij} \bar{x}_i x_j \right) \geq \sum_{i=1}^n \left( a_{ii} |x_i|^2 - \sum_{i \neq j} |a_{ij}| |x_i| |x_j| \right) \\ &\stackrel{\text{AGM}}{\geq} \sum_{i=1}^n a_{ii} |x_i|^2 - \frac{1}{2} \sum_{i \neq j} |a_{ij}| (|x_i|^2 + |x_j|^2) \\ &\geq \frac{1}{2} \left( \sum_{i=1}^n \{a_{ii} |x_i|^2 - \sum_{j \neq i} |a_{ij}| |x_i|^2\} \right) + \frac{1}{2} \left( \sum_{j=1}^n \{a_{jj} |x_j|^2 - \sum_{i \neq j} |a_{ij}| |x_j|^2\} \right) \\ &\geq \sum_{i=1}^n |x_i|^2 \left( a_{ii} - \sum_{j \neq i} |a_{ij}| \right) \geq 0. \end{aligned}$$

□

### §2.8.0.13 (Cholesky decomposition)

**Lemma 2.8.0.14. Cholesky decomposition for s.p.d. matrices**  $\rightarrow$  [Gut09, Sect. 3.4], [Han02, Sect. II.5], [QSS00, Thm. 3.6]

For any s.p.d.  $\mathbf{A} \in \mathbb{K}^{n,n}$ ,  $n \in \mathbb{N}$ , there is a unique upper triangular matrix  $\mathbf{R} \in \mathbb{K}^{n,n}$  with  $r_{ii} > 0$ ,  $i = 1, \dots, n$ , such that  $\mathbf{A} = \mathbf{R}^H \mathbf{R}$  (Cholesky decomposition).

Thm. 2.8.0.11  $\Rightarrow \mathbf{A} = \mathbf{LU}$  (unique *LU*-decomposition of  $\mathbf{A}$ , Lemma 2.3.2.4)

$$\mathbf{A} = \mathbf{LD}\tilde{\mathbf{U}} \quad , \quad \begin{aligned} \mathbf{D} &\triangleq \text{diagonal of } \mathbf{U}, \\ \tilde{\mathbf{U}} &\triangleq \text{normalized upper triangular matrix} \rightarrow \text{Def. 1.1.2.3} \end{aligned}$$

Due to uniqueness of *LU*-decomposition

$$\mathbf{A} = \mathbf{A}^\top \Rightarrow \mathbf{U} = \mathbf{DL}^\top \Rightarrow \boxed{\mathbf{A} = \mathbf{LDL}^\top},$$

with unique  $\mathbf{L}$ ,  $\mathbf{D}$  (diagonal matrix)

$$\mathbf{x}^\top \mathbf{Ax} > 0 \quad \forall \mathbf{x} \neq 0 \Rightarrow \mathbf{y}^\top \mathbf{Dy} > 0 \quad \forall \mathbf{y} \neq 0.$$

►  $\mathbf{D}$  has positive diagonal  $\Rightarrow \mathbf{R} = \sqrt{\mathbf{DL}^\top}$ . □

Formulas analogous to (2.3.2.7)

$$\mathbf{R}^H \mathbf{R} = \mathbf{A} \Rightarrow a_{ik} = \sum_{j=1}^{\min\{i,k\}} \bar{r}_{ji} r_{jk} = \begin{cases} \sum_{j=1}^{i-1} \bar{r}_{ji} r_{jk} + \bar{r}_{ii} r_{ik} & , \text{if } i < k, \\ \sum_{j=1}^{i-1} |r_{ji}|^2 + r_{ii}^2 & , \text{if } i = k. \end{cases} \quad (2.8.0.15)$$

### C++ code 2.8.0.16: Simple Cholesky factorization → GITLAB

```

2 //! simple Cholesky factorization
3 void cholfac(const MatrixXd &A, MatrixXd &R){
4     int n = A.rows();
5     R = A;
6     for(int k = 0; k < n; ++k){
7         for(int j = k+1; j < n; ++j)
8             R.row(j).tail(n-j) -= R.row(k).tail(n-j)*R(k,j)/R(k,k);
9         R.row(k).tail(n-k) /= std::sqrt(R(k,k));
10    }
11    R.triangularView<StrictlyLower>().setZero();
12 }
```

Computational costs (# elementary arithmetic operations) of Cholesky decomposition:  $\frac{1}{6}n^3 + O(n^2)$

(> “half the costs” of LU-factorization, cf. Code in § 2.3.2.6, but this does not mean “twice as fast” in a concrete implementation, because memory access patterns will have a crucial impact, see Rem. 1.4.1.5.)

Gains of efficiency hardly justify the use of Cholesky decomposition in modern numerical algorithms. Savings in memory compared to standard LU-factorization (only one factor  $\mathbf{R}$  has to be stored) offer a stronger reason to prefer the Cholesky decomposition.

The computation of Cholesky-factorization by means of the algorithm of Code 2.8.0.16 is **numerically stable** ( $\rightarrow$  Def. 1.5.5.19)!

Reason: recall Thm. 2.4.0.4: Numerical instability of Gaussian elimination (with any kind of pivoting) manifests itself in massive growth of the entries of intermediate elimination matrices  $\mathbf{A}^{(k)}$ .

Use the relationship between LU-factorization and Cholesky decomposition, which tells us that we only have to monitor the growth of entries of intermediate upper triangular “Cholesky factorization matrices”  $\mathbf{A} = (\mathbf{R}^{(k)})^H \mathbf{R}^{(k)}$ .

Consider: Euclidean vector norm/matrix norm ( $\rightarrow$  Def. 1.5.5.10)  $\|\cdot\|_2$

$$\mathbf{A} = \mathbf{R}^H \mathbf{R} \Rightarrow \|\mathbf{A}\|_2 = \sup_{\|\mathbf{x}\|_2=1} \mathbf{x}^H \mathbf{R}^H \mathbf{R} \mathbf{x} = \sup_{\|\mathbf{x}\|_2=1} (\mathbf{R} \mathbf{x})^H (\mathbf{R} \mathbf{x}) = \|\mathbf{R}\|_2^2.$$

► For all intermediate Cholesky factorization matrices holds:  $\left\|(\mathbf{R}^{(k)})^H\right\|_2 = \left\|\mathbf{R}^{(k)}\right\|_2 = \|\mathbf{A}\|_2^{1/2}$

This rules out a blowup of entries of the  $\mathbf{R}^{(k)}$ .

Computation of the Cholesky decomposition largely agrees with the computation of LU-factorization (without pivoting). Using the latter together with forward and backward substitution ( $\rightarrow$  Sect. 2.3.2) to solve a linear system of equations is algebraically and numerically equivalent to using Gaussian elimination without pivoting.

From these equivalences we conclude:

Solving LSE with s.p.d. system matrix via

Cholesky decomposition + forward & backward substitution

is numerically stable ( $\rightarrow$  Def. 1.5.5.19)



Gaussian elimination *without pivoting* is a *numerically stable* way to solve LSEs with s.p.d. system matrix.

## Learning Outcomes

Principal take-home knowledge and skills from this chapter:

- A clear understanding of the algorithm of Gaussian elimination with and without pivoting (prerequisite knowledge from linear algebra)
- Insight into the relationship between Gaussian elimination and LU-decomposition and the algorithmic relevance of LU-decomposition
- Awareness of the asymptotic complexity of dense Gaussian elimination, LU-decomposition, and elimination for special matrices
- Familiarity with “sparse matrices”: notion, data structures, initialization, benefits
- Insight into the reduced computational complexity of the direct solution of sparse linear systems of equations.

# Bibliography

- [AG11] Uri M. Ascher and Chen Greif. *A first course in numerical methods*. Vol. 7. Computational Science & Engineering. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2011, pp. xxii+552. DOI: [10.1137/1.9780898719987](https://doi.org/10.1137/1.9780898719987) (cit. on pp. 113, 121, 175).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 123, 128, 129, 132, 133, 169).
- [GGK14] W. Gander, M.J. Gander, and F. Kwok. *Scientific Computing*. Vol. 11. Texts in Computational Science and Engineering. Heidelberg: Springer, 2014 (cit. on p. 109).
- [GV89] G.H. Golub and C.F. Van Loan. *Matrix computations*. 2nd. Baltimore, London: John Hopkins University Press, 1989 (cit. on p. 135).
- [Gut09] M.H. Gutknecht. *Lineare Algebra*. Lecture Notes. SAM, ETH Zürich, 2009 (cit. on pp. 107, 113, 114, 121–124, 179).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 121, 122, 132, 179).
- [Hig02] N.J. Higham. *Accuracy and Stability of Numerical Algorithms*. 2nd ed. Philadelphia, PA: SIAM, 2002 (cit. on p. 135).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on pp. 104, 107, 108, 113, 114, 116, 121–123, 128–131).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 105, 107, 109, 113, 121, 124, 170, 177, 179).
- [SST06] A. Sankar, D.A. Spielman, and S.-H. Teng. “Smoothed analysis of the condition numbers and growth factors of matrices”. In: *SIAM J. Matrix Anal. Appl.* 28.2 (2006), pp. 446–476 (cit. on p. 138).
- [SG04] O. Schenk and K. Gärtnner. “Solving Unsymmetric Sparse Systems of Linear Equations with PARDISO”. In: *J. Future Generation Computer Systems* 20.3 (2004), pp. 475–487 (cit. on p. 162).
- [ST96] D.A. Spielman and Shang-Hua Teng. “Spectral partitioning works: planar graphs and finite element meshes”. In: *Foundations of Computer Science, 1996. Proceedings., 37th Annual Symposium on*. Oct. 1996, pp. 96–105. DOI: [10.1109/SFCS.1996.548468](https://doi.org/10.1109/SFCS.1996.548468) (cit. on p. 137).
- [TB97] L.N. Trefethen and D. Bau. *Numerical Linear Algebra*. Philadelphia, PA: SIAM, 1997 (cit. on pp. 135, 137).

# Chapter 3

## Direct Methods for Linear Least Squares Problems

In this chapter we study numerical methods for **overdetermined** linear systems of equations, that is, linear systems with a “tall” rectangular system matrix

$$\mathbf{x} \in \mathbb{R}^n: \quad \text{“}\mathbf{Ax} = \mathbf{b}\text{”}, \quad (3.0.0.1)$$

$\mathbf{b} \in \mathbb{R}^m, \quad \mathbf{A} \in \mathbb{R}^{m,n}, \quad m \geq n.$

$$\left[ \begin{array}{c|c} \mathbf{A} & \mathbf{x} \\ \hline \end{array} \right] = \left[ \begin{array}{c} \mathbf{b} \end{array} \right]$$

We point out that, in contrast to Chapter 1, Chapter 2, we will restrict ourselves to **real** linear systems in this chapter.

Note that the quotation marks in (3.0.0.1) indicate that this is not a well-defined problem in the sense of § 1.5.5.1;  $\mathbf{Ax} = \mathbf{b}$  does not define a mapping  $(\mathbf{A}, \mathbf{b}) \mapsto \mathbf{x}$ , because

- such a vector  $\mathbf{x} \in \mathbb{R}^n$  may not exist,
- and, even if it exists, it may not be unique.

Therefore, first we have to establish a crisp concept of what we mean by a “solution” of (3.0.0.1).

### Contents

---

3.0.1	Overdetermined Linear Systems of Equations: Examples . . . . .	184
<b>3.1</b>	<b>Least Squares Solution Concepts . . . . .</b>	<b>186</b>
3.1.1	Least Squares Solutions . . . . .	187
3.1.2	Normal Equations . . . . .	189
3.1.3	Moore-Penrose Pseudoinverse . . . . .	194
3.1.4	Sensitivity of Least Squares Problem . . . . .	195
<b>3.2</b>	<b>Normal Equation Methods [DR08, Sect. 4.2], [Han02, Ch. 11] . . . . .</b>	<b>197</b>
<b>3.3</b>	<b>Orthogonal Transformation Methods [DR08, Sect. 4.4.2] . . . . .</b>	<b>200</b>
3.3.1	Transformation Idea . . . . .	200
3.3.2	Orthogonal/Unitary Matrices . . . . .	201
3.3.3	QR-Decomposition [Han02, Sect. 13], [Gut09, Sect. 7.3] . . . . .	202
3.3.4	QR-Based Solver for Linear Least Squares Problems . . . . .	214
3.3.5	Modification Techniques for QR-Decomposition . . . . .	218

3.4	Singular Value Decomposition (SVD) . . . . .	224
3.4.1	SVD: Definition and Theory . . . . .	224
3.4.2	SVD in EIGEN . . . . .	227
3.4.3	Solving General Least-Squares Problems by SVD . . . . .	230
3.4.4	SVD-Based Optimization and Approximation . . . . .	233
3.5	Total Least Squares . . . . .	251
3.6	Constrained Least Squares . . . . .	252
3.6.1	Solution via Lagrangian Multipliers . . . . .	253
3.6.2	Solution via SVD . . . . .	254

### 3.0.1 Overdetermined Linear Systems of Equations: Examples

You may think that overdetermined linear systems of equations are exotic, but this is not true. Rather they are very common in mathematical models.

**EXAMPLE 3.0.1.1 (Linear parameter estimation in 1D)** From first principles it is *known* that two physical quantities  $x \in \mathbb{R}$  and  $y \in \mathbb{R}$  (e.g., pressure and density of an ideal gas) are related by a *linear relationship*

$$y = \alpha x + \beta \quad \text{for some unknown coefficients/parameters } \alpha, \beta \in \mathbb{R}. \quad (3.0.1.2)$$

We carry out  $m \in \mathbb{N}$  measurements that yield pairs  $(x_i, y_i) \in \mathbb{R}^2$ ,  $i = 1, \dots, m$ ,  $m \geq 2$ . If the measurements were perfect, we could expect that there exist  $\alpha, \beta \in \mathbb{R}$  such that  $y_i = \alpha x_i + \beta$  for all  $i = 1, \dots, m$ . This is an overdetermined linear system of equations of the form (3.0.0.1):

$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_m & 1 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \Leftrightarrow \mathbf{Ax} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{m,2}, \mathbf{b} \in \mathbb{R}^m, \mathbf{x} \in \mathbb{R}^2. \quad (3.0.1.3)$$

In practice inevitable (“random”) *measurement errors* will thwart the solvability of (3.0.1.3) and for  $m > 2$  the probability that a solution  $\begin{bmatrix} \alpha \\ \beta \end{bmatrix}$  exists is zero, see Rem. 3.1.0.2.  $\square$

**EXAMPLE 3.0.1.4 (Parameter estimation for a linear model)** Ex. 3.0.1.1 can be generalized to higher dimensions:

Given: measured data points  $(\mathbf{x}_i, y_i)$ ,  $\mathbf{x}_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$ ,  $i = 1, \dots, m$ ,  $m \geq n + 1$  ( $y_i$ ,  $\mathbf{x}_i$  affected by measurement errors).

Known: without measurement errors data would satisfy an *affine linear relationship*  $y = \mathbf{a}^\top \mathbf{x} + \beta$ , for some  $\mathbf{a} \in \mathbb{R}^n$ ,  $c \in \mathbb{R}$ .

Plugging in the measured quantities gives  $y_i = \mathbf{a}^\top \mathbf{x}_i + \beta$ ,  $i = 1, \dots, m$ , a linear system of equations of the form

$$\begin{bmatrix} \mathbf{x}_1^\top & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{a} \\ \beta \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \Leftrightarrow \mathbf{Ax} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{m,n+1}, \mathbf{b} \in \mathbb{R}^m, \mathbf{x} \in \mathbb{R}^{n+1}, \quad (3.0.1.5)$$

which is an overdetermined LSE, in case  $m > n + 1$ .  $\square$

**EXAMPLE 3.0.1.6 (Measuring the angles of a triangle [NS02, Sect. 5.1])** We measure the angles of a planar triangle and obtain  $\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}$  (in radians). In the case of perfect measurements the true angles  $\alpha, \beta, \gamma$  would satisfy

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix} = \begin{bmatrix} \tilde{\alpha} \\ \tilde{\beta} \\ \tilde{\gamma} \\ \pi \end{bmatrix}. \quad (3.0.1.7)$$

Measurement errors will inevitably make the measured angles fail to add up to  $\pi$  so that (3.0.1.7) will not have a solution  $[\alpha, \beta, \gamma]^T$ .

Then, why should we add this last equation? This is suggested by a tenet of data science that reads “You cannot afford not to use any piece of information available”. It turns out that solving (3.0.1.7) “in a suitable way” as discussed below in Section 3.1.1 enhances cancellation of measurement errors and gives better estimates for the angles. We will not discuss this here and refer to statistics for an explanation.

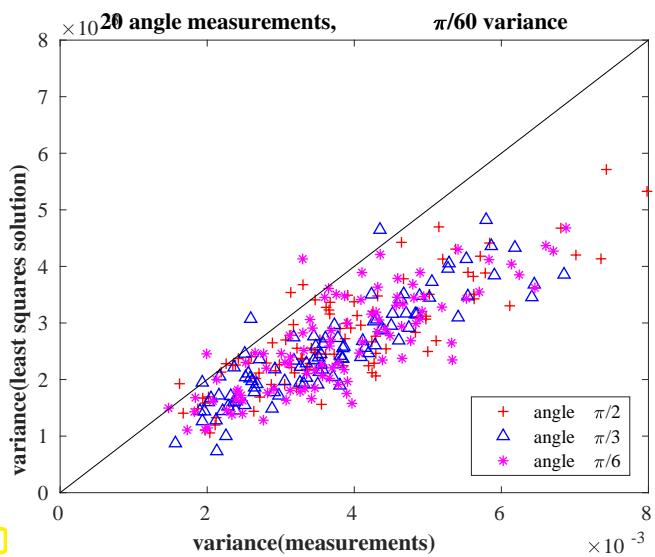


Fig. 69

We observe that in most runs the variance of the estimates from (3.0.1.7) are smaller than those of the raw data. □

### EXAMPLE 3.0.1.8 (Angles in a triangulation)

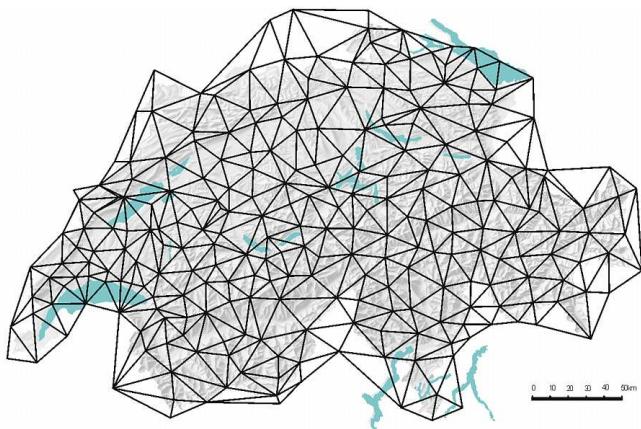


Fig. 70

Die Grundlagen seines Verfahrens hatte Gauss schon 1795 im Alter von 18 Jahren entwickelt. Basis war eine Idee von Pierre-Simon Laplace, die Beträge von Fehlern aufzusummen, so dass sich die Fehler zu Null addieren. Gauss nahm stattdessen die Fehlerquadrate und konnte die künstliche Zusatzanforderung an die Fehler weglassen.

Gauss benutzte dann das Verfahren intensiv bei seiner Vermessung des Königreichs Hannover durch Triangulation. 1821 und 1823 erschien die zweiteilige Arbeit sowie 1826 eine

Here we just report the results of a numerical experiment:

We consider the triangle with angles  $\pi/2$ ,  $\pi/3$ , and  $\pi/6$ . Synthetic “measurement errors” are introduced by adding a normally distributed random perturbation with mean 0 and standard deviation  $\pi/60$  to the exact values of the angles, yielding  $\tilde{\alpha}$ ,  $\tilde{\beta}$ , and  $\tilde{\gamma}$ .

For 20 “measurements” we compute the variance of the raw angles and that of the estimates obtained by solving (3.0.1.7) in least squares sense (→ Section 3.1.1). These variances are plotted for many different “runs”.

In Ex. 2.7.2.3 we learned about the concept and data structures for **planar triangulations** → Def. 2.7.2.4. Such triangulations have been and continue to be of fundamental importance for **geodesy**. In particular before distances could be measured accurately by means of lasers, triangulations were indispensable, because angles could already be determined with high precision. C.F. Gauss pioneered both the use of triangulations in geodesy and the use of the least squares method to deal with measurement errors → [Wikipedia](#).

Ergänzung zur *Theoria combinationis observationum erroribus minimis obnoxiae* (Theorie der den kleinsten Fehlern unterworfenen Kombination der Beobachtungen), in denen Gauss eine Begründung liefern konnte, weshalb sein Verfahren im Vergleich zu den anderen so erfolgreich war: Die Methode der kleinsten Quadrate ist in einer breiten Hinsicht optimal, also besser als andere Methoden.

We now extend Ex. 3.0.1.6 to planar triangulations, for which *measured* values for all internal angles are available. We obtain an overdetermined system of equations by combining the following linear relations:

1. each angle is supposed to be equal to its measured value,
2. the sum of interior angles is  $\pi$  for every triangle,
3. the sum of the angles at an interior node is  $2\pi$ .

If the planar triangulation has  $N_0$  interior vertices and  $M$  cells, then we end up with  $4M + N_0$  equations for the  $3M$  unknown angles.  $\square$

**EXAMPLE 3.0.1.9 ((Relative) point locations from distances [GGK14, Sect. 6.1])** Consider  $n$  points located on the real axis at unknown locations  $x_i \in \mathbb{R}$ ,  $i = 1, \dots, n$ . At least we know that  $x_i < x_{i+1}$ ,  $i = 1, \dots, n-1$ .

We measure the  $m := \binom{n}{2} = \frac{1}{2}n(n-1)$  pairwise distances  $d_{ij} := |x_i - x_j|$ ,  $i, j \in \{1, \dots, n\}$ ,  $i \neq j$ . They are connected to the point positions by the overdetermined linear system of equations

$$\begin{array}{l} x_i - x_j = d_{ij}, \\ 1 \leq j < i \leq n. \end{array} \quad \begin{array}{c} \uparrow \\ \mathbf{Ax} = \mathbf{b} \end{array} \quad \leftrightarrow \quad \left[ \begin{array}{cccccc} -1 & 1 & 0 & \dots & \dots & 0 \\ -1 & 0 & 1 & 0 & & \\ \vdots & & \ddots & \ddots & & \\ \vdots & & & \ddots & \ddots & \\ -1 & \dots & & & 0 & 1 \\ 0 & -1 & 1 & & 0 & \\ \vdots & & & & \vdots & \\ 0 & \dots & & & -1 & 1 \end{array} \right] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} d_{12} \\ d_{13} \\ \vdots \\ d_{1n} \\ d_{23} \\ \vdots \\ d_{n-1,n} \end{bmatrix} \quad (3.0.1.10)$$

Note that we can never expect a unique solution for  $\mathbf{x} \in \mathbb{R}^n$ , because adding a multiple of  $[1, 1, \dots, 1]^\top$  to any solution will again yield a solution, because  $\mathbf{A}$  has a **non-trivial kernel**:  $\mathcal{N}(\mathbf{A}) = [1, 1, \dots, 1]^\top$ . Non-uniqueness can be cured by setting  $x_1 := 0$ , thus removing one component of  $\mathbf{x}$ .

If the measurements were perfect, we could then find  $x_2, \dots, x_n$  from  $d_{i-1,i}$ ,  $i = 2, \dots, n$  by solving a standard (square) linear system of equations. However, as in Ex. 3.0.1.6, using much more information through the overdetermined system (3.0.1.10) helps curb measurement errors.  $\square$

## 3.1 Least Squares Solution Concepts

Throughout we consider the (possibly overdetermined) linear system of equations

$$\mathbf{x} \in \mathbb{R}^n: \quad \mathbf{Ax} = \mathbf{b}, \quad \mathbf{b} \in \mathbb{R}^m, \quad \mathbf{A} \in \mathbb{R}^{m,n}, \quad m \geq n. \quad (3.0.0.1)$$

Recall from linear algebra that  $\mathbf{Ax} = \mathbf{b}$  has a solution, if and only if the right hand side vector  $\mathbf{b}$  lies in the image (range space,  $\rightarrow$  Def. 2.2.1.2) of the matrix  $\mathbf{A}$ :

$$\exists \mathbf{x} \in \mathbb{R}^n: \quad \mathbf{Ax} = \mathbf{b} \iff \mathbf{b} \in \mathcal{R}(\mathbf{A}). \quad (3.1.0.1)$$

☞ Notation for important subspaces associated with a matrix  $\mathbf{A} \in \mathbb{K}^{m,n}$  ( $\rightarrow$  Def. 2.2.1.2)

$$\begin{aligned}\text{image/range: } \mathcal{R}(\mathbf{A}) &:= \{\mathbf{Ax}, \mathbf{x} \in \mathbb{K}^n\} \subset \mathbb{K}^m, \\ \text{kernel/nullspace: } \mathcal{N}(\mathbf{A}) &:= \{\mathbf{x} \in \mathbb{K}^n : \mathbf{Ax} = \mathbf{0}\}.\end{aligned}$$

**Remark 3.1.0.2 (Consistent right hand side vectors are highly improbable)** If  $\mathcal{R}(\mathbf{A}) \neq \mathbb{R}^m$ , then “almost all” perturbations of  $\mathbf{b}$  (e.g., due to measurement errors) will destroy  $\mathbf{b} \in \mathcal{R}(\mathbf{A})$ , because  $\mathcal{R}(\mathbf{A})$  is a “set of measure zero” in  $\mathbb{R}^m$ .  $\square$

### 3.1.1 Least Squares Solutions

#### Definition 3.1.1.1. Least squares solution

For given  $\mathbf{A} \in \mathbb{K}^{m,n}$ ,  $\mathbf{b} \in \mathbb{K}^m$  the vector  $\mathbf{x} \in \mathbb{R}^n$  is a **least squares solution** of the linear system of equations  $\mathbf{Ax} = \mathbf{b}$ , if

$$\begin{aligned}\mathbf{x} \in \operatorname{argmin}_{\mathbf{y} \in \mathbb{K}^n} \|\mathbf{Ay} - \mathbf{b}\|_2^2, \\ \Updownarrow \\ \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \min_{\mathbf{y} \in \mathbb{K}^n} \|\mathbf{Ay} - \mathbf{b}\|_2^2 = \min_{\mathbf{y} \in \mathbb{K}^n} \sum_{i=1}^m \left( \sum_{j=1}^n (\mathbf{A})_{i,j} y_j - b_i \right)^2.\end{aligned}$$

- A least squares solution is any vector  $\mathbf{x}$  that minimizes the Euclidean norm of the **residual**  $\mathbf{r} = \mathbf{b} - \mathbf{Ax}$ , see Def. 2.4.0.1.

We write  $\text{lsq}(\mathbf{A}, \mathbf{b})$  for the set of least squares solutions of the linear system of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ :

$$\text{lsq}(\mathbf{A}, \mathbf{b}) := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \text{ is a least squares solution of } \mathbf{Ax} = \mathbf{b}\} \subset \mathbb{R}^n. \quad (3.1.1.2)$$

**EXAMPLE 3.1.1.3 (linear regression  $\rightarrow$  [DR08, Ex. 4.1])** We consider the problem of parameter estimation for a linear model from Ex. 3.0.1.4:

Given: *measured* data points  $(\mathbf{x}_i, y_i)$ ,  $\mathbf{x}_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$ ,  $i = 1, \dots, m$ ,  $m \geq n + 1$   
 $(y_i, \mathbf{x}_i)$  affected by measurement errors).

Known: without measurement errors data would satisfy **affine linear relationship**

$$y = \mathbf{a}^\top \mathbf{x} + \beta, \quad (3.1.1.4)$$

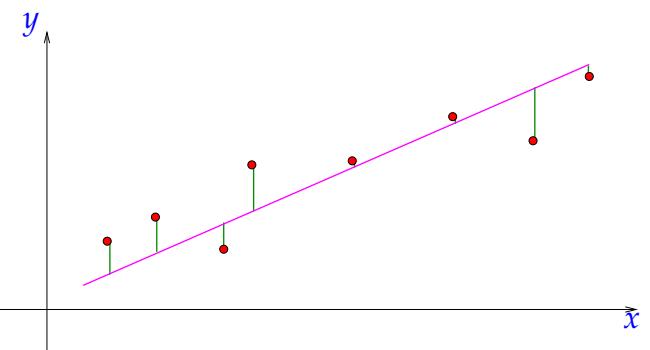
from some parameters  $\mathbf{a} \in \mathbb{R}^n$ ,  $\beta \in \mathbb{R}$ .

Solving the overdetermined linear system of equations in least squares sense we obtain a **least squares** estimate for the parameters  $\mathbf{a}$  and  $\beta$ :

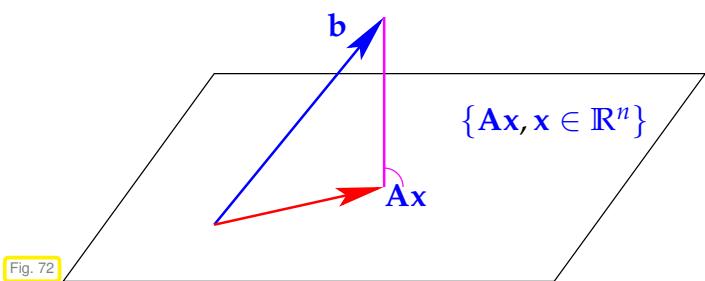
$$(\mathbf{a}, \beta) = \underset{\mathbf{a} \in \mathbb{R}^n, \beta \in \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^m |y_i - \mathbf{a}^\top \mathbf{x}_i - \beta|^2 \quad (3.1.1.5)$$

In statistics, solving (3.1.1.5) is known as **linear regression**

Linear regression for  $n = 1, m = 8$ : “fitting” a line to data points



**§3.1.1.6 (The geometry of least squares problems)** A geometric “proof” for the existence of least squares solutions ( $\mathbb{R} = \mathbb{R}$ )



For a least squares solution  $\mathbf{x} \in \mathbb{R}^n$  the vector  $\mathbf{Ax} \in \mathbb{R}^m$  is the unique **orthogonal projection** of  $\mathbf{b}$  onto

$$\mathcal{R}(\mathbf{A}) = \operatorname{Span}\{(\mathbf{A})_{:,1}, \dots, (\mathbf{A})_{:,n}\},$$

because the orthogonal projection provides the nearest (w.r.t. the Euclidean distance) point to  $\mathbf{b}$  in the subspace (hyperplane)  $\mathcal{R}(\mathbf{A})$ .

From this geometric consideration we conclude that  $\operatorname{lsq}(\mathbf{A}, \mathbf{b})$  is the space of solutions of  $\mathbf{Ax} = \mathbf{b}^*$ , where  $\mathbf{b}^*$  is the orthogonal projection of  $\mathbf{b}$  onto  $\mathcal{R}(\mathbf{A})$ . Since the set of solutions of a linear system of equations invariably is an affine space, this argument teaches that  $\operatorname{lsq}(\mathbf{A}, \mathbf{b})$  is an **affine subspace** of  $\mathbb{R}^n$ !

Geometric intuition yields the following insight:

### Theorem 3.1.1.7. Existence of least squares solutions

For any  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$  a least squares solution of  $\mathbf{Ax} = \mathbf{b}$  ( $\rightarrow$  Def. 3.1.1.1) exists.

*Proof.* The function  $F : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $F(\mathbf{x}) := \|\mathbf{b} - \mathbf{ax}\|_2^2$  is continuous, bounded from below by 0 and  $F(\mathbf{x}) \rightarrow \infty$  for  $\|\mathbf{x}\| \rightarrow \infty$ . Hence, there must be an  $\mathbf{x}^* \in \mathbb{R}^n$  for which it attains its minimum.  $\square$

**§3.1.1.8 (Least squares solution as maximum-likelihood estimator  $\rightarrow$  [DR08, Sect. 4.5]**) Extending the considerations of Ex. 3.0.1.4, a generic **linear parameter estimation** problem seeks to determine the unknown parameter vector  $\mathbf{x} \in \mathbb{R}^n$  from the linear relationship  $\mathbf{Ax} = \mathbf{y}$ , where  $\mathbf{A} \in \mathbb{R}^{m,n}$  is known and  $\mathbf{y} \in \mathbb{R}^m$  is accessible through measurements.

Unfortunately,  $\mathbf{y}$  is affected by measurement errors, Thus we model it as a **random vector**  $\mathbf{y} = \mathbf{y}(\omega)$ ,  $\omega \in \Omega$ ,  $\Omega$  the set of outcomes from a **probability space**.

The measurement errors in different components of  $\mathbf{y}$  are supposed be **unbiased** (expectation = 0), **independent**, and **identically normally distributed** with variance  $\sigma^2$ ,  $\sigma > 0$ , which means

$$\mathbf{y}(\omega) = \mathbf{Ax} + \mathbf{e}(\omega), \quad \omega \in \Omega, \quad (3.1.1.9)$$

where the probability distribution of  $\mathbf{e}$  satisfies

$$\mathbb{P}((\mathbf{e})_\ell \in I_\ell, \ell = 1, \dots, m) = \prod_{\ell=1}^m \frac{1}{\sigma\sqrt{2\pi}} \int_{I_\ell} \exp\left(-\frac{1}{2} \frac{z^2}{\sigma^2}\right) dz, \quad I_\ell \subset \mathbb{R}. \quad (3.1.1.10)$$

From (3.1.1.9) we infer that the probability density of  $\mathbf{y}$  is

$$L(\mathbf{x}; \mathbf{y}) = \prod_{\ell=1}^m \exp\left(-\frac{1}{2}\left(\frac{\mathbf{y}_\ell - (\mathbf{Ax})_\ell}{\sigma}\right)^2\right) = \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{Ax}\|_2^2\right). \quad (3.1.1.11)$$

The last identity follows from  $\exp(x)\exp(y) = \exp(x+y)$ . This probability density function  $\mathbf{y} \mapsto L(\mathbf{x}; \mathbf{y})$  is called the **likelihood** of  $\mathbf{y}$ , and the notation emphasizes its dependence on the parameters  $\mathbf{x}$ .

Assume that we are given a measurement (realization/sample)  $\mathbf{b}$  of  $\mathbf{y}$ . The maximum likelihood principle then suggests that we choose the parameters so that the probability density of  $\mathbf{y}$  becomes maximal at  $\mathbf{b}$ :

$$\mathbf{x}^* \in \mathbb{R}^n \text{ such that } \mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmax}} L(\mathbf{x}; \mathbf{b}) = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmax}} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{b} - \mathbf{Ax}\|_2^2\right).$$

Obviously, due to the monotonicity of  $\xi \mapsto \exp(\xi)$ ,  $\mathbf{x}^*$  is a least squares solution of  $\mathbf{Ax} = \mathbf{b}$  according to Def. 3.1.1.1:

$$\mathbf{x}^* \in \mathbb{R}^n \text{ such that } \mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \|\mathbf{b} - \mathbf{Ax}\|_2^2.$$

□

### 3.1.2 Normal Equations

Appealing to the geometric intuition gleaned from Fig. 72 we infer the orthogonality of  $\mathbf{b} - \mathbf{Ax}$ ,  $\mathbf{x}$  a least squares solution of the overdetermined linear systems of equations  $\mathbf{Ax} = \mathbf{b}$ , to all columns of  $\mathbf{A}$ :

$$\mathbf{b} - \mathbf{Ax} \perp \mathcal{R}(\mathbf{A}) \Leftrightarrow \mathbf{b} - \mathbf{Ax} \perp (\mathbf{A})_{:,j}, \quad j = 1, \dots, n \Leftrightarrow \mathbf{A}^\top(\mathbf{b} - \mathbf{Ax}) = 0.$$

Surprisingly, we have found a square linear system of equations satisfied by the least squares solution. The next theorem gives the formal statement of this discovery. It also completely characterizes  $\operatorname{lsq}(\mathbf{A}, \mathbf{b})$  and reveals a way to compute this set.

#### Theorem 3.1.2.1. Obtaining least squares solutions by solving normal equations

The vector  $\mathbf{x} \in \mathbb{R}^n$  is a least squares solution ( $\rightarrow$  Def. 3.1.1.1) of the linear system of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , if and only if it solves the **normal equations (NEQ)**

$$\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{b}. \quad (3.1.2.2)$$

Note that the normal equations (3.1.2.2) are an  $n \times n$  square linear system of equations with a *symmetric positive semi-definite* coefficient matrix:

$$\begin{bmatrix} \mathbf{A}^\top & & \\ & \mathbf{A} & \\ & & \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{A}^\top & & \\ & \mathbf{A} & \\ & & \mathbf{b} \end{bmatrix},$$

$$\Leftrightarrow \begin{bmatrix} \mathbf{A}^\top \mathbf{A} \end{bmatrix} \begin{bmatrix} \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{A}^\top \\ \vdots \end{bmatrix} \begin{bmatrix} \mathbf{b} \end{bmatrix}.$$

*Proof. (of Thm. 3.1.2.1)*

①: We first show that a least squares solution satisfies the normal equations. Let  $\mathbf{x} \in \mathbb{R}^n$  be a least squares solutions according to Def. 3.1.1.1. Pick an arbitrary  $\mathbf{d} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$  and define the function

$$\varphi_{\mathbf{d}} : \mathbb{R} \rightarrow \mathbb{R}, \quad \varphi_{\mathbf{d}}(\tau) := \|\mathbf{A}(\mathbf{x} + \tau\mathbf{d}) - \mathbf{b}\|_2^2. \quad (3.1.2.3)$$

We find the equivalent expression

$$\varphi_{\mathbf{d}}(\tau) = \tau^2 \mathbf{d}^\top \mathbf{A}^\top \mathbf{A} \mathbf{d} + 2\tau \mathbf{d}^\top \mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) + \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2,$$

which shows that  $\tau \mapsto \varphi_{\mathbf{d}}(\tau)$  is a smooth ( $C^\infty$ ) function.

Moreover, since every  $\mathbf{x} \in \text{lsq}(\mathbf{A}, \mathbf{b})$  is a minimizer of  $\mathbf{y} \mapsto \|\mathbf{A}\mathbf{y} - \mathbf{b}\|_2^2$ , we conclude that  $\tau \mapsto \varphi_{\mathbf{d}}(\tau)$  has a global minimum in  $\tau = 0$ . Necessarily,

$$\frac{d\varphi_{\mathbf{d}}}{d\tau} \Big|_{\tau=0} = 2\mathbf{d}^\top \mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) = 0.$$

Since this holds for *any* vector  $\mathbf{d} \neq 0$ , we conclude (set  $\mathbf{d}$  equal to all the Euclidean unit vectors in  $\mathbb{R}^n$ )

$$\mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) = \mathbf{0},$$

which is equivalent to the normal equations (3.1.2.2).

②: Let  $\mathbf{x}$  be a solution of the normal equations. Then we find by tedious but straightforward computations

$$\begin{aligned} & \|\mathbf{Ay} - \mathbf{b}\|_2^2 - \|\mathbf{Ax} - \mathbf{b}\|_2^2 \\ &= \mathbf{y}^\top \mathbf{A}^\top \mathbf{Ay} - 2\mathbf{y}^\top \mathbf{A}^\top \mathbf{b} + \mathbf{b}^\top \mathbf{b} - \mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} + 2\mathbf{x}^\top \mathbf{A}^\top \mathbf{b} - \mathbf{b}^\top \mathbf{b} \\ &\stackrel{(3.1.2.2)}{=} \mathbf{y}^\top \mathbf{A}^\top \mathbf{Ay} - 2\mathbf{y}^\top \mathbf{A}^\top \mathbf{Ax} + \mathbf{x}^\top \mathbf{A}^\top \mathbf{Ax} \\ &= (\mathbf{y} - \mathbf{x})^\top \mathbf{A}^\top \mathbf{A}(\mathbf{y} - \mathbf{x}) = \|\mathbf{A}(\mathbf{x} - \mathbf{y})\|_2^2 \geq 0. \\ &\implies \|\mathbf{Ay} - \mathbf{b}\| \geq \|\mathbf{Ax} - \mathbf{b}\|. \end{aligned}$$

Since this holds for *any*  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{x}$  must be a global minimizer of  $\mathbf{y} \mapsto \|\mathbf{Ay} - \mathbf{b}\|$ !  $\square$

**EXAMPLE 3.1.2.4 (Normal equations for some examples from Section 3.0.1)** Given  $\mathbf{A}$  and  $\mathbf{b}$  it takes only elementary linear algebra operations to form the normal equations

$$\mathbf{A}^\top \mathbf{A}\mathbf{x} = \mathbf{A}^\top \mathbf{b}. \quad (3.1.2.2)$$

- For Ex. 3.0.1.1,  $\mathbf{A} \in \mathbb{R}^{m,2}$  given in (3.0.1.3) we obtain the normal equations linear system

$$\begin{bmatrix} x_1 & x_2 & \dots & \dots & x_m \\ 1 & 1 & \dots & \dots & 1 \end{bmatrix} \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_m & 1 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} \|\mathbf{x}\|_2^2 & \mathbf{1}^\top \mathbf{x} \\ \mathbf{1}^\top \mathbf{x} & m \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} \mathbf{x}^\top \mathbf{y} \\ \mathbf{1}^\top \mathbf{y} \end{bmatrix},$$

with  $\mathbf{1} = [1, \dots, 1]^\top$ .

- In the case of Ex. 3.0.1.4 and the overdetermined  $m \times (n+1)$  linear system (3.0.1.5), the normal equations read

$$\begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \dots & \mathbf{x}_m \\ 1 & 1 & \dots & \dots & 1 \end{bmatrix} \begin{bmatrix} \mathbf{x}_1^\top & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{a} \\ \beta \end{bmatrix} = \begin{bmatrix} \mathbf{X}\mathbf{X}^\top & \mathbf{X}\mathbf{1} \\ \mathbf{1}^\top \mathbf{X}^\top & m \end{bmatrix} \begin{bmatrix} \mathbf{a} \\ \beta \end{bmatrix} = \begin{bmatrix} \mathbf{X}\mathbf{y} \\ \mathbf{1}^\top \mathbf{y} \end{bmatrix},$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{n,m}$ .

□

**Remark 3.1.2.5 (Normal equations from gradient)** We consider the function

$$J : \mathbb{R}^n \rightarrow \mathbb{R} , \quad J(\mathbf{y}) := \|\mathbf{A}\mathbf{y} - \mathbf{b}\|_2^2. \quad (3.1.2.6)$$

As above, using elementary identities for the Euclidean inner product on  $\mathbb{R}^m$ ,  $J$  can be recast as

$$\begin{aligned} J(\mathbf{y}) &= \mathbf{y}^\top \mathbf{A}^\top \mathbf{A} \mathbf{y} - 2\mathbf{b}^\top \mathbf{A} \mathbf{y} + \mathbf{b}^\top \mathbf{b} \\ &= \sum_{i=1}^n \sum_{j=1}^n (\mathbf{A}^\top \mathbf{A})_{ij} y_i y_j - 2 \sum_{i=1}^m \sum_{j=1}^n b_i (\mathbf{A})_{ij} y_j + \sum_{i=1}^m b_i^2. \end{aligned}$$

Obviously,  $J$  is a **multivariate polynomial** in  $y_1, \dots, y_n$ . As such,  $J$  is an infinitely differentiable function,  $J \in C^\infty(\mathbb{R}^n, \mathbb{R})$ , see [Str09, Bsp. 7.1.5].

According to [Str09, Satz 7.5.3]  $\mathbf{x} \in \mathbb{R}^n$  is a minimizer of  $J$  only if

$$\text{grad } J(\mathbf{x}) = 2\mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{A}^\top \mathbf{b} = \mathbf{0}. \quad (3.1.2.7)$$

This formula for the gradient of  $J$  can easily be confirmed by computing the partial derivatives  $\frac{\partial J}{\partial y_i}$  from the above explicit formula. Observe that (3.1.2.7) is equivalent to the normal equations (3.1.2.2). □

**§3.1.2.8 (The linear least squares problem (→ § 1.5.5.1))** Thm. 3.1.2.1 together with Thm. 3.1.1.7 already confirms that the normal equations will always have a solution and that  $\text{lsq}(\mathbf{A}, \mathbf{b})$  is a subspace of  $\mathbb{R}^n$  parallel to  $\mathcal{N}(\mathbf{A}^\top \mathbf{A})$ . The next theorem gives even more detailed information.

### Theorem 3.1.2.9. Kernel and range of $\mathbf{A}^\top \mathbf{A}$

For  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ , holds

$$\mathcal{N}(\mathbf{A}^\top \mathbf{A}) = \mathcal{N}(\mathbf{A}), \quad (3.1.2.10)$$

$$\mathcal{R}(\mathbf{A}^\top \mathbf{A}) = \mathcal{R}(\mathbf{A}^\top). \quad (3.1.2.11)$$

For the proof we need an basic result from linear algebra:

### Lemma 3.1.2.12. Kernel and range of (Hermitian) transposed matrices

For any matrix  $\mathbf{A} \in \mathbb{K}^{m,n}$  holds

$$\mathcal{N}(\mathbf{A}) = \mathcal{R}(\mathbf{A}^\top)^\perp , \quad \mathcal{N}(\mathbf{A})^\perp = \mathcal{R}(\mathbf{A}^\top). \quad (3.1.2.12)$$

Notation: Orthogonal complement of a subspace  $V \subset \mathbb{K}^k$ :

$$V^\perp := \{\mathbf{x} \in \mathbb{K}^k : \mathbf{x}^H \mathbf{y} = 0 \ \forall \mathbf{y} \in V\}.$$

*Proof. (of Thm. 3.1.2.9)*

①: We first show (3.1.2.10)

$$\begin{aligned} \mathbf{z} \in \mathcal{N}(\mathbf{A}^T \mathbf{A}) &\Leftrightarrow \mathbf{A}^T \mathbf{A} \mathbf{z} = \mathbf{0} \Rightarrow \mathbf{z}^T \mathbf{A}^T \mathbf{A} \mathbf{z} = \|\mathbf{A} \mathbf{z}\|_2^2 = 0 \Leftrightarrow \mathbf{A} \mathbf{z} = \mathbf{0}, \\ \mathbf{A} \mathbf{z} = \mathbf{0} &\Rightarrow \mathbf{A}^T \mathbf{A} \mathbf{z} = \mathbf{0} \Leftrightarrow \mathbf{z} \in \mathcal{N}(\mathbf{A}^T \mathbf{A}). \end{aligned}$$

②: The relationship (3.1.2.11) follows from (3.1.2.10) and Lemma 3.1.2.12:

$$\mathcal{R}(\mathbf{A}^T) \stackrel{\text{Lemma 3.1.2.12}}{=} \mathcal{N}(\mathbf{A})^\perp \stackrel{(3.1.2.10)}{=} \mathcal{N}(\mathbf{A}^T \mathbf{A})^\perp \stackrel{\text{Lemma 3.1.2.12}}{=} \mathcal{R}(\mathbf{A}^T \mathbf{A}).$$

□

### Corollary 3.1.2.13. Uniqueness of least squares solutions

If  $m \geq n$  and  $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$ , then the linear system of equations  $\mathbf{A} \mathbf{x} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , has a unique least squares solution ( $\rightarrow$  3.1.1.1)

$$\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}, \quad (3.1.2.14)$$

that can be obtained by solving the normal equations (3.1.2.2).

Note that  $\mathbf{A}^T \mathbf{A}$  is symmetric positive definite ( $\rightarrow$  Def. 1.1.2.6), if  $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$ .

**Remark 3.1.2.15 (Full-rank condition ( $\rightarrow$  Def. 2.2.1.3))** For a matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$  with  $m \geq n$  is equivalent

$$\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\} \iff \text{rank}(\mathbf{A}) = n. \quad (3.1.2.16)$$

Hence the assumption  $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$  of Cor. 3.1.2.13 is also called a full-rank condition (FRC), because the rank of  $\mathbf{A}$  is maximal.  $\square$

**EXAMPLE 3.1.2.17 (Meaning of full-rank condition for linear models)** We revisit the parameter estimation problem for a linear model.

- For Ex. 3.0.1.1,  $\mathbf{A} \in \mathbb{R}^{m,2}$  given in (3.0.1.3) it is easy to see

$$\text{rank} \left( \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ \vdots & \vdots \\ x_m & 1 \end{bmatrix} \right) = 2 \iff \exists i, j \in \{1, \dots, m\}: x_i \neq x_j,$$

that is, the manifest condition, that the points do not lie on a vertical line.

- In the case of Ex. 3.0.1.4 and the overdetermined  $m \times (n+1)$  linear system (3.0.1.5), we find

$$\text{rank} \begin{pmatrix} \mathbf{x}_1^\top & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^\top & 1 \end{pmatrix} = n+1 \Leftrightarrow \begin{array}{l} \text{There is a subset of } n+1 \\ \text{points } \mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{n+1}} \text{ such that} \\ \{\mathbf{0}, \mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{n+1}}\} \text{ spans a non-} \\ \text{degenerate } n+1\text{-simplex.} \end{array}$$

↓

**Remark 3.1.2.18 (Rank defect in linear least squares problems)** In case the system matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ , of an overdetermined linear system arising from a mathematical fails to have full rank, it hints at **inadequate modelling**:

In this case parameters are redundant, because different sets of parameters yield the same output quantities: the parameters are not “observable”. ↓

**Remark 3.1.2.19 (Hesse matrix of least squares functional)** For the least squares functional

$$J : \mathbb{R}^n \rightarrow \mathbb{R}, \quad J(\mathbf{y}) := \|\mathbf{Ay} - \mathbf{b}\|_2^2. \quad (3.1.2.6)$$

and its explicit form as polynomial in the vector components  $y_j$  we find the Hessian ( $\rightarrow$  Def. 8.4.1.11, [Str09, Satz 7.5.3]) of  $J$ :

$$\mathbf{H}J(\mathbf{y}) = 2\mathbf{A}^\top \mathbf{A}. \quad (3.1.2.20)$$

Thm. 3.1.2.9 implies that  $\mathbf{A}^\top \mathbf{A}$  is positive definite ( $\rightarrow$  Def. 1.1.2.6) if and only if  $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$ .

Therefore, by [Str09, Satz 7.5.3], under the full-rank condition  $J$  has a positive definite Hessian everywhere, and a minimum at every stationary point of its gradient, that is, at every solution of the normal equations.

↓

**Remark 3.1.2.21 (Convex least squares functional)** Another result from analysis tells us that real-valued  $C^1$ -functions on  $\mathbb{R}^n$  whose Hessian has positive eigenvalues uniformly bounded away from zero are **strictly convex**. Hence, if  $\mathbf{A}$  has full rank, the least squares functional  $J$  from (3.1.2.6) is a strictly convex function.

Visualization of a least squares functional  $J : \mathbb{R}^2 \rightarrow \mathbb{R}$  for  $n = 2$

Under the full-rank condition the graph of  $J$  is a paraboloid with  $J(\mathbf{y}) \rightarrow \infty$  for  $\|\mathbf{y}\| \rightarrow \infty$

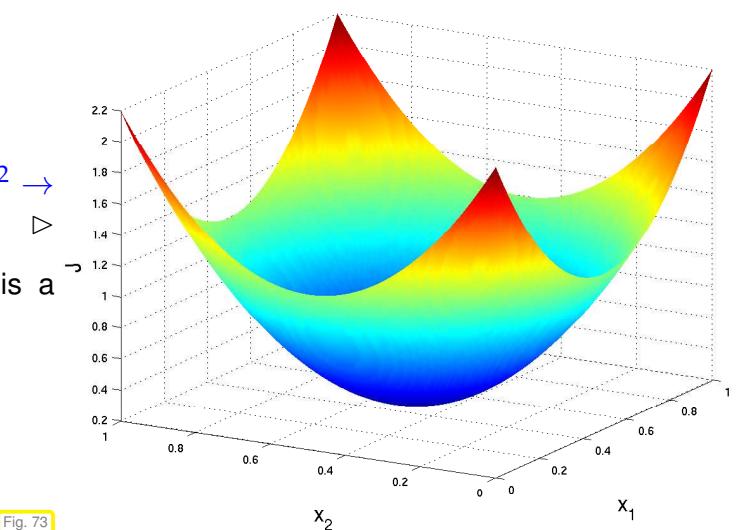


Fig. 73

Now we are in a position to state precisely what we mean by solving an overdetermined ( $m \geq n!$ ) linear system of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , provided that  $\mathbf{A}$  has full (maximal) rank, cf. (3.1.2.16).

(Full rank linear) least squares problem: [DR08, Sect. 4.2]

given:  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m, n \in \mathbb{N}$ ,  $m \geq n$ ,  $\text{rank}(\mathbf{A}) = n$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,

find:  $\mathbf{x} \in \mathbb{R}^n$  such that  $\|\mathbf{Ax} - \mathbf{b}\|_2 = \min\{\|\mathbf{Ay} - \mathbf{b}\|_2 : \mathbf{y} \in \mathbb{R}^n\}$   
 $\Updownarrow$   
 $\mathbf{x} = \underset{\mathbf{y} \in \mathbb{R}^n}{\operatorname{argmin}} \|\mathbf{Ay} - \mathbf{b}\|_2$

☞ A sloppy notation for the minimization problem (3.1.2.22) is  $\|\mathbf{Ax} - \mathbf{b}\|_2 \rightarrow \min$  ↗

### 3.1.3 Moore-Penrose Pseudoinverse

As we have seen in Ex. 3.0.1.9, there can be many least squares solutions of  $\mathbf{Ax} = \mathbf{b}$ , in case  $\mathcal{N}(\mathbf{A}) \neq \{\mathbf{0}\}$ . We can impose another condition to single out a unique element of  $\text{lsq}(\mathbf{A}, \mathbf{b})$ :

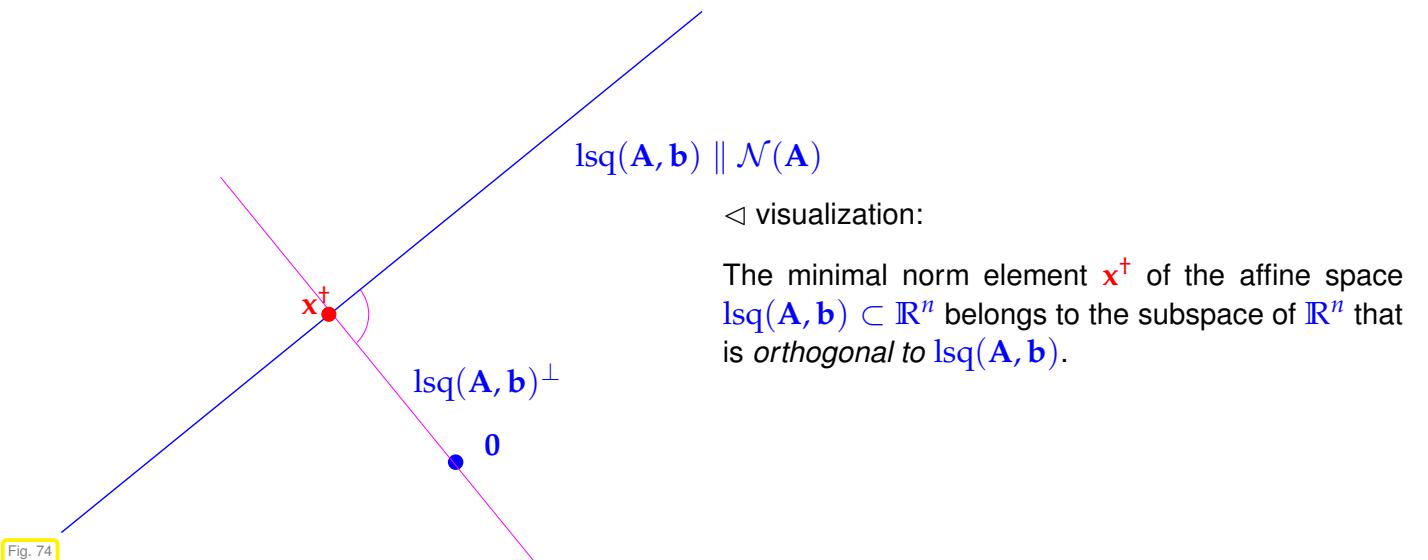
#### Definition 3.1.3.1. Generalized solution of a linear system of equations

The generalized solution  $\mathbf{x}^\dagger \in \mathbb{R}^n$  of a linear system of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , is defined as

$$\mathbf{x}^\dagger := \underset{\mathbf{x} \in \text{lsq}(\mathbf{A}, \mathbf{b})}{\operatorname{argmin}} \|\mathbf{x}\|_2. \quad (3.1.3.2)$$

► The generalized solution is the least squares solution with minimal norm.

**§3.1.3.3 (Reduced normal equations)** Elementary geometry teaches that the minimal norm element of an affine subspace  $L$  (a plane) in Euclidean space is the orthogonal projection of  $\mathbf{0}$  onto  $L$ .



Since the space of least squares solutions of  $\mathbf{Ax} = \mathbf{b}$  is an affine subspace parallel to  $\mathcal{N}(\mathbf{A})$

$$\text{lsq}(\mathbf{A}, \mathbf{b}) = \mathbf{x}^0 + \mathcal{N}(\mathbf{A}), \quad \mathbf{x}^0 \text{ solves normal equations,} \quad (3.1.3.4)$$

the generalized solution  $\mathbf{x}^*$  of  $\mathbf{Ax} = \mathbf{b}$  is contained in  $\mathcal{N}(\mathbf{A})^\perp$ . Therefore, given a basis  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\} \subset \mathbb{R}^n$  of  $\mathcal{N}(\mathbf{A})^\perp$ ,  $k := \dim \mathcal{N}(\mathbf{A})^\perp = n - \dim \mathcal{N}(\mathbf{A})$ , we can find  $\mathbf{y} \in \mathbb{R}^k$  such that  $\mathbf{x}^* = \mathbf{V}\mathbf{y}$ ,  $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_k] \in \mathbb{R}^{n,k}$ . Plugging this representation into the normal equations and multiplying with  $\mathbf{V}^\top$  yields the reduced normal equations

$$\begin{aligned}
 & \mathbf{V}^\top \mathbf{A}^\top \mathbf{A} \mathbf{V} \mathbf{y} = \mathbf{V}^\top \mathbf{A}^\top \mathbf{b} \quad (3.1.3.5) \\
 & \Updownarrow \\
 & \left[ \begin{array}{c|c} \mathbf{V}^\top & \\ \hline & \mathbf{A}^\top \end{array} \right] \left[ \begin{array}{c|c} & \mathbf{A} \\ \hline \mathbf{A} & \end{array} \right] \left[ \begin{array}{c|c} & \mathbf{V} \\ \hline \mathbf{V} & \end{array} \right] \left[ \begin{array}{c} \mathbf{y} \end{array} \right] = \\
 & \left[ \begin{array}{c|c} \mathbf{V}^\top & \\ \hline & \mathbf{A}^\top \end{array} \right] \left[ \begin{array}{c|c} & \mathbf{b} \end{array} \right].
 \end{aligned}$$

The very construction of  $\mathbf{V}$  ensures  $\mathcal{N}(\mathbf{AV}) = \{\mathbf{0}\}$  so that, by Thm. 3.1.2.9 the  $k \times k$  linear system of equations (3.1.3.5) has a unique solution. The next theorem summarizes our insights:

### Theorem 3.1.3.6. Formula for generalized solution

Given  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , the generalized solution  $\mathbf{x}^\dagger$  of the linear system of equations  $\mathbf{Ax} = \mathbf{b}$  is given by

$$\mathbf{x}^\dagger = \mathbf{V}(\mathbf{V}^\top \mathbf{A}^\top \mathbf{A} \mathbf{V})^{-1}(\mathbf{V}^\top \mathbf{A}^\top \mathbf{b}),$$

where  $\mathbf{V}$  is any matrix whose columns form a basis of  $\mathcal{N}(\mathbf{A})^\perp$ .

Terminology: The matrix  $\mathbf{A}^\dagger := \mathbf{V}(\mathbf{V}^\top \mathbf{A}^\top \mathbf{A} \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{A}^\top$  is called the **Moore-Penrose pseudoinverse** of  $\mathbf{A}$ .

☞ notation:  $\mathbf{A}^\dagger \in \mathbb{R}^{n,m} \hat{=} \text{pseudoinverse of } \mathbf{A} \in \mathbb{R}^{m,n}$

Note that the Moore-Penrose pseudoinverse does not depend on the choice of  $\mathbf{V}$ . □

Armed with the concept of generalized solution and the knowledge about its existence and uniqueness we can state the most general linear least squares problem:

(General linear) **least squares problem**:

given:  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m, n \in \mathbb{N}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,

find:  $\mathbf{x}^\dagger \in \mathbb{R}^n$  such that

$$(i) \quad \|\mathbf{Ax}^\dagger - \mathbf{b}\|_2 = \min\{\|\mathbf{Ay} - \mathbf{b}\|_2 : \mathbf{y} \in \mathbb{R}^n\}, \quad (3.1.3.7)$$

(ii)  $\|\mathbf{x}^\dagger\|_2$  is minimal under the condition (i).

### 3.1.4 Sensitivity of Least Squares Problem

Consider the full-rank linear least squares problem introduced in (3.1.2.22):

➤ data  $(\mathbf{A}, \mathbf{b}) \in \mathbb{R}^{m,n} \times \mathbb{R}^m$ , result  $\mathbf{x} = \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^n} \|\mathbf{Ay} - \mathbf{b}\|_2 \in \mathbb{R}^n$

On data space and result space we use the Euclidean norm (2-norm  $\|\cdot\|_2$ ) and the associated matrix norm, see § 1.5.5.3.

Recall Section 2.2.2, where we discussed the **sensitivity** of solutions of square linear systems, that is, the impact of perturbations in the problem data on the result. Now we study how (small) changes in  $\mathbf{A}$  and  $\mathbf{b}$  affect the unique ( $\rightarrow$  Cor. 3.1.2.13) least solution  $\mathbf{x}$  of  $\mathbf{Ax} = \mathbf{b}$  in the case of  $\mathbf{A}$  with full rank ( $\Leftrightarrow \mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$ )

Note: If the matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ , has full rank, then there is a  $c > 0$  such that  $\mathbf{A} + \Delta\mathbf{A}$  still has full rank for all  $\Delta\mathbf{A} \in \mathbb{R}^{m,n}$  with  $\|\Delta\mathbf{A}\|_2 < c$ . Hence, “sufficiently small” perturbations will not destroy the full-rank property of  $\mathbf{A}$ . This is a generalization of the Perturbation Lemma 2.2.2.5.

For square linear systems the condition number of the system matrix ( $\rightarrow$  Def. 2.2.2.7) provided the key gauge of sensitivity. To express the sensitivity of linear least squares problems we also generalize this concept:

#### Definition 3.1.4.1. Generalized condition number of a matrix

Given  $\mathbf{A} \in \mathbb{K}^{m,n}$ ,  $m \geq n$ ,  $\operatorname{rank}(\mathbf{A}) = n$ , we define its **generalized** (Euclidean) **condition number** as

$$\operatorname{cond}_2(\mathbf{A}) := \sqrt{\frac{\lambda_{\max}(\mathbf{A}^H \mathbf{A})}{\lambda_{\min}(\mathbf{A}^H \mathbf{A})}}.$$

☞ notation:  $\lambda_{\min}(\mathbf{A}) \hat{=} \text{smallest (in modulus) eigenvalue of matrix } \mathbf{A}$   
 $\lambda_{\max}(\mathbf{A}) \hat{=} \text{largest (in modulus) eigenvalue of matrix } \mathbf{A}$

For a square regular matrix this agrees with its condition number according to Def. 2.2.2.7, which follows from Cor. 1.5.5.16.

#### Theorem 3.1.4.2. Sensitivity of full-rank linear least squares problem

For  $m \geq n$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\operatorname{rank}(\mathbf{A}) = n$ , let  $\mathbf{x} \in \mathbb{R}^n$  be the solution of the least squares problem  $\|\mathbf{Ax} - \mathbf{b}\| \rightarrow \min$  and  $\hat{\mathbf{x}}$  the solution of the perturbed least squares problem  $\|(\mathbf{A} + \Delta\mathbf{A})\hat{\mathbf{x}} - \mathbf{b}\| \rightarrow \min$ . Then

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_2}{\|\mathbf{x}\|_2} \leq \left( 2 \operatorname{cond}_2(\mathbf{A}) + \operatorname{cond}_2^2(\mathbf{A}) \frac{\|\mathbf{r}\|_2}{\|\mathbf{A}\|_2 \|\mathbf{x}\|_2} \right) \frac{\|\Delta\mathbf{A}\|_2}{\|\mathbf{A}\|_2}$$

holds, where  $\mathbf{r} = \mathbf{Ax} - \mathbf{b}$  is the **residual**.

This means: if  $\|\mathbf{r}\|_2 \ll 1$  ➤ condition of the least squares problem  $\approx \operatorname{cond}_2(\mathbf{A})$   
if  $\|\mathbf{r}\|_2$  “large” ➤ condition of the least squares problem  $\approx \operatorname{cond}_2^2(\mathbf{A})$

For instance, in a linear parameter estimation problem ( $\rightarrow$  Ex. 3.0.1.4) a small residual will be the consequence of small measurement errors.

## 3.2 Normal Equation Methods [DR08, Sect. 4.2], [Han02, Ch. 11]

Given  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ ,  $\text{rank}(\mathbf{A}) = n$ ,  $\mathbf{b} \in \mathbb{R}^m$ , we introduce a first practical numerical method to determine the unique least squares solution ( $\rightarrow$  Def. 3.1.1.1) of the overdetermined linear system of equations  $\mathbf{Ax} = \mathbf{b}$ .

In fact, Cor. 3.1.2.13 suggests a simple algorithm for solving linear least squares problems of the form (3.1.2.22) satisfying the full (maximal) rank condition  $\text{rank}(\mathbf{A}) = n$ : it boils down to solving the **normal equations** (3.1.2.2):

Algorithm: Normal equation method to solve full-rank least squares problem  $\mathbf{Ax} = \mathbf{b}$

- ① Compute **regular** matrix  $\mathbf{C} := \mathbf{A}^\top \mathbf{A} \in \mathbb{R}^{n,n}$ .
- ② Compute right hand side vector  $\mathbf{c} := \mathbf{A}^\top \mathbf{b}$ .
- ③ Solve s.p.d. ( $\rightarrow$  Def. 1.1.2.6) linear system of equations:  $\mathbf{Cx} = \mathbf{c} \rightarrow$  § 2.8.0.13

This can be done in EIGEN in a single line of code:

### C++ code 3.2.0.1: Solving a linear least squares probel via normal equations

```

2 //!
3 //! Solving the overdetermined linear system of equations
4 //! Ax = b by solving normal equations (3.1.2.2)
5 //! The least squares solution is returned by value
6 VectorXd normeqsolve(const MatrixXd &A, const VectorXd &b) {
7     if (b.size() != A.rows()) throw runtime_error("Dimension mismatch");
8     // Cholesky solver
9     VectorXd x = (A.transpose()*A).llt().solve(A.transpose()*b);
10    return x;
}

```

By Thm. 2.8.0.11, for the s.p.d. matrix  $\mathbf{A}^\top \mathbf{A}$  Gaussian elimination remains stable even without pivoting. This is taken into account by requesting the Cholesky decomposition of  $\mathbf{A}^\top \mathbf{A}$  by calling the method `llt()`.

**§3.2.0.2 (Asymptotic complexity of normal equation method)** The problem size parameters for the linear least squares problem (3.1.2.22) are the matrix dimensions  $m, n \in \mathbb{N}$ , where  $n$  small & fixed,  $n \ll m$ , is common.

In Section 1.4.2 and Thm. 2.5.0.2 we discussed the asymptotic complexity of the operations involved in step ①-③ of the normal equation method:

$$\left. \begin{array}{ll} \text{step ①:} & \text{cost } O(mn^2) \\ \text{step ②:} & \text{cost } O(nm) \\ \text{step ③:} & \text{cost } O(n^3) \end{array} \right\} \Rightarrow \text{cost } O(n^2m + n^3) \quad \text{for } m, n \rightarrow \infty.$$

Note that for small fixed  $n$ ,  $n \ll m$ ,  $m \rightarrow \infty$  the computational effort scales linearly with  $m$ .  $\square$

### Remark 3.2.0.3 (Conditioning of normal equations [DR08, pp. 128])

The solution of least squares problems via the normal equation method is *vulnerable to instability*; immediate from Def. 3.1.4.1:



$$\text{cond}_2(\mathbf{A}^\top \mathbf{A}) = \text{cond}_2(\mathbf{A})^2.$$

Recall from Thm. 2.2.2.4:  $\text{cond}_2(\mathbf{A}^\top \mathbf{A})$  governs amplification of (roundoff) errors in  $\mathbf{A}^\top \mathbf{A}$  and  $\mathbf{A}^\top \mathbf{b}$  when solving normal equations (3.1.2.2).

- For fairly ill-conditioned  $\mathbf{A}$  using the normal equations (3.1.2.2) to solve the linear least squares problem from Def. 3.1.1.1 numerically may run the risk of huge amplification of roundoff errors incurred during the computation of the right hand side  $\mathbf{A}^H \mathbf{b}$ : **potential instability** ( $\rightarrow$  Def. 1.5.5.19) of normal equation approach.

↓

**EXAMPLE 3.2.0.4 (Roundoff effects in normal equations  $\rightarrow$  [DR08, Ex. 4.12])** In this example we witness loss of information in the computation of  $\mathbf{A}^H \mathbf{A}$ .



$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ \varepsilon & 0 \\ 0 & \varepsilon \end{bmatrix} \Rightarrow \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 + \varepsilon^2 & 1 & 1 \\ 1 & 1 + \varepsilon^2 & 1 \\ 1 & 1 & 1 + \varepsilon^2 \end{bmatrix}$$

Exp. 1.5.3.14: If  $\varepsilon \approx \sqrt{\text{EPS}}$ , then  $1 + \varepsilon^2 = 1$  in  $\mathbb{M}$  (set of machine numbers, see Def. 1.5.2.4). Hence the computed  $\mathbf{A}^T \mathbf{A}$  will fail to be regular, though  $\text{rank}(\mathbf{A}) = 2$ ,  $\text{cond}_2(\mathbf{A}) \approx \sqrt{\text{EPS}}$ .

#### C++-code 3.2.0.5:

```

2 int main() {
3     MatrixXd A(3,2);
4     // Inquire about machine precision → Ex. 1.5.3.12
5     double eps = std::numeric_limits<double>::epsilon();
6     // « initialization of matrix → § 1.2.1.3
7     A << 1, 1, sqrt(eps), 0, 0, sqrt(eps);
8     // Output rank of A^T A
9     std::cout << "Rank of A: " << A.fullPivLu().rank() << std::endl
10    << "Rank of A^TA: "
11    << (A.transpose() * A).fullPivLu().rank() << std::endl;
12
13    return 0;
}

```

Output:

---

```

1 Rank of A: 2
2 Rank of A^TA: 1

```

---

↓

**Remark 3.2.0.6 (Loss of sparsity when forming normal equations)** Another *reason not to compute  $\mathbf{A}^H \mathbf{A}$* , when both  $m, n$  large:

$\mathbf{A} \text{ sparse} \neq \mathbf{A}^T \mathbf{A} \text{ sparse}$

Example from Rem. 1.3.1.5: “Arrow matrices”

$$\begin{bmatrix} \text{red diagonal} \\ \text{red bottom row} \end{bmatrix} \begin{bmatrix} \text{red diagonal} \\ \text{red bottom row} \end{bmatrix} = \begin{bmatrix} \text{red matrix} \\ \text{red bottom row} \end{bmatrix}.$$

- Consequences for normal equation method, if both  $m, n$  large:
- ◆ Potential memory overflow, when computing  $\mathbf{A}^T \mathbf{A}$
  - ◆ Squanders possibility to use efficient sparse direct elimination techniques, see Section 2.7.4

This situation is faced in Ex. 3.0.1.9, Ex. 3.0.1.8. □

### §3.2.0.7 (Extended normal equations)

A way to avoid the computation of  $\mathbf{A}^H \mathbf{A}$ :

Extend normal equations (3.1.2.2): introduce **residual**  $\mathbf{r} := \mathbf{Ax} - \mathbf{b}$  as new unknown:

$$\mathbf{A}^H \mathbf{Ax} = \mathbf{A}^H \mathbf{b} \Leftrightarrow \mathbf{B} \begin{bmatrix} \mathbf{r} \\ \mathbf{x} \end{bmatrix} := \begin{bmatrix} -\mathbf{I} & \mathbf{A} \\ \mathbf{A}^H & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{r} \\ \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix}. \quad (3.2.0.8)$$

The benefit of using (3.2.0.8) instead of the standard normal equations (3.1.2.2) is that **sparsity is preserved**. However, the conditioning of the system matrix in (3.2.0.8) is not better than that of  $\mathbf{A}^T \mathbf{A}$ .

A more general substitution  $\mathbf{q} := \alpha^{-1}(\mathbf{Ax} - \mathbf{b})$  with  $\alpha > 0$  may improve the conditioning for suitably chosen parameter  $\alpha > 0$ :

$$\mathbf{A}^H \mathbf{Ax} = \mathbf{A}^H \mathbf{b} \Leftrightarrow \mathbf{B}_\alpha \begin{bmatrix} \mathbf{q} \\ \mathbf{x} \end{bmatrix} := \begin{bmatrix} -\alpha \mathbf{I} & \mathbf{A} \\ \mathbf{A}^H & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{q} \\ \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix}. \quad (3.2.0.9)$$

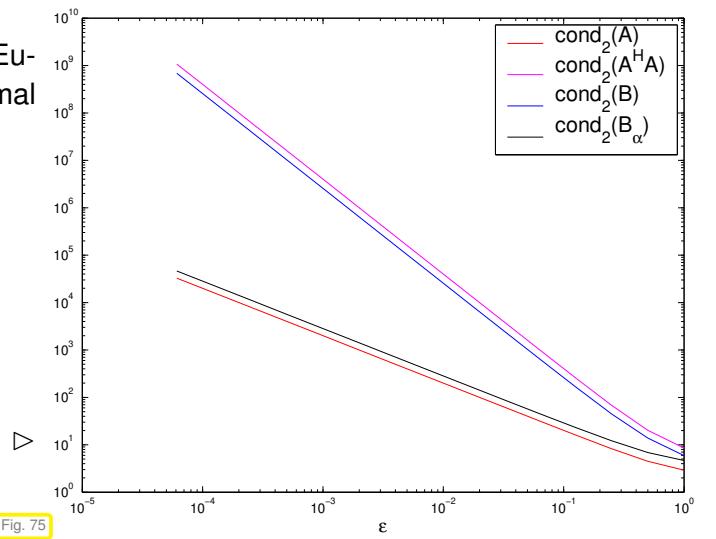
For  $m, n \gg 1$ ,  $\mathbf{A}$  sparse, both (3.2.0.8) and (3.2.0.9) lead to large sparse linear systems of equations, amenable to sparse direct elimination techniques, see Section 2.7.4. □

### EXAMPLE 3.2.0.10 (Conditioning of the extended normal equations)

In this example we explore empirically how the Euclidean condition number of the extended normal equations (3.2.0.9) is influenced by the choice of  $\alpha$ . Consider (3.2.0.8), (3.2.0.9) for

$$\mathbf{A} = \begin{bmatrix} 1+\epsilon & 1 \\ 1-\epsilon & 1 \\ \epsilon & \epsilon \end{bmatrix}.$$

Plot of different condition numbers  
in dependence on  $\epsilon$   
(Here  $\alpha = \epsilon \|\mathbf{A}\|_2 / \sqrt{2}$ )



### 3.3 Orthogonal Transformation Methods [DR08, Sect. 4.4.2]

#### 3.3.1 Transformation Idea

We consider the full-rank linear least squares problem (3.1.2.22)

$$\text{given } \mathbf{A} \in \mathbb{R}^{m,n}, \mathbf{b} \in \mathbb{R}^m \quad \text{find } \mathbf{x} = \underset{\mathbf{y} \in \mathbb{R}^n}{\operatorname{argmin}} \|\mathbf{Ay} - \mathbf{b}\|_2. \quad (3.1.2.22)$$

Setting:  $m \geq n$  and  $\mathbf{A}$  has full (maximum) rank:  $\operatorname{rank}(\mathbf{A}) = n$ .

**§3.3.1.1 (Generalizing the policy underlying Gaussian elimination)** Recall the rationale behind Gaussian elimination ( $\rightarrow$  Section 2.3, Ex. 2.3.1.1)

- By row transformations convert LSE  $\mathbf{Ax} = \mathbf{b}$  to *equivalent* (in terms of set of solutions) LSE  $\mathbf{Ux} = \tilde{\mathbf{b}}$ , which is easier to solve because it has triangular form.

How to adapt this policy to linear least squares problem (3.1.2.22) ?

Two questions: ① What linear least squares problems are “easy to solve” ?

② How can we arrive at them by *equivalent transformations* of (3.1.2.22) ?

Here we call two overdetermined linear systems  $\mathbf{Ax} = \mathbf{b}$  and  $\tilde{\mathbf{A}}\mathbf{x} = \tilde{\mathbf{b}}$  equivalent in the sense of (3.1.2.22), if both have the same set of least squares solutions:  $\operatorname{lsq}(\mathbf{A}, \mathbf{b}) = \operatorname{lsq}(\tilde{\mathbf{A}}, \tilde{\mathbf{b}})$ , see (3.1.1.2).

#### §3.3.1.2 (Triangular linear least squares problems)

The answer to question ① is the same as for LSE/Gaussian elimination:

Linear least squares problems (3.1.2.22) with upper *triangular*  $\mathbf{A}$  are easy to solve!

$$\begin{array}{c}
 \left[ \begin{array}{c|ccccc} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ \text{R} & & & & & \\ \hline & & & & & \\ \text{0} & & & & & \\ \hline & & & & & \end{array} \right] - \left[ \begin{array}{c} x_1 \\ \vdots \\ x_n \end{array} \right] = \left[ \begin{array}{c} b_1 \\ \vdots \\ b_m \end{array} \right]_2 \rightarrow \min \xrightarrow{(*)} \mathbf{x} = \left[ \begin{array}{c|ccccc} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ \text{R} & & & & & & \\ \hline & & & & & & \\ \text{0} & & & & & & \end{array} \right]^{-1} \left[ \begin{array}{c} b_1 \\ \vdots \\ b_n \end{array} \right]
 \end{array}$$

$\mathbf{x} \doteq$  least squares solution

How can we draw the conclusion  $(*)$ ? Obviously, the components  $n+1, \dots, m$  of the vector inside the norm are fixed and do not depend on  $\mathbf{x}$ . All we can do is to make the first components  $1, \dots, n$  vanish, by choosing a suitable  $\mathbf{x}$ , see [DR08, Thm. 4.13]. Obviously,  $\mathbf{x} = \mathbf{R}^{-1}(\mathbf{b})_{1:n}$  accomplishes this.

Note: since  $\mathbf{A}$  has full rank  $n$ , the upper triangular part  $\mathbf{R} \in \mathbb{R}^{n,n}$  of  $\mathbf{A}$  is regular!  $\square$

Answer to question ②:

Idea: If we have a (transformation) matrix  $\mathbf{T} \in \mathbb{R}^{m,m}$  satisfying



$$\|\mathbf{T}\mathbf{y}\|_2 = \|\mathbf{y}\|_2 \quad \forall \mathbf{y} \in \mathbb{R}^m, \quad (3.3.1.3)$$

$$\text{then } \underset{\mathbf{y} \in \mathbb{R}^n}{\operatorname{argmin}} \|\mathbf{A}\mathbf{y} - \mathbf{b}\|_2 = \underset{\mathbf{y} \in \mathbb{R}^n}{\operatorname{argmin}} \|\tilde{\mathbf{A}}\mathbf{y} - \tilde{\mathbf{b}}\|_2,$$

where  $\tilde{\mathbf{A}} = \mathbf{T}\mathbf{A}$  and  $\tilde{\mathbf{b}} = \mathbf{T}\mathbf{b}$ .

The next section will characterize the class of eligible transformation matrices  $\mathbf{T}$ .

### 3.3.2 Orthogonal/Unitary Matrices

**Definition 3.3.2.1. Unitary and orthogonal matrices** → [Gut09, Sect. 2.8]

- $\mathbf{Q} \in \mathbb{K}^{n,n}$ ,  $n \in \mathbb{N}$ , is **unitary**, if  $\mathbf{Q}^{-1} = \mathbf{Q}^H$ .
- $\mathbf{Q} \in \mathbb{R}^{n,n}$ ,  $n \in \mathbb{N}$ , is **orthogonal**, if  $\mathbf{Q}^{-1} = \mathbf{Q}^T$ .

**Theorem 3.3.2.2. Preservation of Euclidean norm**

A matrix is unitary/orthogonal, if and only if the associated linear mapping preserves the 2-norm:

$$\mathbf{Q} \in \mathbb{K}^{n,n} \text{ unitary} \Leftrightarrow \|\mathbf{Q}\mathbf{x}\|_2 = \|\mathbf{x}\|_2 \quad \forall \mathbf{x} \in \mathbb{K}^n.$$

From Thm. 3.3.2.2 we immediately conclude that, if a matrix  $\mathbf{Q} \in \mathbb{K}^{n,n}$  is unitary/orthogonal, then

- all rows/columns (regarded as vectors  $\in \mathbb{K}^n$ ) have Euclidean norm = 1,

- all rows/columns are pairwise orthogonal (w.r.t. Euclidean inner product),
- $|\det \mathbf{Q}| = 1$ ,  $\|\mathbf{Q}\|_2 = 1$ , and all eigenvalues  $\in \{z \in \mathbb{C} : |z| = 1\}$ .
- $\|\mathbf{QA}\|_2 = \|\mathbf{A}\|_2$  for any matrix  $\mathbf{A} \in \mathbb{K}^{n,m}$

### 3.3.3 QR-Decomposition [Han02, Sect. 13], [Gut09, Sect. 7.3]

This section will answer the question whether and how it is possible to find orthogonal transformations that convert any given matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ ,  $\text{rank}(\mathbf{A}) = n$ , to upper triangular form, as required for the application of the “equivalence transformation idea” to full-rank linear least squares problems.

#### 3.3.3.1 Theory

##### §3.3.3.1 (Gram-Schmidt orthogonalisation recalled) → § 1.5.1.1)

Input:  $\{\mathbf{a}^1, \dots, \mathbf{a}^n\} \subset \mathbb{K}^m$

Output:  $\{\mathbf{q}^1, \dots, \mathbf{q}^n\}$  (assuming no premature termination!)

```

1:  $\mathbf{q}^1 := \frac{\mathbf{a}^1}{\|\mathbf{a}^1\|_2}$ ; % 1st output vector
2: for  $j = 2, \dots, n$  do
   { % Orthogonal projection
3:    $\mathbf{q}^j := \mathbf{a}^j$ ;
4:   for  $\ell = 1, 2, \dots, j-1$  do      (GS)
5:     {  $\mathbf{q}^j \leftarrow \mathbf{q}^j - \langle \mathbf{a}^j, \mathbf{q}^\ell \rangle \mathbf{q}^\ell$ ; }
6:   if ( $\mathbf{q}^j = 0$ ) then STOP
7:   else {  $\mathbf{q}^j \leftarrow \frac{\mathbf{q}^j}{\|\mathbf{q}^j\|_2}$ ; }
8:   }

```

##### Theorem 3.3.3.2. Span property of G.S. vectors

If  $\{\mathbf{a}^1, \dots, \mathbf{a}^n\} \subset \mathbb{R}^m$  is linearly independent, then Algorithm (GS) computes orthonormal vectors  $\mathbf{q}^1, \dots, \mathbf{q}^k \in \mathbb{R}^m$  satisfying

$$\text{Span}\{\mathbf{q}^1, \dots, \mathbf{q}^k\} = \text{Span}\{\mathbf{a}^1, \dots, \mathbf{a}^k\}, \quad (1.5.1.2)$$

for all  $\ell \in \{1, \dots, n\}$ .

The span property (1.5.1.2) can be made more explicit in terms of the existence of linear combinations

$$\begin{aligned}
 \mathbf{q}_1 &= t_{11} \mathbf{a}_1 \\
 \mathbf{q}_2 &= t_{12} \mathbf{a}_1 + t_{22} \mathbf{a}_2 \\
 \mathbf{q}_3 &= t_{13} \mathbf{a}_1 + t_{23} \mathbf{a}_2 + t_{33} \mathbf{a}_3 \\
 &\vdots \\
 \mathbf{q}_k &= t_{1n} \mathbf{a}_1 + t_{2n} \mathbf{a}_2 + \cdots + t_{kk} \mathbf{a}_k .
 \end{aligned}
 \quad \Rightarrow \quad \exists \mathbf{T} \in \mathbb{R}^{n,n} \text{ upper triangular: } \mathbf{Q} = \mathbf{AT}, \quad (3.3.3.3)$$

where  $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_n] \in \mathbb{R}^{m,n}$  (with orthonormal columns),  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{R}^{m,n}$ . Note that thanks to the linear independence of  $\{\mathbf{a}^1, \dots, \mathbf{a}^k\}$  and  $\{\mathbf{q}^1, \dots, \mathbf{q}^k\}$ , the matrix  $\mathbf{T} = (t_{ij})_{i,j=1}^k \in \mathbb{R}^{k,k}$  is regular (“non-existent”  $t_{ij}$  are set to zero, of course).

Recall from Lemma 1.3.1.9 that inverses of regular upper triangular matrices are upper triangular

again.

Thus, by (3.3.3.3), we have found an *upper triangular*  $\mathbf{R} := \mathbf{T}^{-1} \in \mathbb{R}^{n,n}$  such that

$$A = QR \leftrightarrow \begin{bmatrix} A \\ \end{bmatrix} = \begin{bmatrix} Q \\ \end{bmatrix} \begin{bmatrix} R \\ \end{bmatrix}$$

Next “augmentation by zero”: add  $m - n$  zero rows at the bottom of  $\mathbf{R}$  and complement columns of  $\mathbf{Q}$  to an orthonormal basis of  $\mathbb{R}^m$ , which yields an orthogonal matrix  $\tilde{\mathbf{Q}} \in \mathbb{R}^{m,m}$ :

$$\blacktriangleright \quad \mathbf{A} = \tilde{\mathbf{Q}} \begin{bmatrix} \mathbf{R} \\ 0 \end{bmatrix}$$

↔

$$\left[ \begin{array}{c|c} \mathbf{A} & \\ \hline \end{array} \right] = \left[ \begin{array}{c|c} \tilde{\mathbf{Q}} & \\ \hline \end{array} \right] \left[ \begin{array}{c|c} \mathbf{R} & \\ \hline 0 & \\ \hline \end{array} \right]$$

$\Leftrightarrow \quad \tilde{\mathbf{Q}}^\top \mathbf{A} = \begin{bmatrix} \mathbf{R} \\ 0 \end{bmatrix}.$

Thus the algorithm of Gram-Schmidt orthonormalization “proves” the following theorem.

**Theorem 3.3.3.4. QR-decomposition → [NS02, Satz 5.2], [Gut07, Sect. 7.3]**

For any matrix  $\mathbf{A} \in \mathbb{K}^{n,k}$  with  $\text{rank}(\mathbf{A}) = k$  there exists

- (i) a unique Matrix  $\mathbf{Q}_0 \in \mathbb{R}^{n,k}$  that satisfies  $\mathbf{Q}_0^H \mathbf{Q}_0 = \mathbf{I}_k$ , and a unique *upper triangular* Matrix  $\mathbf{R}_0 \in \mathbb{K}^{k,k}$  with  $(\mathbf{R}_0)_{i,i} > 0$ ,  $i \in \{1, \dots, k\}$ , such that

$$\mathbf{A} = \mathbf{Q}_0 \cdot \mathbf{R}_0 \quad (\text{"economical" QR-decomposition}),$$

- (ii) a *unitary* Matrix  $\mathbf{Q} \in \mathbb{K}^{n,n}$  and a unique *upper triangular*  $\mathbf{R} \in \mathbb{K}^{n,k}$  with  $(\mathbf{R})_{i,i} > 0$ ,  $i \in \{1, \dots, n\}$ , such that

$$\mathbf{A} = \mathbf{Q} \cdot \mathbf{R} \quad (\text{full QR-decomposition}).$$

If  $\mathbb{K} = \mathbb{R}$  all matrices will be real and  $\mathbf{Q}$  is then *orthogonal*.

Visualisation: “economical” QR-decomposition:  $\mathbf{Q}_0^H \mathbf{Q}_0 = \mathbf{I}_k$  (orthonormal columns),

$$\mathbf{A} = \mathbf{Q}_0 \mathbf{R}_0, \quad \mathbf{Q}_0 \in \mathbb{K}^{n,k}, \quad \mathbf{R}_0 \in \mathbb{K}^{k,k} \quad \text{upper triangular},$$

$$\begin{bmatrix} \mathbf{A} \\ \vdots \end{bmatrix} = \begin{bmatrix} \mathbf{Q}_0 \\ \vdots \end{bmatrix} \begin{bmatrix} \mathbf{R}_0 \\ \vdots \end{bmatrix} \quad \begin{bmatrix} \mathbf{Q}_0 \\ \vdots \end{bmatrix} \text{ orthonormal columns} \quad \begin{bmatrix} \mathbf{R}_0 \\ \vdots \end{bmatrix} \text{ upper triangular}. \quad (3.3.3.5)$$

Visualisation: full QR-decomposition:  $\mathbf{Q}^H \mathbf{Q} = \mathbf{Q} \mathbf{Q}^H = \mathbf{I}_n$  (orthogonal matrix),

$$\mathbf{A} = \mathbf{Q} \mathbf{R}, \quad \mathbf{Q} \in \mathbb{K}^{n,n}, \quad \mathbf{R} \in \mathbb{K}^{n,k},$$

$$\begin{bmatrix} \mathbf{A} \\ \vdots \end{bmatrix} = \begin{bmatrix} \mathbf{Q} \\ \vdots \end{bmatrix} \begin{bmatrix} \mathbf{R} \\ \vdots \end{bmatrix} \quad \begin{bmatrix} \mathbf{Q} \\ \vdots \end{bmatrix} \text{ orthogonal matrix} \quad \begin{bmatrix} \mathbf{R} \\ \vdots \end{bmatrix} \text{ upper triangular}. \quad (3.3.3.6)$$

For square  $\mathbf{A}$ , that is,  $n = k$ , both QR-decompositions coincide.

**Corollary 3.3.3.7. Uniqueness of QR-factorization**

The “economical” QR-factorization (3.3.3.1) of  $\mathbf{A} \in \mathbb{K}^{m,n}$ ,  $m \geq n$ , with  $\text{rank}(\mathbf{A}) = n$  is unique, if we demand  $(\mathbf{R}_0)_{ii} > 0$ ,  $i = 1, \dots, n$ .

*Proof.* We observe that  $\mathbf{R}$  is regular, if  $\mathbf{A}$  has full rank  $n$ . Since the regular upper triangular matrices form a group under multiplication:

$$\mathbf{Q}_1 \mathbf{R}_1 = \mathbf{Q}_2 \mathbf{R}_2 \Rightarrow \mathbf{Q}_1 = \mathbf{Q}_2 \mathbf{R} \quad \text{with upper triangular } \mathbf{R} := \mathbf{R}_2 \mathbf{R}_1^{-1}.$$

►  $\mathbf{I} = \mathbf{Q}_1^H \mathbf{Q}_1 = \mathbf{R}^H \underbrace{\mathbf{Q}_2^H \mathbf{Q}_2}_{=\mathbf{I}} \mathbf{R} = \mathbf{R}^H \mathbf{R}.$

The assertion follows by uniqueness of Cholesky decomposition, Lemma 2.8.0.14.  $\square$

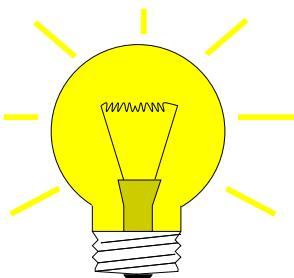
**3.3.3.2 Computation of QR-Decomposition**

In theory, Gram-Schmidt orthogonalization (GS) can be used to compute the QR-factorization of a matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ ,  $\text{rank}(\mathbf{A}) = n$ . However, as we saw in Exp. 1.5.1.5, Gram-Schmidt orthogonalization in the form of Code 1.5.1.3 is *not* a *stable* algorithm.

There is a stable way to compute QR-decompositions, based on the accumulation of orthogonal transformations.

**Corollary 3.3.3.8. Composition of orthogonal transformations**

*The product of two orthogonal/unitary matrices of the same size is again orthogonal/unitary.*

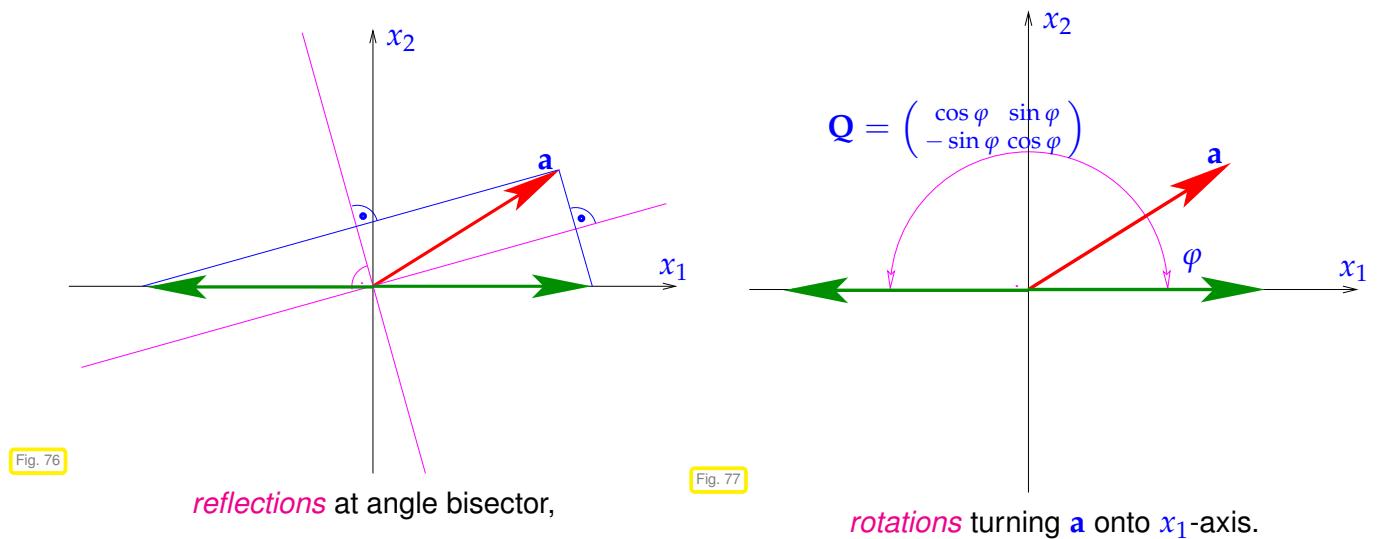


Idea: find simple unitary/orthogonal (row) transformations rendering certain matrix elements zero:

$$\mathbf{Q} \begin{pmatrix} \text{[yellow square]} \\ \text{[yellow square]} \end{pmatrix} = \begin{bmatrix} \text{[blue vertical bar]} \\ 0 \end{bmatrix} \quad \text{with } \mathbf{Q}^H = \mathbf{Q}^{-1}.$$

Recall that this “annihilation of column entries” is the key operation in Gaussian forward elimination, where it is achieved by means of non-unitary row transformations, see Sect. 2.3.2. Now we want to find a counterpart of Gaussian elimination based on unitary row transformations on behalf of numerical stability.

**EXAMPLE 3.3.3.9 (“Annihilating” orthogonal transformations in 2D)** In 2D there are two possible orthogonal transformations make 2nd component of  $\mathbf{a} \in \mathbb{R}^2$  vanish, which, in geometric terms, amounts to mapping the vector onto the  $x_1$ -axis.



Note that in each case we have *two different* length-preserving linear mappings at our disposal. This flexibility will be important for curbing the impact of roundoff. □

Both reflections and rotations are actually used in library routines and both are discussed in the sequel:

**§3.3.3.10 (Householder reflections)** → [GV13, Sect. 5.1.2]) The following so-called **Householder matrices** effect the reflection of a vector into a multiple of the first unit vector with the same length:

$$\mathbf{Q} = \mathbf{H}(\mathbf{v}) := \mathbf{I} - 2 \frac{\mathbf{v}\mathbf{v}^H}{\mathbf{v}^H\mathbf{v}} \quad \text{with} \quad \mathbf{v} = \frac{1}{2}(\mathbf{a} \pm \|\mathbf{a}\|_2 \mathbf{e}_1), \quad (3.3.3.11)$$

where  $\mathbf{e}_1$  is the first Cartesian basis vector. Orthogonality of these matrices can be established by direct computation.

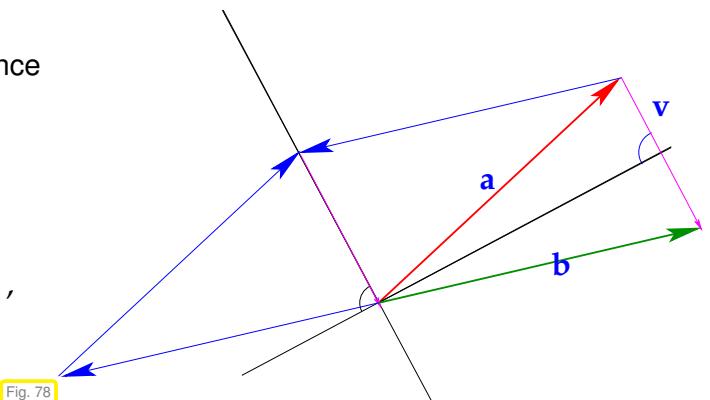
Fig. 78 sketches a “geometric derivation” of Householder reflections:

Given  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$  with  $\|\mathbf{a}\| = \|\mathbf{b}\|$ , the difference vector  $\mathbf{v} = \mathbf{b} - \mathbf{a}$  is orthogonal to the bisector.

$$\begin{aligned} \mathbf{b} &= \mathbf{a} - (\mathbf{a} - \mathbf{b}) = \mathbf{a} - \mathbf{v} \frac{\mathbf{v}^T \mathbf{v}}{\mathbf{v}^T \mathbf{v}} \\ &= \mathbf{a} - 2\mathbf{v} \frac{\mathbf{v}^T \mathbf{a}}{\mathbf{v}^T \mathbf{v}} = \mathbf{a} - 2 \frac{\mathbf{v} \mathbf{v}^T}{\mathbf{v}^T \mathbf{v}} \mathbf{a} = \mathbf{H}(\mathbf{v}) \mathbf{a}, \end{aligned}$$

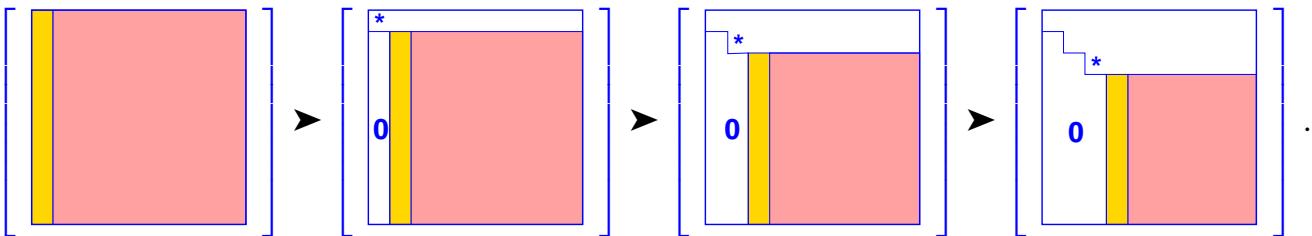
because, due to orthogonality  $(\mathbf{a} - \mathbf{b}) \perp (\mathbf{a} + \mathbf{b})$

$$(\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b}) = (\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b} + \mathbf{a} + \mathbf{b}) = 2(\mathbf{a} - \mathbf{b})^T \mathbf{a}.$$



Suitable successive Householder transformations determined by the leftmost column (“target column”) of shrinking bottom right matrix blocks can be used to achieve upper triangular form  $\mathbf{R}$ . Visualization of

annihilation of lower triangular matrix part for a square matrix:



**Yellow** = "target column **a**" (determines unitary transformation),  
**Red** = modified in course of transformations.

Writing  $\mathbf{Q}_\ell$  for the Householder matrix used in the  $\ell$ -th step we get

**QR-factorization** (QR-decomposition) of  $\mathbf{A} \in \mathbb{C}^{n,n}$ :  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ ,  $\mathbf{Q} := \mathbf{Q}_1^H \cdots \mathbf{Q}_{n-1}^H$  orthogonal matrix ,  
 $\mathbf{R}$  upper triangular matrix .

**Remark 3.3.3.12 (QR-decomposition of “fat” matrices)** We can also apply successive Householder transformation as outlined in § 3.3.3.10 to a matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$  with  $m < n$ . If the first  $m$  columns of  $\mathbf{A}$  are linearly independent, we obtain another variant of the QR-decomposition:

$$A = QR \quad , \quad Q \in \mathbb{R}^{m,m}, \quad R \in \mathbb{R}^{m,n},$$

where  $\mathbf{Q}$  is orthogonal,  $\mathbf{R}$  upper triangular, that is,  $(\mathbf{R})_{i,j} = 0$  for  $i > j$ .

**Remark 3.3.3.13 (Stable implementation of Householder reflections)** In (3.3.3.11) the computation of the vector  $\mathbf{v}$  can be prone to cancellation ( $\rightarrow$  Section 1.5.4), if the vector  $\mathbf{a}$  encloses a very small angle with the first unit vector, because in this case  $\mathbf{v}$  can be very small and beset with a huge relative error. This is a concern, because in the formula for the Householder matrix  $\mathbf{v}$  is normalized to unit length (division by  $\|\mathbf{v}\|_2^2$ ).

Fortunately, two choices for  $v$  are possible in (3.3.3.11) and at most one can be affected by cancellation. The right choice is

$$\mathbf{v} = \begin{cases} \frac{1}{2}(\mathbf{a} + \|\mathbf{a}\|_2 \mathbf{e}_1) & , \text{if } a_1 > 0 \\ \frac{1}{2}(\mathbf{a} - \|\mathbf{a}\|_2 \mathbf{e}_1) & , \text{if } a_1 \leq 0 \end{cases}$$

See [Hig02, Sect. 19.1] and [GV13, Sect. 5.1.3] for a discussion.

**§3.3.3.14 (Givens rotations** → [Han02, Sect. 14], [GOV13]) The 2D rotation displayed in Fig. 77 can be embedded in an identity matrix. Thus, the following orthogonal transformation, a **Givens rotation**, annihilates the  $k$ -th component of a vector  $\mathbf{a} = [a_1, \dots, a_n]^\top \in \mathbb{R}^n$ . Here  $\gamma$  stands for  $\cos(\varphi)$  and  $\sigma$  for

$\sin(\varphi)$ ,  $\varphi$  the angle of rotation, see Fig. 77.

$$\mathbf{G}_{1k}(a_1, a_k) \mathbf{a} := \begin{bmatrix} \gamma & \cdots & \sigma & \cdots & 0 \\ \vdots & \ddots & \vdots & \cdots & \vdots \\ -\sigma & \cdots & \gamma & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \\ a_{nf} \end{bmatrix} = \begin{bmatrix} a_1^{(1)} \\ \vdots \\ 0 \\ \vdots \\ a_n \end{bmatrix}, \text{ if } \begin{array}{l} \gamma = \frac{a_1}{\sqrt{|a_1|^2 + |a_k|^2}}, \\ \sigma = \frac{a_k}{\sqrt{|a_1|^2 + |a_k|^2}}. \end{array} \quad (3.3.3.15)$$

Orthogonality ( $\rightarrow$  Def. 6.2.1.2) of  $\mathbf{G}_{1k}(a_1, a_k)$  is verified immediately. Again, we have two options for an annihilating rotation, see Ex. 3.3.3.9. It will always be possible to choose one that avoids cancellation [GV13, Sect. 5.1.8], see Code 3.3.3.16 for details.

#### C++11 code 3.3.3.16: Stable Givens rotation of a 2D vector, [GV13, Alg. 5.1.3]

```

2 // plane (2D) Givens rotation avoiding cancellation
3 // Computes orthogonal  $\mathbf{G} \in \mathbb{R}^{2,2}$  with  $\mathbf{G}^\top \mathbf{a} = \begin{bmatrix} r \\ 0 \end{bmatrix} =: \mathbf{x}$ ,  $r = \pm \|\mathbf{a}\|_2$ 
4 void planerot(const Vector2d& a, Matrix2d& G, Vector2d& x) {
5     int sign{1};
6     if (a[1] != 0.0) {
7         double t, s, c; //  $s \leftrightarrow \sigma$ ,  $c \leftrightarrow \gamma$ 
8         if (std::abs(a[1]) > std::abs(a[0])) { // Avoid cancellation/overflow
9             t = -a[0]/a[1]; s = 1.0/std::sqrt(1.0+t*t); c = s*t; sign = -1;
10        }
11        else {
12            t = -a[1]/a[0]; c = 1.0/std::sqrt(1.0+t*t); s = c*t;
13        }
14        G << c,s,-s,c; // Form  $2 \times 2$  Givens rotation matrix
15    }
16    else G.setIdentity();
17    x << (sign*a.norm()),0.0;
18}

```

We validate the implementation by straightforward computations using the variable names of Code 3.3.3.16 and  $\mathbf{a} = \begin{bmatrix} a_0 \\ a_1 \end{bmatrix}$ :

- Case  $a_1 \geq a_0$ :  $t = -\frac{a_0}{a_1}$ ,  $s = \frac{1}{\sqrt{1+t^2}}$ ,  $c = st$

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix}^\top \begin{bmatrix} a_0 \\ a_1 \end{bmatrix} = \begin{bmatrix} ca_0 - sa_1 \\ sa_0 + ca_1 \end{bmatrix} = \begin{bmatrix} sta_0 - sa_1 \\ sa_0 = sta_1 \end{bmatrix} = \frac{1}{\sqrt{1+t^2}} \begin{bmatrix} ta_0 - a_1 \\ a_0 + ta_1 \end{bmatrix} = \frac{a_1}{\|\mathbf{a}\|_2} \begin{bmatrix} -\frac{a_0^2}{a_1} - a_1 \\ a_0 - a_0 \end{bmatrix} = \begin{bmatrix} -\|\mathbf{a}\|_2 \\ 0 \end{bmatrix}.$$

- Case  $a_0 > a_1$ :  $t = -\frac{a_1}{a_0}$ ,  $c = \frac{1}{1+t^2}$ ,  $s = ct$

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix}^\top \begin{bmatrix} a_0 \\ a_1 \end{bmatrix} = \begin{bmatrix} ca_0 - cta_1 \\ cta_0 + ca_1 \end{bmatrix} = \frac{1}{\sqrt{1+t^2}} \begin{bmatrix} a_0 - ta_1 \\ ta_0 + a_1 \end{bmatrix} = \frac{a_0}{\|\mathbf{a}\|_2} \begin{bmatrix} a_0 + \frac{a_1^2}{a_0} \\ -a_1 + a_1 \end{bmatrix} = \begin{bmatrix} +\|\mathbf{a}\|_2 \\ 0 \end{bmatrix}.$$

So far, we know how to annihilate a single component of a vector by means of a Givens rotation that targets that component and some other (the first in (3.3.3.15)). However, for the sake of QR-decomposition we aim to map *all* components to zero except for the first.

- ☞ This can be achieved by  $n - 1$  successive Givens rotations, see also Code 3.3.3.18

$$\begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} \xrightarrow{\mathbf{G}_{12}(a_1, a_2)} \begin{bmatrix} a_1^{(1)} \\ 0 \\ a_3 \\ \vdots \\ a_n \end{bmatrix} \xrightarrow{\mathbf{G}_{13}(a_1^{(1)}, a_3)} \begin{bmatrix} a_1^{(2)} \\ 0 \\ 0 \\ a_4 \\ \vdots \\ a_n \end{bmatrix} \xrightarrow{\mathbf{G}_{14}(a_1^{(2)}, a_4)} \dots \xrightarrow{\mathbf{G}_{1n}(a_1^{(n-2)}, a_n)} \begin{bmatrix} a_1^{(n-1)} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (3.3.3.17)$$

- ☞ Notation:  $\mathbf{G}_{ij}(a_1, a_2) \doteq$  Givens rotation (3.3.3.15) aimed at rows  $i$  and  $j$  of the matrix.

#### C++11 code 3.3.3.18: Rotating a vector onto the $x_1$ -axis by successive Givens transformation

```

2 // Orthogonal transformation of a (column) vector into a multiple of
3 // the first unit vector by successive Givens transformations
4 void givenscoltrf(const VectorXd& aln, MatrixXd& Q, VectorXd& aOut){
5     unsigned int n = aln.size();
6     // Assemble rotations in a dense matrix.
7     // For (more efficient) alternatives see Rem. Rem. 3.3.3.20
8     Q.setIdentity();
9     Matrix2d G; Vector2d tmp, xDummy;
10    aOut = aln;
11    for(int j = 1; j < n; ++j) {
12        tmp(0) = aOut(0); tmp(1) = aOut(j);
13        planerot(tmp, G, xDummy); // see Code 3.3.3.16
14        // select 1st and jth element of aOut and use the Map function
15        // to prevent copying; equivalent to aOut([1,j]) in MATLAB
16        Map<VectorXd>, 0, InnerStride<> > aOutMap(aOut.data(), 2, InnerStride<>(j));
17        aOutMap = G * aOutMap;
18        // select 1st and jth column of Q (Q(:,[1,j]) in MATLAB)
19        Map<MatrixXd>, 0, OuterStride<> > QMap(Q.data(), n, 2, OuterStride<>(j*n));
20        QMap = QMap * G.transpose();
21    }
22}

```

Armed with these compound Givens rotations we can proceed as in the case of Householder reflections to accomplish the orthogonal transformation of a full-rank matrix to upper triangular form, see

#### C++11 code 3.3.3.19: QR-decomposition by successive Givens rotations

```

2 //! QR decomposition of square matrix A by successive Givens
3 // transformations
4 void qrgivens(const MatrixXd& A, MatrixXd& Q, MatrixXd& R){
5     unsigned int n = A.rows();
6     // Assemble rotations in a dense matrix.
7     // For (more efficient) alternatives see Rem. Rem. 3.3.3.20
8     Q.setIdentity();
9     Matrix2d G; Vector2d tmp, xDummy;
10    R = A; // In situ transformation
11    for(int i = 0; i < n-1; ++i) {
12        for(int j = n-1; j > i; --j) {
13            tmp(0) = R(j-1, i); tmp(1) = R(j, i);
14            planerot(tmp, G, xDummy); // see Code 3.3.3.16
15            R.block(j-1,0,2,n) = G * R.block(j-1,0,2,n);
16            Q.block(0, j-1, n, 2) = Q.block(0,j-1,n,2) * G.transpose();
17        }
18    }
19}

```

**Remark 3.3.3.20 (Storing orthogonal transformations)** When doing successive orthogonal transformations as in the case of QR-decomposition by means of Householder reflections ( $\rightarrow$  § 3.3.3.10) or Givens rotations ( $\rightarrow$  § 3.3.3.14) it would be prohibitively expensive to assemble and even multiply the transformation matrices!

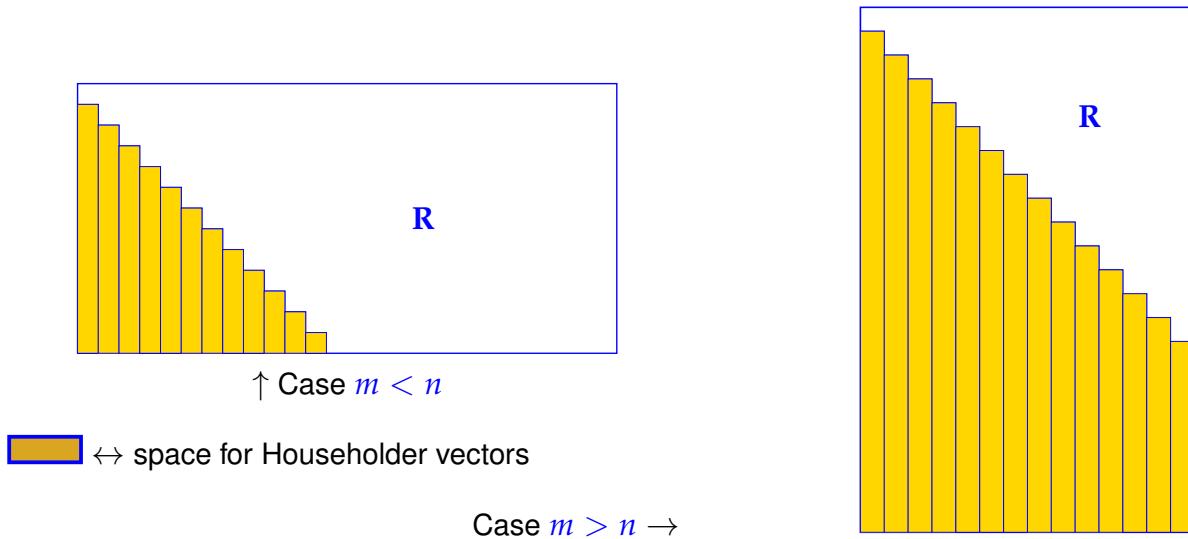
The matrices for the orthogonal transformation are never built in codes!

The transformations are stored in a compressed format.

- ① In the case of **Householder reflections**  $H(\mathbf{v}) \in \mathbb{R}^{m,m}$  (3.3.3.11), see [GOV13],

- store only the last  $n - 1$  components of the normalized vector  $\mathbf{v} \in \mathbb{R}^m$

For QR-decomposition of a matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ , by means of successive Householder reflections  $\mathbf{H}(\mathbf{v}_1) \cdot \dots \cdot \mathbf{H}(\mathbf{v}_k)$ ,  $k := \min\{m, n\}$ , we store the bottom parts of the vectors  $\mathbf{v}_j \in \mathbb{R}^{m-j+1}$ ,  $j = 1, \dots, k$ , whose lengths decrease, in place of the “annihilated” lower triangular part of  $\mathbf{A}$ :



- ② In the case of **Givens rotations**, for a single rotation  $G_{i,j}(a_1, a_2)$ , see [GOV13],

- store row indices  $(i, j)$  and rotation angle [GV13, Sect. 5.1.11],

$$\text{for } \mathbf{G} = \begin{bmatrix} \gamma & \sigma \\ -\sigma & \gamma \end{bmatrix} \Rightarrow \text{store } \rho := \begin{cases} \text{sign}(\sigma) & , \text{ if } \gamma = 0 , \\ \frac{1}{2} \text{sign}(\gamma)\sigma & , \text{ if } |\sigma| < |\gamma| , \\ 2 \text{sign}(\sigma)/\gamma & , \text{ if } |\sigma| \geq |\gamma| , \end{cases}$$

$$\text{which means } \begin{cases} \rho = \pm 1 & \Rightarrow \gamma = 0, \sigma = \pm 1 \\ |\rho| < 1 & \Rightarrow \sigma = 2\rho, \gamma = \sqrt{1 - \sigma^2} \\ |\rho| > 1 & \Rightarrow \gamma = 2/\rho, \sigma = \sqrt{1 - \gamma^2}. \end{cases}$$

Then store  $\mathbf{G}_{ij}(a, b)$  as triple  $(i, j, \rho)$ . The parameter  $\rho$  forgets the sign of the matrix  $\mathbf{G}_{ij}$ , so the signs of the corresponding rows in the transformed matrix  $\mathbf{R}$  have to be changed accordingly. The rationale behind the above convention is to curb the impact of roundoff errors.

**Remark 3.3.3.21 (QR-decomposition of banded matrices)** The advantage of Givens rotations is its **selectivity**, which can be exploited for banded matrices, see Section 2.7.5, Def. 2.7.5.1.

Example: Orthogonal transformation of an  $n \times n$  tridiagonal matrix to upper triangular form, that is, the annihilation of the sub-diagonal, by means of successive Givens rotations:

$$\begin{array}{c} \left[ \begin{array}{ccccccc} * & * & 0 & 0 & 0 & 0 & 0 \\ * & * & 0 & 0 & 0 & 0 & 0 \\ 0 & * & * & 0 & 0 & 0 & 0 \\ 0 & 0 & * & * & 0 & 0 & 0 \\ 0 & 0 & 0 & * & * & 0 & 0 \\ 0 & 0 & 0 & 0 & * & * & 0 \\ 0 & 0 & 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & * \end{array} \right] \xrightarrow{\mathbf{G}_{12}} \left[ \begin{array}{ccccccc} * & * & * & 0 & 0 & 0 & 0 \\ 0 & * & * & 0 & 0 & 0 & 0 \\ 0 & * & * & 0 & 0 & 0 & 0 \\ 0 & 0 & * & * & * & 0 & 0 \\ 0 & 0 & 0 & * & * & * & 0 \\ 0 & 0 & 0 & 0 & * & * & 0 \\ 0 & 0 & 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & * \end{array} \right] \xrightarrow{\mathbf{G}_{23} \dots \mathbf{G}_{n-1,n}} \left[ \begin{array}{ccccccc} * & * & * & 0 & 0 & 0 & 0 \\ 0 & * & * & * & 0 & 0 & 0 \\ 0 & 0 & * & * & * & 0 & 0 \\ 0 & 0 & 0 & * & * & * & 0 \\ 0 & 0 & 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & * \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right] \end{array}$$

\*  $\hat{=}$  entry set to zero by Givens rotation, \*  $\hat{=}$  new non-zero entry (“fill-in”  $\rightarrow$  Def. 2.7.4.3).

This is a manifestation of a more general result, see Def. 2.7.5.1 for notations:

### Theorem 3.3.3.22. QR-decomposition “preserves bandwidth”

If  $\mathbf{A} = \mathbf{QR}$  is the QR-decomposition of a regular matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$ , then  $\text{bw}(\mathbf{R}) \leq \text{bw}(\mathbf{A})$ .

Studying the algorithms sketched above for tridiagonal matrices, we find that a total of at most  $n \cdot \text{bw}(\mathbf{A})$  Givens rotations is required for computing the QR-decomposition. Each of them acts on  $\mathcal{O}(\text{bw}(\mathbf{A}))$  non-zero entries of the matrix, which leads to an asymptotic total computational effort of  $\mathcal{O}(n \cdot \text{bw}(\mathbf{A})^2)$  for  $n, \text{bw}(\mathbf{A}) \rightarrow \infty$ .  $\square$

### 3.3.3.3 QR-Decomposition: Stability

In numerical linear algebra orthogonal transformation methods usually give rise to reliable algorithms, thanks to the norm-preserving property of orthogonal transformations.

**§3.3.3.23 (Stability of unitary/orthogonal transformations)** We consider the mapping (the “transformation” induced by  $\mathbf{Q}$ )

$$F : \mathbb{K}^n \rightarrow \mathbb{K}^n, \quad F(\mathbf{x}) := \mathbf{Q}\mathbf{x}, \quad \mathbf{Q} \in \mathbb{K}^{n,n} \text{ unitary/orthogonal}.$$

We are interested in the sensitivity of  $F$ , that is, the impact of relative errors in the data vector  $\mathbf{x}$  on the output vector  $\mathbf{y} := F(\mathbf{x})$ .

We study the output for a perturbed input vector:

$$\left. \begin{array}{l} \mathbf{Q}\mathbf{x} = \mathbf{y} \Rightarrow \|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 \\ \mathbf{Q}(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{y} + \Delta\mathbf{y} \Rightarrow \mathbf{Q}\Delta\mathbf{x} = \Delta\mathbf{y} \Rightarrow \|\Delta\mathbf{y}\|_2 = \|\Delta\mathbf{x}\|_2 \end{array} \right\} \Rightarrow \frac{\|\Delta\mathbf{y}\|_2}{\|\mathbf{y}\|_2} = \frac{\|\Delta\mathbf{x}\|_2}{\|\mathbf{x}\|_2}.$$

We conclude, that unitary/orthogonal transformations do *not* involve any *amplification of relative errors* in the data vectors.

Of course, this also applies to the “solution” of square linear systems with orthogonal coefficient matrix  $\mathbf{Q} \in \mathbb{R}^{n,n}$ , which, by Def. 6.2.1.2, boils down to multiplication of the right hand side vector with  $\mathbf{Q}^H$ .  $\square$

**Remark 3.3.3.24 (Conditioning of conventional row transformations)** Gaussian elimination as presented in § 2.3.1.3 converts a matrix to upper triangular form by elementary row transformations. Those

add a scalar multiple of a row of the matrix to another matrix and amount to left-multiplication with matrices

$$\mathbf{T} := \mathbf{I}_n + \mu \mathbf{e}_j \mathbf{e}_i^\top, \quad \mu \in \mathbb{K}, \quad i, j \in \{1, \dots, n\}, \quad i \neq j. \quad (3.3.3.25)$$

However, these transformations can lead to a massive amplification of relative errors, which, by virtue of Ex. 2.2.2.1 can be linked to large condition numbers of  $\mathbf{T}$ .

This accounts for fact that the computation of LU-decompositions by means of Gaussian elimination might not be stable, see Ex. 2.4.0.5.  $\square$

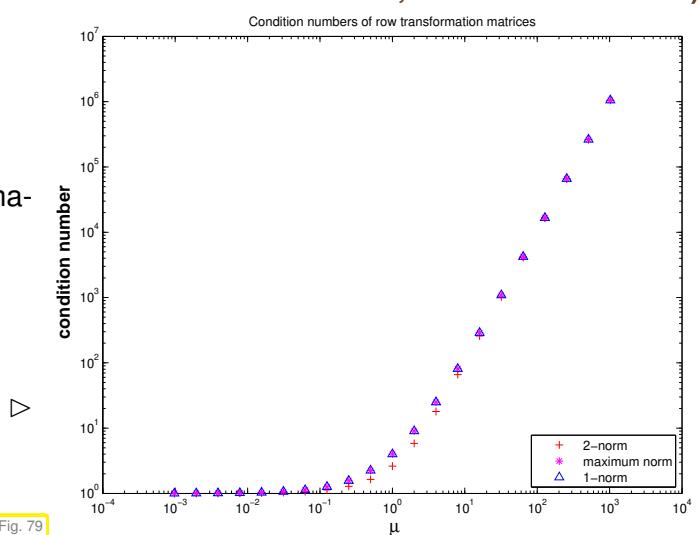
### EXPERIMENT 3.3.3.26 (Conditioning of conventional row transformations, Rem. 3.3.3.24 cnt'd)

Study in 2D:

$2 \times 2$  row transformation matrix, (cf. elimination matrices of Gaussian elimination.

$$\mathbf{T}(\mu) = \begin{pmatrix} 1 & 0 \\ \mu & 1 \end{pmatrix}$$

Euclidean condition numbers of  $\mathbf{T}(\mu)$



$\square$

The perfect conditioning of orthogonal transformation prevents the destructive build-up of roundoff errors.

#### Theorem 3.3.3.27. Stability of Householder QR [Hig02, Thm. 19.4]

Let  $\tilde{\mathbf{R}} \in \mathbb{R}^{m,n}$  be the  $R$ -factor of the QR-decomposition of  $\mathbf{A} \in \mathbb{R}^{m,n}$  computed by means of successive Householder reflections ( $\rightarrow$  § 3.3.3.10). Then there exists an orthogonal  $\mathbf{Q} \in \mathbb{R}^{m,m}$  such that

$$\mathbf{A} + \Delta \mathbf{A} = \mathbf{Q} \tilde{\mathbf{R}} \quad \text{with} \quad \|\Delta \mathbf{A}\|_2 \leq \frac{cmn \text{EPS}}{1 - cmn \text{EPS}} \|\mathbf{A}\|_2, \quad (3.3.3.28)$$

where EPS is the machine precision and  $c > 0$  a small constant independent of  $\mathbf{A}$ .

### 3.3.3.4 QR-Decomposition in EIGEN

EIGEN offers several classes dedicated to computing QR-type decompositions of matrices, for instance **HouseholderQR**. Internally the QR-decomposition is stored in compressed format as explained in Rem. 3.3.3.20. Its computation is triggered by the constructor.

#### C++-code 3.3.3.29: QR-decompositions in EIGEN

```

2 # include <Eigen/QR>
3
4 // Computation of full QR-decomposition (3.3.3.1),
5 // dense matrices built for both QR-factors (expensive!)
6 std::pair<MatrixXd, MatrixXd> qr_decomp_full(const MatrixXd& A) {

```

```

7  Eigen::HouseholderQR<MatrixXd> qr(A);
8  MatrixXd Q = qr.householderQ(); // 
9  MatrixXd R = qr.matrixQR().template triangularView<Eigen::Upper>();
10 return std::pair<MatrixXd, MatrixXd>(Q,R);
11 }

12 // Computation of economical QR-decomposition (3.3.3.1),
13 // dense matrix built for Q-factor (possibly expensive!)
14 std::pair<MatrixXd, MatrixXd> qr_decomp_eco(const MatrixXd& A) {
15     using index_t = MatrixXd::Index;
16     const index_t m = A.rows(), n = A.cols();
17     Eigen::HouseholderQR<MatrixXd> qr(A);
18     MatrixXd Q = (qr.householderQ() * MatrixXd::Identity(m,n)); //
19     MatrixXd R = qr.matrixQR().block(0,0,n,n).template triangularView<Eigen::Upper>();
20     //
21     return std::pair<MatrixXd, MatrixXd>(Q,R);
22 }

```

Note that the method `householderQ` returns the Q-factor in compressed format → Rem. 3.3.3.20. Assignment to a matrix will convert it into a (dense) matrix format, see Line 8; only then the actual computation of the matrix entries is performed. It can also be multiplied with another matrix of suitable size, which is used in Line 19 to extract the Q-factor  $Q_0 \in \mathbb{R}^{m,n}$  of the economical QR-decomposition (3.3.3.1).

The matrix returned by the method `matrixQR()` gives access to a matrix storing the QR-factors in compressed form. Its upper triangular part provides  $\mathbf{R}$ , see Line 20.

A close inspection of the algorithm for the computation of QR-decompositions of  $\mathbf{A} \in \mathbb{R}^{m,n}$  by successive Householder reflections (→ § 3.3.3.10) reveals, that  $n$  transformations costing  $\sim mn$  operations each are required.

### Asymptotic complexity of Householder QR-decomposition

The computational effort for `HouseholderQR()` of  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m > n$ , is  $O(mn^2)$  for  $m, n \rightarrow \infty$ .

**EXPERIMENT 3.3.3.31 (Asymptotic complexity of Householder QR-factorization)** We empirically investigate the (asymptotic) complexity of QR-factorization algorithms in EIGEN through runtime measurements.

### C++-code 3.3.3.32: timing QR-factorizations in EIGEN

```

2  int nruns = 3, minExp = 2, maxExp = 6;
3  MatrixXd tms(maxExp-minExp+1,4);
4  for(int i = 0; i <= maxExp-minExp; ++i) {
5      Timer t1, t2, t3; // timer class
6      int n = std::pow(2, minExp + i); int m = n*n;
7      // Initialization of matrix A
8      MatrixXd A(m,n); A.setZero();
9      A.setIdentity(); A.block(n,0,m-n,n).setOnes();
10     A += VectorXd::LinSpaced(m,1,m) * RowVectorXd::LinSpaced(n,1,n);
11     for(int j = 0; j < nruns; ++j) {
12         // plain QR-factorization in the constructor
13         t1.start(); HouseholderQR<MatrixXd> qr(A); t1.stop();
14         // full decomposition
15         t2.start(); std::pair<MatrixXd, MatrixXd> QR2 = qr_decomp_full(A); t2.stop();
16         // economic decomposition

```

```

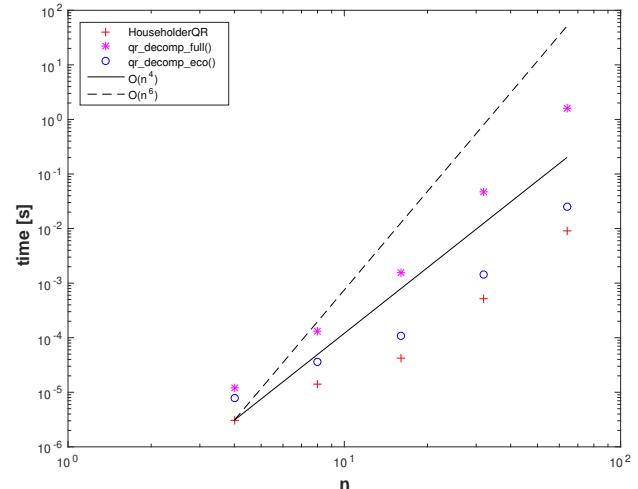
17     t3.start(); std::pair<MatrixXd, MatrixXd> QR3 = qr_decomp_eco(A); t3.stop();
18 }
19 tms(i,0)=n;
20 tms(i,1)=t1.min(); tms(i,2)=t2.min(); tms(i,3)= t3.min();
21 }
```

Timings for

- plain QR-factorization in the constructor of **HouseholderQR**,
- invocation of function `qr_decomp_full()`, see Code 3.4.2.1,
- call to `qr_decomp_eco()` from Code 3.4.2.1.

Platform:

- ◆ ubuntu 14.04 LTS
- ◆ CPU i7-3517U, 1.90GHZ, 4 cores
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB, Mem 8 GB
- ◆ gcc 4.8.4, -O3



The runtimes for the QR-factorization of  $\mathbf{A} \in \mathbb{R}^{n^2, n}$  behave like  $\mathcal{O}(n^2 \cdot n)$  for large  $n$ . □

### 3.3.4 QR-Based Solver for Linear Least Squares Problems

The QR-decomposition introduced in Section 3.3.3, Thm. 3.3.3.4, paves the way for the practical algorithmic realization of the “equivalent orthonormal transformation to upper triangular form”-idea from Section 3.3.1.

We consider the full-rank linear least squares problem Eq. (3.1.2.22): Given  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ ,  $\text{rank}(\mathbf{A}) = n$ ,

$$\text{seek } \mathbf{x} \in \mathbb{R}^n \text{ such that } \|\mathbf{Ax} - \mathbf{b}\|_2 \rightarrow \min .$$

We assume that we are given a

**QR-decomposition**:  $\mathbf{A} = \mathbf{QR}$ ,  $\mathbf{Q} \in \mathbb{R}^{m,m}$  orthogonal,  $\mathbf{R} \in \mathbb{R}^{m,n}$  (regular) upper triangular matrix.

We apply the orthogonal 2-norm preserving ( $\rightarrow$  Thm. 3.3.2.2) transformation encoded in  $\mathbf{Q}$  to  $\mathbf{Ax} - \mathbf{b}$ , the vector inside the 2-norm to be minimized:

$$\|\mathbf{Ax} - \mathbf{b}\|_2 = \left\| \mathbf{Q}(\mathbf{Rx} - \mathbf{Q}^H \mathbf{b}) \right\|_2 = \left\| \mathbf{Rx} - \tilde{\mathbf{b}} \right\|_2 , \quad \tilde{\mathbf{b}} := \mathbf{Q}^H \mathbf{b} .$$

Thus, we have obtained an *equivalent* triangular linear least squares problem:



$$\|Ax - b\|_2 \rightarrow \min \Leftrightarrow \left[ \begin{array}{c|c} \text{yellow blocks} & R_0 \\ \text{white blocks} & 0 \end{array} \right] \left[ \begin{array}{c} x_1 \\ \vdots \\ x_n \end{array} \right] = \left[ \begin{array}{c} \tilde{b}_1 \\ \vdots \\ \tilde{b}_m \end{array} \right]_2 \rightarrow \min .$$

↓

$$x = \left[ \begin{array}{c|c} \text{yellow blocks} & R_0 \\ \text{white blocks} & 0 \end{array} \right]^{-1} \left[ \begin{array}{c} \tilde{b}_1 \\ \vdots \\ \tilde{b}_n \end{array} \right], \quad \text{with residual } r = Q \left[ \begin{array}{c} 0 \\ \vdots \\ 0 \\ \tilde{b}_{n+1} \\ \vdots \\ \tilde{b}_m \end{array} \right].$$

Note: by Thm. 3.3.2.2 the norm of the residual is readily available:  $\|\mathbf{r}\|_2 = \sqrt{\tilde{b}_{n+1}^2 + \dots + \tilde{b}_m^2}$ .

### C++-code 3.3.4.1: QR-based solver for full rank linear least squares problem (3.1.2.22)

```

2 // Solution of linear least squares problem (3.1.2.22) by means of
3 // QR-decomposition
4 // Note:  $\mathbf{A} \in \mathbb{R}^{m,n}$  with  $m > n$ ,  $\text{rank}(\mathbf{A}) = n$  is assumed
5 // Least squares solution returned in x, residual norm as return value
6 double qrLsqSolve(const MatrixXd& A, const VectorXd& b,
7                     VectorXd& x) {
8     const unsigned m = A.rows(), n = A.cols();
9
10    MatrixXd Ab(m, n + 1); Ab << A, b; // Form extended matrix [A, b]
11
12    // QR-decomposition of extended matrix automatically transforms b
13    MatrixXd R = Ab.householderQr().matrixQR().template
14        triangularView<Eigen::Upper>(); //
15
16    MatrixXd R_nn = R.block(0, 0, n, n); // R-factor R_0
17    // Compute least squares solution x = (R) $^{-1}_{1:n, 1:n}$ (Q $^T$ b) $_{1:n}$ 
18    x = R_nn.template triangularView<Eigen::Upper>().solve(R.block(0, n, n, 1));
19    return R(n, n); // residual norm = ||A $\hat{x}$  - b|| $_2$  (why ?)
}

```

Discussion of (some) details of implementation in Code 3.3.4.1:

- The QR-decomposition is computed in a numerically stable way by means of Householder reflections ( $\rightarrow$  § 3.3.3.10) by EIGEN's built-in function `householderQR` available for matrix type. The computational cost of this function when called for an  $m \times n$  matrix is, asymptotically for  $m, n \rightarrow \infty$ ,  $O(n^2m)$ .

- Line 9: We perform the QR-decomposition of the **extended matrix**  $[A, b]$  with  $b$  as rightmost column. Thus, the orthogonal transformations are automatically applied to  $b$ ; the augmented matrix is converted into  $[R, Q^T b]$ , the data of the equivalent upper triangular linear least squares problem. Thus, actually, no information about  $Q$  needs to be stored, if one is interested in the least squares solution  $x$  only.

The idea is borrowed from Gaussian elimination, see Code 2.3.1.4, Line 9.

- Line 13: `MatrixQR()` returns the compressed QR-factorization as a matrix, where the R-factor  $R \in \mathbb{R}^{m,n}$  is contained in the upper triangular part, whose top  $n$  rows give  $R_0$  from see (3.3.3.1).
- Line 18: the components  $(b)_{n+2:m}$  of the vector  $b$  (treated as rightmost column of the augmented matrix) are annihilated when computing the QR-decomposition (by final Householder reflection):  $(Q^H [A, b])_{n+2:m, n} = 0$ . Hence,  $(Q^H [A, b])_{n+1, n+1} = \|(\tilde{b})_{n+1:m}\|_2$ , which gives the norm of the residual.
- A QR-based algorithm is implemented in the `solve()` method available for EIGEN's QR-decomposition, see Code 3.3.4.2.

#### C++ code 3.3.4.2: EIGEN's built-in QR-based linear least squares solver

```

2 // Solving a full-rank least squares problem ||Ax - b||_2 → min in EIGEN
3 double lsqslve_eigen(const MatrixXd& A, const VectorXd& b,
4                      VectorXd& x) {
5     x = A.householderQr().solve(b);
6     return ((A*x-b).norm());
7 }
```

**Remark 3.3.4.3 (QR-based solution of linear systems of equations)** Applying the QR-based algorithm for full-rank linear least squares problems in the case  $m = n$ , that is, to a square linear system of equations  $Ax = b$  with a regular coefficient matrix  $A$ , will compute the solution  $x = A^{-1}b$ . In a sense, the QR-decomposition offers an alternative to Gaussian elimination/LU-decomposition discussed in § 2.3.2.15.

The steps for solving a linear system of equations  $Ax = b$  by means of QR-decomposition are as follows:

- ① QR-decomposition  $A = QR$ , computational costs  $\frac{2}{3}n^3 + O(n^2)$   
(about twice as expensive as LU-decomposition without pivoting)
- $Ax = b$  :      ② orthogonal transformation  $z = Q^H b$ , computational costs  $4n^2 + O(n)$   
(in the case of *compact storage* of reflections/rotations)
- ③ **Backward substitution**, solve  $Rx = z$ , computational costs  $\frac{1}{2}n(n + 1)$

Benefit: we can utterly dispense with any kind of pivoting:

- ⌚ Computing the generalized QR-decomposition  $A = QR$  by means of Householder reflections or Givens rotations is (numerically stable) for any  $A \in \mathbb{C}^{m,n}$ .
- ⌚ For any regular system matrix an LSE can be solved by means of  
QR-decomposition + orthogonal transformation + backward substitution  
in a stable manner.

Drawback: QR-decomposition can hardly ever avoid **massive fill-in** (→ Def. 2.7.4.3) also in situations, where LU-factorization greatly benefits from Thm. 2.7.5.4. ↓

**Remark 3.3.4.4 (QR-based solution of banded LSE)** From Rem. 3.3.3.21, Thm. 3.3.3.22, we know that that particular situation, in which QR-decomposition can **avoid fill-in** (→ Def. 2.7.4.3) is the case of banded

matrices, see Def. 2.7.5.1. For a banded  $n \times n$  linear systems of equations with small fixed bandwidth  $\text{bw}(\mathbf{A}) \leq O(1)$  we incur an

➤ asymptotic computational effort:  $O(n)$  for  $n \rightarrow \infty$

The following code uses QR-decomposition computed by means of selective Givens rotations ( $\rightarrow$  § 3.3.3.14) to solve a tridiagonal linear system of equations  $\mathbf{Ax} = \mathbf{b}$

$$\mathbf{A} = \begin{bmatrix} d_1 & c_1 & 0 & \cdots & 0 \\ e_1 & d_2 & c_2 & & \vdots \\ 0 & e_2 & d_3 & c_3 & \\ \vdots & & \ddots & \ddots & c_{n-1} \\ 0 & \cdots & 0 & e_{n-1} & d_n \end{bmatrix}$$

The matrix is passed in the form of three vectors  $\mathbf{e}, \mathbf{c}, \mathbf{d}$  giving the entries in the non-zero bands.

#### C++ code 3.3.4.5: Solving a tridiagonal system by means of QR-decomposition → [GITLAB](#)

```

2  ///! @brief Solves the tridiagonal system  $\mathbf{Ax} = \mathbf{b}$  with QR-decomposition
3  ///! @param[in] d Vector of dim n; the diagonal elements
4  ///! @param[in] c Vector of dim n-1; the lower diagonal elements
5  ///! @param[in] e Vector of dim n-1; the upper diagonal elements
6  ///! @param[in] b Vector of dim n; the rhs.
7  ///! @param[out] x Vector of dim n
8  VectorXd tridiagqr(VectorXd c, VectorXd d, VectorXd e, VectorXd& b){
9      int n = d.size();
10     // resize the vectors c and d to correct length if needed
11     c.conservativeResize(n); e.conservativeResize(n);
12     double t = d.norm() + e.norm() + c.norm();
13     Matrix2d R; Vector2d z, tmp;
14     for(int k = 0; k < n-1; ++k){
15         tmp(0) = d(k); tmp(1) = e(k);
16         // Use givensrotation to set the entries below the diagonal
17         // to zero
18         planerot(tmp, R, z); // see Code 3.3.3.16
19         if( std::abs(z(0))/t < std::numeric_limits<double>::epsilon() )
20             throw std::runtime_error("A nearly singular");
21         // Update all other entries of the matrix and rhs. which
22         // were affected by the givensrotation
23         d(k) = z(0);
24         b.segment(k,2).applyOnTheLeft(R); // rhs.
25         // Block of the matrix affected by the givensrotation
26         Matrix2d Z;
27         Z << c(k), 0, d(k+1), c(k+1);
28         Z.applyOnTheLeft(R);
29         // Write the transformed block back to the corresponding places
30         c.segment(k,2) = Z.diagonal(); d(k+1) = Z(1,0); e(k) = Z(0,1);
31     }
32     // Note that the e is now above d and c
33     // Backsubstitution acting on upper triangular matrix
34     // with upper bandwidth 2 (stored in vectors).
35     VectorXd x(n);
36     // last row
37     x(n-1) = b(n-1)/d(n-1);
38     if(n >= 2) {
39         // 2nd last row
40         x(n-2) = (b(n-2)-c(n-2)*x(n-1))/d(n-2);
41         // remaining rows
42         for(int i = n-3; i >= 0; --i)
43             x(i) = ( b(i) - c(i) * x(i+1) - e(i)*x(i+2) ) / d(i);
44     }

```

```

45     return x;
46 }
```

**EXAMPLE 3.3.4.6 (Stable solution of LSE by means of QR-decomposition)** Aiming to confirm the claim of superior stability of QR-based approaches ( $\rightarrow$  Rem. 3.3.4.3, § 3.3.3.23) we revisit Wilkinson's counterexample from Ex. 2.4.0.5 for which Gaussian elimination with partial pivoting does not yield an acceptable solution.

Wilkinson matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$

$$(\mathbf{A})_{i,j} := \begin{cases} -1 & \text{for } i > j, j < n, \\ 1 & \text{for } i = j, \\ 0 & \text{for } i < j, j < n, \\ 1 & \text{for } j = n. \end{cases}$$

QR-decomposition produces perfect solution

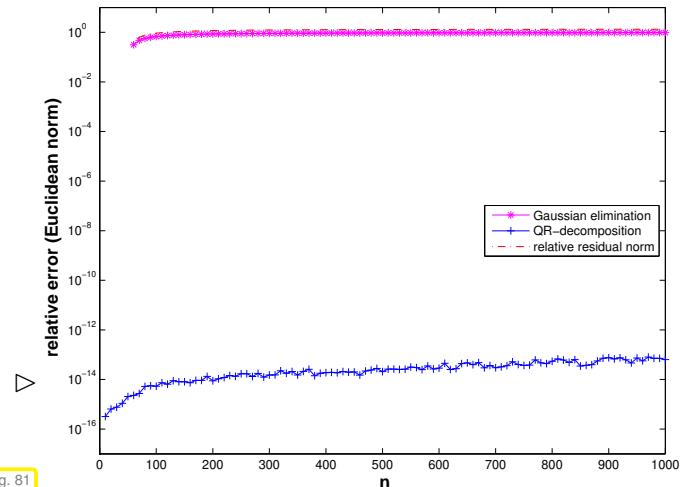


Fig. 8.1

Let us summarize the pros and cons of orthogonal transformation techniques for linear least squares problems:

#### Normal equations vs. orthogonal transformations method

Superior numerical stability ( $\rightarrow$  Def. 1.5.5.19) of orthogonal transformations methods:

- Use orthogonal transformations methods for least squares problems (3.1.3.7), whenever  $\mathbf{A} \in \mathbb{R}^{m,n}$  dense and  $n$  small.

SVD/QR-factorization cannot exploit sparsity:

- Use normal equations in the expanded form (3.2.0.8)/(3.2.0.9), when  $\mathbf{A} \in \mathbb{R}^{m,n}$  sparse ( $\rightarrow$  Notion 2.7.0.1) and  $m, n$  big.

### 3.3.5 Modification Techniques for QR-Decomposition

In § 2.6.0.12 we faced the task of solving a square linear system of equations  $\tilde{\mathbf{A}}\mathbf{x} = \mathbf{b}$  efficiently, whose coefficient matrix  $\tilde{\mathbf{A}}$  was a (rank-1) perturbation of  $\mathbf{A}$ , for which an LU-decomposition was available. Lemma 2.6.0.21 showed a way to reuse the information contained in the LU-decomposition.

A similar task can be posed, if the QR-decomposition of a matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ , has already been computed and then we have to solve a full-rank linear least squares problem  $\|\tilde{\mathbf{A}}\mathbf{x} - \mathbf{b}\|_2 \rightarrow \min$  with  $\tilde{\mathbf{A}} \in \mathbb{R}^{m,n}$  a “slight” perturbation of  $\mathbf{A}$ . If we aim to use orthogonalization techniques it would be desirable to compute the QR-decomposition of  $\tilde{\mathbf{A}}$  with recourse to the QR-decomposition of  $\mathbf{A}$ .

### 3.3.5.1 Rank-1 Modifications

For  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ ,  $\text{rank}(\mathbf{A}) = n$ , we consider the rank-1 modification, cf. Eq. (2.6.0.16),

$$\mathbf{A} \longrightarrow \tilde{\mathbf{A}} := \mathbf{A} + \mathbf{u}\mathbf{v}^\top, \quad \mathbf{u} \in \mathbb{R}^m, \quad \mathbf{v} \in \mathbb{R}^n. \quad (3.3.5.1)$$

We assume that still  $\text{rank}(\tilde{\mathbf{A}}) = n$ .

Given the a full QR-decomposition  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ ,  $\mathbf{Q} \in \mathbb{R}^{m,m}$  orthogonal (stored in some implicit format, see Rem. 3.3.3.20),  $\mathbf{R} \in \mathbb{R}^{m,n}$  upper triangular, according to Thm. 3.3.3.4, the goal is to find an efficient algorithm that yields the (unique) full QR-decomposition of  $\tilde{\mathbf{A}}$ :  $\tilde{\mathbf{A}} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$ .

Step ①: compute  $\mathbf{w} = \mathbf{Q}^\top \mathbf{u}$ , and observe  $\mathbf{A} + \mathbf{u}\mathbf{v}^\top = \mathbf{Q}(\mathbf{R} + \mathbf{w}\mathbf{v}^\top)$ .

➤ Computational effort  $\mathcal{O}(mn)$ , if  $\mathbf{Q}$  stored in suitable (compressed) format.

Step ②: achieve  $\mathbf{w} \rightarrow \|\mathbf{w}\|_2 \mathbf{e}_1$  through applying  $m-1$  Givens rotations, visualization for  $m=n$ :

$$\mathbf{w} = \begin{bmatrix} * \\ * \\ \vdots \\ * \\ * \\ * \end{bmatrix} \xrightarrow{\mathbf{G}_{n-1,n}} \begin{bmatrix} * \\ * \\ \vdots \\ * \\ * \\ 0 \end{bmatrix} \xrightarrow{\mathbf{G}_{n-2,n-1}} \begin{bmatrix} * \\ * \\ \vdots \\ * \\ 0 \\ 0 \end{bmatrix} \xrightarrow{\mathbf{G}_{n-3,n-2}} \dots \xrightarrow{\mathbf{G}_{12}} \begin{bmatrix} * \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Note that this will affect  $\mathbf{R}$  by creating a single non-zero subdiagonal:

$$\begin{aligned} \mathbf{R} &= \begin{bmatrix} * & * & \cdots & * & * & * & * \\ 0 & * & \cdots & * & * & * & * \\ \vdots & \ddots & & & & & \vdots \\ 0 & \cdots & 0 & * & * & * & * \\ 0 & \cdots & 0 & 0 & * & * & * \\ 0 & \cdots & 0 & 0 & 0 & * & * \\ 0 & \cdots & 0 & 0 & 0 & 0 & * \end{bmatrix} \xrightarrow{\mathbf{G}_{n-1,n}} \begin{bmatrix} * & * & \cdots & * & * & * & * \\ 0 & * & \cdots & * & * & * & * \\ \vdots & \ddots & & & & & \vdots \\ 0 & \cdots & 0 & * & * & * & * \\ 0 & \cdots & 0 & 0 & * & * & * \\ 0 & \cdots & 0 & 0 & 0 & * & * \\ 0 & \cdots & 0 & 0 & 0 & 0 & * \end{bmatrix} \xrightarrow{\mathbf{G}_{n-2,n-1}} \\ &\longrightarrow \begin{bmatrix} * & * & \cdots & * & * & * & * \\ 0 & * & \cdots & * & * & * & * \\ \vdots & \ddots & & & & & \vdots \\ 0 & \cdots & 0 & * & * & * & * \\ 0 & \cdots & 0 & 0 & * & * & * \\ 0 & \cdots & 0 & 0 & 0 & * & * \\ 0 & \cdots & 0 & 0 & 0 & 0 & * \end{bmatrix} \xrightarrow{\mathbf{G}_{n-3,n-2}} \dots \xrightarrow{\mathbf{G}_{1,2}} \begin{bmatrix} * & * & \cdots & * & * & * & * \\ * & * & \cdots & * & * & * & * \\ \vdots & \ddots & & & & & \vdots \\ 0 & \cdots & * & * & * & * & * \\ 0 & \cdots & 0 & * & * & * & * \\ 0 & \cdots & 0 & 0 & * & * & * \\ 0 & \cdots & 0 & 0 & 0 & * & * \end{bmatrix} =: \mathbf{R}_1 \end{aligned}$$

We see  $(\mathbf{R}_1)_{i,j} = 0$ , if  $i > j+1$ , it a so-called **upper Hessenberg matrix**. If  $\mathbf{Q}_1$  collects all orthogonal transformations used in Step ②, then

$$\mathbf{A} + \mathbf{u}\mathbf{v}^\top = \mathbf{Q}\mathbf{Q}_1^\top (\underbrace{\mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^\top}_{\text{upper Hessenberg matrix}}) \quad \text{with orthogonal } \mathbf{Q}_1 := \mathbf{G}_{12} \cdots \mathbf{G}_{n-1,n}.$$

➤ Computational effort  $\mathcal{O}(m+n^2)$

Step ③: Convert  $\mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^\top \mapsto$  upper triangular form by successive Givens rotations applied to rows

of the matrix.

$$\begin{aligned}
 \mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^H &= \left[ \begin{array}{ccccccc}
 * & * & \cdots & * & * & * & * \\
 * & * & \cdots & * & * & * & * \\
 \ddots & & & & & & \\
 0 & \cdots & * & * & * & * & * \\
 0 & \cdots & 0 & * & * & * & * \\
 0 & \cdots & 0 & 0 & * & * & * \\
 0 & \cdots & 0 & 0 & 0 & * & *
 \end{array} \right] \xrightarrow{\mathbf{G}_{12}} \left[ \begin{array}{ccccccc}
 * & * & \cdots & * & * & * & * \\
 0 & * & \cdots & * & * & * & * \\
 \ddots & & & & & & \\
 0 & \cdots & * & * & * & * & * \\
 0 & \cdots & 0 & * & * & * & * \\
 0 & \cdots & 0 & 0 & * & * & * \\
 0 & \cdots & 0 & 0 & 0 & * & *
 \end{array} \right] \xrightarrow{\mathbf{G}_{23}} \cdots \\
 &\xrightarrow{\mathbf{G}_{n-2,n-1}} \left[ \begin{array}{ccccccc}
 * & * & \cdots & * & * & * & * \\
 0 & * & \cdots & * & * & * & * \\
 \ddots & & & & & & \\
 0 & \cdots & 0 & * & * & * & * \\
 0 & \cdots & 0 & 0 & * & * & * \\
 0 & \cdots & 0 & 0 & 0 & * & * \\
 0 & \cdots & 0 & 0 & 0 & * & *
 \end{array} \right] \xrightarrow{\mathbf{G}_{n-1,n}} \left[ \begin{array}{ccccccc}
 * & * & \cdots & * & * & * & * \\
 0 & * & \cdots & * & * & * & * \\
 \ddots & & & & & & \\
 0 & \cdots & 0 & * & * & * & * \\
 0 & \cdots & 0 & 0 & * & * & * \\
 0 & \cdots & 0 & 0 & 0 & * & * \\
 0 & \cdots & 0 & 0 & 0 & 0 & *
 \end{array} \right] =: \tilde{\mathbf{R}}.
 \end{aligned}$$

$(\mathbf{G}_{n-1,n} \mathbf{G}_{n-2,n-1} \cdots \cdots \cdots \mathbf{G}_{23} \mathbf{G}_{12}) (\mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^H) = \tilde{\mathbf{R}}$  (upper triangular!). (3.3.5.2)

We need  $n$  Givens rotations acting on matrix rows of length  $n$ :



Computational effort  $\mathcal{O}(n^2)$



$\mathbf{A} + \mathbf{u}\mathbf{v}^H = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$  with  $\tilde{\mathbf{Q}} = \mathbf{Q}\mathbf{Q}_1^H \mathbf{G}_{12}^H \mathbf{G}_{23}^H \cdots \cdots \mathbf{G}_{n-1,n-2}^H \mathbf{G}_{n,n-1}^H$ .



Asymptotic total computational effort  $\mathcal{O}(mn)$  for  $m, n \rightarrow \infty$

For large  $n$  this is much cheaper than the cost  $\mathcal{O}(n^2m)$  for computing the QR-decomposition of  $\tilde{\mathbf{A}}$  from scratch.

### 3.3.5.2 Adding a Column

We obtain an augmented matrix  $\tilde{\mathbf{A}} \in \mathbb{R}^{m,n+1}$  by inserting a column  $\mathbf{v}$  into  $\mathbf{A} \in \mathbb{R}^{m,n}$  at an arbitrary location.

$$\mathbf{A} \in \mathbb{R}^{m,n} \longrightarrow \tilde{\mathbf{A}} = [\mathbf{a}_1, \dots, \mathbf{a}_{k-1}, \mathbf{v}, \mathbf{a}_k, \dots, \mathbf{a}_n], \quad \mathbf{v} \in \mathbb{R}^m, \quad \mathbf{a}_j := (\mathbf{A})_{:,j}. \quad (3.3.5.3)$$

Given: QR-decomposition  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ ,  $\mathbf{Q} \in \mathbb{R}^{m,m}$  orthogonal,  $\mathbf{R} \in \mathbb{R}^{m,n}$  upper triangular, of  $\mathbf{A}$

Sought: QR-decomposition  $\tilde{\mathbf{A}} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$ ,  $\tilde{\mathbf{Q}} \in \mathbb{R}^{m,n}$  orthogonal,  $\tilde{\mathbf{R}} \in \mathbb{R}^{m,n+1}$  upper triangular, of  $\tilde{\mathbf{A}}$  from (3.3.5.3) computed efficiently.



Idea:

Our task is easy in the case  $k = n+1$ !

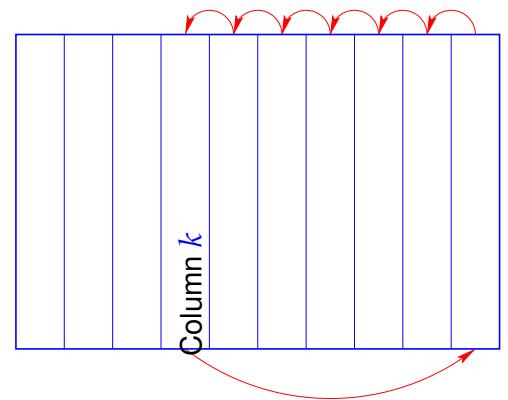
(because we just need target the bottom part of the rightmost column with an orthogonal transformation.)

First move new column to rightmost position

To move columns employ (partial) cyclic permutation of columns

$$k \mapsto n+1, i \mapsto i-1, i = k+1, \dots, n+1,$$

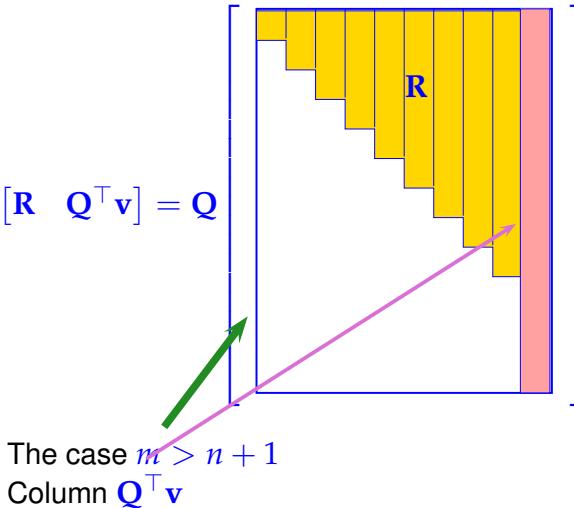
which can be achieved by multiplying from the right with the following (orthogonal!) permutation matrix



$$\mathbf{P} = \begin{bmatrix} 1 & 0 & \cdots & & \cdots & 0 \\ 0 & \ddots & & & & \vdots \\ & & 1 & 0 & & 0 \\ \vdots & & 0 & 0 & & 1 \\ \vdots & & 1 & 0 & \ddots & 0 \\ & & & \ddots & \ddots & \vdots \\ 0 & \cdots & & & 1 & 0 \end{bmatrix} \in \mathbb{R}^{n+1, n+1}$$

This effects the following transformation of  $\tilde{\mathbf{A}}$ :

$$\tilde{\mathbf{A}} \rightarrow \mathbf{A}_1 = \tilde{\mathbf{A}} \mathbf{P}^T = [\mathbf{a}_1, \dots, \mathbf{a}_n, \mathbf{v}] = \mathbf{Q} [\mathbf{R} \quad \mathbf{Q}^T \mathbf{v}] = \mathbf{Q}$$

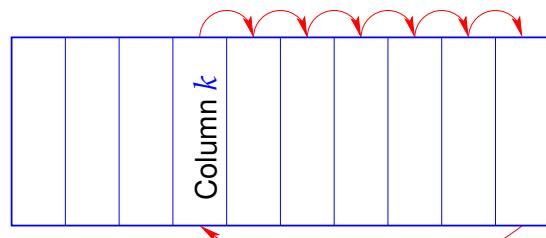


- ① If  $m > n + 1$  there is an orthogonal transformation  $\mathbf{Q}_1 \in \mathbb{R}^{m, m}$ , for instance realized by  $m - n - 1$  Givens rotations or a single Householder reflection, such that

$$\mathbf{Q}_1 \mathbf{Q}^T \mathbf{v} = \begin{bmatrix} * \\ \vdots \\ * \\ * \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \left\{ \begin{array}{l} n+1 \\ m-n-1 \end{array} \right\} \quad \Rightarrow \quad \mathbf{Q}_1 \mathbf{Q}^T \mathbf{A}_1 = \begin{bmatrix} * & \cdots & & \cdots & * \\ 0 & * & & * & \vdots \\ \vdots & \ddots & & * & * \\ & & & * & * \\ 0 & \cdots & & \cdots & 0 \\ \vdots & & & & \vdots \\ 0 & \cdots & & \cdots & 0 \end{bmatrix} \quad \left\{ \begin{array}{l} n+1 \\ m-n-1 \end{array} \right\}$$

&gt;

Computational effort  $\mathcal{O}(m - n)$



To undo the initial moving of the new column we have to apply the inverse permutation:

~ Right-multiply with  $\mathbf{P}^\top$

We can visualize the impact on the structure of the matrix:

$$\mathbf{Q}_1 \mathbf{Q}^\top \tilde{\mathbf{A}} = \mathbf{Q}_1 \mathbf{Q}^\top \mathbf{A}_1 \mathbf{P}^\top = \begin{bmatrix} * & \cdots & * & \cdots & * \\ 0 & * & * & \vdots & \\ \vdots & \ddots & \ddots & \vdots & \\ & \ddots & \ddots & \vdots & \\ & * & * & * & \\ \vdots & * & * & * & \\ 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \\ 0 & \cdots & 0 & \cdots & 0 \end{bmatrix} = \begin{bmatrix} \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} \end{bmatrix}$$

② Apply  $n + 1 - k$  successive Givens rotations to convert  $\mathbf{Q}_1 \mathbf{Q}^\top \tilde{\mathbf{A}}$  into upper triangular form:

$$\mathbf{Q}_1 \mathbf{Q}^\top \mathbf{A}_1 = \begin{bmatrix} * & \cdots & * & \cdots & * \\ 0 & * & * & \vdots & \\ \vdots & \ddots & \ddots & \vdots & \\ & \ddots & \ddots & \vdots & \\ * & * & * & \vdots & \\ * & * & * & \vdots & \\ \vdots & * & 0 & \vdots & \\ 0 & \cdots & 0 & \cdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \\ 0 & \cdots & 0 & \cdots & 0 \end{bmatrix} \xrightarrow{\mathbf{G}_{n,n+1}, \dots, \mathbf{G}_{k,k+1}} \begin{bmatrix} * & \cdots & * & \cdots & * \\ 0 & * & * & \vdots & \\ \vdots & \ddots & \ddots & \vdots & \\ & \ddots & \ddots & \vdots & \\ 0 & 0 & 0 & \ddots & \\ 0 & 0 & 0 & * & \vdots \\ \vdots & 0 & 0 & \cdots & * \\ 0 & \cdots & 0 & \cdots & 0 \end{bmatrix}$$

— target rows of Givens rotations,      ■ new entries  $\neq 0$

➤ Computational effort for this step  $O((n - k)^2)$

➤ Total asymptotic computational cost  $O(n^2 + m)$

### 3.3.5.3 Adding a Row

We are given a matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$  of which a full QR-decomposition ( $\rightarrow$  Thm. 3.3.3.4)  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ ,  $\mathbf{Q} \in \mathbb{R}^{m,m}$  orthogonal,  $\mathbf{R} \in \mathbb{R}^{m,n}$  upper triangular, is already available, maybe only in encoded form ( $\rightarrow$  Rem. 3.3.3.20).

We add another row to the matrix  $\mathbf{A}$  in arbitrary position  $k \in \{1, \dots, m\}$

$$\mathbf{A} \in \mathbb{R}^{m,n} \mapsto \tilde{\mathbf{A}} = \begin{bmatrix} (\mathbf{A})_{1,:} \\ \vdots \\ (\mathbf{A})_{k-1,:} \\ \mathbf{v}^T \\ (\mathbf{A})_{k,:} \\ \vdots \\ (\mathbf{A})_{m,:} \end{bmatrix}, \quad \mathbf{v} \in \mathbb{R}^n. \quad (3.3.5.4)$$

Task: Find an algorithm for the *efficient* computation of the QR-decomposition  $\tilde{\mathbf{A}} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$  of  $\tilde{\mathbf{A}}$  from (3.3.5.4),  $\tilde{\mathbf{Q}} \in \mathbb{R}^{m+1,m+1}$  orthogonal,  $\tilde{\mathbf{R}} \in \mathbb{K}^{m+1,n+1}$  upper triangular.

Step ①: Move new row to the bottom.

Employ partial cyclic permutation of rows of  $\tilde{\mathbf{A}}$ :

$$\text{row } m+1 \leftarrow \text{row } k, \quad \text{row } i \leftarrow \text{row } i+1, \quad i = k, \dots, m.$$

$$\mathbf{P}\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{A} \\ \mathbf{v}^T \end{bmatrix} \xrightarrow{\quad} \begin{bmatrix} \mathbf{Q}^H & 0 \\ 0 & 1 \end{bmatrix} \mathbf{P}\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{R} \\ \mathbf{v}^T \end{bmatrix} = \begin{bmatrix} \mathbf{R} \\ \mathbf{v}^T \end{bmatrix}.$$

Case  $m = n$

Step ②: Restore upper triangular form through Givens rotations ( $\rightarrow$  § 3.3.3.14)

Successively target bottom row and rows from the top to turn leftmost entries of bottom row into zeros.  
Here demonstrated for the case  $m = n$

$$\begin{bmatrix} * & \cdots & & \cdots & * \\ 0 & * & & & \vdots \\ \vdots & 0 & \ddots & & \\ \vdots & \vdots & & * & \vdots \\ 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & * \\ * & \cdots & \cdots & * & * & * \end{bmatrix} \xrightarrow{\mathbf{G}_{1,m+1}} \begin{bmatrix} * & \cdots & & \cdots & * \\ 0 & * & & & \vdots \\ \vdots & 0 & \ddots & & \\ \vdots & \vdots & & * & \vdots \\ 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & * \\ 0 & * & \cdots & * & * & * \end{bmatrix} \xrightarrow{\mathbf{G}_{2,m+1}} \dots$$

$$\dots \xrightarrow{\mathbf{G}_{m-1,m+1}} \begin{bmatrix} * & \cdots & & \cdots & * \\ 0 & * & & & \vdots \\ \vdots & 0 & \ddots & & \\ \vdots & \vdots & & * & \vdots \\ 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & * \\ 0 & \cdots & 0 & * & * \end{bmatrix} \xrightarrow{\mathbf{G}_{m,m+1}} \begin{bmatrix} * & \cdots & & \cdots & * \\ 0 & * & & & \vdots \\ \vdots & 0 & \ddots & & \\ \vdots & \vdots & & * & \vdots \\ 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & * \\ 0 & \cdots & 0 & 0 & 0 \end{bmatrix} := \tilde{\mathbf{R}} \quad (3.3.5.5)$$

Step ③: Incorporate rotations into  $\mathbf{Q}$ :

Setting  $\mathbf{Q}_1 = \mathbf{G}_{m,m+1} \cdots \mathbf{G}_{1,m+1}$  the final QR-decomposition reads

$$\tilde{\mathbf{A}} = \mathbf{P}^T \begin{pmatrix} \mathbf{Q} & 0 \\ 0 & 1 \end{pmatrix} \mathbf{Q}_1^H \tilde{\mathbf{R}} = \tilde{\mathbf{Q}} \tilde{\mathbf{R}} \quad \text{with orthogonal } \tilde{\mathbf{Q}} \in \mathbb{K}^{m+1,m+1},$$

because the product of orthogonal matrices is again orthogonal.

## 3.4 Singular Value Decomposition (SVD)

Beside the QR-decomposition of matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$  there are other factorizations based on orthogonal transformations. The most important among them is the singular value decomposition (SVD), which can be used to tackle linear least squares problems and many other optimization problems beyond, see [Kal96].

### 3.4.1 SVD: Definition and Theory

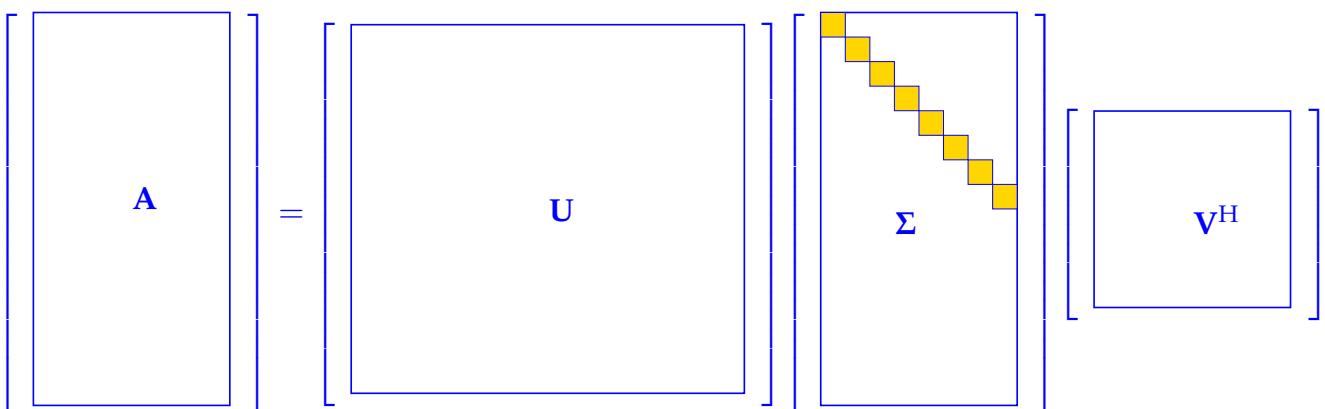
**Theorem 3.4.1.1. singular value decomposition** → [NS02, Thm. 9.6], [Gut09, Thm. 11.1]

For any  $\mathbf{A} \in \mathbb{K}^{m,n}$  there are unitary matrices  $\mathbf{U} \in \mathbb{K}^{m,m}$ ,  $\mathbf{V} \in \mathbb{K}^{n,n}$  and a (generalized) diagonal (\*) matrix  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p) \in \mathbb{R}^{m,n}$ ,  $p := \min\{m, n\}$ ,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  such that

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^H.$$

**Terminology (\*):** A matrix  $\Sigma$  is called a **generalized diagonal matrix**, if  $(\Sigma)_{i,j} = 0$ , if  $i \neq j$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq n$ . We still use the **diag** operator to create it from a vector.

Next, we visualize the structure of the singular value decomposition of a matrix  $\mathbf{A} = \mathbb{K}^{m,n}$ .



$$\begin{bmatrix} \mathbf{A} \\ \vdots \end{bmatrix} = \begin{bmatrix} \mathbf{U} \\ \vdots \end{bmatrix} \begin{bmatrix} \Sigma & \\ & \ddots \end{bmatrix} \begin{bmatrix} \mathbf{V}^H \\ \vdots \end{bmatrix}$$

*Proof. (of Thm. 3.4.1.1, by induction)*

To start the induction note that the assertion of the theorem is immediate for  $n = 1$  or  $m = 1$ .

For the induction step  $(n - 1, m - 1) \Rightarrow (n, m)$  first remember from analysis [Str09, Thm. 4.2.3]: Continuous functions attain extremal values on compact sets (here the unit ball  $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq 1\}$ ). In particular, consider the function  $\mathbf{x} \mapsto \|\mathbf{Ax}\|$ .

$$\blacktriangleright \exists \mathbf{x} \in \mathbb{K}^n, \mathbf{y} \in \mathbb{K}^m, \|\mathbf{x}\| = \|\mathbf{y}\|_2 = 1 : \mathbf{Ax} = \sigma \mathbf{y}, \sigma = \|\mathbf{A}\|_2,$$

where we used the definition of the matrix 2-norm, see Def. 1.5.5.10. By Gram-Schmidt orthogonalization or a similar procedure we can extend the single unit vectors  $\mathbf{x}$  and  $\mathbf{y}$  to orthonormal bases of  $\mathbb{K}^n$  and  $\mathbb{K}^m$ , respectively:  $\exists \tilde{\mathbf{V}} \in \mathbb{K}^{n,n-1}, \tilde{\mathbf{U}} \in \mathbb{K}^{m,m-1}$  such that

$$\mathbf{V} = [\mathbf{x} \ \tilde{\mathbf{V}}] \in \mathbb{K}^{n,n}, \quad \mathbf{U} = [\mathbf{y} \ \tilde{\mathbf{U}}] \in \mathbb{K}^{m,m} \text{ are unitary.}$$

$$\blacktriangleright \mathbf{U}^H \mathbf{A} \mathbf{V} = [\mathbf{y} \ \tilde{\mathbf{U}}]^H \mathbf{A} [\mathbf{x} \ \tilde{\mathbf{V}}] = \left[ \begin{array}{c|c} \mathbf{y}^H \mathbf{A} \mathbf{x} & \mathbf{y}^H \mathbf{A} \tilde{\mathbf{V}} \\ \hline \tilde{\mathbf{U}}^H \mathbf{A} \mathbf{x} & \tilde{\mathbf{U}}^H \mathbf{A} \tilde{\mathbf{V}} \end{array} \right] = \left[ \begin{array}{c|c} \sigma & \mathbf{w}^H \\ \hline \mathbf{0} & \mathbf{B} \end{array} \right] =: \mathbf{A}_1.$$

For the induction argument we have to show that  $\mathbf{w} = \mathbf{0}$ . Since

$$\left\| \mathbf{A}_1 \begin{bmatrix} \sigma \\ \mathbf{w} \end{bmatrix} \right\|_2^2 = \left\| \begin{bmatrix} \sigma^2 + \mathbf{w}^H \mathbf{w} \\ \mathbf{B} \mathbf{w} \end{bmatrix} \right\|_2^2 = (\sigma^2 + \mathbf{w}^H \mathbf{w})^2 + \|\mathbf{B} \mathbf{w}\|_2^2 \geq (\sigma^2 + \mathbf{w}^H \mathbf{w})^2,$$

we conclude

$$\|\mathbf{A}_1\|_2^2 = \sup_{0 \neq \mathbf{x} \in \mathbb{K}^n} \frac{\|\mathbf{A}_1 \mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} \geq \frac{\|\mathbf{A}_1(\sigma)\|_2^2}{\|(\sigma)\|_2^2} \geq \frac{(\sigma^2 + \mathbf{w}^H \mathbf{w})^2}{\sigma^2 + \mathbf{w}^H \mathbf{w}} = \sigma^2 + \mathbf{w}^H \mathbf{w}. \quad (3.4.1.2)$$

$$\sigma^2 = \|\mathbf{A}\|_2^2 = \left\| \mathbf{U}^H \mathbf{A} \mathbf{V} \right\|_2^2 = \|\mathbf{A}_1\|_2^2 \stackrel{(3.4.1.2)}{\implies} \|\mathbf{A}_1\|_2^2 = \|\mathbf{A}_1\|_2^2 + \|\mathbf{w}\|_2^2 \Rightarrow \mathbf{w} = \mathbf{0}.$$

$$\blacktriangleright \mathbf{A}_1 = \left[ \begin{array}{c|c} \sigma & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{B} \end{array} \right].$$

Then apply the induction argument to  $\mathbf{B}$ . □

### Definition 3.4.1.3. Singular value decomposition (SVD)

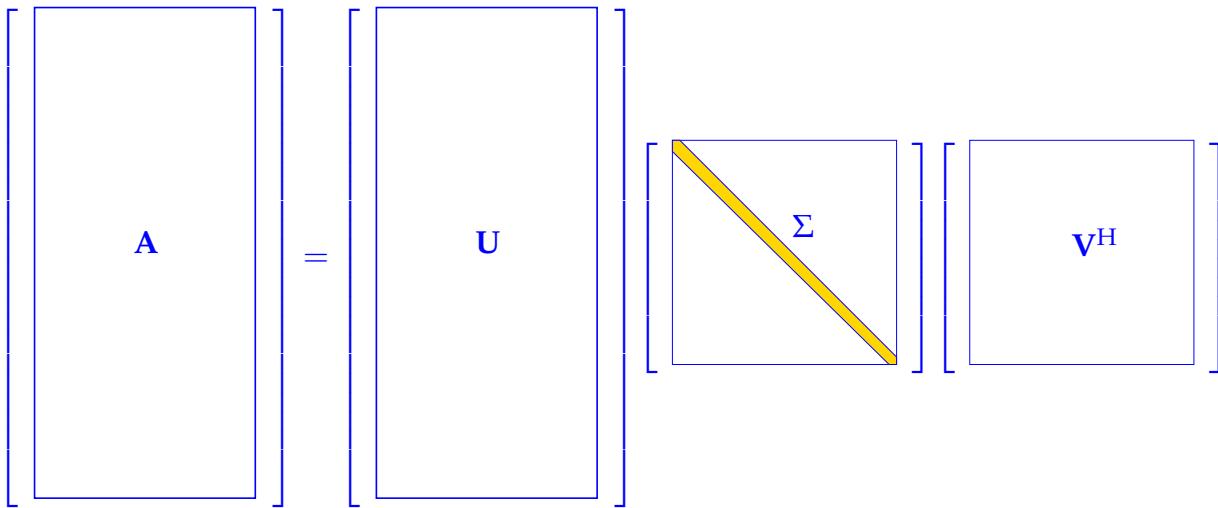
The decomposition  $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^H$  of Thm. 3.4.1.1 is called **singular value decomposition (SVD)** of  $\mathbf{A}$ . The diagonal entries  $\sigma_i$  of  $\Sigma$  are the **singular values** of  $\mathbf{A}$ . The columns of  $\mathbf{U}/\mathbf{V}$  are the left/right **singular vectors** of  $\mathbf{A}$ .

As in the case of the QR-decomposition, compare (3.3.3.1) and (3.3.3.1), we can also drop the bottom zero rows of  $\Sigma$  and the corresponding columns of  $\mathbf{U}$  in the case of  $m > n$ . Thus we end up with an "economical" singular value decomposition of  $\mathbf{A} \in \mathbb{K}^{m,n}$ :

$$\begin{aligned} m \geq n: \quad \mathbf{A} &= \mathbf{U}\Sigma\mathbf{V}^H, \quad \mathbf{U} \in \mathbb{K}^{m,n}, \quad \Sigma \in \mathbb{K}^{n,n}, \quad \mathbf{V} \in \mathbb{K}^{n,n}, \quad \mathbf{U}^H\mathbf{U} = \mathbf{I}_n, \mathbf{V} \text{ unitary,} \\ &\quad (\mathbf{U} \text{ orthonormal columns}) \\ m < n: \quad \mathbf{A} &= \mathbf{U}\Sigma\mathbf{V}^H, \quad \mathbf{U} \in \mathbb{K}^{m,m}, \quad \Sigma \in \mathbb{K}^{m,m}, \quad \mathbf{V} \in \mathbb{K}^{n,m}, \quad \mathbf{U} \text{ unitary,} \mathbf{V}^H\mathbf{V} = \mathbf{I}_m. \\ &\quad (\mathbf{V} \text{ orthonormal columns}) \end{aligned} \tag{3.4.1.4}$$

with true diagonal matrices  $\Sigma$ , whose diagonals contain the singular values of  $\mathbf{A}$ .

Visualization of economical SVD for  $m > 0$ :



### Lemma 3.4.1.5.

The squares  $\sigma_i^2$  of the non-zero singular values of  $\mathbf{A}$  are the non-zero eigenvalues of  $\mathbf{A}^H\mathbf{A}$ ,  $\mathbf{A}\mathbf{A}^H$  with associated eigenvectors  $(\mathbf{V})_{:,1}, \dots, (\mathbf{V})_{:,p}$ ,  $(\mathbf{U})_{:,1}, \dots, (\mathbf{U})_{:,p}$ , respectively.

*Proof.*  $\mathbf{A}\mathbf{A}^H$  and  $\mathbf{A}^H\mathbf{A}$  are similar ( $\rightarrow$  Lemma 9.1.0.6) to diagonal matrices with non-zero diagonal entries  $\sigma_i^2$  ( $\sigma_i \neq 0$ ), e.g.,

$$\mathbf{A}\mathbf{A}^H = \mathbf{U}\Sigma\mathbf{V}^H\mathbf{V}\Sigma^H\mathbf{U}^H = \mathbf{U} \underbrace{\Sigma\Sigma^H}_{\text{diagonal matrix}} \mathbf{U}^H.$$

□

□

### Remark 3.4.1.6 (SVD and additive rank-1 decomposition $\rightarrow$ [Gut09, Cor. 11.2], [NS02, Thm. 9.8])

Recall from linear algebra: rank-1 matrices are tensor products of vectors

$$\mathbf{A} \in \mathbb{K}^{m,n} \quad \text{and} \quad \text{rank}(\mathbf{A}) = 1 \iff \exists \mathbf{u} \in \mathbb{K}^m, \mathbf{v} \in \mathbb{K}^n: \quad \mathbf{A} = \mathbf{u}\mathbf{v}^H, \tag{3.4.1.7}$$

because  $\text{rank}(\mathbf{A}) = 1$  means that  $\mathbf{Ax} = \mu(\mathbf{x})\mathbf{u}$  for some  $\mathbf{u} \in \mathbb{K}^m$  and linear form  $\mathbf{x} \mapsto \mu(\mathbf{x})$ . By the Riesz representation theorem the latter can be written as  $\mu(\mathbf{x}) = \mathbf{v}^H\mathbf{x}$ .

► The Singular value decomposition provides an additive decomposition into rank-1 matrices:

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H = \sum_{j=1}^p \sigma_j (\mathbf{U})_{:,j} (\mathbf{V})_{:,j}^H. \tag{3.4.1.8}$$

Since the columns of  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal, we immediately conclude:

$$\mathbf{A}(\mathbf{V})_{:,j} = \sigma_j(\mathbf{U})_j \quad , \quad \mathbf{A}^H(\mathbf{U})_j = \sigma_j(\mathbf{V})_{:,j} \quad , \quad j \in \{1, \dots, p\} . \quad (3.4.1.9)$$

□

### Remark 3.4.1.10 (Uniqueness of SVD)

The SVD from Def. 3.4.1.3 is not (necessarily) unique, but the singular values are.

*Proof.* Proof by contradiction: assume that  $\mathbf{A}$  has two singular value decompositions

$$\mathbf{A} = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^H = \mathbf{U}_2 \Sigma_2 \mathbf{V}_2^H \Rightarrow \mathbf{U}_1 \underbrace{\Sigma_1 \Sigma_1^H}_{=\text{diag}(\sigma_1^2, \dots, \sigma_m^2)} \mathbf{U}_1^H = \mathbf{A} \mathbf{A}^H = \mathbf{U}_2 \underbrace{\Sigma_2 \Sigma_2^H}_{=\text{diag}(\sigma_1^2, \dots, \sigma_m^2)} \mathbf{U}_2^H .$$

The two diagonal matrices are similar, which implies that they have the same eigenvalues, which agree with their diagonal entries. Since the latter are sorted, the diagonals must agree. □

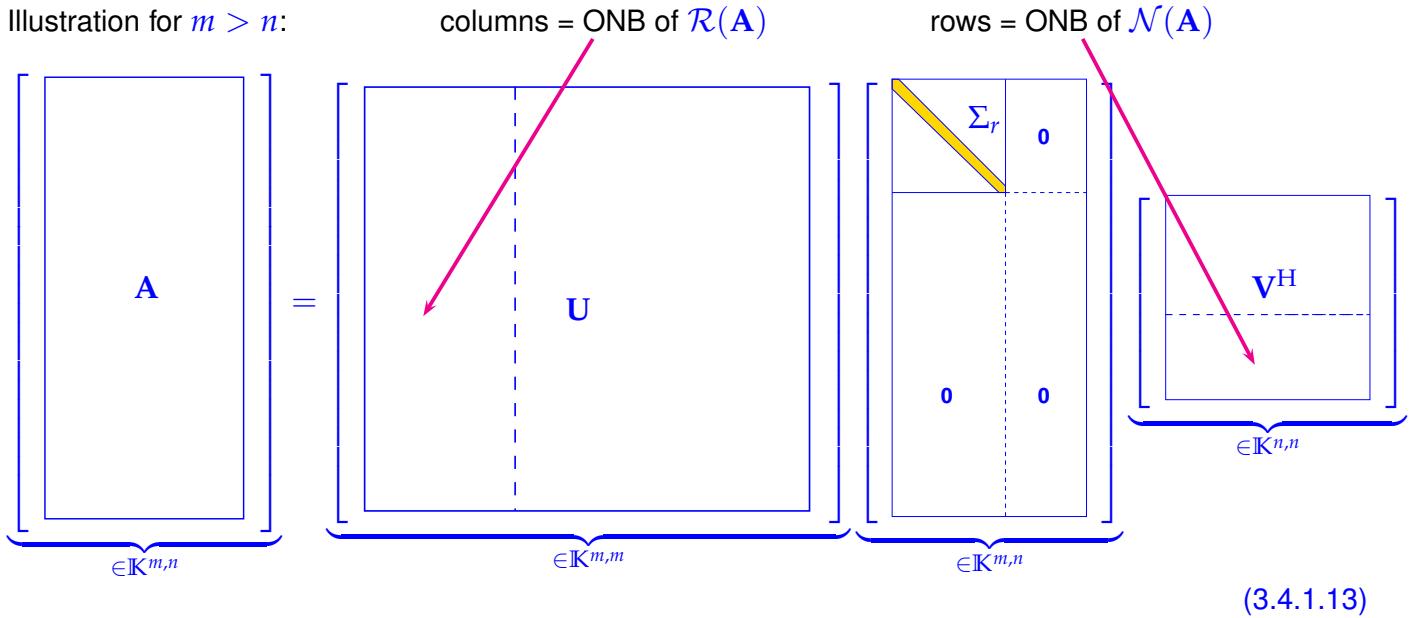
**§3.4.1.11 (SVD, nullspace, and image space)** The SVD give complete information about all crucial subspaces associated with a matrix:

**Lemma 3.4.1.12. SVD and rank of a matrix** → [NS02, Cor. 9.7]

If, for some  $1 \leq r \leq p := \min\{m, n\}$ , the singular values of  $\mathbf{A} \in \mathbb{K}^{m,n}$  satisfy  $\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0$ , then

- $\text{rank}(\mathbf{A}) = r$  (no. of non-zero singular values),
- $\mathcal{N}(\mathbf{A}) = \text{Span}\{(\mathbf{V})_{:,r+1}, \dots, (\mathbf{V})_{:,n}\}$ ,
- $\mathcal{R}(\mathbf{A}) = \text{Span}\{(\mathbf{U})_{:,1}, \dots, (\mathbf{U})_{:,r}\}$ .

Illustration for  $m > n$ :



(3.4.1.13)

□

## 3.4.2 SVD in EIGEN

The EIGEN class **JacobiSVD** is constructed from a matrix data type, computes the SVD of its argument during construction and offers access methods `MatrixU()`, `singularValues()`, and `MatrixV()`

to request the SVD-factors and singular values.

#### C++-code 3.4.2.1: Computing SVDs in EIGEN

```

2 #include <Eigen/SVD>
3
4 // Computation of (full) SVD  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$  → Thm. 3.4.1.1
5 // SVD factors are returned as dense matrices in natural order
6 std::tuple<MatrixXd,MatrixXd,MatrixXd> svd_full(const MatrixXd& A) {
7     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeFullU | Eigen::ComputeFullV);
8     MatrixXd U = svd.matrixU(); // get unitary (square) matrix  $\mathbf{U}$ 
9     MatrixXd V = svd.matrixV(); // get unitary (square) matrix  $\mathbf{V}$ 
10    VectorXd sv = svd.singularValues(); // get singular values as vector
11    MatrixXd Sigma = MatrixXd::Zero(A.rows(), A.cols());
12    const unsigned p = sv.size(); // no. of singular values
13    Sigma.block(0,0,p,p) = sv.asDiagonal(); // set diagonal block of  $\Sigma$ 
14    return std::tuple<MatrixXd,MatrixXd,MatrixXd>(U,Sigma,V);
15}
16
17 // Computation of economical (thin) SVD  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$ , see (3.4.1.4)
18 // SVD factors are returned as dense matrices in natural order
19 std::tuple<MatrixXd,MatrixXd,MatrixXd> svd_eco(const MatrixXd& A) {
20     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeThinU | Eigen::ComputeThinV);
21     MatrixXd U = svd.matrixU(); // get matrix  $\mathbf{U}$  with orthonormal columns
22     MatrixXd V = svd.matrixV(); // get matrix  $\mathbf{V}$  with orthonormal columns
23     VectorXd sv = svd.singularValues(); // get singular values as vector
24     MatrixXd Sigma = sv.asDiagonal(); // build diagonal matrix  $\Sigma$ 
25     return std::tuple<MatrixXd,MatrixXd,MatrixXd>(U,Sigma,V);
26}

```

The second argument in the constructor of **JacobiSVD** determines, whether the methods `matrixU()` and `matrixV()` return the factor for the full SVD of Def. 3.4.1.3 or of the economical (thin) SVD (3.4.1.4): `Eigen::ComputeFull*` will select the full versions, whereas `Eigen::ComputeThin*` picks the economical versions → [documentation](#).

Internally, the computation of the SVD is done by a sophisticated algorithm, for which key steps rely on orthogonal/unitary transformations. Also there we reap the benefit of the exceptional stability brought about by norm-preserving transformations → § 3.3.3.23.

EIGEN's algorithm for computing SVD is (numerically) stable → Def. 1.5.5.19

**§3.4.2.2 (Computational cost of computing the SVD)** According to EIGEN's [documentation](#) the SVD of a general dense matrix involves the following asymptotic complexity:

$$\text{cost(economical SVD of } \mathbf{A} \in \mathbb{K}^{m,n}) = O(\min\{m, n\}^2 \max\{m, n\})$$

► The computational effort is (asymptotically) *linear in the larger matrix dimension*. ↴

**EXAMPLE 3.4.2.3 (SVD-based computation of the rank of a matrix)** Based on Lemma 3.4.1.12, the SVD is the main tool for the stable computation of the **rank** of a matrix (→ Def. 2.2.1.3)

However, theory as reflected in Lemma 3.4.1.12 entails identifying zero singular values, which must rely

on a **threshold condition** in a numerical code, recall Rem. 1.5.3.15. Given the SVD  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$ ,  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min\{m,n\}})$ , of a matrix  $\mathbf{A} \in \mathbb{K}^{m,n}$ ,  $\mathbf{A} \neq \mathbf{0}$  and a **tolerance**  $\text{tol} > 0$ , we define the **numerical rank**

$$r := \#\left\{\sigma_i : |\sigma_i| \geq \text{tol} \max_j\{|\sigma_j|\}\right\}. \quad (3.4.2.4)$$

The following code implements this rule.

#### C++-code 3.4.2.5: Computing rank of a matrix through SVD

```

2 // Computation of the numerical rank of a non-zero matrix by means of
3 // singular value decomposition, cf. (3.4.2.4).
4 MatrixXd::Index rank_by_svd(const MatrixXd &A, double tol = EPS) {
5     if (A.norm() == 0) return MatrixXd::Index(0);
6     Eigen::JacobiSVD<MatrixXd> svd(A);
7     const VectorXd sv = svd.singularValues(); // Get sorted singular values as
8     // vector
9     MatrixXd::Index n = sv.size();
10    MatrixXd::Index r = 0;
11    // Test relative size of singular values
12    while ((r < n) && (sv(r) >= sv(0)*tol)) r++;
13    return r;
}
```

EIGEN offers an equivalent built-in method `rank()` for objects representing singular value decompositions:

#### C++-code 3.4.2.6: Using `rank()` in EIGEN

```

2 // Computation of the numerical rank of a matrix by means of SVD
3 MatrixXd::Index rank_eigen(const MatrixXd &A, double tol = EPS) {
4     return A.jacobiSvd().setThreshold(tol).rank();
5 }
```

The method `setThreshold()` passes `tol` from (3.4.2.4) to `rank()`. □

**EXAMPLE 3.4.2.7 (Computation of nullspace and image space of matrices)** “Computing” a subspace of  $\mathbb{R}^k$  amounts to making available a (stable) **basis** of that subspace, ideally an **orthonormal basis**.

Lemma 3.4.1.12 taught us how to glean orthonormal bases of  $\mathcal{N}(\mathbf{A})$  and  $\mathcal{R}(\mathbf{A})$  from the SVD of a matrix  $\mathbf{A}$ . This immediately gives a numerical method and its implementation is given in the next two codes.

#### C++-code 3.4.2.8: ONB of $\mathcal{N}(\mathbf{A})$ through SVD

```

2 // Computation of an ONB of the kernel of a matrix
3 MatrixXd nullspace(const MatrixXd &A, double tol = EPS) {
4     using index_t = MatrixXd::Index;
5     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeFullV);
6     index_t r = svd.setThreshold(tol).rank();
7     // Rightmost columns of V provide ONB of N(A)
8     MatrixXd Z = svd.matrixV().rightCols(A.cols() - r);
9     return Z;
10 }
```

### C++-code 3.4.2.9: ONB of $\mathcal{R}(\mathbf{A})$ through SVD

```

2 // Computation of an ONB of the image space of a matrix
3 MatrixXd rangespace (const MatrixXd &A, double tol = EPS) {
4     using index_t = MatrixXd::Index;
5     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeFullV);
6     index_t r = svd.setThreshold(tol).rank();
7     // r left columns of U provide ONB of  $\mathcal{R}(\mathbf{A})$ 
8     return svd.matrixU().leftCols(r);
9 }
```

### 3.4.3 Solving General Least-Squares Problems by SVD

In a similar fashion as explained for QR-decomposition in Section 3.3.4, the singular valued decomposition (SVD, → Def. 3.4.1.3) can be used to transform **general** linear least squares problems (3.1.3.7) into a simpler form. In the case of SVD-based orthogonal transformation methods this simpler form involves merely a diagonal matrix.

Here we consider the most general setting

$$\mathbf{Ax} = \mathbf{b} \in \mathbb{K}^m \quad \text{with} \quad \mathbf{A} \in \mathbb{K}^{m,n}, \quad \text{rank}(\mathbf{A}) = r \leq \min\{m, n\}.$$

In particular, we drop the assumption of full rank of  $\mathbf{A}$ . This means that the minimum norm condition (ii) in the definition (3.1.3.7) of a linear least squares problem may be required for singling out a unique solution.

We recall the (full) SVD of  $\mathbf{A} \in \mathbb{K}^{m,n}$ :

$$\begin{array}{ccccccc}
 \mathbf{A} & = & [\mathbf{U}_1 & \mathbf{U}_2] & \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix} & \begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix} \\
 & = & \begin{bmatrix} \mathbf{U}_1 & \mathbf{U}_2 \end{bmatrix} & & \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix} & \begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix} \\
 & & \in \mathbb{K}^{m,m} & & \in \mathbb{K}^{m,n} & \in \mathbb{K}^{n,n} \\
 & & & & & & (3.4.3.1)
 \end{array}$$

with  $\mathbf{U}_1 \in \mathbb{K}^{m,r}$ ,  $\mathbf{U}_2 \in \mathbb{K}^{m,m-r}$  with *orthonormal* columns,  $\mathbf{U} := [\mathbf{u}_1, \mathbf{u}_2]$  *unitary*,  
 $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{K}^{r,r}$  (*singular values*, Def. 3.4.1.3),  
 $\mathbf{V}_1 \in \mathbb{K}^{n,r}$ ,  $\mathbf{V}_2 \in \mathbb{K}^{n,n-r}$  with *orthonormal* columns,  $\mathbf{V} := [\mathbf{v}_1, \mathbf{v}_2]$  *unitary*.

We can proceed in two different ways, both of which we elaborate in the sequel:

**Approach ①:** We can use the invariance of the 2-norm of a vector with respect to multiplication with  $\mathbf{U} := [\mathbf{U}_1, \mathbf{U}_2]$ , see Thm. 3.3.2.2, together with the fact that  $\mathbf{U}$  is unitary,  $\mathbf{U}^{-1} = \mathbf{U}^H$ , see Def. 6.2.1.2:

$$[\mathbf{U}_1, \mathbf{U}_2] \cdot \begin{bmatrix} (\mathbf{U}_1^H) \\ (\mathbf{U}_2^H) \end{bmatrix} = \mathbf{I}.$$

►  $\|\mathbf{Ax} - \mathbf{b}\|_2 = \left\| [\mathbf{U}_1 \ \mathbf{U}_2] \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix} \mathbf{x} - \mathbf{b} \right\|_2 = \left\| \begin{bmatrix} \Sigma_r \mathbf{V}_1^H \mathbf{x} \\ 0 \end{bmatrix} - \begin{bmatrix} \mathbf{U}_1^H \mathbf{b} \\ \mathbf{U}_2^H \mathbf{b} \end{bmatrix} \right\|_2 \quad (3.4.3.2)$

We follow the same strategy as in the case of QR-based solvers for full-rank linear least squares problems.

We choose  $\mathbf{x}$  such that the first  $r$  components of  $\begin{bmatrix} \Sigma_r \mathbf{V}_1^H \mathbf{x} \\ 0 \end{bmatrix} - \begin{bmatrix} \mathbf{U}_1^H \mathbf{b} \\ \mathbf{U}_2^H \mathbf{b} \end{bmatrix}$  vanish:

➤ (possibly *underdetermined*)  $r \times n$  linear system  $\Sigma_r \mathbf{V}_1^H \mathbf{x} = \mathbf{U}_1^H \mathbf{b}$ . (3.4.3.3)

To fix a unique solution in the case  $r < n$  we appeal to the **minimal norm condition** in (3.1.3.7): appealing to the considerations of § 3.1.3.3, the solution  $\mathbf{x}$  of (3.4.3.3) is unique up to contributions from

$$\mathcal{N}(\mathbf{V}_1^H) \stackrel{\text{Lemma 3.1.2.12}}{=} \mathcal{R}(\mathbf{V}_1)^\perp \stackrel{\text{orthonormality}}{=} \mathcal{R}(\mathbf{V}_2). \quad (3.4.3.4)$$

Since  $\mathbf{V}$  is unitary, the minimal norm solution is obtained by setting contributions from  $\mathcal{R}(\mathbf{V}_2)$  to zero, which amounts to choosing  $\mathbf{x} \in \mathcal{R}(\mathbf{V}_1)$ . This converts (3.4.3.3) into

$$\Sigma_r \underbrace{\mathbf{V}_1^H \mathbf{V}_1}_{=I} \mathbf{z} = \mathbf{U}_1^H \mathbf{b} \Rightarrow \mathbf{z} = \Sigma_r^{-1} \mathbf{U}_1^H \mathbf{b}.$$

► **generalized solution** → Def. 3.1.3.1  $\mathbf{x}^\dagger = \mathbf{V}_1 \Sigma_r^{-1} \mathbf{U}_1^H \mathbf{b}$ ,  $\|\mathbf{r}\|_2 = \left\| \mathbf{U}_2^H \mathbf{b} \right\|_2$ . (3.4.3.5)

**Approach ②:** From Thm. 3.1.2.1 we know that the generalized least-squares solution  $\mathbf{x}^\dagger$  of  $\mathbf{Ax} = \mathbf{b}$  solves the normal equations (3.1.2.2), and in § 3.1.3.3 we saw that  $\mathbf{x}^\dagger$  lies in the orthogonal complement of  $\mathcal{N}(\mathbf{A})$ :

$$\mathbf{A}^H \mathbf{A} \mathbf{x}^\dagger = \mathbf{A}^H \mathbf{b}, \quad \mathbf{x}^\dagger \in \mathcal{N}(\mathbf{A})^\perp. \quad (3.4.3.6)$$

By Lemma 3.4.1.12 and using the notations from (3.4.3.1) together with the fact that the columns of  $\mathbf{V}$  form an orthonormal basis of  $\mathbb{K}^n$ :

$$\mathcal{N}(\mathbf{A}) = \mathcal{R}(\mathbf{V}_2) \Leftrightarrow \mathcal{N}(\mathbf{A})^\perp = \mathcal{R}(\mathbf{V}_1). \quad (3.4.3.7)$$

Hence, we can write  $\mathbf{x}^\dagger = \mathbf{V}_1 \mathbf{y}$  for some  $\mathbf{y} \in \mathbb{K}^r$ . We plug this representation into the normal equations and also multiply with  $\mathbf{V}_1^H$ , similar to what we did in § 3.1.3.3:

►  $\mathbf{V}_1^H \mathbf{A}^H \mathbf{A} \mathbf{V}_1 \mathbf{y} = \mathbf{V}_1^H \mathbf{b}$ . (3.4.3.8)

Next, insert the SVD of  $\mathbf{A}$ :

$$\mathbf{V}_1^H \mathbf{V} \Sigma \underbrace{\mathbf{U}^H \mathbf{U}}_{=I} \Sigma \mathbf{V} \mathbf{V}_1^H \mathbf{y} = \mathbf{V}_1^H \mathbf{V} \Sigma \mathbf{U}^H \mathbf{b}. \quad (3.4.3.9)$$

Then we switch to the block-partitioned form as in (3.4.3.1):

$$\underbrace{\mathbf{V}_1^H [\mathbf{V}_1, \mathbf{V}_2]}_{\in \mathbb{K}^{r,n}} \underbrace{\begin{bmatrix} \Sigma_r & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}}_{\in \mathbb{K}^{n,n}}^H \underbrace{\begin{bmatrix} \Sigma_r & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}}_{\in \mathbb{K}^{n,n}} \underbrace{\begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix}}_{\in \mathbb{K}^{n,r}} \mathbf{V}_1 \mathbf{y} = \mathbf{V}_1^H [\mathbf{V}_1, \mathbf{V}_2] \underbrace{\begin{bmatrix} \Sigma_r & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}}_{\in \mathbb{K}^{n,n}}^H \underbrace{\begin{bmatrix} \mathbf{U}_1^H \\ \mathbf{U}_2^H \end{bmatrix}}_{\in \mathbb{K}^{n,r}} \mathbf{b} \quad (3.4.3.10)$$

$$\begin{aligned} & \uparrow\downarrow \\ [\mathbf{I}, \mathbf{O}] \left[ \begin{array}{c|c} \Sigma_r^2 & \mathbf{O} \\ \hline \mathbf{O} & \mathbf{O} \end{array} \right] \left[ \begin{array}{c} \mathbf{I} \\ \mathbf{O} \end{array} \right] \mathbf{y} &= [\mathbf{I}, \mathbf{O}] \left[ \begin{array}{c|c} \Sigma_r & \mathbf{O} \\ \hline \mathbf{O} & \mathbf{O} \end{array} \right]^H \left[ \begin{array}{c} \mathbf{U}_1^H \\ \mathbf{U}_2^H \end{array} \right] \mathbf{b} \quad (3.4.3.11) \\ & \uparrow\downarrow \\ \Sigma_r^2 \mathbf{y} &= \Sigma_r \mathbf{U}_1^H \mathbf{b}. \quad (3.4.3.12) \end{aligned}$$

Cancelling the invertible matrix  $\Sigma_r$  on both sides yields formula (3.4.3.5).

We remind that in a practical implementation, as in Code 3.4.2.5, one has to resort to the numerical rank from (3.4.2.4):

$$r = \max\{i : \sigma_i / \sigma_1 > \text{tol}\},$$

where we have assumed that the singular values  $\sigma_j$  are sorted according to decreasing modulus.

#### C++-code 3.4.3.13: Computing generalized solution of $\mathbf{Ax} = \mathbf{b}$ via SVD

```

2 #include <Eigen/SVD>
3
4 VectorXd lsqsvd(const MatrixXd &A, const VectorXd &b) {
5     // Compute economical SVD, compare Code 3.4.2.1
6     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeThinU | Eigen::ComputeThinV);
7     VectorXd sv = svd.singularValues();
8     unsigned int r = svd.rank(); // Numerical rank, default tolerance
9     MatrixXd U = svd.matrixU(), V = svd.matrixV();
10    //  $\mathbf{x}^\dagger = \mathbf{V}_1 \Sigma_r^{-1} \mathbf{U}_1^H \mathbf{b}$ , see
11    // (3.4.3.5)
12    return V.leftCols(r) * (sv.head(r).cwiseInverse().asDiagonal() *
13                             (U.leftCols(r).adjoint() * b));
14}

```

The `solve()` method directly returns the generalized solution

#### C++-code 3.4.3.14: Computing generalized solution of $\mathbf{Ax} = \mathbf{b}$ via SVD

```

2 VectorXd lsqsvd_eigen(const MatrixXd &A, const VectorXd &b) {
3     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeThinU | Eigen::ComputeThinV);
4     return svd.solve(b);
5 }

```

**Remark 3.4.3.15 (Pseudoinverse and SVD → [Han02, Ch. 12], [DR08, Sect. 4.7])** From Thm. 3.1.3.6 we could conclude a general formula for the Moore-Penrose pseudoinverse of any matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ . Now, the solution formula (3.4.3.5) directly yields a concrete incarnation of the pseudoinverse  $\mathbf{A}^+$ .

#### Theorem 3.4.3.16. Pseudoinverse and SVD

If  $\mathbf{A} \in \mathbb{K}^{m,n}$  has the SVD decomposition  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$  partitioned as in (3.4.3.1), then its Moore-Penrose pseudoinverse ( $\rightarrow$  Thm. 3.1.3.6) is given by  $\mathbf{A}^+ = \mathbf{V}_1 \Sigma_r^{-1} \mathbf{U}_1^H$ .

### 3.4.4 SVD-Based Optimization and Approximation

For the general least squares problem (3.1.3.7) we have seen the use of SVD for its numerical solution in Section 3.4.3. The the SVD was a powerful tool for solving a minimization problem for a 2-norm. In many other contexts the SVD is also a key component in numerical optimization.

#### 3.4.4.1 Norm-Constrained Extrema of Quadratic Forms

We consider the following problem of finding the extrema of quadratic forms on the Euclidean unit sphere  $\{\mathbf{x} \in \mathbb{K}^n : \|\mathbf{x}\|_2 = 1\}$ :

$$\text{given } \mathbf{A} \in \mathbb{K}^{m,n}, m \geq n, \text{ find } \mathbf{x} \in \mathbb{K}^n, \|\mathbf{x}\|_2 = 1, \|\mathbf{Ax}\|_2 \rightarrow \min. \quad (3.4.4.1)$$

Use that multiplication with orthogonal/unitary matrices preserves the 2-norm ( $\rightarrow$  Thm. 3.3.2.2) and resort to the (full) singular value decomposition  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$  ( $\rightarrow$  Def. 3.4.1.3):

$$\begin{aligned} \min_{\|\mathbf{x}\|_2=1} \|\mathbf{Ax}\|_2^2 &= \min_{\|\mathbf{x}\|_2=1} \left\| \mathbf{U}\Sigma\mathbf{V}^H\mathbf{x} \right\|_2^2 = \min_{\|\mathbf{V}^H\mathbf{x}\|_2=1} \left\| \mathbf{U}\Sigma(\mathbf{V}^H\mathbf{x}) \right\|_2^2 \\ [\mathbf{y} = \mathbf{V}^H\mathbf{x}] &= \min_{\|\mathbf{y}\|_2=1} \|\Sigma\mathbf{y}\|_2^2 = \min_{\|\mathbf{y}\|_2=1} (\sigma_1^2 y_1^2 + \dots + \sigma_n^2 y_n^2) \geq \sigma_n^2. \end{aligned}$$

Since the singular values are assumed to be sorted as  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ , the minimum with value  $\sigma_n^2$  is attained for  $y_n^2 = 1$  and  $y_1 = \dots = y_{n-1} = 0$ , that is,  $\mathbf{V}^H\mathbf{x} = \mathbf{y} = \mathbf{e}_n$  ( $\hat{n}$ -th Cartesian basis vector in  $\mathbb{R}^n$ ).  $\Rightarrow$  minimizer  $\mathbf{x}^* = \mathbf{Ve}_n = (\mathbf{V})_{:,n}$ , minimal value  $\|\mathbf{Ax}^*\|_2 = \sigma_n$ .

#### C++ code 3.4.4.2: Solving (3.4.4.1) with EIGEN $\rightarrow$ GITLAB

```

2 // EIGEN based function for solving (3.4.4.1);
3 // minimizer returned in x, minimum as return value
4 double minconst(VectorXd &x, const MatrixXd &A) {
5     MatrixXd::Index m=A.rows(), n=A.cols();
6     if (m < n) throw std::runtime_error("A must be tall matrix");
7     // SVD factor U is not computed!
8     Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeThinV);
9     x.resize(n); x.setZero(); x(n-1) = 1.0; // e_n
10    x = svd.matrixV() * x;
11    return (svd.singularValues()(n-1));
12}

```

By similar arguments we can solve the corresponding norm constrained maximization problem

$$\text{given } \mathbf{A} \in \mathbb{K}^{m,n}, m \geq n, \text{ find } \mathbf{x} \in \mathbb{K}^n, \|\mathbf{x}\|_2 = 1, \|\mathbf{Ax}\|_2 \rightarrow \max,$$

and obtain the solution based on the SVD  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$  of  $\mathbf{A}$ :

$$\sigma_1 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{Ax}\|_2, (\mathbf{V})_{:,1} = \underset{\|\mathbf{x}\|_2=1}{\operatorname{argmax}} \|\mathbf{Ax}\|_2. \quad (3.4.4.3)$$

Recall: The Euclidean matrix norm (2-norm) of the matrix  $\mathbf{A}$  ( $\rightarrow$  Def. 1.5.5.10) is defined as the maximum in (3.4.4.3). Thus we have proved the following theorem:

#### Lemma 3.4.4.4. SVD and Euclidean matrix norm

If  $\mathbf{A} \in \mathbb{K}^{m,n}$  has singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ ,  $p := \min\{m, n\}$ , then its Euclidean matrix norm is given by  $\|\mathbf{A}\|_2 = \sigma_1(\mathbf{A})$ .

If  $m = n$  and  $\mathbf{A}$  is regular/invertible, then its 2-norm condition number is  $\operatorname{cond}_2(\mathbf{A}) = \sigma_1/\sigma_n$ .

**EXAMPLE 3.4.4.5 (Fit of hyperplanes)** For an important application from computational geometry, this example studies the power and versatility of orthogonal transformations in the context of (generalized) least squares minimization problems.

From school recall the **Hesse normal form** of a hyperplane  $\mathcal{H}$  (= affine subspace of dimension  $d - 1$ ) in  $\mathbb{R}^d$ :

$$\mathcal{H} = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{c} + \mathbf{n}^\top \mathbf{x} = 0\}, \quad \|\mathbf{n}\|_2 = 1. \quad (3.4.4.6)$$

where  $\mathbf{n}$  is the unit normal to  $\mathcal{H}$  and  $|\mathbf{c}|$  gives the distance of  $\mathcal{H}$  from  $\mathbf{0}$ . The Hesse normal form is convenient for computing the distance of points from  $\mathcal{H}$ , because the

$$\text{Euclidean distance of } \mathbf{y} \in \mathbb{R}^d \text{ from the plane is } \text{dist}(\mathcal{H}, \mathbf{y}) = |\mathbf{c} + \mathbf{n}^\top \mathbf{y}|, \quad (3.4.4.7)$$

Goal: given the point coordinate vectors  $\mathbf{y}_1, \dots, \mathbf{y}_m \in \mathbb{R}^d$ ,  $m > d$ , find  $\mathcal{H} \leftrightarrow \{\mathbf{c} \in \mathbb{R}, \mathbf{n} \in \mathbb{R}^d, \|\mathbf{n}\|_2 = 1\}$ , such that

$$\sum_{j=1}^m \text{dist}(\mathcal{H}, \mathbf{y}_j)^2 = \sum_{j=1}^m |\mathbf{c} + \mathbf{n}^\top \mathbf{y}_j|^2 \rightarrow \min. \quad (3.4.4.8)$$

Note that (3.4.4.8) is **not** a linear least squares problem due to the constraint  $\|\mathbf{n}\|_2 = 1$ . However, it turns out to be a minimization problem with *almost* the structure of (3.4.4.1) ( $\mathbf{y}_{k,\ell} := (\mathbf{y}_k)_\ell$ ):

$$(3.4.4.8) \Leftrightarrow \left\| \begin{bmatrix} 1 & y_{1,1} & \cdots & y_{1,d} \\ 1 & y_{2,1} & \cdots & y_{2,d} \\ \vdots & \vdots & & \vdots \\ 1 & y_{m,1} & \cdots & y_{m,d} \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{n}_1 \\ \vdots \\ \mathbf{n}_d \end{bmatrix} \right\|_2 \rightarrow \min \quad \text{under the constraint } \|\mathbf{n}\|_2 = 1.$$

Note that the solution component  $\mathbf{c}$  is not subject to the constraint. One is tempted to use this freedom to make one component of  $\mathbf{Ax}$  vanish, but which one is not clear. This is why we need another preparatory step.

Step ①: To convert the minimization problem into the form (3.4.4.1) we start with a QR-decomposition ( $\rightarrow$  Section 3.3.3)

$$\mathbf{A} := \begin{bmatrix} 1 & y_{1,1} & \cdots & y_{1,d} \\ 1 & y_{2,1} & \cdots & y_{2,d} \\ \vdots & \vdots & & \vdots \\ 1 & y_{m,1} & \cdots & y_{m,d} \end{bmatrix} = \mathbf{QR}, \quad \mathbf{R} := \begin{bmatrix} r_{11} & r_{12} & \cdots & \cdots & r_{1,d+1} \\ 0 & r_{22} & \cdots & \cdots & r_{2,d+1} \\ \vdots & \ddots & & & \vdots \\ 0 & & & & r_{d+1,d+1} \\ 0 & \cdots & \cdots & & 0 \\ \vdots & & & & \vdots \\ 0 & \cdots & \cdots & & 0 \end{bmatrix} \in \mathbb{R}^{m,d+1}.$$

$$\|\mathbf{Ax}\|_2 \rightarrow \min \Leftrightarrow \|\mathbf{Rx}\|_2 = \left\| \begin{bmatrix} r_{11} & r_{12} & \cdots & \cdots & r_{1,d+1} \\ 0 & r_{22} & \cdots & \cdots & r_{2,d+1} \\ \vdots & \ddots & & & \vdots \\ 0 & & & & r_{d+1,d+1} \\ 0 & \cdots & \cdots & & 0 \\ \vdots & & & & \vdots \\ 0 & \cdots & \cdots & & 0 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{n}_1 \\ \vdots \\ \mathbf{n}_d \end{bmatrix} \right\|_2 \rightarrow \min. \quad (3.4.4.9)$$

Step ② Note that, if  $\mathbf{n}$  is the solution of (3.4.4.8), then necessarily (why?)

$$c \cdot r_{11} + n_1 \cdot r_{12} + \cdots + r_{1,d+1} \cdot n_d = 0.$$

This insight converts (3.4.4.9) to

$$\left\| \begin{bmatrix} r_{22} & r_{23} & \cdots & \cdots & r_{2,d+1} \\ 0 & r_{33} & \cdots & \cdots & r_{3,d+1} \\ \vdots & \ddots & & & \vdots \\ 0 & & & & r_{d+1,d+1} \end{bmatrix} \begin{bmatrix} n_1 \\ \vdots \\ n_d \end{bmatrix} \right\|_2 \rightarrow \min, \quad \|\mathbf{n}\|_2 = 1. \quad (3.4.4.10)$$

(3.4.4.10) is now a problem of type (3.4.4.1), minimization on the Euclidean sphere. Hence, (3.4.4.10) can be solved using the SVD-based algorithm implemented in Code 3.4.4.2.

Note: Since  $r_{11} = \|(\mathbf{A})_{:,1}\|_2 = \sqrt{m} \neq 0$ ,  $c = -r_{11}^{-1} \sum_{j=1}^d r_{1,j+1} n_j$  can always be computed.

This algorithm is implemented as case  $p==dim+1$  in the following code, making heavy use of EIGEN's block access operations and the built-in QR-decomposition and SVD factorization.

#### C++-code 3.4.4.11: (Generalized) distance fitting of a hyperplane: solution of (3.4.4.12)

```

2 // Fitting of a hyperplane, returning Hesse normal form
3 void clsq(const MatrixXd& A, const unsigned dim,
4     double& c, VectorXd& n) {
5     unsigned p = A.cols(), m = A.rows();
6     if (p < dim + 1) { cerr << "not enough unknowns\n"; return; }
7     if (m < dim) { cerr << "not enough equations\n"; return; }
8
9     m = std::min(m, p);
10    // First step: orthogonal transformation, see Code 3.3.4.1
11    MatrixXd R = A.householderQr().matrixQR().template triangularView<Eigen::Upper>();
12    // compute matrix V from SVD composition of R, solve (3.4.4.10)
13    MatrixXd V = R.block(p - dim, p - dim, m + dim - p, dim)
14        .jacobiSvd(Eigen::ComputeFullV).matrixV();
15    n = V.col(dim - 1);
16
17    MatrixXd R_topleft = R.topLeftCorner(p - dim, p - dim);
18    c = -(R_topleft.template triangularView<Eigen::Upper>()
19        .solve(R.block(0, p - dim, p - dim, dim)) * n)(0);
20}

```

Note that Code 3.4.4.11 solves the general problem: For  $\mathbf{A} \in \mathbb{K}^{m,n}$  find  $\mathbf{n} \in \mathbb{R}^d$ ,  $\mathbf{c} \in \mathbb{R}^{n-d}$  such that

$$\left\| \mathbf{A} \begin{bmatrix} \mathbf{c} \\ \mathbf{n} \end{bmatrix} \right\|_2 \rightarrow \min \quad \text{with constraint} \quad \|\mathbf{n}\|_2 = 1. \quad (3.4.4.12)$$

#### 3.4.4.2 Best Low-Rank Approximation

**§3.4.4.13 (Low-Rank matrix compression)** **Matrix compression** addresses the problem of approximating a given “generic” matrix (of a certain class) by means of matrix, whose “information content”, that is, the number of reals needed to store it, is significantly lower than the information content of the original matrix.

Sparse matrices ( $\rightarrow$  Notion 2.7.0.1) are a prominent class of matrices with “low information content”. Unfortunately, they cannot approximate dense matrices very well. Another type of matrices that enjoy “low information content”, also called **data sparse**, are low-rank matrices.

**Lemma 3.4.4.14.**

If  $\mathbf{A} \in \mathbb{R}^{m,n}$  has rank  $r \leq \min\{m, n\}$  ( $\rightarrow$  Def. 2.2.1.3), then there exist  $\mathbf{U} \in \mathbb{R}^{m,r}$  and  $\mathbf{V} \in \mathbb{R}^{n,r}$ , such that  $\mathbf{A} = \mathbf{U}\mathbf{V}^\top$ .

*Proof.* The lemma is a straightforward consequence of Lemma 3.4.1.12 and (3.4.1.13): choose  $\mathbf{U}, \mathbf{V}$  as first  $r$  columns of the corresponding SVD-factors.  $\square$

None of the columns of  $\mathbf{U}$  and  $\mathbf{V}$  can vanish. Hence, in addition, we may assume that the columns of  $\mathbf{U}$  are normalized:  $\|(\mathbf{U})_{:,j}\|_2 = 1, j = 1, \dots, r$ .

It takes only  $r(m + n - r)$  real numbers to store  $\mathbf{A} \in \mathbb{R}^{m,n}$  with  $\text{rank}(\mathbf{A}) = r$ .

Thus approximating a given matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$  with a rank- $r$  matrix,  $r \ll \min\{m, n\}$ , can be regarded as an instance of matrix compression. The approximation error with respect to some matrix norm  $\|\cdot\|$  will be minimal if we choose the **best approximation**

$$\mathbf{A}_r := \operatorname{argmin}\{\|\mathbf{A} - \mathbf{B}\| : \mathbf{B} \in \mathbb{R}^{m,n}, \text{rank}(\mathbf{B}) = r\}, \quad 1 \leq r \leq \min\{m, n\}. \quad (3.4.4.15)$$

$\square$

Here we explore low-rank best approximation of general matrices with respect to the Euclidean matrix norm  $\|\cdot\|_2$  induced by the 2-norm for vectors ( $\rightarrow$  Def. 1.5.5.10), and the Frobenius norm  $\|\cdot\|_F$ .

**Definition 3.4.4.16. Frobenius norm**

The **Frobenius norm** of  $\mathbf{A} \in \mathbb{K}^{m,n}$  is defined as

$$\|\mathbf{A}\|_F^2 := \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2.$$

It should be obvious that  $\|\mathbf{A}\|_F$  invariant under orthogonal/unitary transformations of  $\mathbf{A}$ . Thus the Frobenius norm of a matrix  $\mathbf{A}$ ,  $\text{rank}(\mathbf{A}) = r$ , can be expressed through its singular values  $\sigma_j$ :

Frobenius norm and SVD:  $\|\mathbf{A}\|_F^2 = \sum_{j=1}^r \sigma_j^2 \quad (3.4.4.17)$

☞ notation:  $\mathcal{R}_r(m, n) := \{\mathbf{A} \in \mathbb{K}^{m,n} : \text{rank}(\mathbf{A}) \leq r\}, m, n, r \in \mathbb{N}$

The next profound result links best approximation in  $\mathcal{R}_r(m, n)$  and the singular value decomposition ( $\rightarrow$  Def. 3.4.1.3).

**Theorem 3.4.4.18. Best low rank approximation** → [Gut09, Thm. 11.6]

Let  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$  be the SVD of  $\mathbf{A} \in \mathbb{K}^{m,n}$  ( $\rightarrow$  Thm. 3.4.1.1). For  $1 \leq k \leq \text{rank}(\mathbf{A})$  set

$$\mathbf{A}_k := \mathbf{U}_k \Sigma_k \mathbf{V}_k^H = \sum_{\ell=1}^k \sigma_\ell (\mathbf{U})_{:, \ell} (\mathbf{V})_{:, \ell}^H \quad \text{with} \quad \begin{aligned} \mathbf{U}_k &:= \left[ (\mathbf{U})_{:, 1}, \dots, (\mathbf{U})_{:, k} \right] \in \mathbb{K}^{m, k}, \\ \mathbf{V}_k &:= \left[ (\mathbf{V})_{:, 1}, \dots, (\mathbf{V})_{:, k} \right] \in \mathbb{K}^{n, k}, \\ \Sigma_k &:= \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbb{K}^{k, k}. \end{aligned}$$

Then, for  $\|\cdot\| = \|\cdot\|_F$  and  $\|\cdot\| = \|\cdot\|_2$ , holds true

$$\|\mathbf{A} - \mathbf{A}_k\| \leq \|\mathbf{A} - \mathbf{F}\| \quad \forall \mathbf{F} \in \mathcal{R}_k(m, n),$$

that is,  $\mathbf{A}_k$  is the rank- $k$  best approximation of  $\mathbf{A}$  in the matrix norms  $\|\cdot\|_F$  and  $\|\cdot\|_2$ .

This theorem teaches us that the rank- $k$ -matrix that is closest to  $\mathbf{A}$  (rank- $k$  best approximation) in both the Euclidean matrix norm and the Frobeniusnorm ( $\rightarrow$  Def. 3.4.4.16) can be obtained by truncating the rank-1 sum expansion (3.4.1.8) obtained from the SVD of  $\mathbf{A}$  after  $k$  terms.

*Proof.* (of Thm. 3.4.4.18) Write  $\mathbf{A}_k = \mathbf{U}_k \Sigma_k \mathbf{V}_k^H$ . Obviously, with  $r = \text{rank } \mathbf{A}$ ,

$$\text{rank } \mathbf{A}_k = k \quad \text{and} \quad \|\mathbf{A} - \mathbf{A}_k\| = \|\Sigma - \Sigma_k\| = \begin{cases} \sigma_{k+1} & , \text{for } \|\cdot\| = \|\cdot\|_2, \\ \sqrt{\sigma_{k+1}^2 + \dots + \sigma_r^2} & , \text{for } \|\cdot\| = \|\cdot\|_F. \end{cases}$$

① Pick  $\mathbf{B} \in \mathbb{K}^{n,n}$ ,  $\text{rank } \mathbf{B} = k$ .

$$\Rightarrow \dim \mathcal{N}(\mathbf{B}) = n - k \Rightarrow \mathcal{N}(\mathbf{B}) \cap \text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_{k+1}\} \neq \{0\},$$

where  $\mathbf{v}_i$ ,  $i = 1, \dots, n$  are the columns of  $\mathbf{V}$ . For  $\mathbf{x} \in \mathcal{N}(\mathbf{B}) \cap \text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_{k+1}\}$ ,  $\|\mathbf{x}\|_2 = 1$

$$\begin{aligned} \mathbf{x} &= \sum_{j=1}^{k+1} (\mathbf{v}_j^H \mathbf{x}) \mathbf{v}_j, \\ \|\mathbf{A} - \mathbf{B}\|_2^2 &\geq \|(\mathbf{A} - \mathbf{B})\mathbf{x}\|_2^2 = \|\mathbf{A}\mathbf{x}\|_2^2 = \left\| \sum_{j=1}^{k+1} \sigma_j (\mathbf{v}_j^H \mathbf{x}) \mathbf{u}_j \right\|_2^2 = \sum_{j=1}^{k+1} \sigma_j^2 (\mathbf{v}_j^H \mathbf{x})^2 \geq \sigma_{k+1}^2, \end{aligned}$$

because  $\sum_{j=1}^{k+1} (\mathbf{v}_j^H \mathbf{x})^2 = \|\mathbf{x}\|_2^2 = 1$ .

② Find ONB  $\{\mathbf{z}_1, \dots, \mathbf{z}_{n-k}\}$  of  $\mathcal{N}(\mathbf{B})$  and assemble it into a matrix  $\mathbf{Z} = [\mathbf{z}_1 \dots \mathbf{z}_{n-k}] \in \mathbb{K}^{n, n-k}$

$$\|\mathbf{A} - \mathbf{B}\|_F^2 \geq \|(\mathbf{A} - \mathbf{B})\mathbf{Z}\|_F^2 = \|\mathbf{A}\mathbf{Z}\|_F^2 = \sum_{i=1}^{n-k} \|\mathbf{A}\mathbf{z}_i\|_2^2 = \sum_{i=1}^{n-k} \sum_{j=1}^r \sigma_j^2 (\mathbf{v}_j^H \mathbf{z}_i)^2$$

□

□

The following code computes the low-rank best approximation of a dense matrix in EIGEN.

**C++ code 3.4.4.19: SVD-based low-rank matrix compression, → GITLAB**

```
2 MatrixXd lowrankbestapprox(const Eigen::MatrixXd &A, unsigned int k) {
3 // Compute economical SVD, compare Code 3.4.2.1
```

```

4   const Eigen::JacobiSVD<MatrixXd> svd(A, Eigen::ComputeThinU |
5                                             Eigen::ComputeThinV);
6   // Form matrix product  $\mathbf{U}_k \Sigma_k \mathbf{V}^H$ .
7   // Extract  $\Sigma_k$  as diagonal matrix of largest  $k$  singular values.
8   // EIGEN provides singular values in decreasing order!
9   return (svd.matrixU().leftCols(k)) *
10      (svd.singularValues().head(k).asDiagonal()) *
11      (svd.matrixV().leftCols(k).transpose());
12 }

```

**§3.4.4.20 (Error of low-rank best approximation of a matrix)** Since the matrix norms  $\|\cdot\|_2$  and  $\|\cdot\|_F$  are invariant under multiplication with orthogonal (unitary) matrices, we immediately obtain expressions for the norms of the best approximation error:

$$\text{In Euclidean matrix norm: } \left\| \mathbf{A} - \mathbf{U}_k \Sigma_k \mathbf{V}_k^H \right\|_2 = \sigma_{k+1}, \quad (3.4.4.21)$$

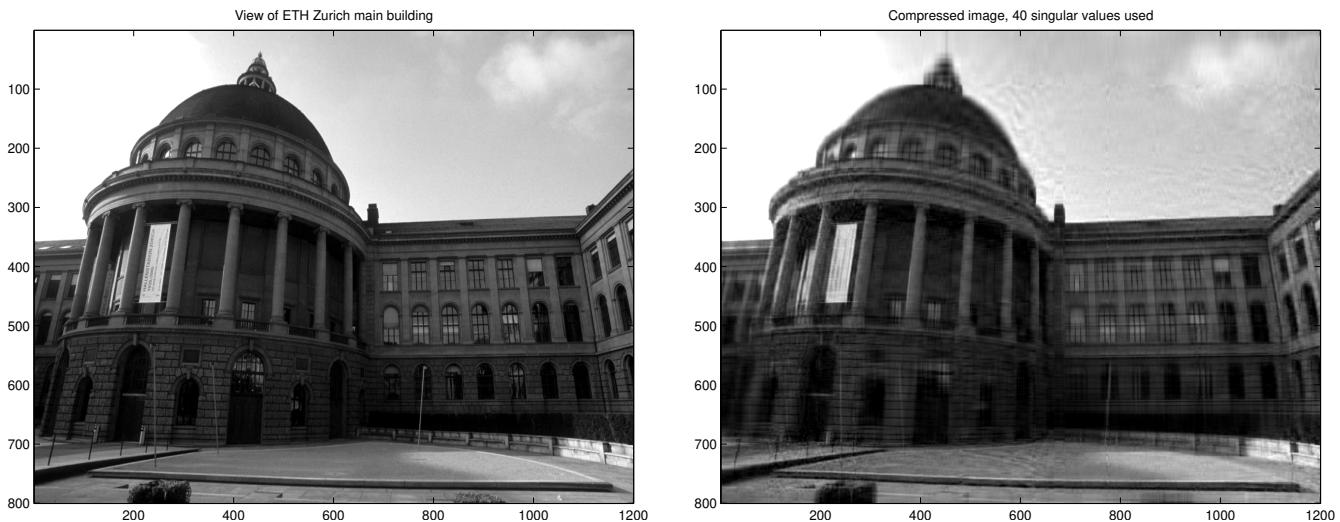
$$\text{in Frobenius norm: } \left\| \mathbf{A} - \mathbf{U}_k \Sigma_k \mathbf{V}_k^H \right\|_F^2 = \sum_{j=k+1}^{\min\{m,n\}} \sigma_j^2. \quad (3.4.4.22)$$

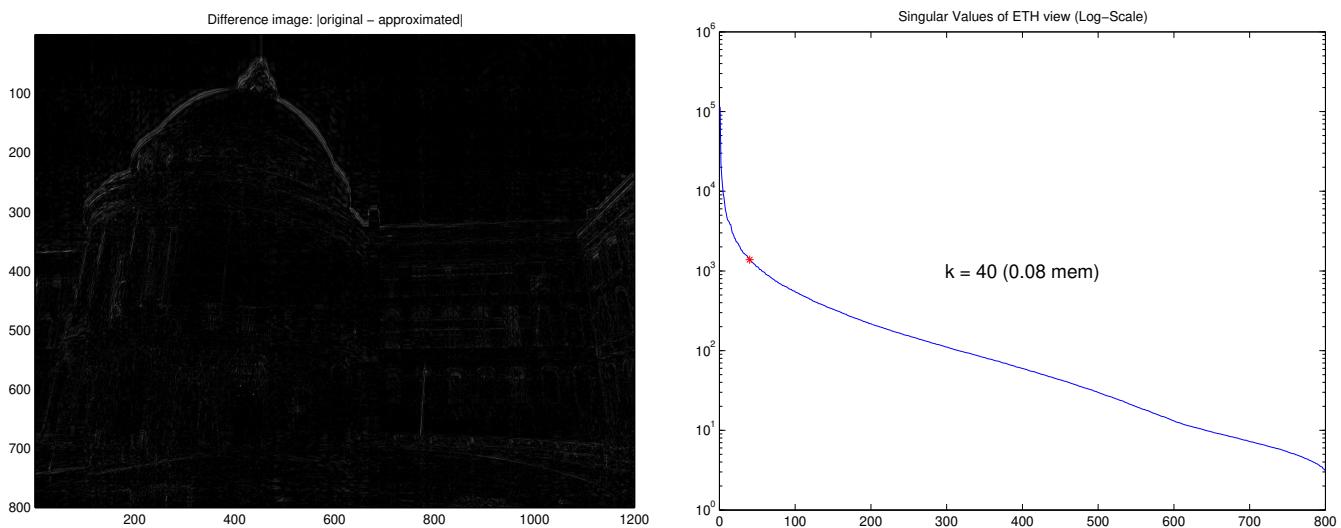
This provides precise information about the best approximation error for rank- $k$  matrices. In particular, the decay of the singular values of the matrix governs the convergence of the rank- $k$  best approximation error as  $k$  increases.  $\square$

**EXAMPLE 3.4.4.23 (Image compression)** A rectangular greyscale image composed of  $m \times n$  pixels (greyscale, BMP format) can be regarded as a matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $(\mathbf{A})_{i,j} \in \{0, \dots, 255\}$ , cf. Ex. 9.3.2.1. Thus low-rank approximation of the image matrix is a way to compress the image.

► Thm. 3.4.4.18  $\Rightarrow$  best rank- $k$  approximation of image:  $\tilde{\mathbf{A}} = \mathbf{U}_k \Sigma_k \mathbf{V}^T$

Of course, the matrices  $\mathbf{U}_k$ ,  $\mathbf{V}_k$ , and  $\Sigma_k$  are available from the economical (thin) SVD (3.4.1.4) of  $\mathbf{A}$ .

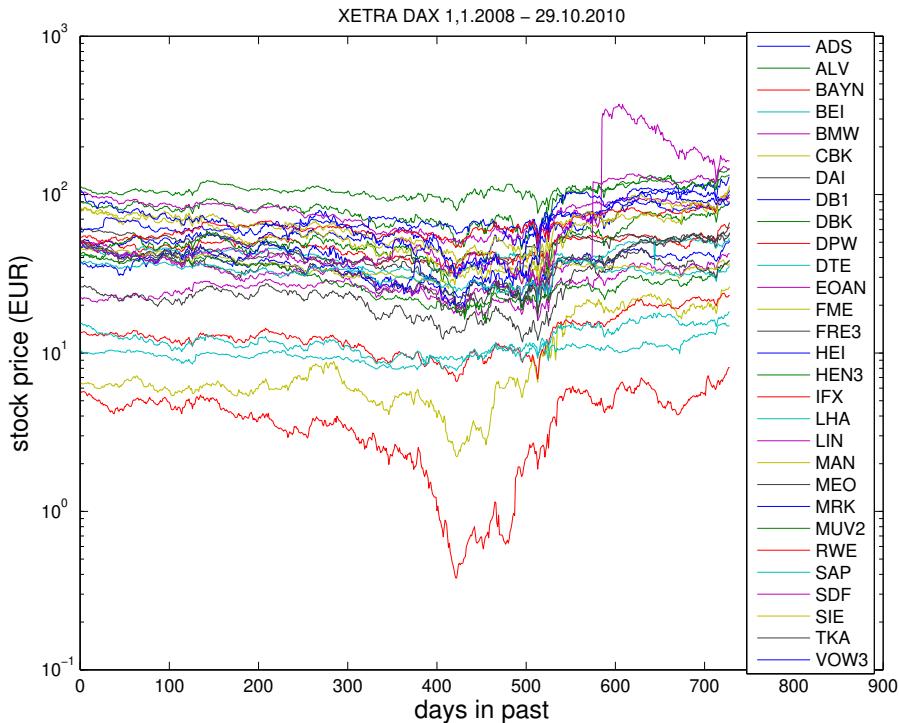




Note that there are better and faster ways to compress images than SVD (JPEG, Wavelets, etc.)

### 3.4.4.3 Principal Component Data Analysis (PCA)

**EXAMPLE 3.4.4.24 (Trend analysis)** The objective is to extract information in the form of hidden “trends” from data.



We are given time series data:

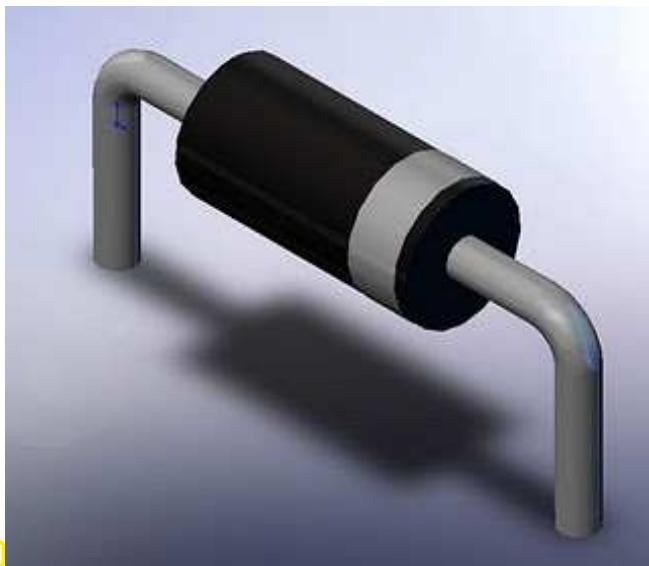
◁ (end of day) stock prizes  
 $\hat{=} n$  data vectors  $\in \mathbb{R}^m$

Rephrased in the language of linear algebra:

Are there underlying governing trends ?

Are there a few vectors  $\mathbf{u}_1, \dots, \mathbf{u}_p$ ,  $p \ll n$ , such that, *approximately* all other data vectors  $\in \text{Span}\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ .

**EXAMPLE 3.4.4.25 (Classification from measured data)** Data vectors belong to different classes, where those in the same class are “qualitatively similar” in the sense that they are small (random) perturbations of a typical data vector. The task is to tease out the typical data patterns and tell which class every data vector belongs to.



Given: measured  $U$ - $I$  characteristics of  $n$  diodes in a box  
 (data points  $(U_j, I_j^{(k)})$ ,  $j = 1, \dots, m$ ,  $k = 1, \dots, n$ )

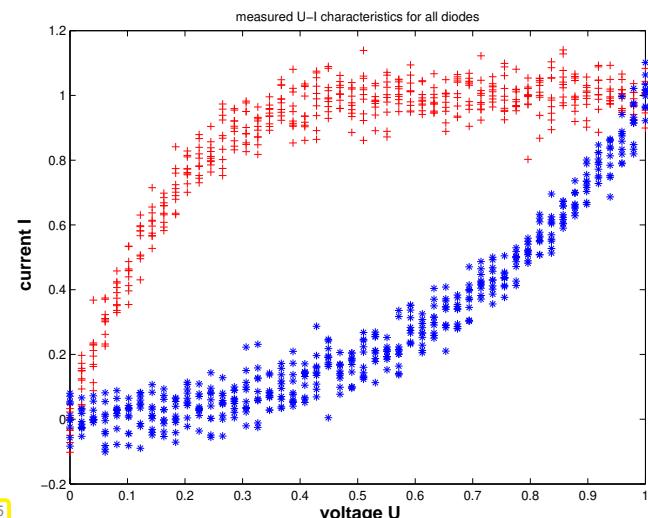
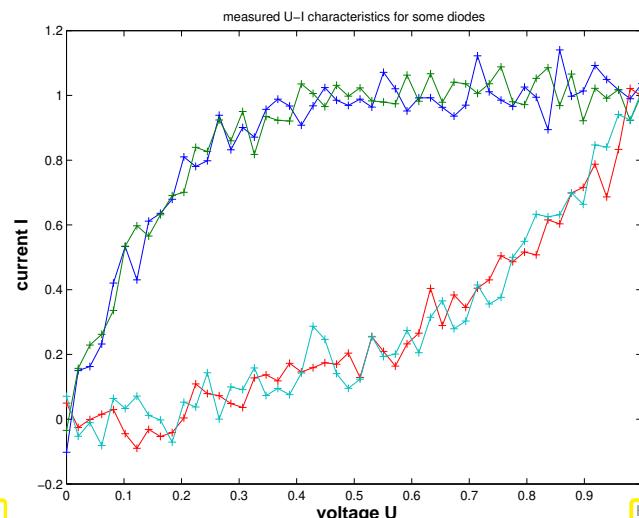
Classification problem: find out

- how many different types of diodes in box,
- the  $U$ - $I$  characteristic of each type.



Measurement errors !  
 Manufacturing tolerances !

The following plots display possible (“synthetic”) measured data for two types of diodes; measurement errors and manufacturing tolerances taken into account by additive (Gaussian) random perturbations (noise).



Ex. 3.4.4.24 and Ex. 3.4.4.25 present typical tasks that can be tackled by **principal component analysis**. Now we give an abstract description as a problem of linear algebra.

Given:  $n$  data points  $\mathbf{a}_j \in \mathbb{R}^m$ ,  $j = 1, \dots, n$ , in  $m$ -dimensional (feature) space  
 (e.g.,  $\mathbf{a}_j$  may represent a finite *time series* or a *measured relationship* of physical quantities)

In Ex. 3.4.4.24:  $n \hat{=} \text{number of stocks}$ ,  
 $m \hat{=} \text{number of days, for which stock prices are recorded}$

- ♦ Extreme case: all stocks follow exactly *one* trend

$$\leftrightarrow \quad \mathbf{a}_j \in \text{Span}\{\mathbf{u}\} \quad \forall j = 1, \dots, n,$$

for a **trend vector**  $\mathbf{u} \in \mathbb{R}^m$ ,  $\|\mathbf{u}\|_2 = 1$ .

- ♦ Unlikely case: all stocks prices are governed by  $p < n$  trends:

$$\leftrightarrow \quad \mathbf{a}_j \in \text{Span}\{\mathbf{u}_1, \dots, \mathbf{u}_p\} \quad \forall j = 1, \dots, m, \quad (3.4.4.26)$$

with **orthonormal trend vectors**  $\mathbf{u}_i \in \mathbb{R}^m$ ,  $i = 1, \dots, p$ .

Why unlikely ? Small *random fluctuations* will be present in each stock prize

Why orthonormal ? Trends should be as “independent as possible” (minimally correlated)

Expressed using the terminology linear algebra:

$$(3.4.4.26) \quad \Leftrightarrow \quad \begin{aligned} \text{rank}(\mathbf{A}) = p \text{ for } \mathbf{A} := [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{R}^{m,n}, \\ \mathcal{R}(\mathbf{A}) = \text{Span}\{\mathbf{u}_1, \dots, \mathbf{u}_p\} \end{aligned} \quad (3.4.4.27)$$

- ◆ Realistic: stock prizes *approximately* follow a few trends

$$\mathbf{a}_j \in \text{Span}\{\mathbf{u}_1, \dots, \mathbf{u}_p\} + \text{“small perturbations” } \forall j = 1, \dots, m,$$

with orthonormal trend vectors  $\mathbf{u}_i, i = 1, \dots, p$ .

Task (PCA): determine (minimal)  $p$  and orthonormal trend vectors  $\mathbf{u}_i, i = 1, \dots, p$

Now **singular value decomposition** (SVD) according to Def. 3.4.1.3 comes into play, because Lemma 3.4.1.12 tells us that it can supply an orthonormal basis of the image space of a matrix, cf. Code 3.4.2.9.

Issue: how to deal with (small, random) **perturbations** ?

Recall Rem. 3.4.1.6, (3.4.1.8): If  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$  is the SVD of  $\mathbf{A} \in \mathbb{R}^{m,n}$ , then  $(\mathbf{u}_j) \hat{=} \text{columns of } \mathbf{U}$ ,  $\mathbf{v}_j \hat{=} \text{columns of } \mathbf{V}$

$$\left[ \begin{array}{c|c} & \mathbf{A} \\ \hline & \end{array} \right] = \sigma_1 \left[ \begin{array}{c|c} & \mathbf{u}_1 \\ \hline & \end{array} \right] \left[ \begin{array}{c|c} & \mathbf{v}_1^T \\ \hline & \end{array} \right] + \sigma_2 \left[ \begin{array}{c|c} & \mathbf{u}_2 \\ \hline & \end{array} \right] \left[ \begin{array}{c|c} & \mathbf{v}_2^T \\ \hline & \end{array} \right] + \dots$$

This already captures the case (3.4.4.26) and we see that the columns of  $\mathbf{U}$  supply the trend vectors we are looking for!

- ① no perturbations:

$$\text{SVD: } \mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H \text{ satisfies } \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p > \sigma_{p+1} = \dots = \sigma_{\min\{m,n\}} = 0, \\ \text{orthonormal trend vectors } (\mathbf{U})_{:,1}, \dots, (\mathbf{U})_{:,p}.$$

- ② with perturbations:

$$\text{SVD: } \mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H \text{ satisfies } \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \gg \sigma_{p+1} \approx \dots \approx \sigma_{\min\{m,n\}} \approx 0, \\ \text{orthonormal trend vectors } (\mathbf{U})_{:,1}, \dots, (\mathbf{U})_{:,p}.$$

If there is a pronounced gap in distribution of the singular values, which separates  $p$  large from  $\min\{m, n\} - p$  relatively small singular values, this hints that  $\mathcal{R}\mathbf{A}$  has essentially dimension  $p$ . It depends on the application what one accepts as a “pronounced gap”.

Frequently used criterion:

$$p = \min \left\{ q: \sum_{j=1}^q \sigma_j^2 \geq (1 - \tau) \sum_{j=1}^{\min\{m,n\}} \sigma_j^2 \right\} \text{ for } \tau \ll 1.$$

What is the Information carried by  $\mathbf{V}$  in PCA context ?

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} = \sigma_1 \begin{bmatrix} \mathbf{u}_1 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^\top \end{bmatrix} + \sigma_2 \begin{bmatrix} \mathbf{u}_2 \end{bmatrix} \begin{bmatrix} \mathbf{v}_2^\top \end{bmatrix} + \dots$$

$j$ -th data set ( $\leftrightarrow$  time series # $j$ ) in  $j$ -th column of  $\mathbf{A}$

$$(3.4.1.8) \Rightarrow (\mathbf{A})_{:,j} = \sigma_1 \mathbf{u}_1(\mathbf{v}_1)_j + \sigma_2 \mathbf{u}_2(\mathbf{v}_2)_j + \dots$$

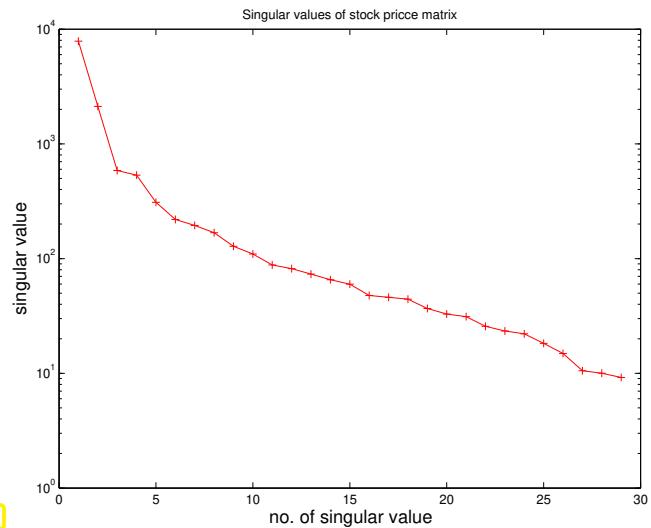
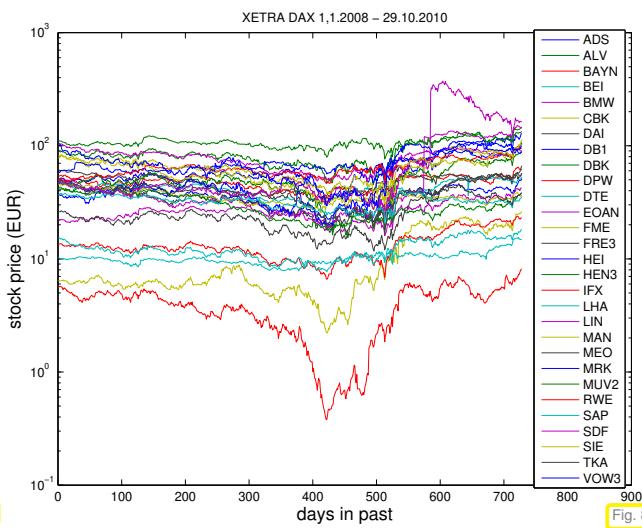
► The  $j$ -th row of  $\mathbf{V}$  (up to the  $p$ -th component) gives the **weights** with which the  $p$  identified trends contribute to data set  $j$ .

**EXAMPLE 3.4.4.28 (PCA of stock prices → Ex. 3.4.4.24cnt'd)** Stock prices are given as a large matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ :

columns of  $\mathbf{A}$  → time series of end of day stock prices of individual stocks  
 rows of  $\mathbf{A}$  → closing prices of DAX stocks on a particular day

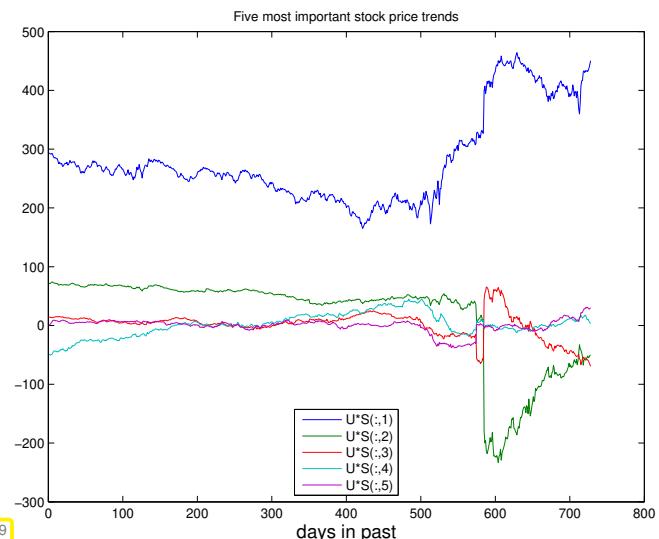
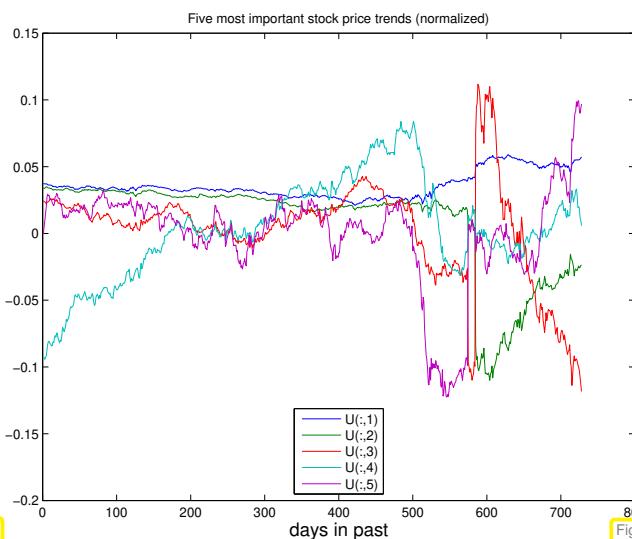
The data were obtained from Yahoo Finance in 2016:

```
#!/bin/csh
foreach i (ADS ALV BAYN BEI BMW CBK DAI DBK DB1 LHA DPW DTE EOAN FRE3 \
           FME HEI HEN3 IFX SDF LIN MAN MRK MEO MUV2 RWE SAP SIE TKA VOW3)
  wget -O "i".csv "http://ichart.finance.yahoo.com/table.csv?s=i.DE&a=00&b=1&
  c=2008&d=09&e=30&f=2010&g=d&ignore=.csv"
  sed -i -e 's/-/_/g' "i".csv
end
```



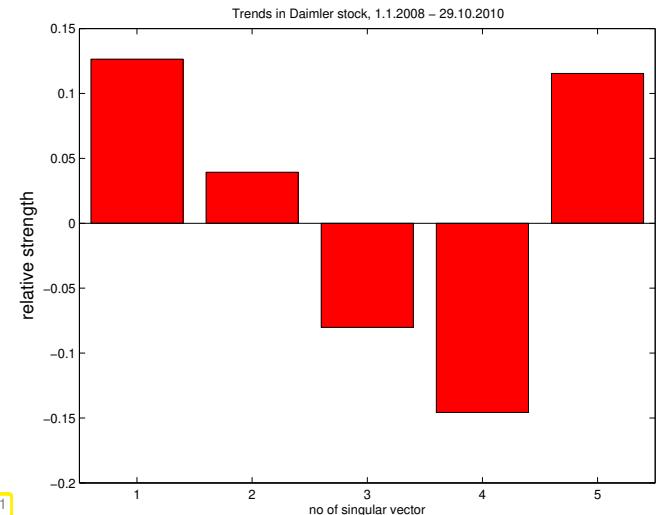
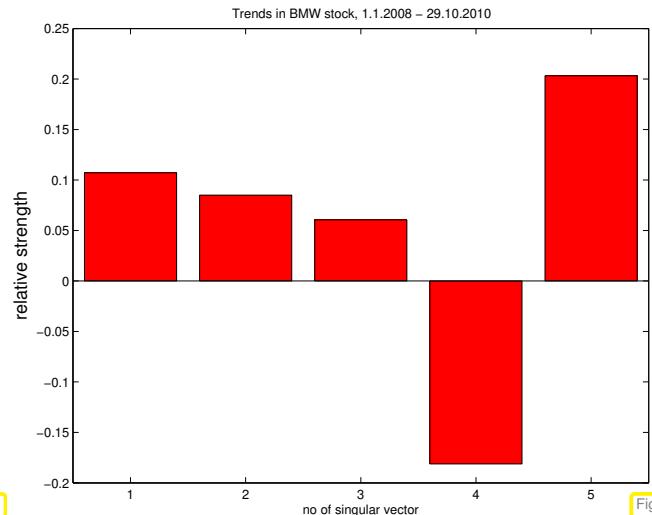
We observe a pronounced decay of the singular values of  $\mathbf{A}$ . The plot of Fig. 87 is given in linear-logarithmic scale. The neat alignment of larger singular values indicates approximate **exponential decay** of the singular values.

➤ a few trends (corresponding to a few of the largest singular values) govern the time series.



Columns of  $\mathbf{U}$  ( $\rightarrow$  Fig. 88) in SVD  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$  provide trend vectors, cf. Ex. 3.4.4.24 & Ex. 3.4.4.29.

When weighted with the corresponding singular value, the importance of a trend contribution emerges, see Fig. 89

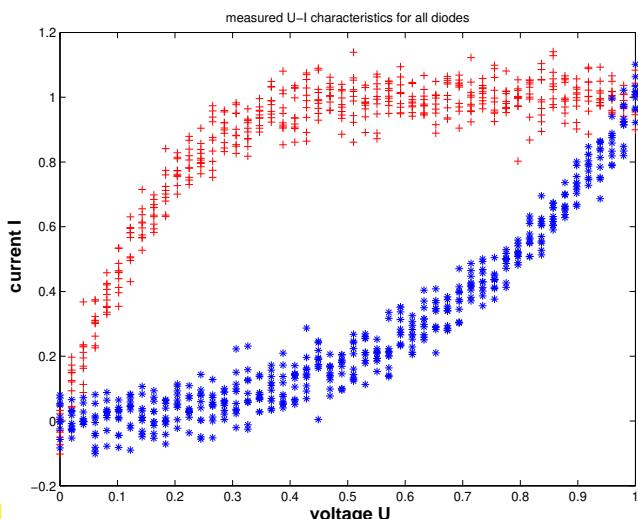
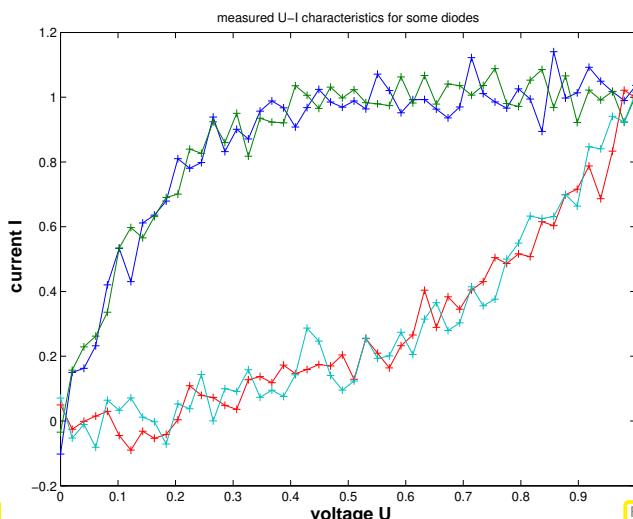


Stocks of companies from the same sector of the economy should display similar contributions of major trend vectors, because their prices can be expected to be more closely correlated than stock prices in general. This is evident in 90 and Fig. 91 for two car makers.

### EXAMPLE 3.4.4.29 (Principal component analysis for data classification $\rightarrow$ Ex. 3.4.4.25 cnt'd)

Given: measured  $U$ - $I$  characteristics of  $n = 20$  unknown diodes,  $I(U)$  available for  $m = 50$  voltages.

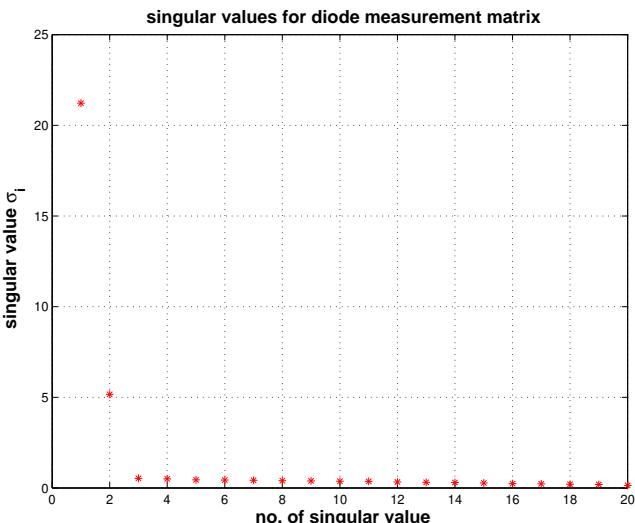
Sought: Number of different types of diodes in batch and reconstructed  $U$ - $I$  characteristic for each type.



Data matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \gg n$ :

Rows  $\mathbf{A}$  → series of measurements for different diodes (times/locations etc.),  
Columns of  $\mathbf{A}$  → measured values corresponding to one diode (time/location etc.).

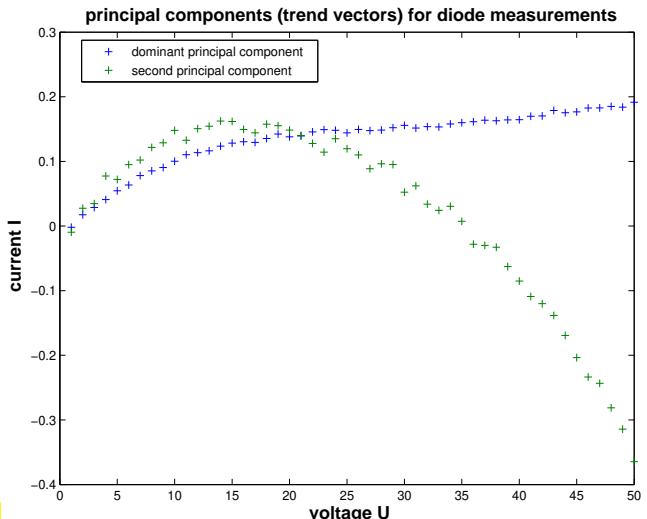
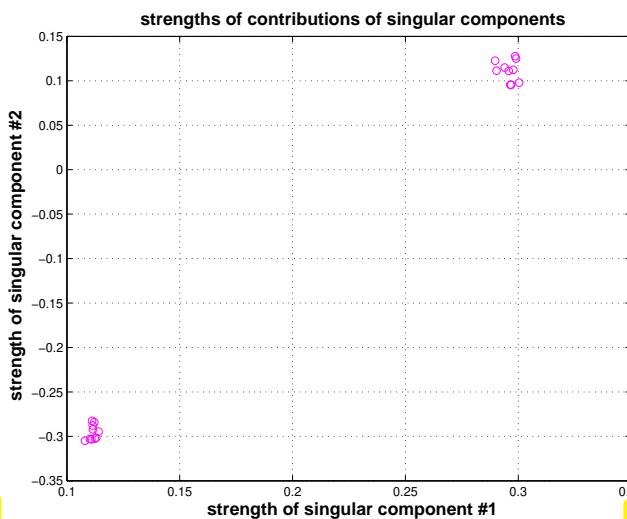
Goal of PCA: detect linear correlations between columns of  $\mathbf{A}$



← distribution of singular values of matrix  
two dominant singular values !

measurements display linear correlation with **two principal components**

**two types of diodes in batch**



Observations:

- ◆ First two rows of  $\mathbf{V}$ -matrix specify strength of contribution of the two leading principal components to each measurement

- Points  $(\mathbf{V})_{:,1:2}$ , which correspond to different diodes are neatly clustered in  $\mathbb{R}^2$ . To determine the type of diode  $i$ , we have to identify the cluster to which the point  $((\mathbf{V})_{i,1}, \mathbf{V}_{i,2})$  belongs ( $\rightarrow$  **cluster analysis**, course “machine learning”, see Rem. 3.4.4.40 below).
- ◆ The principal components themselves do not carry much useful information in this example.

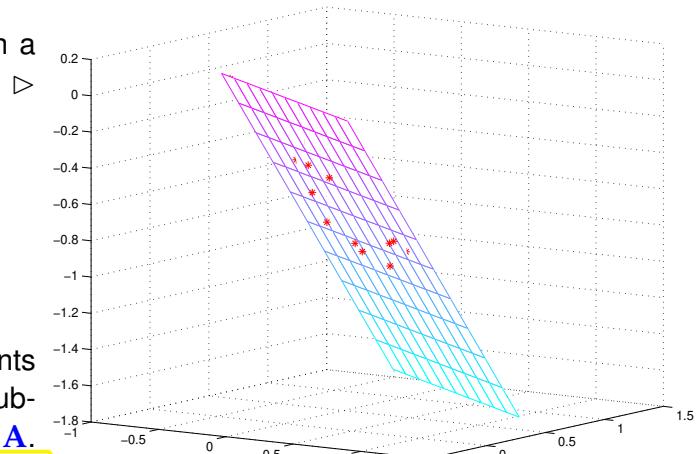
**EXAMPLE 3.4.4.30 (Data points confined to a subspace)** More abstractly, above we tried to identify a subspace to which all data points  $\mathbf{a}_i$  were “close”. We saw that the SVD of

Data points  $\bullet \leftrightarrow \mathbf{a}_i \in \mathbb{R}^3$  “almost” located on a plane

Non-zero singular values of  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ :

$$\begin{array}{r} 3.1378 \\ 1.8092 \\ \hline 0.1792 \end{array}$$

The third singular value is much smaller, which hints that the data points approximately lie in a 2D subspace spanned by the two first singular vectors of  $\mathbf{A}$ .



**§3.4.4.31 (Proper orthogonal decomposition (POD))** In the previous Ex. 3.4.4.30 we saw that the singular values of a matrix whose columns represent data points  $\in \mathbb{R}^m$  tell us whether these points are all “approximately located” in a lower-dimensional subspace  $V \subset \mathbb{R}^m$ . This is linked to the following problem:

Problem of proper orthogonal decomposition (POD):

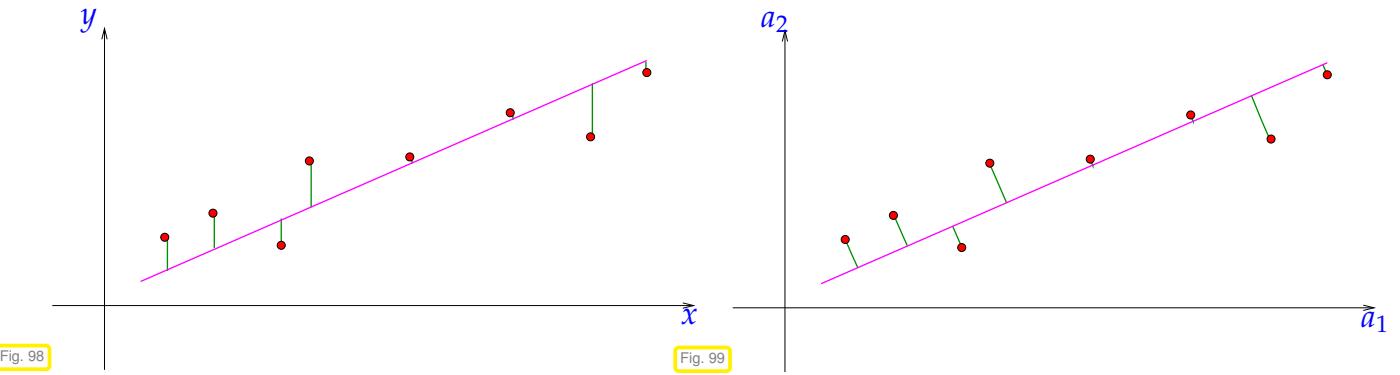
Given: Data points  $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^m$ ,  $m, n \in \mathbb{N}$

Sought: For  $k \leq \min\{m, n\}$ , find a **subspace**  $U_k \subset \mathbb{R}^m$  such that

$$U_k = \underset{W \subset \mathbb{R}^m, \dim W=k}{\operatorname{argmin}} \sum_{j=1}^n \inf_{w \in W} \|\mathbf{a}_j - w\|_2^2, \quad (3.4.4.32)$$

that is, we seek that  $k$ -dimensional subspace  $U_k$  of  $\mathbb{R}^m$  for which the sum of squared distances of the data points to  $U_k$  is minimal.

We have already seen a similar problem in Ex. 3.4.4.5. For  $m = 2$  we want to point out the difference to linear regression:



Linear regression: Minimize the sum of squares of *vertical distances*.

POD: Minimize the sum of squares of (*minimal*) *distances*.

By finding a  $k$ -dimensional subspace we mean finding a, preferably orthonormal, basis of that subspace. Let us assume that  $\{\mathbf{w}_1, \dots, \mathbf{w}_k\}$  is an orthonormal basis (ONB) of a  $k$ -dimensional subspace  $W \subset \mathbb{R}^m$ . Then the **orthogonal projection**  $P_W \mathbf{x}$  of a point  $\mathbf{x} \in \mathbb{R}^m$  onto  $W$  is given by

$$P_W \mathbf{x} = \sum_{j=1}^k (\mathbf{w}_j^\top \mathbf{x}) \mathbf{w}_j = \mathbf{W} \mathbf{W}^\top \mathbf{x}, \quad (3.4.4.33)$$

where  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k] \in \mathbb{R}^{m,k}$ . This formula is closely related to the normal equations for a linear least squares problem, see Thm. 3.1.2.1 and § 3.1.1.6 for a visualization of an orthogonal projection.  $\|\mathbf{x} - P_W \mathbf{x}\|_2$  is the (minimal) distance of  $\mathbf{x}$  to  $W$ .

Hence, again writing  $\mathbf{W} \in \mathbb{R}^{m,k}$  for the matrix whose columns form an ONB ( $\Rightarrow \mathbf{W}^\top \mathbf{W} = \mathbf{I}$ ) of  $W \subset \mathbb{R}^m$ , we have

$$\sum_{j=1}^n \inf_{w \in W} \|a_j - w\|_2^2 = \sum_{j=1}^n \|a_j - \mathbf{W} \mathbf{W}^\top a_j\|_2^2 = \|\mathbf{A} - \mathbf{W} \mathbf{W}^\top \mathbf{A}\|_F^2, \quad (3.4.4.34)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm of a matrix, see Def. 3.4.4.16, and  $\mathbf{A} = [a_1, \dots, a_n] \in \mathbb{R}^{m,n}$ .

Note that  $\text{rank}(\mathbf{W} \mathbf{W}^\top \mathbf{A}) \leq k$ . Let us write  $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top$  for the SVD of  $\mathbf{A}$ , and  $\mathbf{U}_k \in \mathbb{R}^{m,k}$ ,  $\mathbf{V}_k \in \mathbb{R}^{n,k}$ , and  $\Sigma_k \in \mathbb{R}^{k,k}$  for the truncated SVD-factors of  $\mathbf{A}$  as introduced in Thm. 3.4.4.18. Then the rank- $k$  best approximation result of that theorem implies

$$\|\mathbf{A} - \mathbf{W} \mathbf{W}^\top \mathbf{A}\|_F \geq \|\mathbf{A} - \mathbf{U}_k \Sigma_k \mathbf{V}_k\|_F \quad \forall \mathbf{W} \in \mathbb{R}^{m,k}, \quad \mathbf{W}^\top \mathbf{W} = \mathbf{I}.$$

In fact, we can find  $\mathbf{W} \in \mathbb{R}^{m,k}$  with orthonormal columns that realizes the minimum: just choose  $\mathbf{W} := \mathbf{U}_k$  and verify

$$\mathbf{W} \mathbf{W}^\top \mathbf{A} = \mathbf{U}_k \mathbf{U}_k^\top \mathbf{U} \Sigma \mathbf{V}^\top = \mathbf{U}_k [\mathbf{I}_k \quad \mathbf{O}] \Sigma \mathbf{V}^\top = \mathbf{U}_k \Sigma_k \mathbf{V}_k^\top.$$

### Theorem 3.4.4.35. Solution of POD problem

The subspace  $U_k$  spanned by the first  $k$  left singular vectors of  $\mathbf{A} = [a_1, \dots, a_n] \in \mathbb{R}^{m,n}$  solves the **POD problem** (3.4.4.32):

$$U_k = \mathcal{R}((\mathbf{U})_{:,1:k}) \quad \text{with} \quad \mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top \quad \text{the SVD of } \mathbf{A}.$$

Appealing to (3.4.4.22), the sum of the squared distances can be obtained as the sum of the squares of the remaining singular values  $\sigma_{k+1}, \dots, \sigma_p$ ,  $p := \min\{m, n\}$ , of  $\mathbf{A}$ :

$$\sum_{j=1}^n \inf_{w \in U_k} \|a_j - w\|_2^2 = \sum_{\ell=k+1}^p \sigma_\ell^2. \quad (3.4.4.36)$$

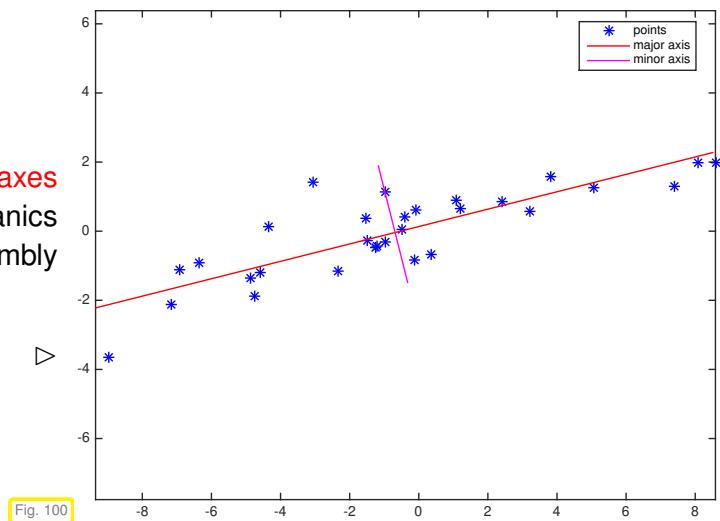
As a consequence, the decay of the singular values again predicts how close the data points are to the POD subspaces  $U_k, k = 1, \dots, p - 1$ .  $\square$

**EXAMPLE 3.4.4.37 (Principal axis of a point cloud)** Given  $m > 2$  points  $\mathbf{x}_j \in \mathbb{R}^k, j = 1, \dots, m$ , in  $k$ -dimensional space, we ask what is the “longest” and “shortest” diameter  $\mathbf{d}_+$  and  $\mathbf{d}_-$ . This question can be stated rigorously in several different ways: here we ask for directions for which the point cloud will have maximal/minimal variance, when projected onto that direction:

$$\begin{aligned} \mathbf{d}_+ &:= \underset{\|\mathbf{v}\|=1}{\operatorname{argmax}} Q(\mathbf{v}), \\ \mathbf{d}_- &:= \underset{\|\mathbf{v}\|=1}{\operatorname{argmin}} Q(\mathbf{v}), \quad Q(\mathbf{v}) := \sum_{j=1}^m |(\mathbf{x}_j - \mathbf{c})^\top \mathbf{v}|^2, \quad \mathbf{c} = \frac{1}{m} \sum_{j=1}^m \mathbf{x}_j. \end{aligned} \quad (3.4.4.38)$$

The directions  $\mathbf{d}_+, \mathbf{d}_-$  are called the **principal axes** of the point cloud, a term borrowed from mechanics and connected with the axes of inertia of an assembly of point masses.

Principal axes of a point cloud in 2D



$\mathbf{d}_+, \mathbf{d}_-$  can be computed by computing the extremizers of  $\mathbf{x} \mapsto \|\mathbf{Ax}\|_2$  with

$$\mathbf{A} = \begin{bmatrix} (\mathbf{x}_1 - \mathbf{c})^\top \\ \vdots \\ (\mathbf{x}^m - \mathbf{c})^\top \end{bmatrix} \in \mathbb{R}^{m,k} \Rightarrow \begin{aligned} \mathbf{d}_+ &= \underset{\|\mathbf{v}\|=1}{\operatorname{argmax}} \|\mathbf{Ax}\|_2, \\ \mathbf{d}_- &= \underset{\|\mathbf{v}\|=1}{\operatorname{argmin}} \|\mathbf{Ax}\|_2, \end{aligned} \quad (3.4.4.39)$$

on  $\{\mathbf{x} \in \mathbb{R}^k : \|\mathbf{x}\|_2 = 1\}$ , using the SVD-based method presented in Section 3.4.4.1.  $\square$

**Remark 3.4.4.40 (Algorithm for cluster analysis)** In data classification as presented in Ex. 3.4.4.29, after we have identified the  $p$  main trends ( $\leftrightarrow$  singular values and left singular vectors) and how much they contribute to every data point ( $\leftrightarrow$  right singular vectors), the last step is to perform a cluster analysis based on the right singular vectors  $\mathbf{v}_1, \dots, \mathbf{v}_p$ .

Now we study the abstract problem of cluster analysis:

Given:  $\blacklozenge N$  data points  $\mathbf{x}_i \in \mathbb{R}^k, i = 1, \dots, N$ ,

$\blacklozenge$  Assume: number  $n$  of desired clusters is known in advance.

Sought: Partitioning of index set  $\{1, \dots, N\} = I_1 \cup \dots \cup I_n$ , achieving minimal **mean least squares error**

$$\text{mlse} := \sum_{l=1}^n \sum_{i \in I_l} \|\mathbf{x}_i - \mathbf{m}_l\|_2^2, \quad \mathbf{m}_l = \frac{1}{\#I_l} \sum_{i \in I_l} \mathbf{x}_i. \quad (3.4.4.41)$$

The subsets  $\{\mathbf{x}_i : i \in I_l\}$  are called the **clusters**. The points  $\mathbf{m}_l$  are their **centers of gravity**.

The Algorithm involves two components:

- ❶ Splitting of a cluster by separation along its **principal axis**, see Ex. 3.4.4.37 and Code 3.4.4.45:

$$\mathbf{a}_l := \underset{\|\mathbf{v}\|_2=1}{\operatorname{argmax}} \left\{ \sum_{i \in I_l} |(\mathbf{x}_i - \mathbf{m}_l)^\top \mathbf{v}|^2 \right\}. \quad (3.4.4.42)$$

Relies on the algorithm from Ex. 3.4.4.37, see Code 3.4.4.45.

- ❷ Improvement of clusters using the **Lloyd-Max algorithm**, see Code 3.4.4.46. It involves two steps in turns:

- (a) Given centers of gravity  $\mathbf{m}_l$  redistribute points according to

$$I_l := \{i \in \{1, \dots, N\} : \|\mathbf{x}_i - \mathbf{m}_l\|_2 \leq \|\mathbf{x}_i - \mathbf{m}_k\|_2 \forall k \neq l\}, \quad (3.4.4.43)$$

that is, we assign each point to the nearest center of gravity, see Code 3.4.4.46.

- (b) Recompute centers of gravity

$$\mathbf{m}_l = \frac{1}{\#I_l} \sum_{i \in I_l} \mathbf{x}_i. \quad (3.4.4.44)$$

We start with a single cluster, and then do repeated splitting (❶) and cluster rearrangement (❷) until we have reached the desired final number  $n$  of clusters, see Code 3.4.4.47.

#### C++-code 3.4.4.45: Principal axis point set separation

```

2 // Separation of a set of points whose coordinates are stored in the
3 // columns of X according to their location w.r.t. the principal axis
4 std::pair<VectorXi, VectorXi> princaxissep(const MatrixXd & X){
5     int N = X.cols(); // no. of points
6     VectorXd g = X.rowwise().sum() / N; // Center of gravity, cf. (3.4.4.44)
7     MatrixXd Y = X - g.replicate(1,N); // Normalize point coordinates.
8     // Compute principal axes, cf. (3.4.4.42) and (3.4.4.43). Note that the
9     // SVD of a symmetric matrix is available through an orthonormal
10    // basis of eigenvectors.
11    SelfAdjointEigenSolver<MatrixXd> es(Y*Y.transpose());
12    // Major principal axis
13    Eigen::VectorXd a = es.eigenvectors().rightCols<1>();
14    // Coordinates of points w.r.t. to major principal axis
15    Eigen::VectorXd c = a.transpose()*Y;
16    // Split point set according to locations of projections on principal
17    // axis
18    // std::vector with indices to prevent resizing of matrices
19    std::vector<int> i1, i2;
20    for(int i = 0; i < c.size(); ++i){
21        if(c(i) >= 0)
22            i1.push_back(i);
23        else
24            i2.push_back(i);
25    }
26    // return the mapped std::vector as Eigen::VectorXd
27    return std::pair<VectorXi, VectorXi>(VectorXi::Map(i1.data(), i1.size()),
28                                         VectorXi::Map(i2.data(), i2.size()));
}
```

#### C++-code 3.4.4.46: Lloyd-Max algorithm for cluster identification

```

2 template <class Derived>
3 std::tuple<double, VectorXi, VectorXd> distcomp(const MatrixXd & X, const
```

```

    MatrixBase<Derived> & C){
4   // Compute squared distances
5   // d.row(j) = squared distances from all points in X to cluster j
6   MatrixXd d(C.cols(), X.cols());
7   for(int j = 0; j < C.cols(); ++j){
8     MatrixXd Dv = X - C.col(j).replicate(1, X.cols());
9     d.row(j) = Dv.array().square().colwise().sum();
10 }
11 // Compute minimum distance point association and sum of minimal
12 // squared distances
12 VectorXi idx(d.cols()); VectorXd mx(d.cols());
13 for(int j = 0; j < d.cols(); ++j){
14   // mx(j) tells the minimal squared distance of point j to the
15   // nearest cluster
16   // idx(j) tells to which cluster point j belongs
17   mx(j) = d.col(j).minCoeff(&idx(j));
18 }
19 double sumd = mx.sum(); // sum of all squared distances
20 // Computer sum of squared distances within each cluster
21 VectorXd cds(C.cols()); cds.setZero();
22 for(int j = 0; j < idx.size(); ++j) // loop over all points
23   cds(idx(j)) += mx(j);
24 return std::make_tuple(sumd, idx, cds);
25 }

26 // Lloyd-Max iterative vector quantization algorithm for discrete point
27 // sets; the columns of X contain the points  $x_i$ , the columns of
28 // C initial approximations for the centers of the clusters. The final
29 // centers are returned in C, the index vector idx specifies
30 // the association of points with centers.
31 template <class Derived>
32 void lloydmax(const MatrixXd & X, MatrixBase<Derived> & C, VectorXi & idx, VectorXd
33   & cds, const double tol = 0.0001){
34   int k = X.rows(); // dimension of space
35   int N = X.cols(); // no. of points
36   int n = C.cols(); // no. of clusters
37   if(k != C.rows())
38     throw std::logic_error("dimension mismatch");
39   double sd_old = std::numeric_limits<double>::max();
40   double sd;
41   std::tie(sd, idx, cds) = distcomp(X,C);
42   // Terminate, if sum of squared minimal distances has not changed much
43   while( (sd_old-sd)/sd > tol ){
44     // Compute new centers of gravity according to (3.4.4.44)
45     MatrixXd Ctmp(C.rows(),C.cols()); Ctmp.setZero();
46     // number of points in cluster for normalization
47     VectorXi nj(n); nj.setZero();
48     for(int j = 0; j < N; ++j){ // loop over all points
49       Ctmp.col(idx(j)) += X.col(j);
50       ++nj(idx(j)); // count associated points for normalization
51     }
52     for(int i = 0; i < Ctmp.cols(); ++i){
53       if(nj(i) > 0)
54         C.col(i) = Ctmp.col(i)/nj(i); // normalization
55     }
56     sd_old = sd;
57     // Get new minimum association of the points to cluster points
58     // for next iteration
59     std::tie(sd, idx, cds) = distcomp(X,C);
60   }
}

```

```

61 // Note: this function is needed to allow a call with an rvalue
62 // && stands for an rvalue reference and allows rvalue arguments
63 // such as C.leftCols(nc) to be passed by reference (C++ 11 feature)
64 template <class Derived>
65 void lloydmax(const MatrixXd & X, MatrixBase<Derived> && C, VectorXi & idx, VectorXd
66   & cds, const double tol = 0.0001){
67   lloydmax(X, C, idx, cds);
68 }
```

### C++-code 3.4.4.47: Clustering of point set

```

2 // n-quantization of point set in k-dimensional space based on
3 // minimizing the mean square error of Euclidean distances. The
4 // columns of the matrix X contain the point coordinates, n specifies
5 // the desired number of clusters.
6 std::pair<MatrixXd, VectorXi> pointcluster(const MatrixXd & X, const int n){
7   int N = X.cols(); // no. of points
8   int k = X.rows(); // dimension of space
9   // Start with two clusters obtained by principal axis separation
10  int nc = 1; // Current number of clusters
11  // Initial single cluster encompassing all points
12  VectorXi Ibig = VectorXi::LinSpaced(N,0,N-1);
13  int nbig = 0; // Index of largest cluster
14  MatrixXd C(X.rows(), n); // matrix for cluster midpoints
15  C.col(0) = X.rowwise().sum()/N; // center of gravity
16  VectorXi idx(N); idx.setOnes();
17  // Split largest cluster into two using the principal axis separation
18  // algorithm
19  while(nc < n){
20    VectorXi i1, i2;
21    MatrixXd Xbig(k, Ibig.size());
22    for(int i = 0; i < Ibig.size(); ++i) // slicing
23      Xbig.col(i) = X.col(Ibig(i));
24    // separate Xbig into two clusters, i1 and i2 are index vectors
25    std::tie(i1, i2) = princaxissep(Xbig);
26    // new cluster centers of gravity
27    VectorXd c1(k), c2(k); c1.setZero(); c2.setZero();
28    for(int i = 0; i < i1.size(); ++i)
29      c1 += X.col(Ibig(i1(i)));
30    for(int i = 0; i < i2.size(); ++i)
31      c2 += X.col(Ibig(i2(i)));
32    c1 /= i1.size(); c2 /= i2.size(); // normalization
33    C.col(nbig) = c1;
34    C.col(nbig+1) = c2;
35    ++nc; // Increase number of clusters
36    // Improve clusters by Lloyd-Max iteration
37    VectorXd cds; // saves mean square error of clusters
38    // Note C.leftCols(nc) is passed as rvalue reference (C++ 11)
39    lloydmax(X,C.leftCols(nc),idx,cds);
40    // Identify cluster with biggest contribution to mean square error
41    cds.maxCoeff(&nbig);
42    int counter = 0;
43    // update Ibig with indices of points in cluster with biggest
44    // contribution
45    for(int i = 0; i < idx.size(); ++i){
46      if(idx(i) == nbigr)
47        Ibig(counter) = i;
48        ++counter;
49    }
50 }
```

```

49     }
50     lbig.conservativeResize(counter);
51   }
52   return std::make_pair(C, idx);
53 }
```

### 3.5 Total Least Squares

In the examples of Section 3.0.1 we generally considered overdetermined linear systems of equations  $\mathbf{Ax} = \mathbf{b}$ , for which only the right hand side vector  $\mathbf{b}$  was affected by measurement errors. However, also the entries of the coefficient matrix  $\mathbf{A}$  may have been obtained by measurement. This is the case, for instance, in the nodal analysis of electric circuits → Ex. 2.1.0.3. Then, it may be legitimate to seek a “better” matrix based on information contained in the whole linear system. This is the gist of the total least squares approach.

Given: overdetermined linear system of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m > n$ .

Known: LSE solvable  $\Leftrightarrow \mathbf{b} \in \text{Im}(\mathbf{A})$ , if  $\mathbf{A}, \mathbf{b}$  were not perturbed,  
but  $\mathbf{A}, \mathbf{b}$  are perturbed (measurement errors).

Sought: Solvable overdetermined system of equations  $\widehat{\mathbf{Ax}} = \widehat{\mathbf{b}}$ ,  $\widehat{\mathbf{A}} \in \mathbb{R}^{m,n}$ ,  $\widehat{\mathbf{b}} \in \mathbb{R}^m$ ,  
“nearest” to  $\mathbf{Ax} = \mathbf{b}$ .

☞ least squares problem “turned upside down”: now we are allowed to tamper with system matrix and right hand side vector!

**Total least squares problem:**

Given:  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m > n$ ,  $\text{rank}(\mathbf{A}) = n$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,

find:  $\widehat{\mathbf{A}} \in \mathbb{R}^{m,n}$ ,  $\widehat{\mathbf{b}} \in \mathbb{R}^m$  with

$$\left\| [\mathbf{A} \quad \mathbf{b}] - [\widehat{\mathbf{A}} \quad \widehat{\mathbf{b}}] \right\|_F \rightarrow \min , \quad \widehat{\mathbf{b}} \in \mathcal{R}(\widehat{\mathbf{A}}) . \quad (3.5.0.1)$$

$$\widehat{\mathbf{b}} \in \mathcal{R}(\widehat{\mathbf{A}}) \Rightarrow \text{rank}([\widehat{\mathbf{A}} \quad \widehat{\mathbf{b}}]) = n \quad \blacktriangleright \quad (3.5.0.1) \Rightarrow [\widehat{\mathbf{A}} \quad \widehat{\mathbf{b}}] = \underset{\text{rank}(\widehat{\mathbf{X}})=n}{\text{argmin}} \left\| [\mathbf{A} \quad \mathbf{b}] - \widehat{\mathbf{X}} \right\|_F .$$



$[\widehat{\mathbf{A}} \quad \widehat{\mathbf{b}}]$  is the rank- $n$  best approximation of  $[\mathbf{A} \quad \mathbf{b}]$ !

We face the problem to compute the best rank- $n$  approximation of the given matrix  $[\mathbf{A} \quad \mathbf{b}]$ , a problem already treated in Section 3.4.4.2: Thm. 3.4.4.18 tells us how to use the SVD of  $[\mathbf{A} \quad \mathbf{b}]$

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top , \quad \mathbf{U} \in \mathbb{R}^{m,n+1} , \quad \Sigma \in \mathbb{R}^{n+1,n+1} , \quad \mathbf{V} \in \mathbb{R}^{n+1,n+1} , \quad (3.5.0.2)$$

to construct  $\begin{bmatrix} \hat{\mathbf{A}} & \hat{\mathbf{b}} \end{bmatrix}$ :

$$[\mathbf{A} \quad \mathbf{b}] = \mathbf{U}\Sigma\mathbf{V}^\top = \sum_{j=1}^{n+1} \sigma_j(\mathbf{U})_{:,j}(\mathbf{V})_{:,j}^\top \stackrel{\text{Thm. 3.4.4.18}}{\implies} \begin{bmatrix} \hat{\mathbf{A}} & \hat{\mathbf{b}} \end{bmatrix} = \sum_{j=1}^n \sigma_j(\mathbf{U})_{:,j}(\mathbf{V})_{:,j}^\top. \quad (3.5.0.3)$$

$$\stackrel{\mathbf{V} \text{ orthogonal}}{\implies} [\hat{\mathbf{A}} \quad \hat{\mathbf{b}}](\mathbf{V})_{:,n+1} = \hat{\mathbf{A}}(\mathbf{V})_{1:n,n+1} + \hat{\mathbf{b}}(\mathbf{V})_{n+1,n+1} = 0. \quad (3.5.0.4)$$

► (3.5.0.4) also provides the solution  $\mathbf{x}$  of  $\hat{\mathbf{S}}\mathbf{x} = \hat{\mathbf{b}}$ ,

$$\mathbf{x} := -\hat{\mathbf{A}}^{-1}\hat{\mathbf{b}} = -(\mathbf{V})_{1:n,n+1}/(\mathbf{V})_{n+1,n+1}, \quad (3.5.0.5)$$

if  $(\mathbf{V})_{n+1,n+1} \neq 0$  (in numerical sense, of course).

#### C++-code 3.5.0.6: Total least squares via SVD

```

2 // computes only solution x of fitted consistent LSE
3 VectorXd lsqtotal(const MatrixXd& A, const VectorXd& b) {
4     const unsigned m = A.rows(), n = A.cols();
5     MatrixXd C(m, n + 1); C << A, b; // C = [A, b]
6     // We need only the SVD-factor V, see (3.5.0.3)
7     MatrixXd V = C.jacobiSvd(Eigen::ComputeThinU | Eigen::ComputeThinV).matrixV();
8
9     // Compute solution according to (3.5.0.5);
10    double s = V(n, n);
11    if (std::abs(s) < 1.0E-15) { cerr << "No solution!\n"; exit(1); }
12    return (-V.col(n).head(n) / s);
13 }
```

## 3.6 Constrained Least Squares

In the examples of Section 3.0.1 we expected all components of the right hand side vectors to be possibly affected by measurement errors. However, it might happen that some data are very reliable and in this case we would like the corresponding equation to be satisfied exactly.

► linear least squares problem with **linear constraint** defined as follows:

*Linear least squares problem with linear constraint:*

Given:  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m \geq n$ ,  $\text{rank}(\mathbf{A}) = n$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  
 $\mathbf{C} \in \mathbb{R}^{p,n}$ ,  $p < n$ ,  $\text{rank}(\mathbf{C}) = p$ ,  $\mathbf{d} \in \mathbb{R}^p$

Find:  $\mathbf{x} \in \mathbb{R}^n$  such that  $\|\mathbf{Ax} - \mathbf{b}\|_2 \rightarrow \min$  and  $\boxed{\mathbf{Cx} = \mathbf{d}}.$

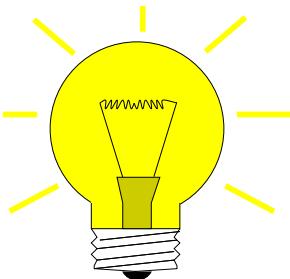
Linear constraint

(3.6.0.1)

Here the constraint matrix  $\mathbf{C}$  collects all the coefficients of those  $p$  equations that are to be satisfied exactly, and the vector  $\mathbf{d}$  the corresponding components of the right hand side vector. Conversely, the  $m$  equations of the (overdetermined) LSE  $\mathbf{Ax} = \mathbf{b}$  cannot be satisfied and are treated in a least squares sense.

### 3.6.1 Solution via Lagrangian Multipliers

**§3.6.1.1 (A saddle point problem)** Recall important technique from multidimensional calculus for tackling constrained minimization problems: [Lagrange multipliers](#), see [Str09, Sect. 7.9].



Idea: coupling the constraint using the [Lagrange multiplier](#)  $\mathbf{m} \in \mathbb{R}^p$

$$\mathbf{x} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \sup_{\mathbf{m} \in \mathbb{R}^p} L(\mathbf{y}, \mathbf{m}), \quad (3.6.1.2)$$

$$L(\mathbf{y}, \mathbf{m}) := \frac{1}{2} \|\mathbf{Ay} - \mathbf{b}\|^2 + \mathbf{m}^\top (\mathbf{Cy} - \mathbf{d}). \quad (3.6.1.3)$$

$L$  as defined in (3.6.1.3) is called a [Lagrange function](#). The simple heuristics behind Lagrange multipliers is the observation:

$$\sup_{\mathbf{m} \in \mathbb{R}^p} L(\mathbf{y}, \mathbf{m}) = \infty, \quad \text{in case } \mathbf{Cx} \neq \mathbf{d}!$$

- ▶ A minimum in (3.6.1.2) can only attained, if the constraint is satisfied!

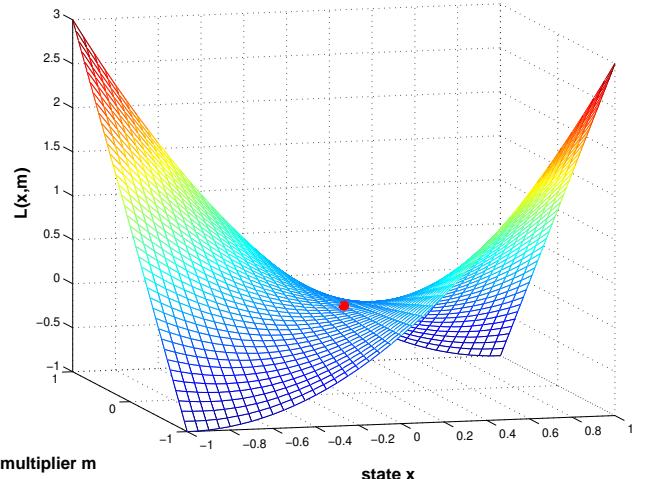
(3.6.1.2) is called a [saddle point problem](#).

Solution of min-max problem like (3.6.1.2) is called a [saddle point](#).

Saddle point of  $F(x, m) = x^2 - 2xm$

▷

Note that the function is “flat” in the saddle point •, that is, both the derivative with respect to  $\mathbf{x}$  and with respect to  $\mathbf{m}$  has to vanish.



↓

**§3.6.1.4 (Augmented normal equations)** In a saddle point the Lagrange function is “flat”, that is, all its partial derivatives have to vanish there.

- ▶ Necessary (and sufficient) conditions for the solution  $\mathbf{x}$  of (3.6.1.2) and a saddle point in  $\mathbf{x}, \mathbf{q}$ : (For a similar technique employing multi-dimensional calculus see Rem. 3.1.2.5)
- ▶

$$\frac{\partial L}{\partial \mathbf{x}}(\mathbf{x}, \mathbf{q}) = \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}) + \mathbf{C}^\top \mathbf{q} \stackrel{!}{=} 0, \quad (3.6.1.5a)$$

$$\frac{\partial L}{\partial \mathbf{m}}(\mathbf{x}, \mathbf{q}) = \mathbf{Cy} - \mathbf{d} \stackrel{!}{=} 0. \quad (3.6.1.5b)$$

$$\begin{bmatrix} \mathbf{A}^\top \mathbf{A} & \mathbf{C}^\top \\ \mathbf{C} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{q} \end{bmatrix} = \begin{bmatrix} \mathbf{A}^\top \mathbf{b} \\ \mathbf{d} \end{bmatrix} \quad \text{Augmented normal equations (matrix saddle point problem)} \quad (3.6.1.6)$$

As we know, a direct elimination solution algorithm for (3.6.1.6) amounts to finding an LU-decomposition of the coefficient matrix. Here we opt for its symmetric variant, the **Cholesky decomposition**, see Section 2.8. On the block-matrix level it can be found by considering the equation

$$\begin{bmatrix} \mathbf{A}^\top \mathbf{A} & \mathbf{C}^\top \\ \mathbf{C} & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{R}^\top & 0 \\ \mathbf{G} & -\mathbf{S}^\top \end{bmatrix} \begin{bmatrix} \mathbf{R} & \mathbf{G}^\top \\ 0 & \mathbf{S} \end{bmatrix}, \quad \mathbf{R}, \mathbf{S} \in \mathbb{R}^{n,n} \text{ upper triangular matrices, } \mathbf{G} \in \mathbb{R}^{p,n}.$$

Thus the blocks of the Cholesky factors of the coefficient matrix of the linear system (3.6.1.6) can be determined in four steps.

- ① Compute  $\mathbf{R}$  from  $\mathbf{R}^\top \mathbf{R} = \mathbf{A}^\top \mathbf{A} \rightarrow$  Cholesky decomposition  $\rightarrow$  Section 2.8,
- ② Compute  $\mathbf{G}$  from  $\mathbf{R}^\top \mathbf{G}^\top = \mathbf{C}^\top \rightarrow$   $n$  forward substitutions  $\rightarrow$  Section 2.3.2,
- ③ Compute  $\mathbf{S}$  from  $\mathbf{S}^\top \mathbf{S} = \mathbf{G} \mathbf{G}^\top \rightarrow$  Cholesky decomposition  $\rightarrow$  Section 2.8.

**§3.6.1.7 (Extended augmented normal equations)** The same caveats as those discussed for the regular normal equations in Rem. 3.2.0.3, Ex. 3.2.0.4, and Rem. 3.2.0.6, apply to the direct use of the augmented normal equations (3.6.1.6):

1. their condition number can be much bigger than that of the matrix  $\mathbf{A}$ ,
2. forming  $\mathbf{A}^\top \mathbf{A}$  may be vulnerable to roundoff,
3. the matrix  $\mathbf{A}^\top \mathbf{A}$  may not be sparse, though  $\mathbf{A}$  is.

As in § 3.2.0.7 also in the case of the augmented normal equations (3.6.1.6) switching to an extended version by introducing the residual  $\mathbf{r} = \mathbf{Ax} - \mathbf{b}$  as a new unknown is a remedy, cf. (3.2.0.8). This leads to the following linear system of equations.

$$\begin{bmatrix} -\mathbf{I} & \mathbf{A} & 0 \\ \mathbf{A}^\top & 0 & \mathbf{C}^\top \\ 0 & \mathbf{C} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{r} \\ \mathbf{x} \\ \mathbf{m} \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ 0 \\ \mathbf{d} \end{bmatrix} \quad \hat{=} \quad \text{Extended augmented normal equations} \quad (3.6.1.8)$$

## 3.6.2 Solution via SVD



Idea: Identify the subspace in which the solution can vary without violating the constraint. Since  $\mathbf{C}$  has *full rank*, this subspace agrees with the nullspace/kernel of  $\mathbf{C}$ .

From Lemma 3.4.1.12 and Ex. 3.4.2.7 we have learned that the **SVD** can be used to compute (an orthonormal basis of) the nullspace  $\mathcal{N}(\mathbf{C})$ . This suggests the following method for solving the constrained linear least squares problem (3.6.0.1).

- ① Compute an orthonormal basis of  $\mathcal{N}(\mathbf{C})$  using SVD ( $\rightarrow$  Lemma 3.4.1.12, (3.4.3.1)):

$$\mathbf{C} = \mathbf{U} [\Sigma \ 0] \begin{bmatrix} \mathbf{V}_1^\top \\ \mathbf{V}_2^\top \end{bmatrix}, \quad \mathbf{U} \in \mathbb{R}^{p,p}, \Sigma \in \mathbb{R}^{p,p}, \mathbf{V}_1 \in \mathbb{R}^{n,p}, \mathbf{V}_2 \in \mathbb{R}^{n,n-p}$$

►  $\mathcal{N}(\mathbf{C}) = \mathcal{R}(\mathbf{V}_2)$  .

and the **particular solution** of the constraint equation

$$\mathbf{x}_0 := \mathbf{V}_1 \Sigma^{-1} \mathbf{U}^\top \mathbf{d} .$$

► Representation of the solution  $\mathbf{x}$  of (3.6.0.1):  $\mathbf{x} = \mathbf{x}_0 + \mathbf{V}_2 \mathbf{y}$ ,  $\mathbf{y} \in \mathbb{R}^{n-p}$ .

- ② Insert this representation in (3.6.0.1). This yields a standard linear least squares problem with coefficient matrix  $\mathbf{A}\mathbf{V}_2 \in \mathbb{R}^{m,n-p}$  and right hand side vector  $\mathbf{b} - \mathbf{Ax}_0 \in \mathbb{R}^m$

$$\|\mathbf{A}(\mathbf{x}_0 + \mathbf{V}_2 \mathbf{y}) - \mathbf{b}\|_2 \rightarrow \min \Leftrightarrow \|\mathbf{A}\mathbf{V}_2 \mathbf{y} - (\mathbf{b} - \mathbf{Ax}_0)\| \rightarrow \min .$$

## Learning Outcomes

After having studied the contents of this chapter you should be able to

- give a rigorous definition of the least squares solution of an (overdetermined) linear system of equations,
- state the (extended) normal equations for any overdetermined linear system of equations,
- tell conditions for uniqueness and existence of solutions of the normal equations,
- define (economical) QR-decomposition and SVD of a matrix,
- know the asymptotic computational effort of computing economical QR and SVD factorizations,
- explain the use of QR-decomposition and, in particular, Givens rotations, for solving (overdetermined) linear systems of equations (in least squares sense),
- use SVD to solve least squares, (constrained) optimization, and low-rank best approximation problems
- explain the ideas underlying principal component analysis (PCA) and proper orthogonal decomposition (POD),
- formulate the augmented (extended) normal equations for a linearly constrained least squares problem.

# Bibliography

- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 187, 188, 194, 197–223, 232).
- [GGK14] W. Gander, M.J. Gander, and F. Kwok. *Scientific Computing*. Vol. 11. Texts in Computational Science and Engineering. Heidelberg: Springer, 2014 (cit. on p. 186).
- [GV13] Gene H. Golub and Charles F. Van Loan. *Matrix computations*. Fourth. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, 2013, pp. xiv+756 (cit. on pp. 206–208, 210).
- [Gut07] M. Gutknecht. “Linear Algebra”. 2007 (cit. on p. 204).
- [Gut09] M.H. Gutknecht. *Lineare Algebra*. Lecture Notes. SAM, ETH Zürich, 2009 (cit. on pp. 201, 202, 224, 226, 237).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 197–199, 202, 207, 232).
- [Hig02] N.J. Higham. *Accuracy and Stability of Numerical Algorithms*. 2nd ed. Philadelphia, PA: SIAM, 2002 (cit. on pp. 207, 212).
- [Kal96] D. Kalman. “A singularly valuable decomposition: The SVD of a matrix”. In: *The College Mathematics Journal* 27 (1996), pp. 2–23 (cit. on p. 224).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on pp. 185, 204, 224, 226, 227).
- [QMN16] Alfio Quarteroni, Andrea Manzoni, and Federico Negri. *Reduced basis methods for partial differential equations*. Vol. 92. Unitext. Springer, Cham, 2016, pp. xi+296.
- [Str19] D. Strang. *Linear Algebra and Learning from Data*. Cambridge University Press, 2019.
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on pp. 191, 193, 225, 253).

# Chapter 4

## Filtering Algorithms

This chapter continues the theme of *numerical linear algebra*, also covered in Chapter 1, 2, 10. We will come across very special linear transformations ( $\leftrightarrow$  matrices) and related algorithms. Surprisingly, these form the basis of a host of very important numerical methods for **signal processing**.

**§4.0.0.1 (Time-discrete signals and sampling)** From the perspective of **signal processing** we can identify

$$\text{vector } \mathbf{x} \in \mathbb{R}^n \leftrightarrow \text{finite discrete (= sampled) signal.}$$

**Sampling** converts a time-continuous signal, represented by some real-valued physical quantity (pressure, voltage, power, etc.) into a time-discrete signal:

$X = X(t) \triangleq \text{time-continuous signal}, 0 \leq t \leq T,$

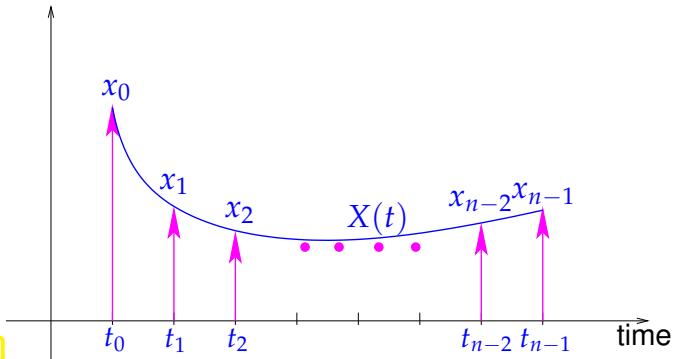
“sampling”:  $x_j = X(j\Delta t), j = 0, \dots, n-1, n \in \mathbb{N}, n\Delta t \leq T.$

$\Delta t > 0 \triangleq \text{time between samples.}$

As already indicated by the indexing the sampled values can be arranged in a vector  $\mathbf{x} \stackrel{\text{Fig. 102}}{=} [x_0, \dots, x_{n-1}]^\top \in \mathbb{R}^n$ .

Note that in this chapter, as is customary in signal processing, we adopt a C++-style indexing from 0: the components of a vector with length  $n$  carry indices  $\in \{0, \dots, n-1\}$ .

As an idealization one sometimes considers a signal of infinite duration  $X = X(t), -\infty < t < \infty$ . In this case sampling yields a **bi-infinite** time-discrete signal, represented by a sequence  $(x_k)_{k \in \mathbb{Z}}$ . If this sequence has a finite number of non-zero terms only, then we write  $(0, \dots, x_\ell, x_{\ell+1}, \dots, x_{n-1}, x_n, 0, \dots)$ .  $\square$



### Contents

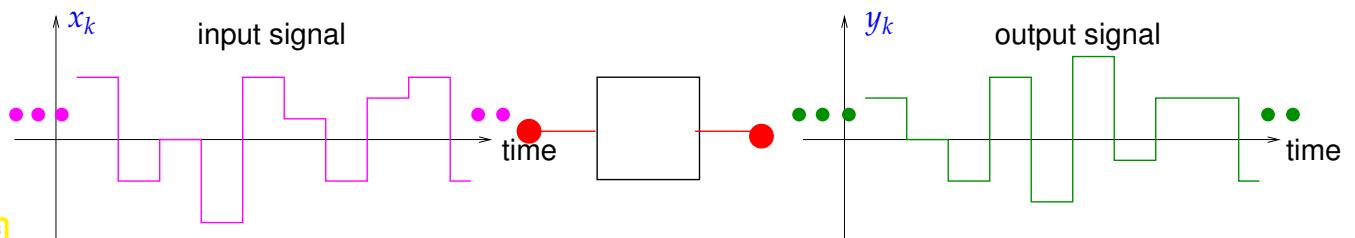
4.1 Discrete Convolutions . . . . .	258
4.2 Discrete Fourier Transform (DFT) . . . . .	267
4.2.1 Discrete Convolution via DFT . . . . .	273
4.2.2 Frequency filtering via DFT . . . . .	274
4.2.3 Real DFT . . . . .	281
4.2.4 Two-dimensional DFT . . . . .	283
4.2.5 Semi-discrete Fourier Transform [QSS00, Sect. 10.11] . . . . .	288

<b>4.3 Fast Fourier Transform (FFT)</b>	296
<b>4.4 Trigonometric transformations</b>	304
4.4.1 Sine transform	305
4.4.2 Cosine transform	311
<b>4.5 Toeplitz Matrix Techniques</b>	312
4.5.1 Toeplitz Matrix Arithmetic	315
4.5.2 The Levinson Algorithm	316

## 4.1 Discrete Convolutions

### §4.1.0.1 (Discrete finite linear time-invariant causal channels/filters)

Now we study a finite linear time-invariant causal channel (filter), which is widely used model for digital communication channels, e.g. in wireless communication theory. Mathematically speaking, a (discrete) channel/filter is a mapping  $F : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z})$  from the vector space  $\ell^\infty(\mathbb{Z})$  of bounded input sequences  $\{x_j\}_{j \in \mathbb{Z}}$  to bounded output sequences  $\{y_j\}_{j \in \mathbb{Z}}$ .



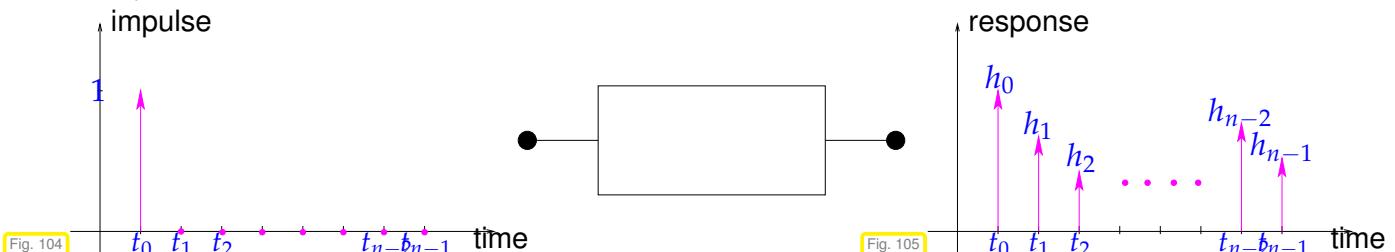
$$\text{Channel/filter: } F : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z}) , \quad (y_j)_{j \in \mathbb{Z}} = F((x_j)_{j \in \mathbb{Z}}) . \quad (4.1.0.2)$$

For the description of filters we rely on special input signals, analogous to the description of a linear mapping  $\mathbb{R}^n \mapsto \mathbb{R}^m$  through a matrix, that is, its action on unit vectors.

#### Definition 4.1.0.3. Impulse response

The **impulse response** of a channel/filter is the output for a single unit impulse at  $t = 0$  as input, that is, the input signal is  $x_j = \delta_{j,0} := \begin{cases} 1 & \text{if } j = 0 \\ 0 & \text{else} \end{cases}$  (Kronecker symbol).

Visualization of the (finite!) impulse response  $(\dots, 0, h_0, \dots, h_{n-1}, 0, \dots)$  of a (causal, see Def. 4.1.0.11 below) channel/filter:



In order to link digital filters to linear algebra, we have to assume certain properties that are indicated by the attributes “finite”, “linear”, “time-invariant” and “causal”:

**Definition 4.1.0.4. Finite channel/filter**

A filter  $F : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z})$  is called **finite**, if every input signal of finite duration produces an output signal of finite duration,

$$\{ \exists M \in \mathbb{N}: |j| > M \Rightarrow x_j = 0 \} \Rightarrow \exists N \in \mathbb{N}: |k| > N \Rightarrow (F((x_j)_{j \in \mathbb{Z}}))_k = 0 , \quad (4.1.0.5)$$

➤ The impulse response of a finite filter can be described by a vector  $\mathbf{h}$  of finite length  $n$ .

It should not matter when exactly signal is fed into the channel. To express this intuition more rigorously we introduce the **time shift operator** for signals: for  $m \in \mathbb{Z}$

$$S_m : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z}) , \quad S_m((x_j)_{j \in \mathbb{Z}}) = (x_{j-m})_{j \in \mathbb{Z}} . \quad (4.1.0.6)$$

**Definition 4.1.0.7. Time-invariant channel/filter**

A filter  $F : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z})$  is called **time-invariant**, if shifting the input in time leads to the same output shifted in time by the same amount; it *commutes with the time shift operator* from (4.1.0.6):

$$\forall (x_j)_{j \in \mathbb{Z}} \in \ell^\infty(\mathbb{Z}), \forall m \in \mathbb{Z}: F(S_m((x_j)_{j \in \mathbb{Z}})) = S_m(F((x_j)_{j \in \mathbb{Z}})) . \quad (4.1.0.8)$$

**Definition 4.1.0.9. Linear channel/filter**

A filter  $F : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z})$  is called **linear**, if  $F$  is a linear mapping:

$$F(\alpha(x_j)_{j \in \mathbb{Z}} + \beta(y_j)_{j \in \mathbb{Z}}) = \alpha F((x_j)_{j \in \mathbb{Z}}) + \beta F((y_j)_{j \in \mathbb{Z}}) \quad \forall (x_j)_{j \in \mathbb{Z}}, (y_j)_{j \in \mathbb{Z}} \in \ell^\infty(\mathbb{Z}), \alpha, \beta \in \mathbb{R} . \quad (4.1.0.10)$$

Slightly rewritten, this means that for all scaling factors  $\alpha, \beta \in \mathbb{R}$

$$\text{output}(\alpha \cdot \text{signal 1} + \beta \cdot \text{signal 2}) = \alpha \cdot \text{output(signal 1)} + \beta \cdot \text{output(signal 2)} .$$

Of course, a signal should not trigger an output before it arrives at the filter; output may depend only on past and present inputs, not on the future.

**Definition 4.1.0.11. Causal channel/filter**

A filter  $F : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z})$  is called **causal** (or physical, or nonanticipative), if the output does not start before the input

$$\forall M \in \mathbb{N}: (x_j)_{j \in \mathbb{Z}} \in \ell^\infty(\mathbb{Z}), x_j = 0 \quad \forall j \leq M \Rightarrow F((x_j)_{j \in \mathbb{Z}})_k = 0 \quad \forall k \leq M . \quad (4.1.0.12)$$

Acronym: **LT-FIR**  $\hat{=}$  finite ( $\rightarrow$  Def. 4.1.0.4), linear ( $\rightarrow$  Def. 4.1.0.9), time-invariant ( $\rightarrow$  Def. 4.1.0.7), and causal ( $\rightarrow$  Def. 4.1.0.11) filter  $F : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z})$

The input response of a finite and causal filter is a sequence of the form  $(\dots, 0, h_0, h_1, \dots, h_{n-1}, 0, \dots)$ ,  $n \in \mathbb{N}$ . Such an impulse response is depicted in Fig. 105.  $\square$

**§4.1.0.13 (Transmission through finite, time-invariant, linear, causal filters)** Let  $(\dots, 0, h_0, h_1, \dots, h_{n-1}, 0, \dots)$ ,  $n \in \mathbb{N}$ , be the impulse response ( $\rightarrow$  4.1.0.3) of a finite ( $\rightarrow$  Def. 4.1.0.4), linear ( $\rightarrow$  Def. 4.1.0.9), time-invariant ( $\rightarrow$  Def. 4.1.0.7), and causal ( $\rightarrow$  Def. 4.1.0.11) filter (LT-FIR)  $F : \ell^\infty(\mathbb{Z}) \rightarrow \ell^\infty(\mathbb{Z})$ :

$$F((\delta_{j,0})_{j \in \mathbb{Z}}) = (\dots, 0, h_0, h_1, \dots, h_{n-1}, 0, \dots).$$

Owing to **time-invariance** we already know the response to a shifted unit pulse

$$F((\delta_{j,k})_{j \in \mathbb{Z}}) = (h_{j-k})_{j \in \mathbb{Z}} = \left( \dots \ 0 \quad \begin{matrix} h_0 & h_1 & \dots & h_{n-1} & 0 & \dots \end{matrix} \right). \quad \begin{matrix} \uparrow \\ t = k\Delta t \end{matrix} \quad \begin{matrix} \uparrow \\ t = (k+n-1)\Delta t \end{matrix}$$

Every finite input signal  $(\dots, 0, x_0, x_1, \dots, x_{m-1}, 0, \dots) \in \ell^\infty(\mathbb{Z})$  can be written as the superposition of scaled unit impulses, which, in turns, are time-shifted copies of a unit pulse at  $t = 0$ :

$$(x_j)_{j \in \mathbb{Z}} = \sum_{k=0}^{m-1} x_k (\delta_{j,k})_{j \in \mathbb{Z}} = \sum_{k=0}^{m-1} x_k S_k ((\delta_{j,0})_{j \in \mathbb{Z}}),$$

where  $S_k$  is the time-shift operator from (4.1.0.6). Applying the filter on both sides of this equation and using **linearity** leads to the general formula for the output signal  $(y_j)_{j \in \mathbb{Z}}$

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \\ \vdots \\ y_{m+n-3} \\ y_{m+n-2} \end{bmatrix} = x_0 \begin{bmatrix} h_0 \\ \vdots \\ h_{n-1} \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + x_1 \begin{bmatrix} 0 \\ h_0 \\ \vdots \\ h_{n-1} \\ 0 \\ \vdots \\ 0 \end{bmatrix} + x_2 \begin{bmatrix} 0 \\ 0 \\ h_0 \\ \vdots \\ h_{n-1} \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \dots + x_{m-1} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ h_0 \\ \vdots \\ 0 \\ \vdots \\ h_{n-1} \end{bmatrix}.$$

Thus, in compact notation we can write the non-zero components of the output signal  $(y_j)_{j \in \mathbb{Z}}$  as  
channel is causal and finite!

$$y_k = \sum_{j=0}^{m+n-1} h_{k-j} x_j, \quad k = 0, \dots, m+n-2 \quad (h_j := 0 \text{ for } j < 0 \text{ and } j \geq n). \quad (4.1.0.14)$$

Summary of the above considerations:

### Superposition of impulse responses

The output  $(\dots, 0, y_0, y_1, y_2, \dots)$  of a finite, time-invariant, linear, and causal channel for finite length input  $x = (\dots, 0, x_0, \dots, x_{n-1}, 0, \dots) \in \ell^\infty(\mathbb{Z})$  is  
a superposition of  $x_j$ -weighted  $j\Delta t$  time-shifted impulse responses.

**EXAMPLE 4.1.0.16 (Visualization: superposition of impulse responses)** The following diagrams give a visual display of considerations of § 4.1.0.13, namely of the superposition of impulse responses for a

particular finite, time-invariant, linear, and causal filter (LT-FIR), and an input signal of duration  $3\Delta t$ ,  $\Delta t \triangleq$  time between samples.

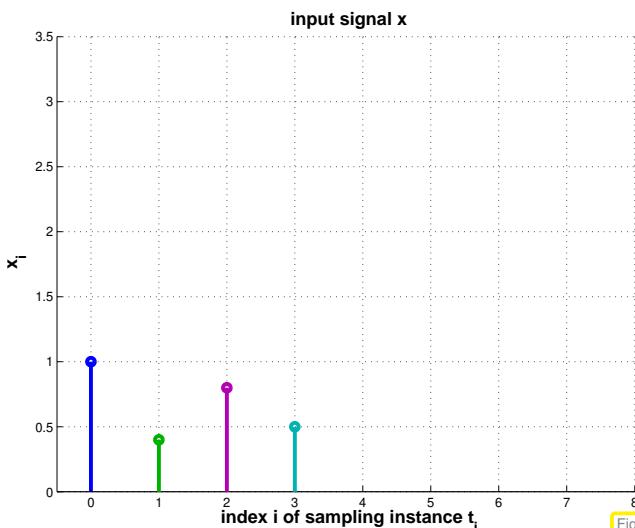


Fig. 106

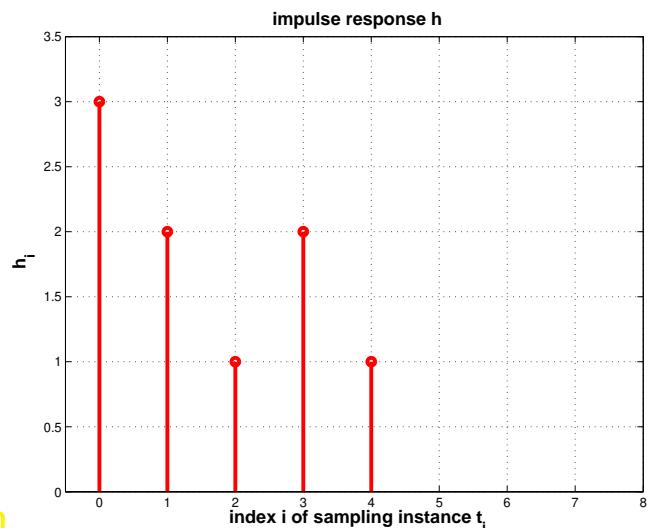


Fig. 107

Output = linear superposition of impulse responses:

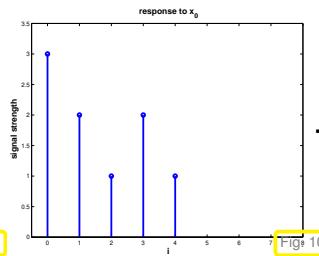


Fig. 108

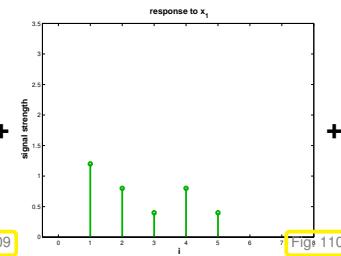


Fig. 109

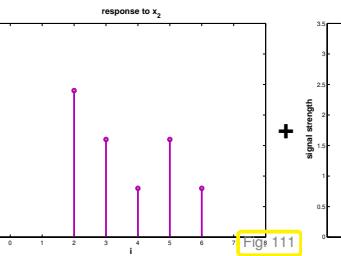


Fig. 110

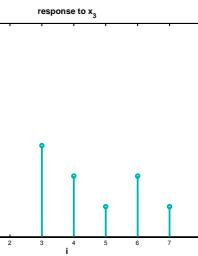


Fig. 111

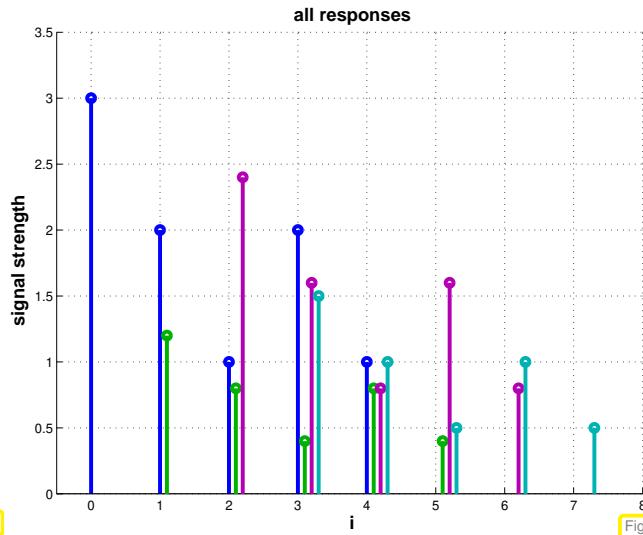
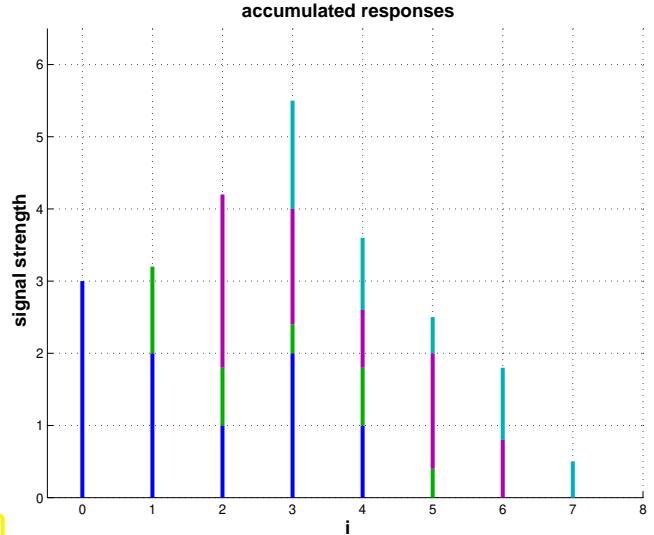


Fig. 112

Fig. 113



**Remark 4.1.0.17 (Reduction to finite signals)** We consider a finite ( $\rightarrow$  Def. 4.1.0.4), linear ( $\rightarrow$  Def. 4.1.0.9), time-invariant ( $\rightarrow$  Def. 4.1.0.7), and causal ( $\rightarrow$  Def. 4.1.0.11) filter (LT-FIR) with impulse response  $(\dots, 0, h_0, h_1, \dots, h_{n-1}, 0, \dots)$ ,  $n \in \mathbb{N}$ . From (4.1.0.14) we learn that

$$\text{duration(output signal)} \leq \text{duration(input signal)} + \text{duration(impulse response)} .$$

Therefore, if we know that all input signals are of the form  $(\dots, x_0, x_1, \dots, x_{m-1}, 0, \dots)$ , we can model them as vectors  $\mathbf{x} = [x_0, \dots, x_{m-1}]^\top \in \mathbb{R}^m$ , cf. § 4.0.0.1, and the filter can be viewed as a linear mapping  $F : \mathbb{R}^m \rightarrow \mathbb{R}^{m+n-1}$ , which takes us to the realm of linear algebra.

Thus, for the filter we have a matrix representation of (4.1.0.14). Writing  $\mathbf{y} = [y_0, \dots, y_{2n-2}]^\top \in \mathbb{R}^{2n-1}$  for the vector of the output signal we find in the case  $m = n$

$$\begin{bmatrix} y_0 \\ \vdots \\ y_{2n-2} \end{bmatrix} = \begin{bmatrix} h_0 & 0 & \cdots & 0 \\ h_1 & \ddots & & \\ \vdots & & \ddots & 0 \\ h_{n-1} & 0 & \cdots & h_1 & h_0 \\ 0 & & & & \\ \vdots & & & & \\ 0 & & \cdots & 0 & h_{n-1} \end{bmatrix} \begin{bmatrix} x_0 \\ \vdots \\ x_{n-1} \end{bmatrix}. \quad (4.1.0.18)$$

□

**EXAMPLE 4.1.0.19 (Multiplication of polynomials)** “Surprisingly” the bilinear operation (4.1.0.14) that takes two input vectors and produces an output vector with double the number of entries (—1) also governs the multiplication of polynomials:

$$p(z) = \sum_{k=0}^{n-1} a_k z^k, \quad q(z) = \sum_{k=0}^{n-1} b_k z^k \quad \Rightarrow \quad (pq)(z) = \sum_{k=0}^{2n-2} \underbrace{\left( \sum_{j=0}^k a_j b_{k-j} \right)}_{=:c_k} z^k \quad (4.1.0.20)$$

Here the roles of  $h_k$ ,  $x_k$  are played by the  $a_k$  and  $b_k$ .

➤ the coefficients of the product polynomial can be obtained through an operation similar to finite, time-invariant, linear, and causal (LT-FIR) filtering! □

**§4.1.0.21 (Discrete convolution)** Both in (4.1.0.14) and (4.1.0.20) we recognize the same pattern of a particular *bi-linear* combination of

- discrete signals in § 4.1.0.13,
- polynomial coefficient sequences in Ex. 4.1.0.19.

### Definition 4.1.0.22. Discrete convolution

Given  $\mathbf{x} = [x_0, \dots, x_{n-1}]^\top \in \mathbb{K}^n$ ,  $\mathbf{h} = [h_0, \dots, h_{n-1}]^\top \in \mathbb{K}^n$  their **discrete convolution** is the vector  $\mathbf{y} \in \mathbb{K}^{2n-1}$  with components

$$y_k = \sum_{j=0}^{n-1} h_{k-j} x_j, \quad k = 0, \dots, 2n-2 \quad (h_j := 0 \text{ for } j < 0). \quad (4.1.0.23)$$

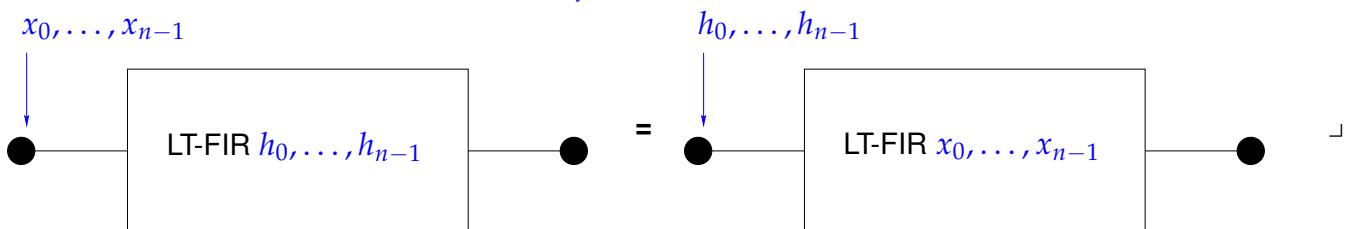
☞ Notation for discrete convolution (4.1.0.23):

$$\mathbf{y} = \mathbf{h} * \mathbf{x}.$$

Defining  $x_j := 0$  for  $j < 0$ , we find that **discrete convolution is commutative** :

$$y_k = \sum_{j=0}^{n-1} h_{k-j} x_j = \sum_{l=0}^{n-1} h_l x_{k-l}, \quad k = 0, \dots, 2n-2, \quad (\text{that is, } \mathbf{h} * \mathbf{x} = \mathbf{x} * \mathbf{h}),$$

obtained by index transformation  $l \leftarrow k - j$ .



**Remark 4.1.0.24 (Multiplication of polynomials)** The formula (4.1.0.23) for the discrete convolution also occurs in a context completely detached from signal processing.

Consider two polynomials in  $t$  of degree  $n-1$ ,  $n \in \mathbb{N}$ ,

$$p(t) = \sum_{j=0}^{n-1} a_j t^j, \quad q(t) = \sum_{j=0}^{n-1} b_j t^j, \quad a_j, b_j \in \mathbb{K}.$$

Their product  $pq$  will be a polynomial of degree  $2n-2$ :

$$(pq)(t) = \sum_{j=0}^{2n-2} c_j t^j, \quad c_j = \sum_{\ell=\max\{0, j-(n-1)\}}^{\min\{j, n-1\}} a_\ell b_{j-\ell}, \quad j = 0, \dots, 2n-2. \quad (4.1.0.25)$$

Let us introduce dummy coefficients for  $p(t)$  and  $q(t)$ ,  $a_j, b_j$ ,  $j = 2n, \dots, 2n-2$ , all set to 0. This can be easily done in a computer code by resizing the coefficient vectors of  $p$  and  $q$  and filling the new entries with zeros. The above formula for  $c_j$  can then be rewritten as

$$c_j = \sum_{\ell=0}^j a_\ell b_{j-\ell}, \quad j = 0, \dots, 2n-2. \quad (4.1.0.26)$$

Hence, the coefficients of the product polynomial can be obtained as the discrete convolution of the coefficient vectors of  $p$  and  $q$ :

$$[c_0 \ c_1 \ \dots \ c_{2n-2}]^\top = \mathbf{a} * \mathbf{b} \quad ! \quad (4.1.0.27)$$

Moreover, this provides another proof for the commutativity of discrete convolution.  $\square$

**Remark 4.1.0.28 (Convolution of sequences)** The notion of a discrete convolution of Def. 4.1.0.22 naturally extends to sequences  $\in \ell^\infty(\mathbb{N}_0)$ , that is, bounded mappings  $\mathbb{N}_0 \mapsto \mathbb{K}$ : the (discrete) convolution of two sequences  $(x_j)_{j \in \mathbb{N}_0}$ ,  $(y_j)_{j \in \mathbb{N}_0}$  is the sequence  $(z_j)_{j \in \mathbb{N}_0}$  defined by

$$z_k := \sum_{j=0}^k x_{k-j} y_j = \sum_{j=0}^k x_j y_{k-j}, \quad k \in \mathbb{N}_0.$$

In this context recall the product formula for power series, **Cauchy product**, which can be viewed as a multiplication rule for “infinite polynomials” = power series.  $\square$

**Remark 4.1.0.29 (Linear filtering of periodic signals)** An  $n$ -periodic signal ( $n \in \mathbb{N}$ ) is a sequence  $(x_j)_{j \in \mathbb{Z}}$  satisfying  $x_{j+n} = x_j \quad \forall j \in \mathbb{Z}$ .

➤ An  $n$ -periodic signal  $(x_j)_{j \in \mathbb{Z}}$  uniquely determined by the values  $x_0, \dots, x_{n-1}$  and can be associated with a vector  $\mathbf{x} = [x_0, \dots, x_{n-1}]^\top \in \mathbb{R}^n$ .

Whenever the input signal of a finite, time-invariant filter is  $n$ -periodic, so will be the output signal. Thus, in the  $n$ -periodic setting, a causal, linear, and time-invariant filter (LT-FIR) will give rise to a linear mapping  $\mathbb{R}^n \mapsto \mathbb{R}^n$  according to

$$y_k = \sum_{j=0}^{n-1} p_{k-j} x_j \quad \text{for some } p_0, \dots, p_{n-1} \in \mathbb{R}, \quad p_k := p_{k-n} \text{ for all } k \in \mathbb{Z}. \quad (4.1.0.30)$$

In matrix notation (4.1.0.30) reads

$$\begin{bmatrix} y_0 \\ \vdots \\ y_{n-1} \end{bmatrix} = \underbrace{\begin{bmatrix} p_0 & p_{n-1} & p_{n-2} & \cdots & \cdots & p_1 \\ p_1 & p_0 & p_{n-1} & & & \vdots \\ p_2 & p_1 & p_0 & \ddots & & \\ \vdots & & \ddots & \ddots & \ddots & \\ & & & \ddots & \ddots & \\ \vdots & & & & \ddots & p_{n-1} \\ p_{n-1} & \cdots & & & p_1 & p_0 \end{bmatrix}}_{=: \mathbf{P}} \begin{bmatrix} x_0 \\ \vdots \\ x_{n-1} \end{bmatrix}. \quad (4.1.0.31)$$

where

$$(\mathbf{P})_{ij} = p_{i-j}, \quad 1 \leq i, j \leq n, \quad \text{with } p_j := p_{j+n} \text{ for } 1 - n \leq j < 0.$$

Note that the coefficients  $p_0, \dots, p_{n-1}$ , which can be regarded as **periodic impulse response**, do **not** agree with the impulse response (→ Def. 4.1.0.3) of the filter. However, they can easily be deduced from it, keeping in mind that the  $n$ -periodic sequence is the output signal for the input  $\sum_{k \in \mathbb{Z}} (\delta_{nk,j})_{j \in \mathbb{Z}}$ . If  $(\dots, h_0, h_1, \dots, h_m, 0, \dots)$  is the impulse response of LT-FIR filter, then

$$p_j = \sum_{\ell \in \mathbb{Z}} h_{j+\ell n}, \quad j \in \{0, \dots, n-1\}. \quad (4.1.0.32)$$

The following special variant of a discrete convolution operation is motivated by the preceding Rem. 4.1.0.29.

### Definition 4.1.0.33. Discrete periodic convolution

The **discrete periodic convolution** of two  $n$ -periodic sequences  $(p_k)_{k \in \mathbb{Z}}, (x_k)_{k \in \mathbb{Z}}$  yields the  $n$ -periodic sequence

$$(y_k) := (p_k) *_n (x_k), \quad y_k := \sum_{j=0}^{n-1} p_{k-j} x_j = \sum_{j=0}^{n-1} x_{k-j} p_j, \quad k \in \mathbb{Z}.$$

☞ notation for discrete periodic convolution:  $(p_k) *_n (x_k)$

Since  $n$ -periodic sequences can be identified with vectors in  $\mathbb{K}^n$  (see above), we can also introduce the discrete periodic convolution of vectors:

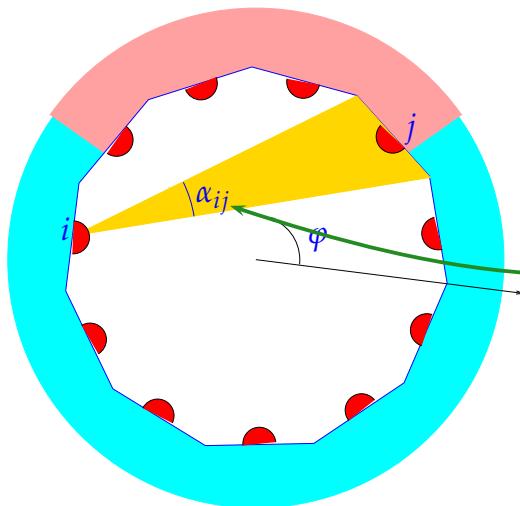
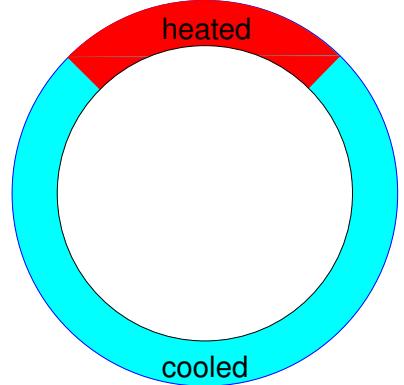
Def. 4.1.0.33  $\Rightarrow$  discrete periodic convolution of vectors:  $\mathbf{y} = \mathbf{p} *_n \mathbf{x} \in \mathbb{K}^n, \quad \mathbf{p}, \mathbf{x} \in \mathbb{K}^n$ .

**EXAMPLE 4.1.0.34 (Radiative heat transfer)** Beyond signal processing discrete periodic convolutions occur in many mathematical models:

An engineering problem:

- ◆ cylindrical pipe,
- ◆ heated on part  $\Gamma_H$  of its perimeter ( $\rightarrow$  prescribed heat flux),
- ◆ cooled on remaining perimeter  $\Gamma_K$  ( $\rightarrow$  constant heat flux).

Task: compute local heat fluxes.



**Modeling** (discretization):

- approximation by regular  $n$ -polygon, edges  $\Gamma_j$ ,
- isotropic radiation of each edge  $\Gamma_j$  (power  $I_j$ ),

$$\text{radiative heat flow } \Gamma_j \rightarrow \Gamma_i: P_{ji} := \frac{\alpha_{ij}}{\pi} I_j,$$

$$\text{opening angle: } \alpha_{ij} = \pi \gamma_{|i-j|}, \quad 1 \leq i, j \leq n,$$

$$\text{power balance: } \underbrace{\sum_{i=1, i \neq j}^n P_{ji}}_{=I_j} - \sum_{i=1, i \neq j}^n P_{ij} = Q_j. \quad (4.1.0.35)$$

$Q_j \triangleq$  heat flux through  $\Gamma_j$ , satisfies

$$Q_j := \int_{\frac{2\pi}{n}(j-1)}^{\frac{2\pi}{n}j} q(\varphi) d\varphi, \quad q(\varphi) := \begin{cases} \text{local heating} & , \text{ if } \varphi \in \Gamma_H, \\ -\frac{1}{|\Gamma_K|} \int_{\Gamma_H} q(\varphi) d\varphi & (\text{const.}), \text{ if } \varphi \in \Gamma_K. \end{cases}$$

$$(4.1.0.35) \quad \Rightarrow \quad \text{LSE: } I_j - \sum_{i=1, i \neq j}^n \frac{\alpha_{ij}}{\pi} I_i = Q_j, \quad j = 1, \dots, n.$$

e.g.  $n = 8$ :

$$\begin{bmatrix} 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 \\ -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 \\ -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 \\ -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 & -\gamma_2 & -\gamma_4 \\ -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 \\ -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 \\ -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 \\ -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 \end{bmatrix} \begin{bmatrix} I_1 \\ I_2 \\ I_3 \\ I_4 \\ I_5 \\ I_6 \\ I_7 \\ I_8 \end{bmatrix} = \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \\ Q_4 \\ Q_5 \\ Q_6 \\ Q_7 \\ Q_8 \end{bmatrix}. \quad (4.1.0.36)$$

This is a linear system of equations with symmetric, singular, and (by Lemma 9.1.0.5,  $\sum \gamma_i \leq 1$ ) positive semidefinite ( $\rightarrow$  Def. 1.1.2.6) system matrix.

Note that the matrices from (4.1.0.31) and (4.1.0.36) have the same structure!

Also observe that the LSE from (4.1.0.36) can be written by means of the discrete periodic convolution ( $\rightarrow$  Def. 4.1.0.33) of vectors  $\mathbf{y} = (1, -\gamma_1, -\gamma_2, -\gamma_3, -\gamma_4, -\gamma_3, -\gamma_2, -\gamma_1)$ ,  $\mathbf{x} = (I_1, \dots, I_8)$

$$(4.1.0.36) \leftrightarrow \mathbf{y} *_8 \mathbf{x} = [Q_1, \dots, Q_8]^\top.$$

□

**§4.1.0.37 (Circulant matrices)** In Ex. 4.1.0.34 we have already seen a coefficient matrix of a special form, which is common enough to warrant giving it a particular name:

**Definition 4.1.0.38. Circulant matrix** → [Han02, Sect. 54]

A matrix  $\mathbf{C} = [c_{ij}]_{i,j=1}^n \in \mathbb{K}^{n,n}$  is **circulant**  
 $\Leftrightarrow \exists (p_k)_{k \in \mathbb{Z}} \text{ } n\text{-periodic sequence: } c_{ij} = p_{j-i}, 1 \leq i, j \leq n.$

☞ Notation: We write  $\text{circul}(\mathbf{p}) \in \mathbb{K}^{n,n}$  for the circulant matrix generated by the periodic sequence/vector  $\mathbf{p} = [p_0, \dots, p_{n-1}]^\top \in \mathbb{K}^n$

- ☞ Circulant matrix has constant (main, sub- and super-) diagonals (for which indices  $j - i = \text{const.}$ ).
- ☞ columns/rows arise by *cyclic permutation* from first column/row.

Similar to the case of banded matrices ( $\rightarrow$  Section 2.7.5):

“information content” of circulant matrix  $\mathbf{C} \in \mathbb{K}^{n,n} = n$  numbers  $\in \mathbb{K}$ .  
(Obviously, one vector  $\mathbf{u} \in \mathbb{K}^n$  enough to define circulant matrix  $\mathbf{C} \in \mathbb{K}^{n,n}$ )

Structure of a generic circulant matrix (“constant diagonals”):

$$\text{circul}(\mathbf{p}) = \begin{bmatrix} p_0 & p_1 & p_2 & \cdots & & \cdots & p_{n-1} \\ p_{n-1} & p_0 & & & & & p_{n-2} \\ p_{n-2} & & p_0 & & & & \vdots \\ \vdots & & & p_0 & & & \vdots \\ \vdots & & & & p_0 & & \vdots \\ p_2 & & p_1 & & \cdots & & p_1 \\ p_1 & & p_2 & \cdots & & & p_0 \end{bmatrix}$$

Write  $\mathbf{Z}((u_k)) \in \mathbb{K}^{n,n}$  for the circulant matrix generated by the  $n$ -periodic sequence  $(u_k)_{k \in \mathbb{Z}}$ . Denote by  $\mathbf{y} := (y_0, \dots, y_{n-1})^\top$ ,  $\mathbf{x} = (x_0, \dots, x_{n-1})^\top$  the vectors associated to  $n$ -periodic sequences.

Then the commutativity of the discrete periodic convolution ( $\rightarrow$  Def. 4.1.0.33) involves

$$\text{circul}(\mathbf{x})\mathbf{y} = \text{circul}(\mathbf{y})\mathbf{x}. \quad (4.1.0.39)$$

□

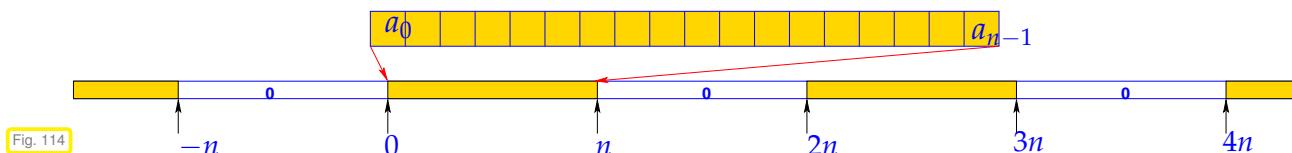
**Remark 4.1.0.40 (Reduction to periodic convolution)** Recall discrete convolution ( $\rightarrow$  Def. 4.1.0.22) of two vectors  $\mathbf{a} = (a_0, \dots, a_{n-1})^\top \in \mathbb{K}^n$ ,  $\mathbf{b} = (b_0, \dots, b_{n-1})^\top \in \mathbb{K}^n$ .

$$z_k := (\mathbf{a} * \mathbf{b})_k = \sum_{j=0}^{n-1} a_j b_{k-j}, \quad k = 0, \dots, 2n-2.$$

Now expand  $a_0, \dots, a_{n-1}$  and  $b_0, \dots, b_{n-1}$  to  $2n - 1$ -periodic sequences by zero padding.

$$x_k := \begin{cases} a_k & , \text{ if } 0 \leq k < n , \\ 0 & , \text{ if } n \leq k < 2n - 1 \end{cases} , \quad y_k := \begin{cases} b_k & , \text{ if } 0 \leq k < n , \\ 0 & , \text{ if } n \leq k < 2n - 1 , \end{cases} \quad (4.1.0.41)$$

and periodic extension:  $x_k = x_{2n-1+k}$ ,  $y_k = y_{2n-1+k}$  for all  $k \in \mathbb{Z}$ . The zero components prevent interaction of different periods:



$$\mathbf{a} * \mathbf{b})_k = (\mathbf{x} *_{2n-1} \mathbf{y})_k, \quad k = 0, \dots, 2n-2. \quad (4.1.0.42)$$

In the spirit of (4.1.0.18) we can switch to a matrix view of the reduction to periodic convolution:

$$\begin{bmatrix} z_0 \\ \vdots \\ z_{2n-2} \end{bmatrix} = \underbrace{\begin{bmatrix} y_0 & 0 & & & & & \\ y_1 & & & & & & \\ & \ddots & & & & & \\ & & 0 & & & & \\ & & & 0 & & & \\ & & & & \ddots & & \\ & & & & & 0 & y_{n-1} \\ & & & & & & 0 \\ y_{n-1} & & & y_1 & y_0 & 0 & & \\ & 0 & & & y_1 y_0 & & & \\ & & \ddots & & & & & \\ & & & 0 & & & & \\ & & & & \ddots & & & \\ & & & & & 0 & y_0 & \\ 0 & & & 0 & y_{n-1} & & & \\ & & & & & & & \\ & & & & & & & \end{bmatrix}}_{\text{a } (2n-1) \times (2n-1) \text{ circulant matrix!}} \begin{bmatrix} x_0 \\ \vdots \\ x_{n-1} \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \quad (4.1.0.43)$$

Discrete convolution can be realized by multiplication with a circulant matrix ( $\rightarrow$  4.1.0.38)

## 4.2 Discrete Fourier Transform (DFT)

Algorithms dealing with circulant matrices make use of their very special spectral properties. Full understanding requires familiarity with the theory of eigenvalues and eigenvectors of matrices from linear

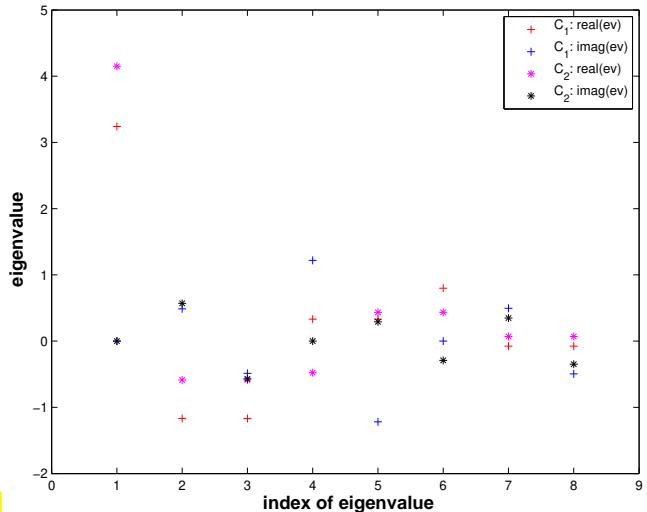
algebra, see [NS02, Ch. 7], [Gut09, Ch. 9].

**EXPERIMENT 4.2.0.1 (Eigenvectors of circulant matrices)** Now we are about to discover a very deep truth ...

Experimentally, we examine the eigenvalues and eigenvectors of two random  $8 \times 8$  circulant matrices  $\mathbf{C}_1, \mathbf{C}_2$  ( $\rightarrow$  Def. 4.1.0.38), generated from random vectors with entries even distributed in  $[0, 1]$ , `VectorXd::Random(n)`.

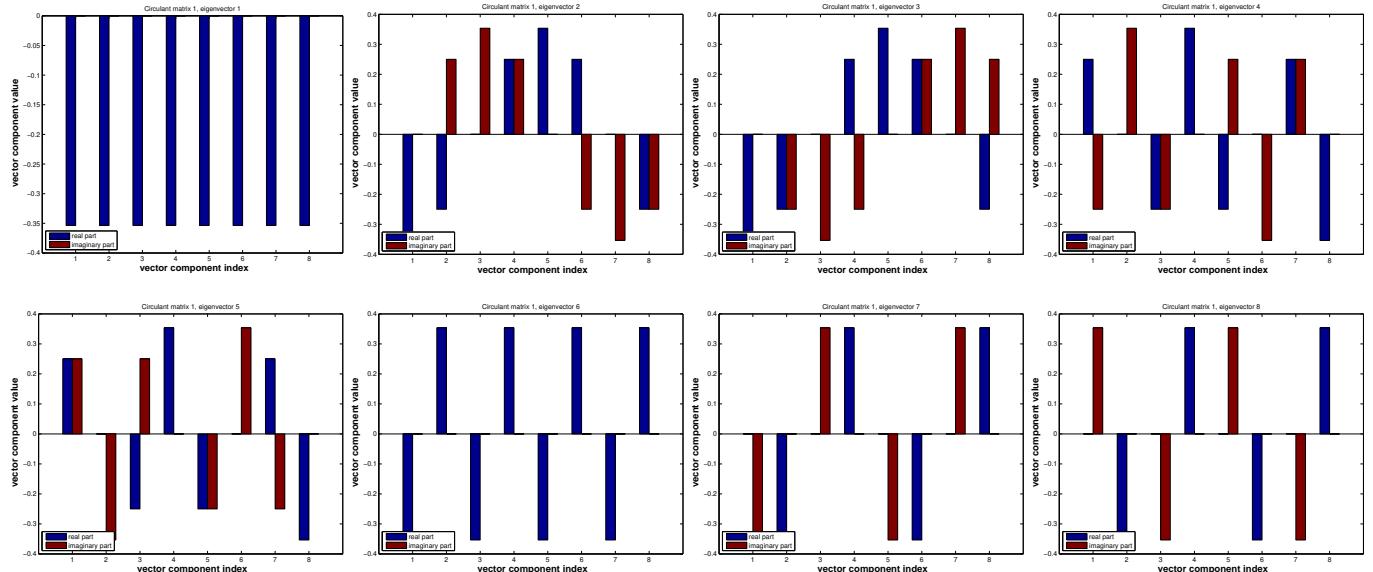
eigenvalues (real part)  $\triangleright$

Little relationship between (complex!) eigenvalues can be observed, as can be expected from random matrices with entries  $\in [0, 1]$ .

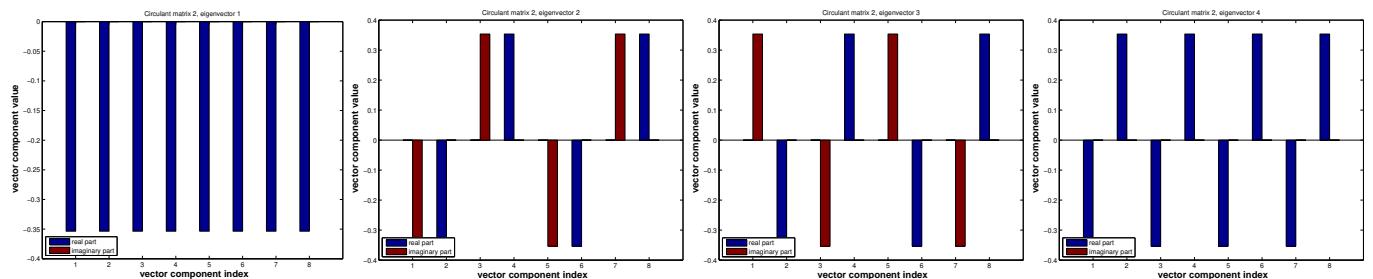


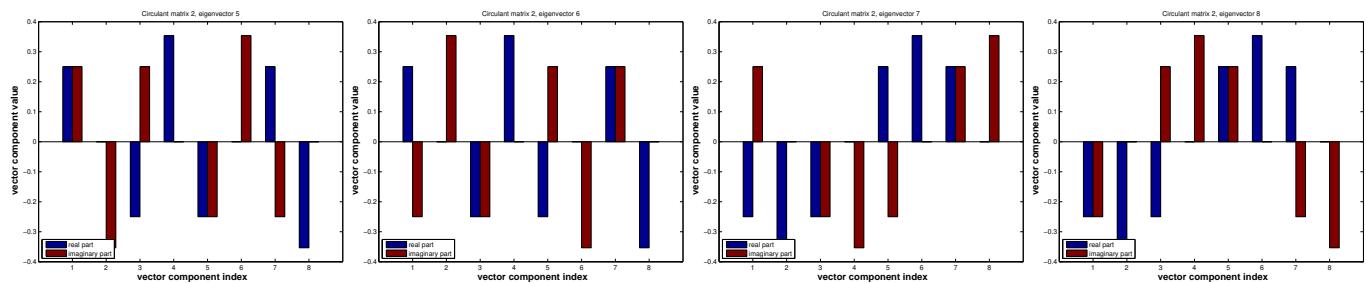
Now: the **surprise** ...

Eigenvectors of matrix  $\mathbf{C}_1$ , visualized through the size of the real and imaginary parts of their components.



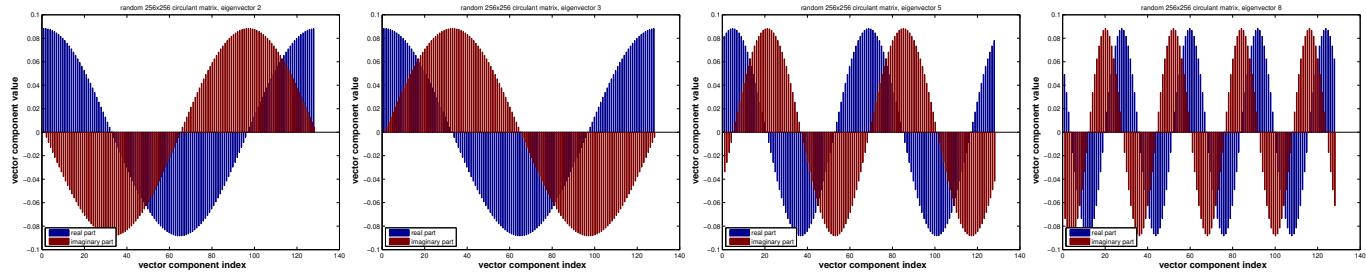
Eigenvectors of matrix  $\mathbf{C}_2$





Observation: different random circulant matrices have the same eigenvectors!

Eigenvectors of  $C = \text{gallery}(\text{'circul'}, (1:128))$ :



The eigenvectors remind us of sampled trigonometric functions  $\cos(k/n)$ ,  $\sin(k/n)$ ,  $k = 0, \dots, n - 1$ !  $\square$

**Remark 4.2.0.2 (Eigenvectors of commuting matrices)** An abstract result from linear algebra puts the surprising observation made in Exp. 4.2.0.1 in a wider context.

### Theorem 4.2.0.3. Commuting matrices have the same eigenvectors

If  $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{n,n}$  commute, that is,  $\mathbf{AB} = \mathbf{BA}$ , and  $\mathbf{A}$  has  $n$  distinct eigenvalues, then the eigenspaces of  $\mathbf{A}$  and  $\mathbf{B}$  coincide.

*Proof.* Let  $\mathbf{v} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  be an eigenvector of  $\mathbf{A}$  with eigenvalue  $\lambda$ . Then

$$(\mathbf{A} - \lambda \mathbf{I})\mathbf{v} = 0 \Rightarrow \mathbf{B}(\mathbf{A} - \lambda \mathbf{I})\mathbf{v} = 0 \stackrel{\mathbf{BA}=\mathbf{AB}}{\Rightarrow} (\mathbf{A} - \lambda \mathbf{I})\mathbf{B}\mathbf{v} = 0.$$

Since in the case of  $n$  distinct eigenvalues  $\dim \mathcal{N}(\mathbf{A} - \lambda \mathbf{I}) = 1$ , we conclude that there is  $\xi \in \mathbb{K}$ :  $\mathbf{B}\mathbf{v} = \xi \mathbf{v}$ ,  $\mathbf{v}$  is an eigenvector of  $\mathbf{B}$ . Since the eigenvector of  $\mathbf{A}$  span  $\mathbb{K}^n$ , there cannot be eigenvectors of  $\mathbf{B}$  that are not eigenvectors of  $\mathbf{A}$ .

Moreover, there is a basis of  $\mathbb{K}^n$  consisting of eigenvectors of  $\mathbf{B}$ ;  $\mathbf{B}$  can be diagonalized.  $\square$

Next, by straightforward calculation one verifies that every circulant matrix commutes with the unitary and circulant **cyclic permutation matrix**

$$\mathbf{S} = \begin{bmatrix} 0 & 0 & 0 & \cdots & \cdots & 0 & 1 \\ 1 & 0 & 0 & & & \vdots & 0 \\ 0 & 1 & 0 & \ddots & & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & & \vdots \\ & & & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & & \cdots & 0 & 1 & 0 \end{bmatrix}. \quad (4.2.0.4)$$

As a unitary matrix  $\mathbf{S}$  can be diagonalized. Observe that  $\mathbf{S}^n - \mathbf{I} = \mathbf{O}$ , that is the minimal polynomial of  $\mathbf{S}$  is  $\xi \mapsto \xi^n - 1$ , which is irreducible, because it has  $n$  distinct roots (of unity). Therefore, by Thm. 4.2.0.3,  $\mathbf{S}$  has  $n$  different eigenvalues and every eigenvector of  $\mathbf{S}$  is also eigenvector of any circulant matrix.  $\square$

**Remark 4.2.0.5 (Why using  $\mathbb{K} = \mathbb{C}$ ?)** In Exp. 4.2.0.1 we saw that we get *complex* eigenvalues/eigenvectors for general circulant matrices. More generally, in many cases real matrices can be diagonalized only in  $\mathbb{C}$ , which is the ultimate reason for the importance of complex numbers.

Complex numbers also allow an elegant handling on trigonometric functions: recall from analysis the unified treatment of trigonometric functions via the *complex exponential function*

$$\exp(it) = \cos(t) + i \sin(t), \quad t \in \mathbb{R}.$$

**C!** The field of complex numbers  $\mathbb{C}$  is the *natural framework* for the analysis of linear, time-invariant filters, and the development of algorithms for circulant matrices.  $\square$

**§4.2.0.6 (Eigenvectors of circulant matrices)** Now we verify by direct computations that circulant matrices all have a particular set of eigenvectors. This will entail computing in  $\mathbb{C}$ , cf. Rem. 4.2.0.15.

notation: *n*th root of unity  $\omega_n := \exp(-2\pi i/n) = \cos(2\pi/n) - i \sin(2\pi/n)$ ,  $n \in \mathbb{N}$

satisfies  $\bar{\omega}_n = \omega_n^{-1}$  ,  $\boxed{\omega_n^n = 1}$  ,  $\omega_n^{n/2} = -1$  ,  $\omega_n^k = \omega_n^{k+n} \quad \forall k \in \mathbb{Z}$  , (4.2.0.7)

$$\sum_{k=0}^{n-1} \omega_n^{kj} = \begin{cases} n & , \text{if } j = 0 \pmod{n}, \\ 0 & , \text{if } j \neq 0 \pmod{n}. \end{cases} \quad (4.2.0.8)$$

(4.2.0.8) is a simple consequence of the geometric sum formula

$$\begin{aligned} \sum_{k=0}^{n-1} q^k &= \frac{1-q^n}{1-q} \quad \forall q \in \mathbb{C} \setminus \{1\}, \quad n \in \mathbb{N}. \\ \Rightarrow \sum_{k=0}^{n-1} \omega_n^{kj} &= \frac{1-\omega_n^{nj}}{1-\omega_n^j} = \frac{1-\exp(-2\pi ij)}{1-\exp(-2\pi ij/n)} = 0, \end{aligned} \quad (4.2.0.9)$$

because  $\exp(-2\pi ij) = \omega_n^{nj} = (\omega_n^j)^j = 1$  for all  $j \in \mathbb{Z}$ .



In expressions like  $\omega_n^{kl}$  the term “*kl*” will always designate an exponent and will never play the role of a superscript.

Now we want to confirm the conjecture gleaned from Exp. 4.2.0.1 that vectors with powers of roots of unity are eigenvectors for any circulant matrix. We do this by simple and straightforward computations:

Consider: circulant matrix  $\mathbf{C} \in \mathbb{C}^{n,n}$  circulant matrix ( $\rightarrow$  Def. 4.1.0.38), with  $c_{ij} = u_{i-j}$ , for a  $n$ -periodic sequence  $(u_k)_{k \in \mathbb{Z}}, u_k \in \mathbb{C}$

We “guess” an eigenvector:  $\mathbf{v}_k \in \mathbb{C}^n$  with  $\mathbf{v}_k := [\omega_n^{jk}]_{j=0}^{n-1} \in \mathbb{C}^n$ ,  $k \in \{0, \dots, n-1\}$ .  
 $(u_{j-l}\omega_n^{lk})_{l \in \mathbb{Z}}$  is  $n$ -periodic!

$$\begin{aligned} (\mathbf{C}\mathbf{v}_k)_j &= \sum_{l=0}^{n-1} u_{j-l}\omega_n^{lk} = \sum_{l=j-n+1}^j u_{j-l}\omega_n^{lk} \\ &= \sum_{l=0}^{n-1} u_l\omega_n^{(j-l)k} = \omega_n^{jk} \sum_{l=0}^{n-1} u_l\omega_n^{-lk} = \lambda_k \cdot \omega_n^{jk} = \lambda_k \cdot (\mathbf{v}_k)_j. \end{aligned} \quad (4.2.0.10)$$

↑  
change of summation index      ↗  
independent of  $j$ !



$\mathbf{v}_k$  is eigenvector of  $\mathbf{C}$  for eigenvalue  $\lambda_k = \sum_{l=0}^{n-1} u_l \omega_n^{-lk}$ .

The set  $\{\mathbf{v}_0, \dots, \mathbf{v}_{n-1}\} \subset \mathbb{C}^n$  provides the so-called *orthogonal trigonometric basis* of  $\mathbb{C}^n$  = *eigen-vector basis* for circulant matrices

$$\{\mathbf{v}_0, \dots, \mathbf{v}_{n-1}\} = \left\{ \begin{bmatrix} \omega_n^0 \\ \vdots \\ \vdots \\ \omega_n^0 \end{bmatrix}, \begin{bmatrix} \omega_n^0 \\ \omega_n^1 \\ \vdots \\ \vdots \\ \omega_n^{n-1} \end{bmatrix}, \dots, \begin{bmatrix} \omega_n^0 \\ \omega_n^{n-2} \\ \omega_n^{2(n-2)} \\ \vdots \\ \vdots \\ \omega_n^{(n-1)(n-2)} \end{bmatrix}, \begin{bmatrix} \omega_n^0 \\ \omega_n^{n-1} \\ \omega_n^{2(n-1)} \\ \vdots \\ \vdots \\ \omega_n^{(n-1)^2} \end{bmatrix} \right\}. \quad (4.2.0.11)$$

From (4.2.0.8) we can conclude orthogonality of the basis vectors by straightforward computations:

$$\mathbf{v}_k := (\omega_n^{jk})_{j=0}^{n-1} \in \mathbb{C}^n: \quad \mathbf{v}_k^H \mathbf{v}_m = \sum_{j=0}^{n-1} \omega_n^{-jk} \omega_n^{jm} = \sum_{j=0}^{n-1} \omega_n^{(m-k)j} \stackrel{(4.2.0.8)}{=} 0, \text{ if } k \neq m. \quad (4.2.0.12)$$

The matrix effecting the change of basis  $\text{trigonometrical basis} \rightarrow \text{standard basis}$  is called the **Fourier-matrix**

$$\mathbf{F}_n = \begin{bmatrix} \omega_n^0 & \omega_n^0 & \cdots & \omega_n^0 \\ \omega_n^0 & \omega_n^1 & \cdots & \omega_n^{n-1} \\ \omega_n^0 & \omega_n^2 & \cdots & \omega_n^{2n-2} \\ \vdots & \vdots & & \vdots \\ \omega_n^0 & \omega_n^{n-1} & \cdots & \omega_n^{(n-1)^2} \end{bmatrix} = [\omega_n^{lj}]_{l,j=0}^{n-1} \in \mathbb{C}^{n,n}. \quad (4.2.0.13)$$

#### Lemma 4.2.0.14. Properties of Fourier matrix

The scaled Fourier-matrix  $\frac{1}{\sqrt{n}} \mathbf{F}_n$  is unitary ( $\rightarrow$  Def. 6.2.1.2):  $\mathbf{F}_n^{-1} = \frac{1}{n} \mathbf{F}_n^H = \frac{1}{n} \bar{\mathbf{F}}_n$ .

*Proof.* The lemma is immediate from (4.2.0.12) and (4.2.0.8), because

$$(\mathbf{F}_n \mathbf{F}_n^H)_{l,j} = \sum_{k=0}^{n-1} \omega_n^{(l-1)k} \overline{\omega_n^{(j-1)k}} = \sum_{k=0}^{n-1} \omega_n^{(l-1)k} \omega_n^{-(j-1)k} = \sum_{k=0}^{n-1} \omega_n^{k(l-j)}, \quad 1 \leq l, j \leq n.$$

□

**Remark 4.2.0.15 (Spectrum of Fourier matrix)** We draw a conclusion from the properties stated in Lemma 4.2.0.14:

$$\frac{1}{n^2} \mathbf{F}_n^4 = I \Rightarrow \sigma\left(\frac{1}{\sqrt{n}} \mathbf{F}_n\right) \subset \{1, -1, i, -i\},$$

because, if  $\lambda \in \mathbb{C}$  is an eigenvalue of  $\mathbf{F}_n$ , then there is an eigenvector  $\mathbf{x} \in \mathbb{C}^n \setminus \{0\}$  such that  $\mathbf{F}_n \mathbf{x} = \lambda \mathbf{x}$ , see Def. 9.1.0.1. □

#### Lemma 4.2.0.16. Diagonalization of circulant matrices ( $\rightarrow$ Def. 4.1.0.38)

For any circulant matrix  $\mathbf{C} \in \mathbb{K}^{n,n}$ ,  $c_{ij} = u_{i-j}$ ,  $(u_k)_{k \in \mathbb{Z}}$   $n$ -periodic sequence, holds true

$$\mathbf{C} \bar{\mathbf{F}}_n = \bar{\mathbf{F}}_n \operatorname{diag}(d_1, \dots, d_n), \quad \mathbf{d} = \mathbf{F}_n [u_0, \dots, u_{n-1}]^\top.$$

*Proof.* Straightforward computation, see (4.2.0.10). □

Conclusion (from  $\bar{\mathbf{F}}_n = n\mathbf{F}_n^{-1}$ ):

$$\mathbf{C} = \mathbf{F}_n^{-1} \operatorname{diag}(d_1, \dots, d_n) \mathbf{F}_n . \quad (4.2.0.17)$$

Lemma 4.2.0.16, (4.2.0.17)  $\Rightarrow$  multiplication with Fourier-matrix will be crucial operation in algorithms for circulant matrices and discrete convolutions.

Therefore this operation has been given a special name:

#### Definition 4.2.0.18. Discrete Fourier transform (DFT)

The linear map  $\mathcal{F}_n : \mathbb{C}^n \mapsto \mathbb{C}^n$ ,  $\mathcal{F}_n(\mathbf{y}) := \mathbf{F}_n \mathbf{y}$ ,  $\mathbf{y} \in \mathbb{C}^n$ , is called **discrete Fourier transform (DFT)**, i.e. for  $\mathbf{c} := \mathcal{F}_n(\mathbf{y})$

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{kj} , \quad k = 0, \dots, n-1 . \quad (4.2.0.19)$$

Recall the convention also adopted for the discussion of the DFT: vector indexes range from 0 to  $n - 1$ !

Terminology:  $\mathbf{c} = \mathbf{F}_n \mathbf{y}$  is also called the (discrete) Fourier transform of  $\mathbf{y}$

From  $\mathbf{F}_n^{-1} = \frac{1}{n} \bar{\mathbf{F}}_n$  ( $\rightarrow$  Lemma 4.2.0.14) we find the **inverse discrete Fourier transform**:

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{kj} \Leftrightarrow y_j = \frac{1}{n} \sum_{k=0}^{n-1} c_k \omega_n^{-kj} \quad (4.2.0.20)$$

#### §4.2.0.21 (Discrete Fourier transform in EIGEN and MATLAB)

- EIGEN-functions for discrete Fourier transform (and its inverse):

DFT:  $\mathbf{c} = \text{fft.fwd}(\mathbf{y}) \leftrightarrow$  inverse DFT:  $\mathbf{y} = \text{fft.inv}(\mathbf{c})$ ;

Before using `fft`, remember to:

1. **# include** <unsupported/Eigen/FFT>
2. Instantiate helper class `Eigen::FFT<double> fft;`  
(The template argument should always be **double**.)

#### C++ code 4.2.0.22: Demo: discrete Fourier transform in EIGEN → GITLAB

```

2 int main() {
3     using Comp = complex<double>;
4     const VectorXcd::Index n = 5;
5     VectorXcd y(n), c(n), x(n);
6     y << Comp(1,0),Comp(2,1),Comp(3,2),Comp(4,3),Comp(5,4);
7     FFT<double> fft; // DFT transform object
8     c = fft.fwd(y); // DTF of y, see Def. 4.2.0.18
9     x = fft.inv(c); // inverse DFT of c, see (4.2.0.20)
10

```

```

11   cout << "y = " << y.transpose() << endl
12   << "c = " << c.transpose() << endl
13   << "x = " << x.transpose() << endl;
14   return 0;
15 }
```

- MATLAB-functions for discrete Fourier transform (and its inverse):

DFT: `c=fft(y)`  $\leftrightarrow$  inverse DFT: `y=ifft(c);`

↳

### 4.2.1 Discrete Convolution via DFT

Coding the formula from Def. 4.1.0.33 one would code discrete periodic convolution as follows:

**C++11 code 4.2.1.1: Discrete periodic convolution: straightforward implementation**  
[→ GITLAB](#)

```

2 VectorXcd pconv(const VectorXcd& u, const VectorXcd& x) {
3     using idx_t = VectorXcd::Index; // may be unsigned !
4     const idx_t n = x.size();
5     VectorXcd z = VectorXcd::Zero(n);
6     // Need signed indices when differences are formed
7     for (long k = 0; k < n; ++k) {
8         for (long j = 0; j < n; ++j) {
9             long ind = (k - j < 0 ? n + k - j : k - j);
10            z(k) += u(ind)*x(j);
11        }
12    }
13 }
```

In § 4.1.0.37 we have seen that periodic convolution ( $\rightarrow$  Def. 4.1.0.33) amounts to multiplication with a circulant matrix. In addition, (4.2.0.17) reduces multiplication with a circulant matrix to two multiplications with the Fourier matrix  $\mathbf{F}_n$  (= DFT) and (componentwise) scaling operations.

Summary:

$$\text{discrete periodic convolution } z_k = \sum_{j=0}^{n-1} u_{k-j} x_j \quad (\rightarrow \text{Def. 4.1.0.33}), \quad k = 0, \dots, n-1$$

↑

$$\text{multiplication with circulant matrix } (\rightarrow \text{Def. 4.1.0.38}) \quad \mathbf{z} = \mathbf{Cx}, \quad \mathbf{C} := [u_{i-j}]_{i,j=1}^n.$$

Idea:

$$(4.2.0.17) \quad \Rightarrow \quad \mathbf{z} = \mathbf{F}_n^{-1} \operatorname{diag}(\mathbf{F}_n \mathbf{u}) \mathbf{F}_n \mathbf{x}$$

This formula is usually referred to as *convolution theorem*:

### Theorem 4.2.1.2. Convolution theorem

The discrete periodic convolution  $*_n$  between  $n$ -dimensional vectors  $\mathbf{u}$  and  $\mathbf{x}$  is equal to the inverse DFT of the component-wise product between the DFTs of  $\mathbf{u}$  and  $\mathbf{x}$ ; i.e.:

$$(\mathbf{u}) *_n (\mathbf{x}) := \sum_{j=0}^{n-1} u_{k-j} x_j = \mathbf{F}_n^{-1} [(\mathbf{F}_n \mathbf{u})_j (\mathbf{F}_n \mathbf{x})_j]_{j=1}^n.$$

### C++11 code 4.2.1.3: Discrete periodic convolution: DFT implementation → GITLAB

```
2 VectorXcd pconvfft(const VectorXcd& u, const VectorXcd& x) {
3     Eigen::FFT<double> fft;
4     VectorXcd tmp = (fft.fwd(u)).cwiseProduct(fft.fwd(x));
5     return fft.inv(tmp);
6 }
```

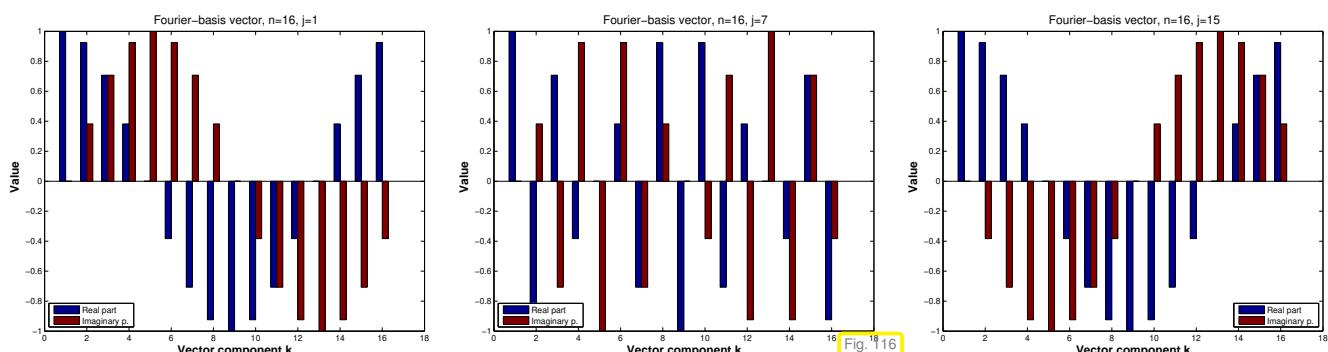
In Rem. 4.1.0.40 we learned that the discrete convolution of  $n$ -vectors (→ Def. 4.1.0.22) can be accomplished by the *periodic* discrete convolution of  $2n - 1$ -vectors (obtained by zero padding, see Rem. 4.1.0.40):

### C++11 code 4.2.1.4: Implementation of discrete convolution (→ Def. 4.1.0.22) based on periodic discrete convolution → GITLAB

```
2 VectorXcd myconv(const VectorXcd& h, const VectorXcd& x) {
3     const long n = h.size();
4     // Zero padding, cf. (4.1.0.41)
5     VectorXcd hp(2*n - 1), xp(2*n - 1);
6     hp << h, VectorXcd::Zero(n - 1);
7     xp << x, VectorXcd::Zero(n - 1);
8     // Periodic discrete convolution of length 2n - 1, ???
9     return pconvfft(hp, xp);
10 }
```

## 4.2.2 Frequency filtering via DFT

The trigonometric basis vectors, when interpreted as time-periodic signals, represent harmonic oscillations. This is illustrated when plotting some vectors of the trigonometric basis ( $n = 16$ ):



“slow oscillation/low frequency”

“fast oscillation/high frequency”

“slow oscillation/low frequency”

► Dominant coefficients of a signal after transformation to trigonometric basis indicate dominant frequency components.

Terminology: coefficients of a signal w.r.t. trigonometric basis = signal in **frequency domain**, original signal = **time domain**.

Recall: DFT (4.2.0.19) and inverse DFT (4.2.0.20)

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{kj} \quad \Leftrightarrow \quad y_j = \frac{1}{n} \sum_{k=0}^{n-1} c_k \omega_n^{-kj} \quad (4.2.0.20)$$

Consider  $y_k \in \mathbb{R}$   $\Rightarrow c_k = \bar{c}_{n-k}$ , because  $\omega_n^{kj} = \bar{\omega}_n^{(n-k)j}$ , and  $n = 2m + 1$

$$\begin{aligned} \blacktriangleright ny_j &= c_0 + \sum_{k=1}^m c_k \omega_n^{-kj} + \sum_{k=m+1}^{2m} c_k \omega_n^{-kj} = c_0 + \sum_{k=1}^m c_k \omega_n^{-kj} + c_{n-k} \omega_n^{(k-n)j} \\ &= c_0 + 2 \sum_{k=1}^m \operatorname{Re}(c_k) \cos(2\pi kj/n) + \operatorname{Im}(c_k) \sin(2\pi kj/n), \end{aligned}$$

since  $\omega_n^\ell = \cos(2\pi \ell/n) + i \sin(2\pi \ell/n)$ .

$\Rightarrow |c_k|, |c_{n-k}|$  measures the strength with which an oscillation with frequency  $k$  is represented in the signal,  $0 \leq k \leq \lfloor \frac{n}{2} \rfloor$ .

**EXAMPLE 4.2.2.1 (Frequency identification with DFT)** Extraction of characteristic frequencies from a distorted discrete periodical signal, generated by the following C++ code:

#### C++ code 4.2.2.2: Generation of noisy sinusoidal signal → [GITLAB](#)

```

2 VectorXd signalgen() {
3     const int N = 64;
4     const ArrayXd t = ArrayXd::LinSpaced(N, 0, N);
5     const VectorXd x = ((2*M_PI/N*t).sin() + (14*M_PI/N*t).sin()).matrix();
6     return x + VectorXd::Random(N);
7 }
```

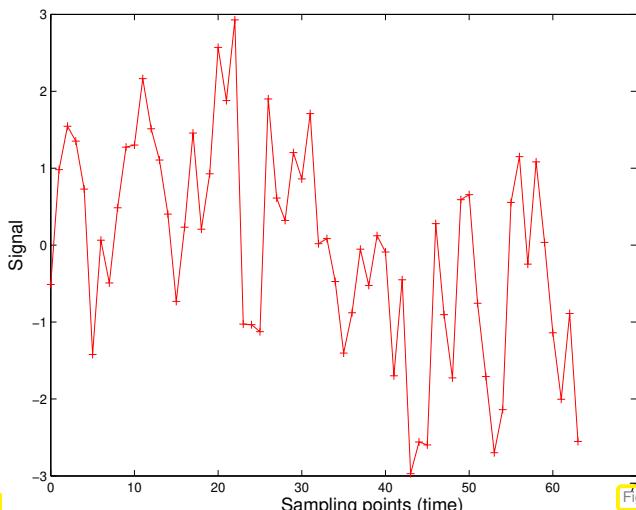


Fig. 117

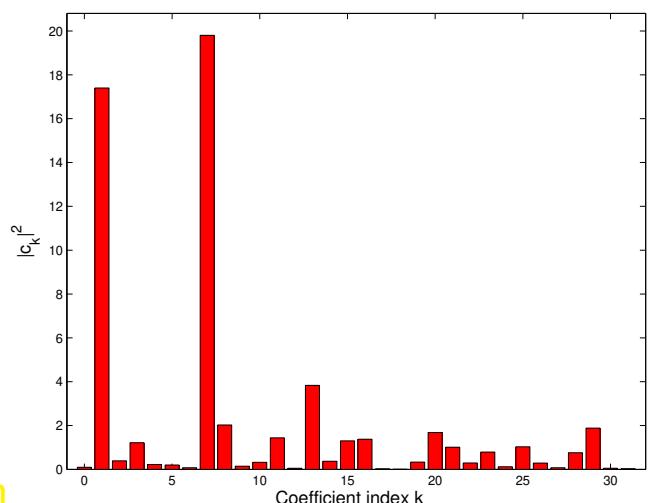


Fig. 118

Fig. 118 was generated by the following MATLAB code. A corresponding C++ code is also available → [GITLAB](#).

#### MATLAB code 4.2.2.3: Frequency extraction → Fig. 118

```
1 c = fft(y); p = (c.*conj(c))/64;
```

```

2
3 figure('Name','power spectrum');
4 bar(0:31,p(1:32),'r');
5 set(gca,'fontsize',14);
6 axis([-1 32 0 max(p)+1]);
7 xlabel('{\bf index k of Fourier coefficient}', 'Fontsize', 14);
8 ylabel('{\bf |c_k|^2}', 'Fontsize', 14);

```

We observe that frequencies present in unperturbed signal become evident in frequency domain, whereas it is hard to tell them from the time-domain signal. □

#### EXAMPLE 4.2.2.4 (Detecting periodicity in data)

DFT: a computer's eye for periodic patterns in data

The following C++ code processes actual web search data and performs a frequency analysis using DFT:

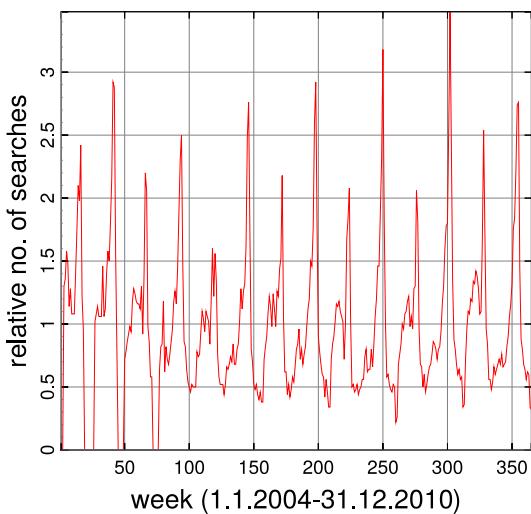
#### C++ code 4.2.2.5: Extraction of periodic patterns by DFT → [GITLAB](#)

```

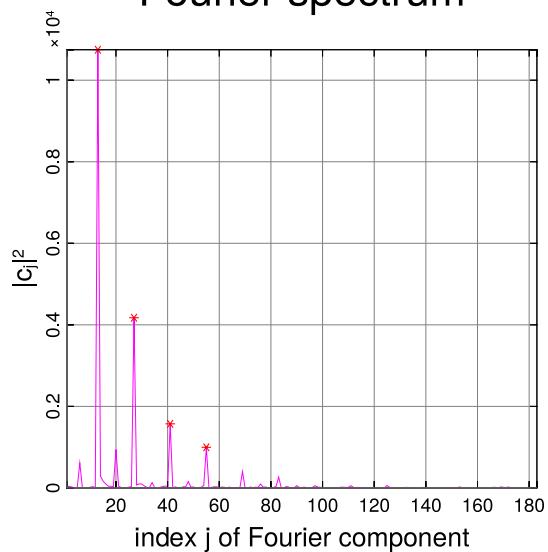
2 VectorXd x = read("trend.dat");
3 const int n = x.size();
4
5 mgl::Figure trend;
6 trend.plot(x, "r");
7 trend.title("Google: 'Vorlesungsverzeichnis'");
8 trend.grid();
9 trend.xlabel("week (1.1.2004–31.12.2010)");
10 trend.ylabel("relative no. of searches");
11 trend.save("searchdata");
12
13 Eigen::FFT<double> fft;
14 VectorXcd c = fft.fwd(x);
15 VectorXd p = c.cwiseProduct(c.conjugate()).real()
16     .segment(2, std::floor((n+1)/2));
17
18 // plot power spectrum
19 mgl::Figure fig; fig.plot(p, "m"); fig.grid();
20 fig.title("Fourier spectrum");
21 fig.xlabel("index j of Fourier component");
22 fig.ylabel("|c_j|^2");
23
24 // mark 4 highest peaks with red star
25 VectorXd p_sorted = p;
26 // sort descending
27 std::sort(p_sorted.data(), p_sorted.data() + p.size(), std::greater<double>());
28 for (int i = 0; i < 4; ++i) {
29     // get position of p_sorted[i]
30     int idx = std::find(p.data(), p.data() + p.size(), p_sorted[i]) - p.data();
31     std::vector<double> ind = { double(idx + 1) }; // save in vectors
32     std::vector<double> val = { p_sorted[i] }; // to be able to plot
33     fig.plot(ind, val, "r*");
34 }
35 fig.save("fourierdata");

```

## Google: 'Vorlesungsverzeichnis'

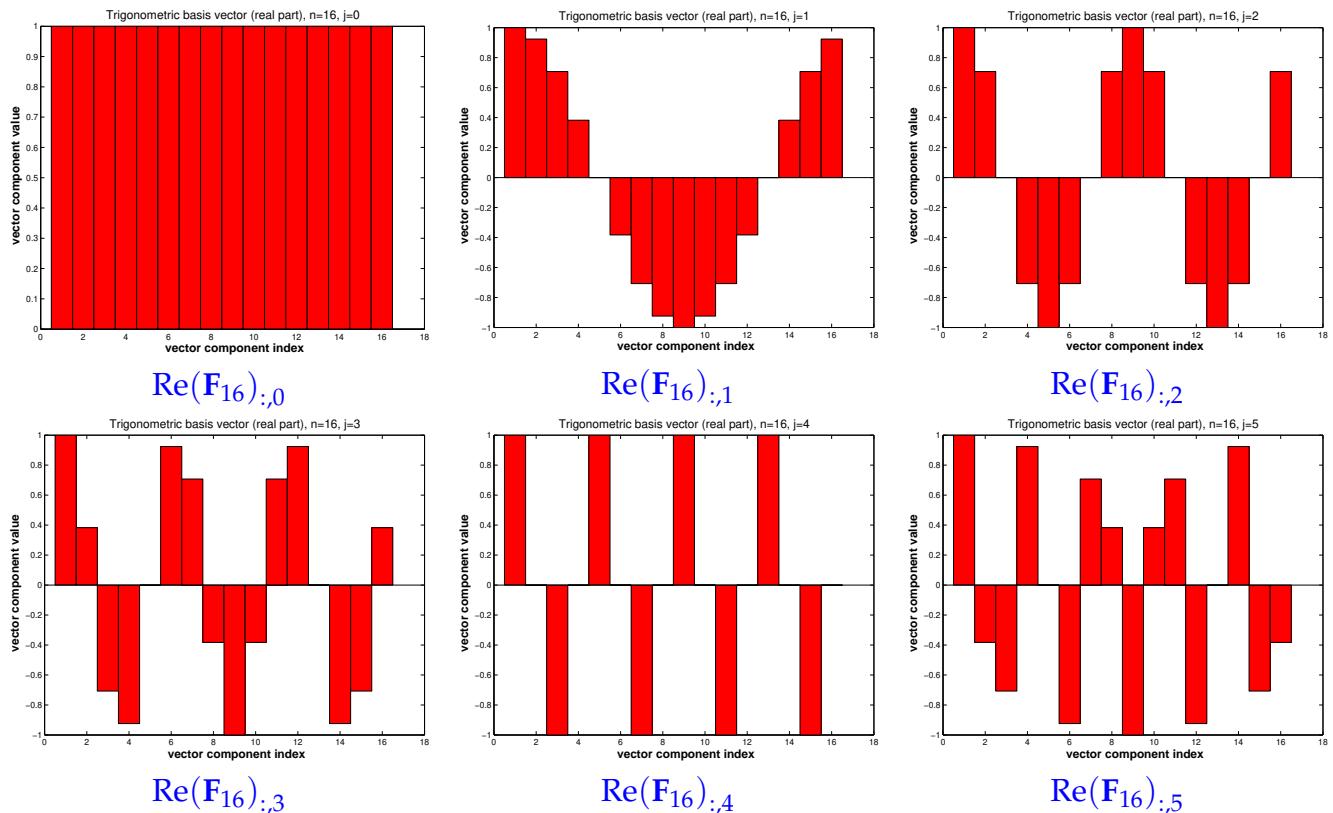


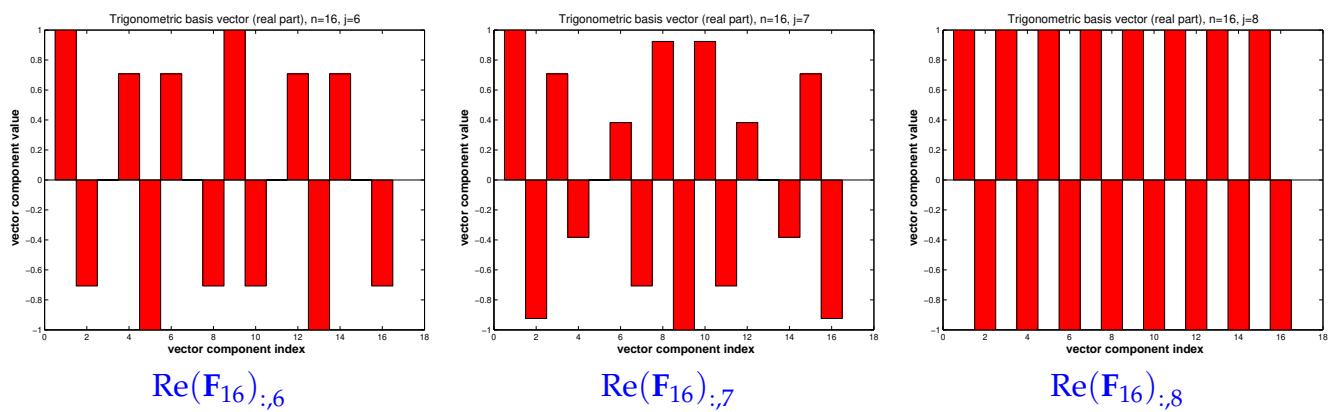
## Fourier spectrum



Pronounced peaks in the power spectrum point to periodic structure of the data. Location of peaks tells lengths of dominant periods.

**Remark 4.2.2.6 (“Low” and “high” frequencies)** Plots of real parts of trigonometric basis vectors  $(\mathbf{F}_n)_{:,j}$  (= columns of Fourier matrix),  $n = 16$ .



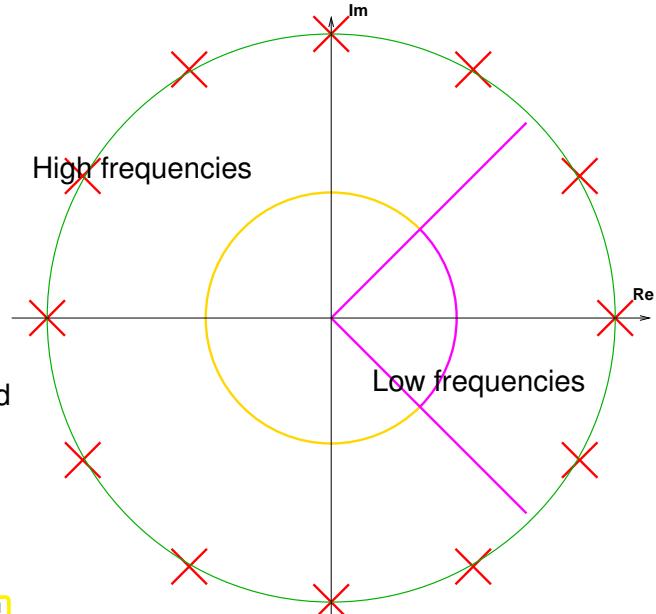


By elementary trigonometric identities:

$$\begin{aligned}\text{Re}(\mathbf{F}_n)_{:,j} &= \left( \text{Re } \omega_n^{(j-1)k} \right)_{k=0}^{n-1} \\ &= (\text{Re } \exp(2\pi i(j-1)k/n))_{k=0}^{n-1} \\ &= (\cos(2\pi(j-1)x))_{x=0,\frac{1}{n},\dots,1-\frac{1}{n}}.\end{aligned}$$

Slow oscillations/low frequencies  $\leftrightarrow j \approx 1$  and  $j \approx n$ .

Fast oscillations/high frequencies  $\leftrightarrow j \approx n/2$ .



► Frequency filtering of real discrete periodic signals by suppressing certain “Fourier coefficients”.

#### C++11-code 4.2.2.7: DFT-based frequency filtering → GITLAB

```

2 void freqfilter(const VectorXd& y, int k,
3                 VectorXd& low, VectorXd& high) {
4     const VectorXd::Index n = y.size();
5     if (n%2 != 0)
6         throw std::runtime_error("Even vector length required!");
7     const VectorXd::Index m = y.size()/2;
8
9     Eigen::FFT<double> fft; // DFT helper object
10    VectorXcd c = fft.fwd(y); // Perform DFT of input vector
11
12    VectorXcd clow = c;
13    // Set high frequency coefficients to zero, Fig. 121
14    for (int j = -k; j <= +k; ++j) clow(m+j) = 0;
15    // (Complementary) vector of high frequency coefficients
16    VectorXcd chigh = c - clow;
17
18    // Recover filtered time-domain signals
19    low = fft.inv(clow).real();
20    high = fft.inv(chigh).real();
21}

```

(It can be optimised exploiting  $y_j \in \mathbb{R}$  and  $c_{n/2-k} = \bar{c}_{n/2+k}$ )

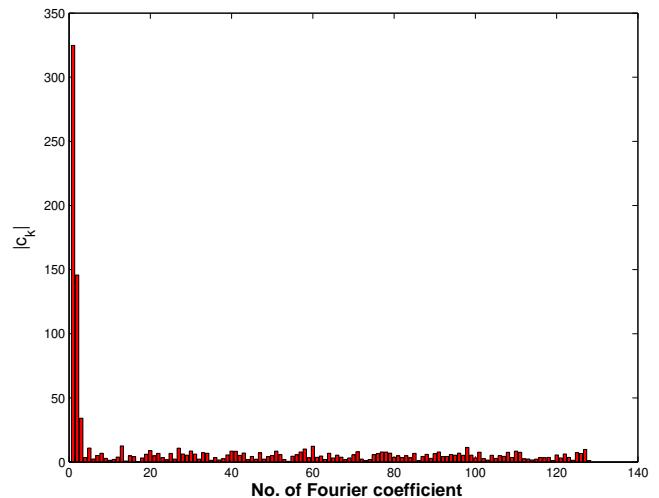
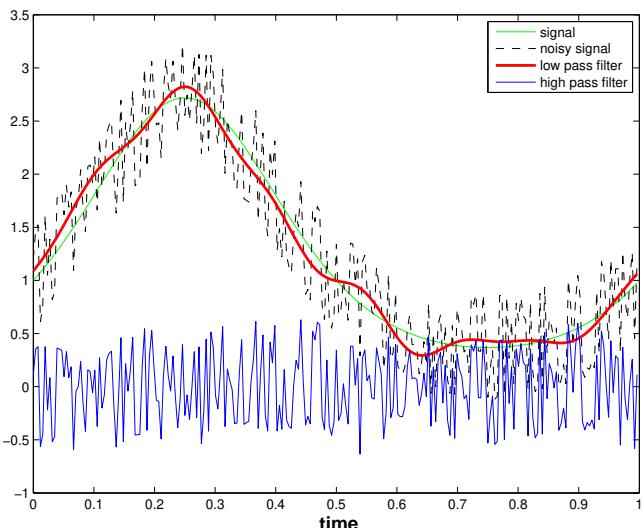
Map  $y \mapsto \text{low}$  (in Code 4.2.2.7)  $\hat{=}$  **low pass filter**.

Map  $y \mapsto \text{high}$  (in Code 4.2.2.7)  $\hat{=}$  **high pass filter**.  $\downarrow$

#### EXAMPLE 4.2.2.8 (Frequency filtering by DFT) Noisy signal:

```
n = 256; y = exp(sin(2*pi*((0:n-1)')/n)) + 0.5*sin(exp(1:n)');
```

Frequency filtering by Code 4.2.2.7 with  $k = 120$ .



Low pass filtering can be used for *denoising*, that is, the removal of high frequency perturbations of a signal.  $\downarrow$

#### EXAMPLE 4.2.2.9 (Sound filtering by DFT) Frequency filtering is ubiquitous in sound processing. Here we demonstrate it in MATLAB, which offers tools for audio processing.

##### MATLAB-code 4.2.2.10: DFT based sound compression

```

1 % Low pass sound filtering by DFT
2
3 % Read sound data
4 [y,Fs,nbits] = wavread('hello.wav');
5 sound(y,Fs);
6
7 n = length(y);
8 fprintf('Read wav File: %d samples, rate = %d/s, nbits = %d\n',
9     n,Fs,nbits);
10 k = 1; s{k} = y; leg{k} = 'Sampled signal';
11
12 C = fft(y);
13
14 figure('name','sound signal');
15 plot((22000:44000)/Fs,s{1}(22000:44000),'r-');
16 title('samples sound signal','fontsize',14);
17 xlabel('{\bf time[s]}','fontsize',14);
18 ylabel('{\bf sound pressure}','fontsize',14);
19 grid on;
```

```

19 % print -depsc2 '../PICTURES/soundsignal.eps';
20
21
22 figure('name','sound frequencies');
23 plot(1:n,abs(c).^2,'m-');
24 title('power spectrum of sound signal','fontsize',14);
25 xlabel('{\bf index k of Fourier coefficient}','fontsize',14);
26 ylabel('{\bf |c_k|^2}','fontsize',14);
27 grid on;
28
29 % print -depsc2 '../PICTURES/soundpower.eps';
30
31 figure('name','sound frequencies');
32 plot(1:3000,abs(c(1:3000)).^2,'b-');
33 title('low frequency power spectrum','fontsize',14);
34 xlabel('{\bf index k of Fourier coefficient}','fontsize',14);
35 ylabel('{\bf |c_k|^2}','fontsize',14);
36 grid on;
37
38 % print -depsc2 '../PICTURES/soundlowpower.eps';
39
40 for m=[1000,3000,5000]
41
42 % Low pass filtering
43 cf = zeros(1,n);
44 cf(1:m) = c(1:m); cf(n-m+1:end) = c(n-m+1:end);
45
46 % Reconstruct filtered signal
47 yf = ifft(cf);
48 % No idea why this is necessary
49 wavwrite(yf,Fs,nbits,sprintf('hellof%d.wav',m));
50 cy = wavread(sprintf('hellof%d.wav',m));
51 sound(cy,Fs,nbits);
52
53 k = k+1;
54 s{k} = real(yf);
55 leg{k} = sprintf('cutt-off = %d',m);
56 end
57
58 % Plot original signal and filtered signals
59 figure('name','sound filtering');
60 plot((30000:32000)/Fs,s{1}(30000:32000),'r-',...
61 (30000:32000)/Fs,s{2}(30000:32000),'b--',...
62 (30000:32000)/Fs,s{3}(30000:32000),'m--',...
63 (30000:32000)/Fs,s{2}(30000:32000),'k--');
64 xlabel('{\bf time[s]}','fontsize',14);
65 ylabel('{\bf sound pressure}','fontsize',14);
66 legend(leg,'location','southeast');
67
68 % print -depsc2 '../PICTURES/soundfiltered.eps';

```

### MATLAB-code 4.2.2.11: DFT based low pass frequency filtering of sound

```

1 [y, sf, nb] = wavread('hello.wav');
2 C = fft(y); C(m+1:end-m) = 0;
3 wavwrite(ifft(C), sf, nb, 'filtered.wav');
```

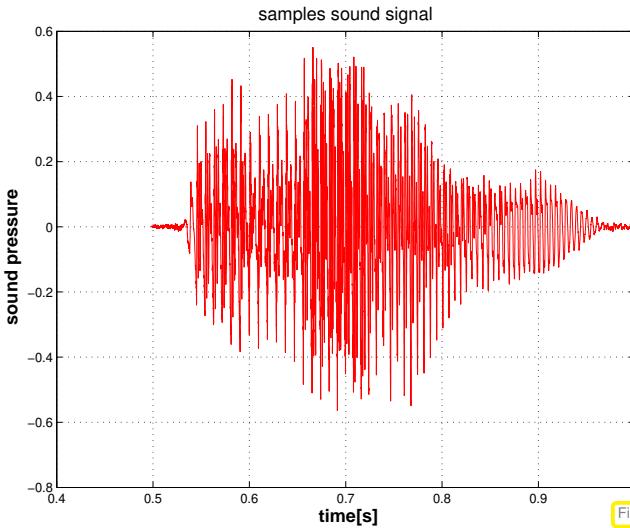


Fig. 122

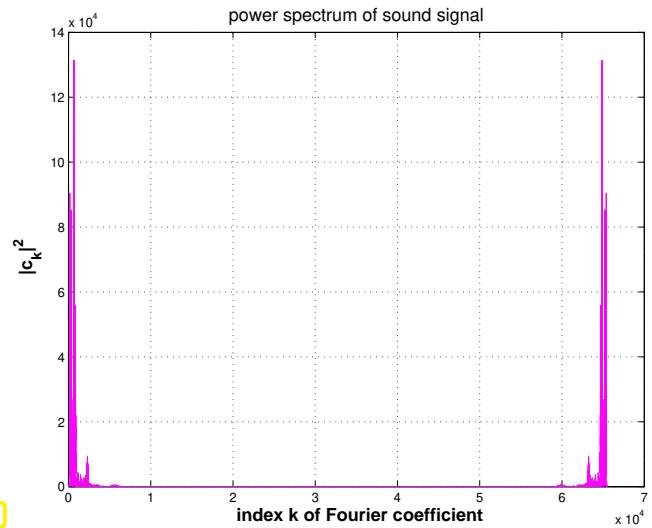


Fig. 123

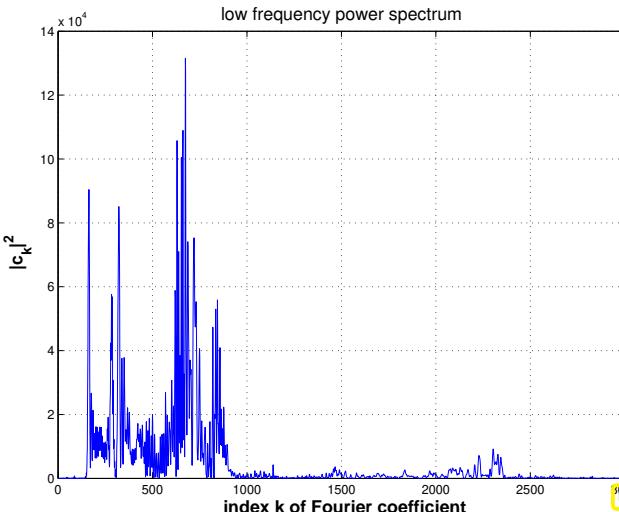


Fig. 124

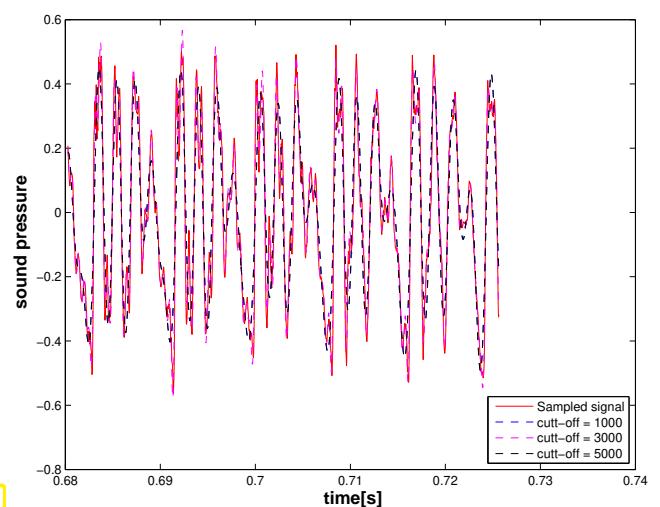


Fig. 125

The **power spectrum** of a signal  $\mathbf{y} \in \mathbb{C}^n$  is the vector  $(|c_j|^2)_{j=0}^{n-1}$ , where  $\mathbf{c} = \mathbf{F}_n \mathbf{y}$  is the discrete Fourier transform of  $\mathbf{y}$ . □

### 4.2.3 Real DFT

Every time-discrete signal obtained from sampling a time-dependent physical quantity will yield a **real** vector. Of course, a real vector contains only half the information compared to complex vector of the same length. We aim to exploit this for a more efficient implementation of DFT.

Task: Efficient implementation of DFT (Def. 4.2.0.18)  $(c_0, \dots, c_{n-1})$  for *real* coefficients  $(y_0, \dots, y_{n-1})^\top \in \mathbb{R}^n$ ,  $n = 2m$ ,  $m \in \mathbb{N}$ .

If  $y_j \in \mathbb{R}$  in DFT formula (4.2.0.19), we obtain redundant output

$$\omega_n^{(n-k)j} = \bar{\omega}_n^{kj}, \quad k = 0, \dots, n-1,$$

$$\Rightarrow \bar{c}_k = \sum_{j=0}^{n-1} y_j \bar{\omega}_n^{kj} = \sum_{j=0}^{n-1} y_j \omega_n^{(n-k)j} = c_{n-k}, \quad k = 1, \dots, n-1.$$

Idea: map  $\mathbf{y} \in \mathbb{R}^n$ , to  $\mathbb{C}^m$  and use DFT of length  $m$ .

$$h_k = \sum_{j=0}^{m-1} (y_{2j} + iy_{2j+1}) \omega_m^{jk} = \boxed{\sum_{j=0}^{m-1} y_{2j} \omega_m^{jk}} + i \cdot \boxed{\sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}}, \quad (4.2.3.1)$$

$$\bar{h}_{m-k} = \sum_{j=0}^{m-1} \overline{y_{2j} + iy_{2j+1}} \bar{\omega}_m^{j(m-k)} = \boxed{\sum_{j=0}^{m-1} y_{2j} \omega_m^{jk}} - i \cdot \boxed{\sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}}. \quad (4.2.3.2)$$

$$\Rightarrow \boxed{\sum_{j=0}^{m-1} y_{2j} \omega_m^{jk}} = \frac{1}{2}(h_k + \bar{h}_{m-k}), \quad \boxed{\sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}} = -\frac{1}{2}i(h_k - \bar{h}_{m-k}).$$

Use simple identities for roots of unity:

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{jk} = \boxed{\sum_{j=0}^{m-1} y_{2j} \omega_m^{jk}} + \omega_n^k \cdot \boxed{\sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}}. \quad (4.2.3.3)$$

$$\Rightarrow \begin{cases} c_k = \frac{1}{2}(h_k + \bar{h}_{m-k}) - \frac{1}{2}i\omega_n^k(h_k - \bar{h}_{m-k}), & k = 0, \dots, m-1, \\ c_m = \operatorname{Re}\{h_0\} - \operatorname{Im}\{h_0\}, \\ c_k = \bar{c}_{n-k}, & k = m+1, \dots, n-1. \end{cases} \quad (4.2.3.4)$$

### C++11 code 4.2.3.5: DFT of real vectors of length $n/2 \rightarrow$ GITLAB

```

2 // Perform fft on a real vector y of even
3 // length and return (complex) coefficients in c
4 // Note: EIGEN's DFT method fwd() has this already implemented and
5 // we could also just call: c = fft.fwd(y);
6 void fftreal(const VectorXd& y, VectorXcd& c) {
7     const unsigned n = y.size(), m = n/2;
8     if (n % 2 != 0) { std::cout << "n must be even!\n"; return; }
9
10    // Step I: compute h from (4.2.3.1), (4.2.3.2)
11    std::complex<double> i(0,1); // Imaginary unit
12    VectorXcd yc(m);
13    for (unsigned j = 0; j < m; ++j)
14        yc(j) = y(2*j) + i*y(2*j + 1);
15
16    Eigen::FFT<double> fft;
17    VectorXcd d = fft.fwd(yc), h(m + 1);
18    h << d, d(0);
19
20    c.resize(n);
21    // Step II: implementation of (4.2.3.4)
22    for (unsigned k = 0; k < m; ++k) {
23        c(k) = (h(k) + std::conj(h(m-k)))/2. - i/2.*std::exp(-2.*k/n*M_PI*i)*(h(k) -
24            std::conj(h(m-k)));
25    }
26    c(m) = std::real(h(0)) - std::imag(h(0));

```

```

26   for (unsigned k = m+1; k < n; ++k) c(k) = std::conj(c(n-k));
27 }
```

#### 4.2.4 Two-dimensional DFT

In this we study the frequency decomposition of matrices. Due to the natural analogy

one-dimensional data (“signal”)  $\longleftrightarrow$  vector  $\mathbf{y} \in \mathbb{C}^n$ ,

two-dimensional data (“image”)  $\longleftrightarrow$  matrix.  $\mathbf{Y} \in \mathbb{C}^{m,n}$ ,

these techniques are of fundamental importance for [image processing](#).

**§4.2.4.1 (Matrix Fourier modes)** A two-dimensional trigonometric basis of  $\mathbb{C}^{m,n}$  is given by the [tensor product matrices](#)

$$\left\{ (\mathbf{F}_m)_{:,j} (\mathbf{F}_n)_{:,\ell}^\top, 1 \leq j \leq m, 1 \leq \ell \leq n \right\} \subset \mathbb{C}^{m,n}. \quad (4.2.4.2)$$

Let a matrix  $\mathbf{C} \in \mathbb{C}^{m,n}$  be given as a linear combination of these basis matrices with coefficients  $y_{j_1,j_2} \in \mathbb{C}$ ,  $0 \leq j_1 < m$ ,  $0 \leq j_2 < n$ :

$$\mathbf{C} = \sum_{j_1=0}^{m-1} \sum_{j_2=0}^{n-1} y_{j_1,j_2} (\mathbf{F}_m)_{:,j_1} (\mathbf{F}_n)_{:,\ell}^\top. \quad (4.2.4.3)$$

Then the entries of  $\mathbf{C}$  can be computed by two [nested](#) discrete Fourier transforms:

$$(\mathbf{C})_{k_1,k_2} = \sum_{j_1=0}^{m-1} \sum_{j_2=0}^{n-1} y_{j_1,j_2} \omega_m^{j_1 k_1} \omega_n^{j_2 k_2} = \sum_{j_1=0}^{m-1} \omega_m^{j_1 k_1} \left( \sum_{j_2=0}^{n-1} \omega_n^{j_2 k_2} y_{j_1,j_2} \right), \quad 0 \leq k_1 < m, 0 \leq k_2 < n.$$

□

The coefficients can also be regarded as entries of a matrix  $\mathbf{Y} \in \mathbb{C}^{m,n}$ . Thus we can rewrite the above expressions: for all  $0 \leq k_1 < m$ ,  $0 \leq k_2 < n$

$$(\mathbf{C})_{k_1,k_2} = \sum_{j_1=0}^{m-1} (\mathbf{F}_n(\mathbf{Y})_{j_1,:})_{k_2} \omega_m^{j_1 k_1} \quad \Rightarrow \quad \boxed{\mathbf{C} = \mathbf{F}_m (\mathbf{F}_n \mathbf{Y}^\top)^\top = \mathbf{F}_m \mathbf{Y} \mathbf{F}_n}. \quad (4.2.4.4)$$

This formula defines the [two-dimensional discrete Fourier transform](#) of the matrix  $\mathbf{Y} \in \mathbb{C}^{m,n}$ . By Lemma 4.2.0.14 we immediately get the inversion formula:

$$\mathbf{C} = \sum_{j_1=0}^{m-1} \sum_{j_2=0}^{n-1} y_{j_1,j_2} (\mathbf{F}_m)_{:,j_1} (\mathbf{F}_n)_{:,\ell}^\top \Rightarrow \boxed{\mathbf{Y} = \mathbf{F}_m^{-1} \mathbf{C} \mathbf{F}_n^{-1} = \frac{1}{mn} \bar{\mathbf{F}}_m \mathbf{C} \bar{\mathbf{F}}_n}. \quad (4.2.4.5)$$

#### C++ code 4.2.4.6: Two-dimensional discrete Fourier transform → [GITLAB](#)

```

2 template <typename Scalar>
3 void fft2(Eigen::MatrixXcd &C, const Eigen::MatrixBase<Scalar> &Y) {
4     using idx_t = Eigen::MatrixXcd::Index;
```

```

5   const idx_t m = Y.rows(), n=Y.cols();
6   C.resize(m,n);
7   Eigen::MatrixXcd tmp(m,n);
8
9   Eigen::FFT<double> fft; // Helper class for DFT
10  // Transform rows of matrix Y
11  for (idx_t k=0;k<m;k++) {
12    Eigen::VectorXcd tv(Y.row(k));
13    tmp.row(k) = fft.fwd(tv).transpose();
14  }
15
16  // Transform columns of temporary matrix
17  for (idx_t k=0;k<n;k++) {
18    Eigen::VectorXcd tv(tmp.col(k));
19    C.col(k) = fft.fwd(tv);
20  }
21

```

#### C++ code 4.2.4.7: Inverse two-dimensional discrete Fourier transform → [GITLAB](#)

```

2 template <typename Scalar>
3 void ifft2(Eigen::MatrixXcd &C, const Eigen::MatrixBase<Scalar> &Y) {
4   using idx_t = Eigen::MatrixXcd::Index;
5   const idx_t m = Y.rows(), n=Y.cols();
6   fft2(C,Y.conjugate()); C = C.conjugate() / (m*n);
7

```

**Remark 4.2.4.8 (Two-dimensional DFT in MATLAB)** The two-dimensional DFT is provided by the MATLAB function: `fft2(Y)`.

In MATLAB the two-dimensional DFT is defined through *nested one-dimensional DFTs* (4.2.0.19) as follows:

$$\text{fft2}(Y) = \text{fft}(\text{fft}(Y) \cdot ') \cdot '$$

Here: `.'` simply transposes the matrix (no complex conjugation) □

**§4.2.4.9 (Periodic convolution of matrices)** We consider the following bilinear mapping  $B : \mathbb{C}^{m,n} \times \mathbb{C}^{m,n} \rightarrow \mathbb{C}^{m,n}$ :

$$(B(X, Y))_{k,\ell} = \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} (X)_{i,j} (Y)_{k-i \mod m, \ell-j \mod n}, \quad (4.2.4.10)$$

which defines the **two-dimensional discrete periodic convolution**, cf. Def. 4.1.0.33.

#### C++ code 4.2.4.11: Straightforward implementation of 2D discrete periodic convolution → [GITLAB](#)

```

2 // Straightforward implementation of 2D periodic convolution
3 template <typename Scalar1, typename Scalar2, class EigenMatrix>
4 void pmconv_basic(const Eigen::MatrixBase<Scalar1> &X, const
5   Eigen::MatrixBase<Scalar2> &Y,
6   EigenMatrix &Z) {
7   using idx_t = typename EigenMatrix::Index;
8   using val_t = typename EigenMatrix::Scalar;
9   const idx_t n=X.cols(), m=X.rows();
10  if ((m!=Y.rows()) || (n!=Y.cols())) throw std::runtime_error("pmconv: size

```

```

    mismatch");
10 Z.resize(m,n);
11 // Implementation of (4.2.4.10)
12 auto idxwrap = [] (const idx_t L, int i)
13 { if (i>=L) i -= L; else if (i<0) i += L; return i; };
14 for (int i=0;i<m;i++) for (int j=0;j<n;j++) {
15     val_t s=0;
16     for (int k=0;k<m;k++) for (int l=0;l<n;l++) {
17         s += X(k,l)*Y(idxwrap(m,i-k),idxwrap(n,j-l));
18     Z(i,j) = s;
19 }
20 }
```

The 2D discrete periodic convolution admits a diagonalization by switching to the trigonometric basis of  $\mathbb{C}^{m,n}$ , analogous to (4.2.0.17), see Section 4.2.1.

In (4.2.4.10) set  $\mathbf{Y} = (\mathbf{F}_m)_{:,r}(\mathbf{F}_n)_{s,:} \in \mathbb{C}^{m,n} \leftrightarrow (\mathbf{Y})_{i,j} = \omega_m^{ri}\omega_n^{sj}, 0 \leq i < m, 0 \leq j < n$ :

$$\begin{aligned}
(\mathbf{B}(\mathbf{X}, \mathbf{Y}))_{k,\ell} &= \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} (\mathbf{X})_{i,j} (\mathbf{Y})_{k-i \bmod m, \ell-j \bmod n} \\
&= \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} (\mathbf{X})_{i,j} \omega_m^{r(k-i)} \omega_n^{s(\ell-j)} \\
&= \left( \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} (\mathbf{X})_{i,j} \overline{\omega_m^{ri} \omega_n^{sj}} \right) \cdot \omega_m^{rk} \omega_n^{s\ell}. \\
\Rightarrow \mathbf{B}(\mathbf{X}, \underbrace{(\mathbf{F}_m)_{:,r}(\mathbf{F}_n)_{s,:}}_{\text{"eigenvector"}}) &= \left( \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} (\mathbf{X})_{i,j} \overline{\omega_m^{ri} \omega_n^{sj}} \right) (\mathbf{F}_m)_{:,r}(\mathbf{F}_n)_{s,:}. \quad (4.2.4.12) \\
&\text{see Eq. (4.2.4.3)}
\end{aligned}$$

Hence, the (complex conjugated) two-dimensional discrete Fourier transform of  $\bar{\mathbf{X}}$  according to (4.2.4.3) provides the eigenvalues of the anti-linear mapping  $\mathbf{Y} \mapsto \mathbf{B}(\mathbf{X}, \mathbf{Y}), \mathbf{X} \in \mathbb{C}^{m,n}$  fixed.

This suggests the following DFT-based algorithm for evaluating the periodic convolution of matrices:

- ① Compute  $\hat{\mathbf{Y}}$  by inverse 2D DFT of  $\bar{\mathbf{Y}}$ , see Code 4.2.4.7
- ② Compute  $\hat{\mathbf{X}}$  by 2D DFT of  $\bar{\mathbf{X}}$ , see Code 4.2.4.6.
- ③ Component-wise multiplication of  $\hat{\mathbf{X}}$  and  $\hat{\mathbf{Y}}$ :  $\hat{\mathbf{Z}} = \hat{\mathbf{X}} \cdot \hat{\mathbf{Y}}$ .
- ④ Compute  $\mathbf{Z}$  through inverse 2D DFT of  $\hat{\mathbf{Z}}$ .

#### C++ code 4.2.4.13: DFT-based 2D discrete periodic convolution → GITLAB

```

2 // DFT based implementation of 2D periodic convolution
3 template <typename Scalar1,typename Scalar2,class EigenMatrix>
4 void pmconv(const Eigen::MatrixBase<Scalar1> &X,const Eigen::MatrixBase<Scalar2> &Y,
5             EigenMatrix &Z) {
6     using Comp = std::complex<double>;
7     using idx_t = typename EigenMatrix::Index;
8     using val_t = typename EigenMatrix::Scalar;
9     const idx_t n=X.cols(), m=X.rows();
10    if ((m!=Y.rows()) || (n!=Y.cols())) throw std::runtime_error("pmconv: size
11 mismatch");
12    Z.resize(m,n); Eigen::MatrixXcd Xh(m,n), Yh(m,n);
```

```

12 // Step ❶: 2D DFT of Y
13 fft2(Yh,(Y.template cast<Comp>()) );
14 // Step ❷: 2D DFT of X
15 fft2(Xh,(X.template cast<Comp>()) );
16 // Steps ❸, ❹: inverse DFT of component-wise product
17 ifft2(Z,Xh.cwiseProduct(Yh) );
18 }

```

**EXAMPLE 4.2.4.14 (Deblurring by DFT)** 2D discrete convolutions are important for image processing. Let a Gray-scale pixel image be stored in the matrix  $\mathbf{P} \in \mathbb{R}^{m,n}$ , actually  $\mathbf{P} \in \{0, \dots, 255\}^{m,n}$ , see also Ex. 9.3.2.1.

Write  $(p_{l,k})_{l,k \in \mathbb{Z}}$  for the periodically extended image:

$$p_{l,j} = (\mathbf{P})_{l+1,j+1} \quad \text{for } 1 \leq l \leq m, 1 \leq j \leq n, \quad p_{l,j} = p_{l+m,j+n} \quad \forall l, k \in \mathbb{Z}.$$

**Blurring** = pixel values get replaced by weighted averages of near-by pixel values  
(effect of distortion in optical transmission systems)

$$c_{l,j} = \sum_{k=-L}^L \sum_{q=-L}^L s_{k,q} p_{l+k,j+q}, \quad 0 \leq l < m, \quad 0 \leq j < n, \quad L \in \{1, \dots, \min\{m, n\}\}. \quad (4.2.4.15)$$

blurred image      point spread function

Does this ring a bell? Hidden in (4.2.4.15) is a 2D discrete periodic convolution, see Eq. (4.2.4.10). In light of the algorithm implemented in Code 4.2.4.13 it is hardly surprising that DFT comes handy for reversing the effect of the blurring!

Note that usually:  $L$  small,  $s_{k,m} \geq 0$ ,  $\sum_{k=-L}^L \sum_{q=-L}^L s_{k,q} = 1$  (an averaging)

In the experiments we used:  $L = 5$  and the PSF  $s_{k,q} = \frac{1}{1 + k^2 + q^2}$ .

#### C++11-code 4.2.4.16: Point spread function (PSF) → GITLAB

```

2 void psf(const long L, MatrixXd& S) {
3     VectorXd x = VectorXd::LinSpaced(2*L+1, -L, L);
4     MatrixXd X,Y;
5     meshgrid(x, x, X, Y);
6     MatrixXd E = MatrixXd::Ones(2*L+1, 2*L+1);
7     S = E.cwiseQuotient(E + X.cwiseProduct(X) + Y.cwiseProduct(Y));
8     S /= S.sum();
9 }

```

Original

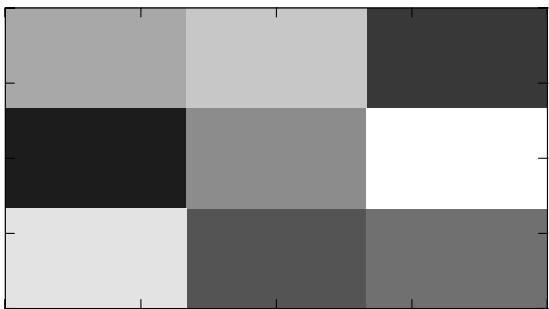


Fig. 126

Blurred image

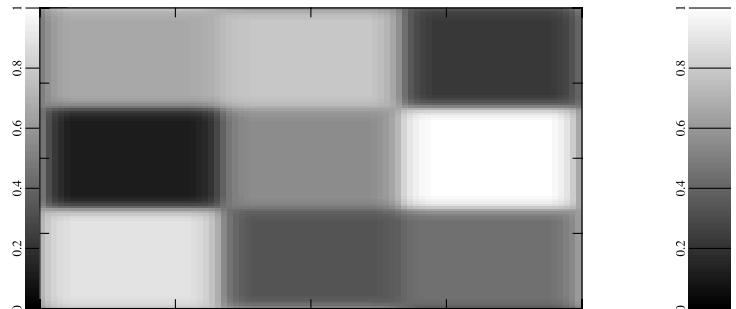


Fig. 127

Note: (4.2.4.15) defines a linear operator  $\mathcal{B} : \mathbb{R}^{m,n} \mapsto \mathbb{R}^{m,n}$  ("blurring operator")

#### C++11-code 4.2.4.17: Blurring operator → [GITLAB](#)

```

2  MatrixXd blur(const MatrixXd& P, const MatrixXd& S) {
3      const long m = P.rows(), n = P.cols(),
4          M = S.rows(), N = S.cols(), L = (M-1)/2;
5      if (M != N) {
6          std::cout << "Error: S not quadratic!\n";
7      }
8
9      MatrixXd C(m,n);
10     for (long i = 1; i <= m; ++i) {
11         for (long j = 1; j <= n; ++j) {
12             double s = 0;
13             for (long k = 1; k <= (2*L+1); ++k) {
14                 for (long q = 1; q <= (2*L+1); ++q) {
15                     double ki = i + k - L - 1;
16                     if (ki < 1) ki += m;
17                     else if (ki > m) ki -= m;
18                     double jq = j + q - L - 1;
19                     if (jq < 1) jq += n;
20                     else if (jq > n) jq -= n;
21                     s += P(ki-1, jq-1)*S(k-1,q-1);
22                 }
23             }
24             C(i-1,j-1) = s;
25         }
26     }
27     return C;
28 }
```

Now we revisit the considerations of § 4.2.4.1 and recall the derivation of (4.2.0.10) and Lemma 4.2.0.16.

$$\left( \mathcal{B}(\left( \omega_m^{\nu k} \omega_n^{\mu q} \right)_{k,q \in \mathbb{Z}}) \right)_{l,j} = \sum_{k=-L}^L \sum_{q=-L}^L s_{k,q} \omega_m^{\nu(l+k)} \omega_n^{\mu(j+q)} = \omega_m^{\nu l} \omega_n^{\mu j} \sum_{k=-L}^L \sum_{q=-L}^L s_{k,q} \omega_m^{\nu k} \omega_n^{\mu q}.$$

►  $\mathbf{V}_{\nu,\mu} := (\omega_m^{\nu k} \omega_n^{\mu q})_{k,q \in \mathbb{Z}}, 0 \leq \mu < m, 0 \leq \nu < n$  are the eigenvectors of  $\mathcal{B}$ :

$$\mathcal{B}\mathbf{V}_{\nu,\mu} = \lambda_{\nu,\mu} \mathbf{V}_{\nu,\mu} \quad , \quad \text{eigenvalue} \quad \lambda_{\nu,\mu} = \underbrace{\sum_{k=-L}^L \sum_{q=-L}^L s_{k,q} \omega_m^{\nu k} \omega_n^{\mu q}}_{\text{2-dimensional DFT of point spread function!}} \quad (4.2.4.18)$$

Thus the inversion of the blurring operator boils down to componentwise scaling in “Fourier domain”, see See also Code 4.2.4.13 for the same idea.

#### C++11-code 4.2.4.19: DFT based deblurring → GITLAB

```

2  MatrixXd deblur(const MatrixXd& C, const MatrixXd& S, const double tol=1e-3) {
3      const long m = C.rows(), n = C.cols(),
4          M = S.rows(), N = S.cols(), L = (M-1)/2;
5      if (M != N) {
6          std::cerr << "Error: S not quadratic!\n";
7      }
8
9      MatrixXd Spad = MatrixXd::Zero(m,n);
10     // Zero padding
11     Spad.block(0, 0, L+1, L+1) = S.block(L, L, L+1, L+1);
12     Spad.block(m-L, n-L, L, L) = S.block(0, 0, L, L);
13     Spad.block(0, n-L, L+1, L) = S.block(L, 0, L+1, L);
14     Spad.block(m-L, 0, L, L+1) = S.block(0, L, L, L+1);
15     // Inverse of blurring operator (fft2 expects a complex matrix)
16     MatrixXcd SF = fft2(Spad.cast<complex>());
17     // Test for invertibility
18     if (SF.cwiseAbs().minCoeff() < tol*SF.cwiseAbs().maxCoeff()) {
19         std::cerr << "Error: Deblurring impossible!\n";
20     }
21     // DFT based deblurring
22     return fft2(ifft2(C.cast<complex>()).cwiseQuotient(SF)).real();
23 }
```

#### 4.2.5 Semi-discrete Fourier Transform [QSS00, Sect. 10.11]

Starting from Rem. 4.1.0.29 we mainly looked at time-discrete  $n$ -periodic signals, which can be mapped to vectors  $\in \mathbb{R}^n$ . This led to discrete periodic convolution (→ Def. 4.1.0.33) and the discrete Fourier transform (DFT) (→ Def. 4.2.0.18) as (bi-)linear mappings in  $\mathbb{C}^n$ .

In this section we are concerned with **non-periodic** signals of **infinite duration** as introduced in § 4.0.0.1.



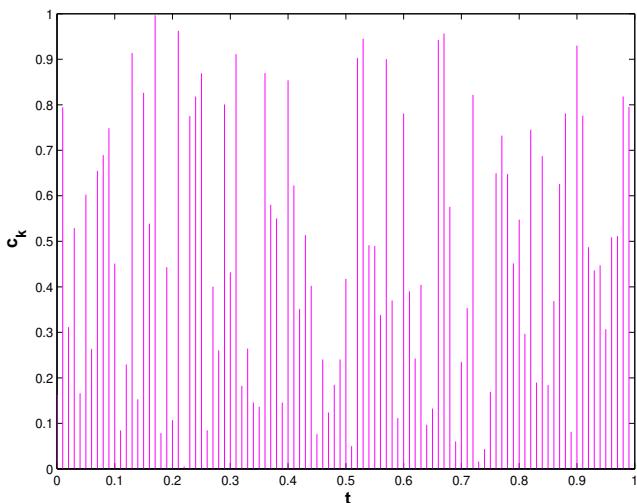
Idea: Study the limit  $n \rightarrow \infty$  for the  $n$ -periodic setting and DFT.

Let  $(y_k)_{k \in \mathbb{Z}}$  be an  $n$ -periodic sequence (signal),  $n = 2m + 1$ ,  $m \in \mathbb{N}$ . Thanks to periodicity we can rewrite

$$\text{DFT} \rightarrow \text{Def. 4.2.0.18: } c_k = \sum_{j=-m}^m y_j \exp\left(-2\pi i \frac{kj}{n}\right), \quad k = 0, \dots, n-1. \quad (4.2.5.1)$$

Now we associate a point  $t_k \in [0, 1[$  with each index  $k$  for the components of the transformed signal  $(c_k)_{k=0}^{n-1}$ :

$$k \in \{0, \dots, n-1\} \iff t_k := \frac{k}{n}. \quad (4.2.5.2)$$



▷ “Squeezing” a vector  $\in \mathbb{R}^n$  into  $[0, 1]$ .

Thus we can read the values  $c_k$  as sampled values of a function defined on  $[0, 1]$

$$c_k \leftrightarrow c(t_k),$$

$$t_k = \frac{k}{n}, \quad k = 0, \dots, n - 1.$$

This makes it possible to pass from a discrete finite signal to a continuous signal.

Thus, formally, we can rewrite (4.2.5.1) as

$$\text{DFT: } c(t_k) := c_k = \sum_{j=-m}^m y_j \exp(-2\pi i j t_k), \quad k = 0, \dots, n - 1. \quad (4.2.5.3)$$

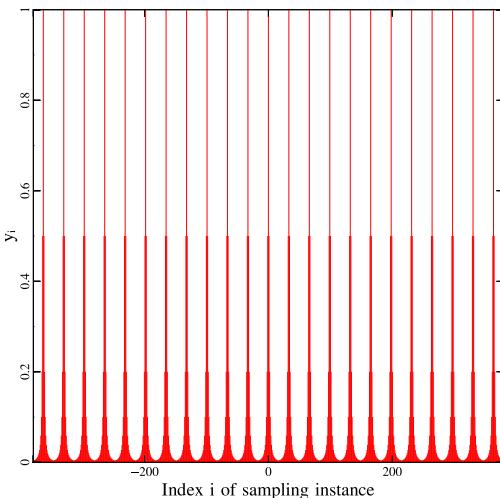
The notation indicates that we read  $c_k$  as the value of a *function*  $c : [0, 1] \mapsto \mathbb{C}$  for argument  $t_k$ .

#### EXAMPLE 4.2.5.4 (“Squeezed” DFT of a periodically truncated signal)

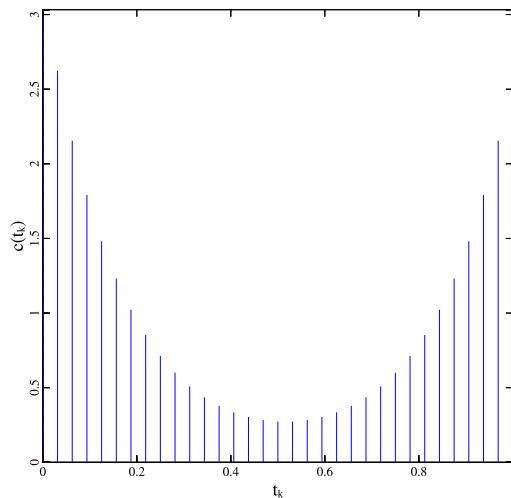
Bi-infinite discrete signal, “concentrated around 0”:  $y_j = \frac{1}{1+j^2}$ ,  $j \in \mathbb{Z}$ .

We examine the DFT of the  $2m + 1$ -periodic signal obtained by periodic extension of  $(y_k)_{k=-m}^m$ .

Period of signal  $y_i = 33$



DFT of period 33 signal



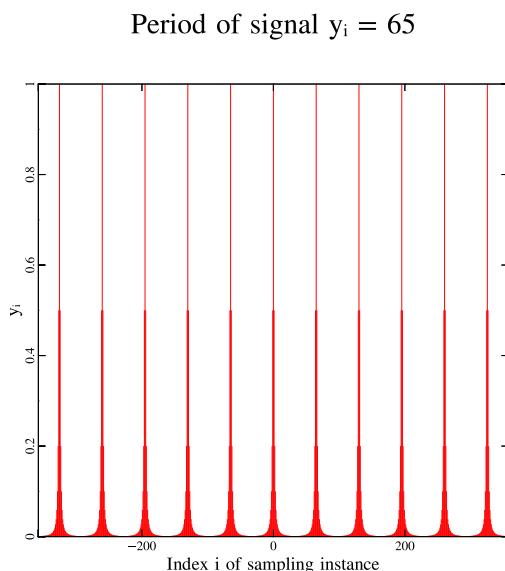


Fig. 131

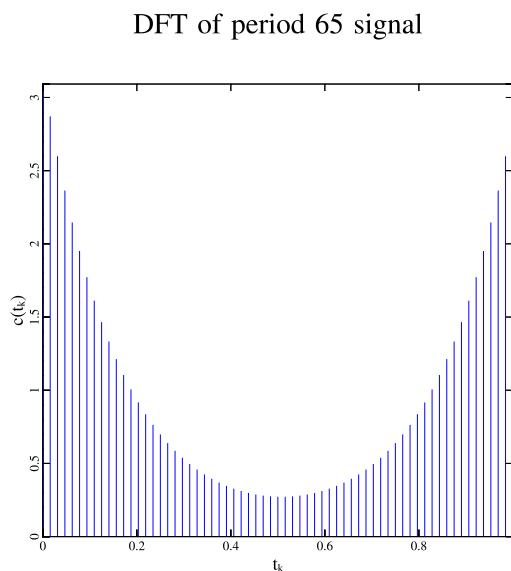


Fig. 132

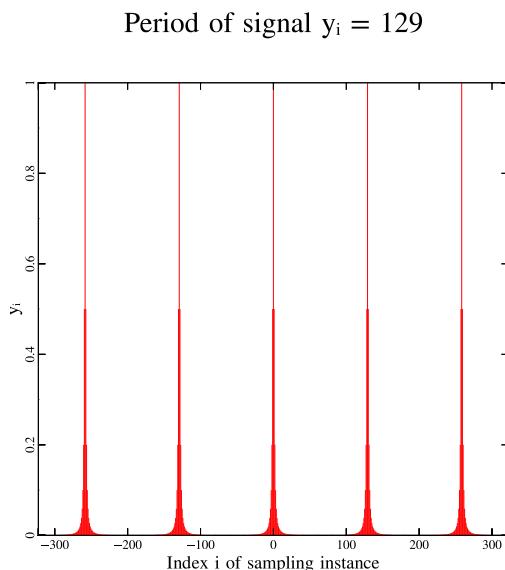


Fig. 133

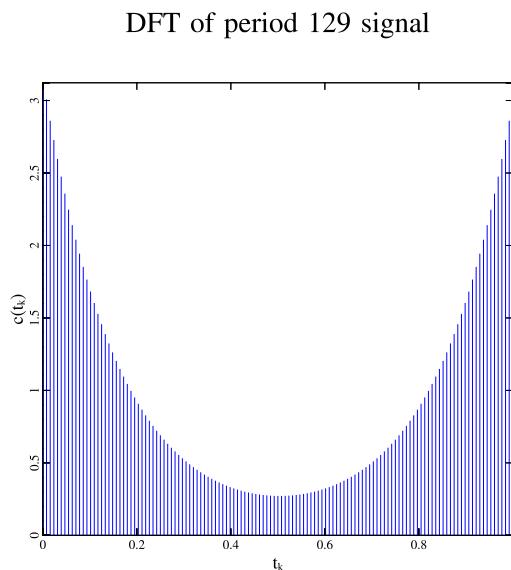


Fig. 134

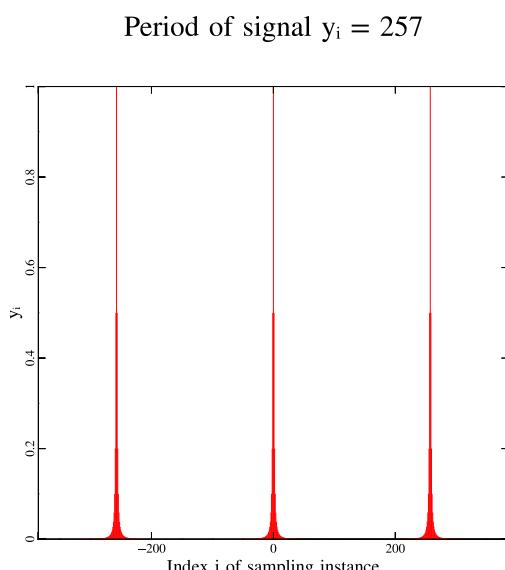


Fig. 135

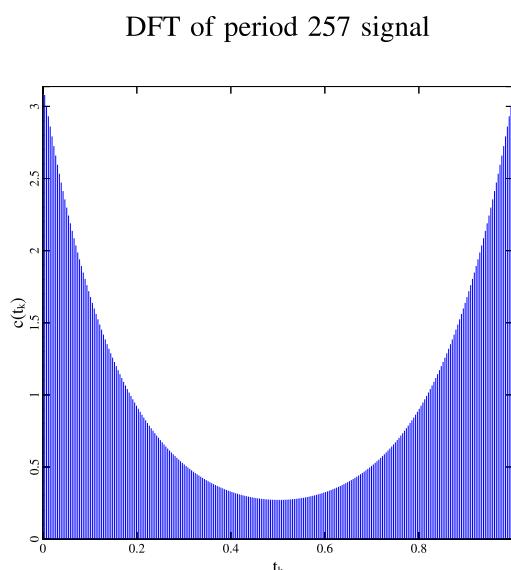


Fig. 136

Visual impression:  $c(t_k)$  “converge” to a function  $c : [0, 1] \mapsto \mathbb{R}$ .

**C++ code 4.2.5.5: Plotting a periodically truncated signal and its DFT** → [GITLAB](#)

```

2 // Visualize limit  $m \rightarrow \infty$  for a  $2m+1$ -periodic signal and
3 // its discrete Fourier transform "squeezed" into [0,1].
4 int main() {
5     // range of plot for visualization of discrete signal
6     const int N = 3*257;
7     // function defining discrete signal
8     auto yfn = [] (VectorXd k) {
9         return (1/(1 + k.array()*k.array())) . matrix();
10    };
11    // loop over different periods  $2^l+1$ 
12    for (int mpow = 4; mpow <= 7; ++mpow) {
13        int m = std::pow(2, mpow);
14        VectorXd ybas = yfn(VectorXd::LinSpaced(2*m+1, -m, m));
15        const int Ncp = std::floor(N/(2*m+1));
16        const int n = ybas.size();
17
18        VectorXd y(n*Ncp);
19        for (int i = 0; i < Ncp; ++i)
20            y.segment(n*i, n) = ybas;
21
22        // Stem plots vertical lines from 0 to the y-value
23        Stem s1;
24        s1.title = "Period of signal  $y_i =$  " + std::to_string(2*m+1);
25        s1.xlabel = "Index  $i$  of sampling instance";
26        s1.ylabel = " $y_i$ ";
27        s1.file = "persig" + std::to_string(mpow); // where to save
28        s1.plot(VectorXd::LinSpaced(n*Ncp, -n*Ncp/2., n*Ncp/2.), y, "r");
29
30        Eigen::FFT<double> fft;
31        // DFT of wrapped signal (one period)
32        VectorXd yconc(n); yconc << ybas.tail(m+1), ybas.head(m);
33        VectorXcd c = fft.fwd(yconc);
34        Stem s2;
35        s2.title = "DFT of period " + std::to_string(2*m+1) + " signal";
36        s2.xlabel = " $t_k$ ";
37        s2.ylabel = " $c(t_k)$ ";
38        s2.file = "persigdft" + std::to_string(mpow);
39        s2.plot(VectorXd::LinSpaced(2*m+1, 0, 1), c.real(), "b");
40    }
41    return 0;
42 }
```

Now we pass to the limit  $m \rightarrow \infty$  (and keep the function perspective  $c_k = c(t_k)$ )

Note: passing to the limit amounts to dropping the assumption of periodicity!



$$c(t) = \sum_{k \in \mathbb{Z}} y_k \exp(-2\pi i k t) .$$

(4.2.5.6)

**Terminology:** The series (= infinite sum) on the right hand side of (4.2.5.6) is called a **Fourier series** ([link](#))

The function  $c : [0, 1] \mapsto \mathbb{C}$  defined by (4.2.5.6) is called the **Fourier transform** of the sequence  $(y_k)_{k \in \mathbb{Z}}$  (, if the series converges).

### Corollary 4.2.5.7. Periodicity of Fourier transforms

Fourier transforms  $t \mapsto c(t)$  are **1-periodic** functions  $\mathbb{R} \rightarrow \mathbb{C}$ .

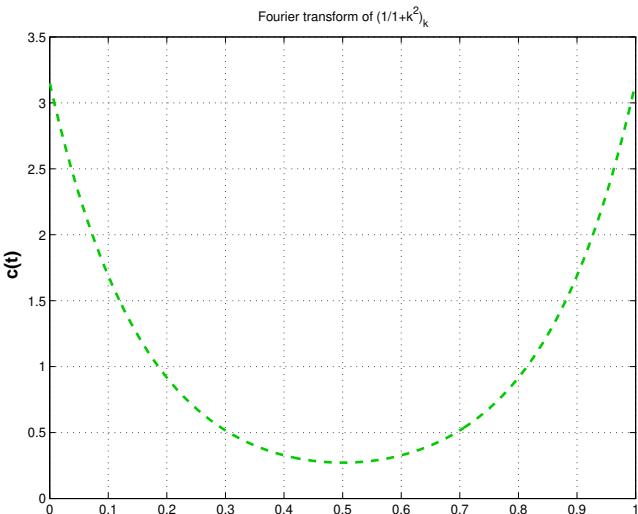


Fig. 137

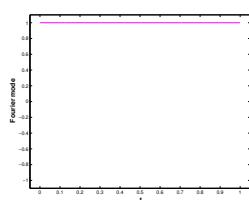
Thus, the limit we “saw” in Ex. 4.2.5.4 is actually the Fourier transform of the sequence  $(y_k)_{k \in \mathbb{Z}}$ !

▷ Fourier transform of  $y_k := \frac{1}{1+k^2}$

From (4.2.5.6) we conclude:

Fourier transform  
=

weighted sum of Fourier modes  $t \mapsto \exp(-2\pi i k t)$ ,  
 $k \in \mathbb{Z}$



→ related [animation](#) on Wikipedia.

It is possible to derive a closed-form expression for the function displayed in Fig. 138:

$$c(t) = \sum_{k \in \mathbb{Z}} \frac{1}{1+k^2} \exp(-2\pi i k t) = \frac{\pi}{e^\pi - e^{-\pi}} \cdot (e^{\pi - 2\pi t} + e^{2\pi t - \pi}) \in C^\infty([0, 1]).$$

Note that when considered as a 1-periodic function on  $\mathbb{R}$ , this  $c(t)$  is merely continuous.

**Remark 4.2.5.8 (Decay conditions for bi-infinite signals)** The considerations above were based on

- ◆ truncation of  $(y_k)_{k \in \mathbb{Z}}$  to  $(y_k)_{k=-m}^m$  and
- ◆ periodic continuation to an  $2m+1$ -periodic signal.

Obviously, only if the signal is *concentrated around  $k = 0$*  this procedure will not lose essential information contained in the signal.

$$\text{Minimal requirement: } \lim_{k \rightarrow \infty} |y_k| = 0, \quad (4.2.5.9)$$

$$\text{Stronger requirement: } \sum_{k \in \mathbb{Z}} |y_k| < \infty. \quad (4.2.5.10)$$

$$(4.2.5.10) \quad \Rightarrow \quad \text{Fourier series (4.2.5.6) ([link](#)) converges uniformly [Str09, Def. 4.8.1]} \\ \Rightarrow \quad c : [0, 1] \mapsto \mathbb{C} \text{ is continuous [Str09, Thm. 4.8.1].}$$

**Remark 4.2.5.11 (Numerical summation of Fourier series)**

Assuming sufficiently fast decay of the signal  $(y_k)_{k \in \mathbb{Z}}$  for  $k \rightarrow \infty$  ( $\rightarrow$  Rem. 4.2.5.8), we can *approximate* the Fourier series (4.2.5.6) by a **Fourier sum**

$$c(t) \approx c_M(t) := \sum_{k=-M}^M y_k \exp(-2\pi i k t), \quad M \gg 1. \quad (4.2.5.12)$$

Task: Approximate evaluation of  $c(t)$  at  $N$  equidistant points  $t_j := \frac{j}{N}, j = 0, \dots, N$  (e.g., for plotting it).

$$c(t_j) = \lim_{M \rightarrow \infty} \sum_{k=-M}^M y_k \exp(-2\pi i k t_j) \approx \sum_{k=-M}^M y_k \exp(-2\pi i \frac{kj}{N}), \quad (4.2.5.13)$$

for  $j = 0, \dots, N - 1$ .

Note: If  $N = M \gg$  (4.2.5.13) is a **discrete Fourier transform** ( $\rightarrow$  Def. 4.2.0.18).

#### C++ code 4.2.5.14: DFT-based evaluation of Fourier sum at equidistant points $\rightarrow$ GITLAB

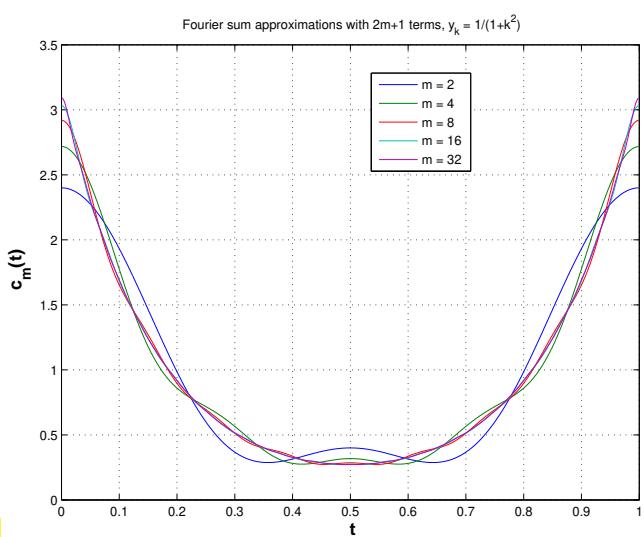
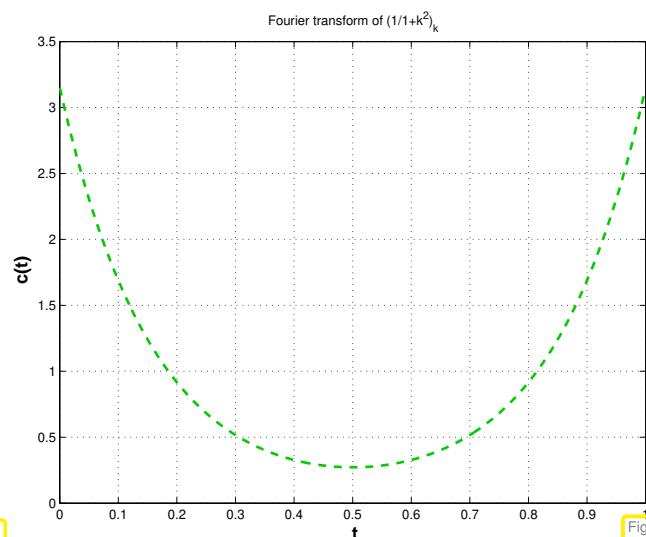
```

2 # include "feval.hpp" // evaluate scalar function with a vector
3 // DFT based approximate evaluation of Fourier series
4 // signal is a handle to a function providing the y_k
5 // M specifies truncation of series according to (4.2.5.12)
6 // N is the number of equidistant evaluation points for c in [0,1].
7 template <class Function>
8 VectorXcd foursum(const Function& signal, int M, int N) {
9     const int m = 2*M+1; // length of the signal
10    // sample signal
11    VectorXd y = feval(signal, VectorXd::LinSpaced(m, -M, M));
12    // Ensure that there are more sampling points than terms in series
13    int l; if (m > N) { l = ceil(double(m)/N); N *= l; } else l = 1;
14    // Zero padding and wrapping of signal, see Code 4.2.2.7
15    VectorXd y_ext = VectorXd::Zero(N);
16    y_ext.head(M+1) = y.tail(M+1);
17    y_ext.tail(M) = y.head(M);
18    // Perform DFT and decimate output vector
19    Eigen::FFT<double> fft;
20    Eigen::VectorXcd k = fft.fwd(y_ext), c(N/l);
21    for (int i = 0; i < N/l; ++i) c(i) = k(i*l);
22    return c;
23 }
```

#### EXAMPLE 4.2.5.15 (Convergence of Fourier sums)

Infinite signal, satisfying decay condition (4.2.5.10):  $y_k = \frac{1}{1+k^2}$ , see Ex. 4.2.5.4.

Monitored: approximation of Fourier transform  $c(t)$  by Fourier sums  $c_m(t)$ , see (4.2.5.12).



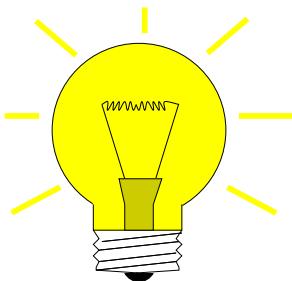
Observation: Convergence of Fourier sums in “eyeball norm”; quantitative statements about convergence can be deduced from Thm. 4.2.5.28.

Same limit considerations as above for the inverse DFT (4.2.0.20),  $n = 2m + 1$ ,

$$y_j = \frac{1}{n} \sum_{k=0}^{n-1} c_k \exp(2\pi i \frac{jk}{n}), \quad j = -m, \dots, m. \quad (4.2.5.16)$$

Adopt function perspective as before:  $c_k \leftrightarrow c(t_k)$ ,  $t_k = \frac{k}{n}$ , cf. (4.2.5.2).

$$y_j = \frac{1}{n} \sum_{k=0}^{n-1} c(t_k) \exp(2\pi i j t_k), \quad j = -m, \dots, m. \quad (4.2.5.17)$$



Now: pass to the limit  $m \rightarrow \infty$  in (4.2.5.17)

Idea: right hand side of (4.2.5.17) = **Riemann sum**, cf. [Str09, Sect. 6.2]



in the limit  $m \rightarrow \infty$  the sum becomes an **integral**!

$$y_j = \int_0^1 c(t) \exp(2\pi i j t) dt. \quad (4.2.5.18)$$

This formula is the inversion of the summation of a Fourier series (4.2.5.6)!

The formula (4.2.5.18) allows to recover the signal  $(y_k)_{k \in \mathbb{Z}}$  from its Fourier transform  $c(t)$ .

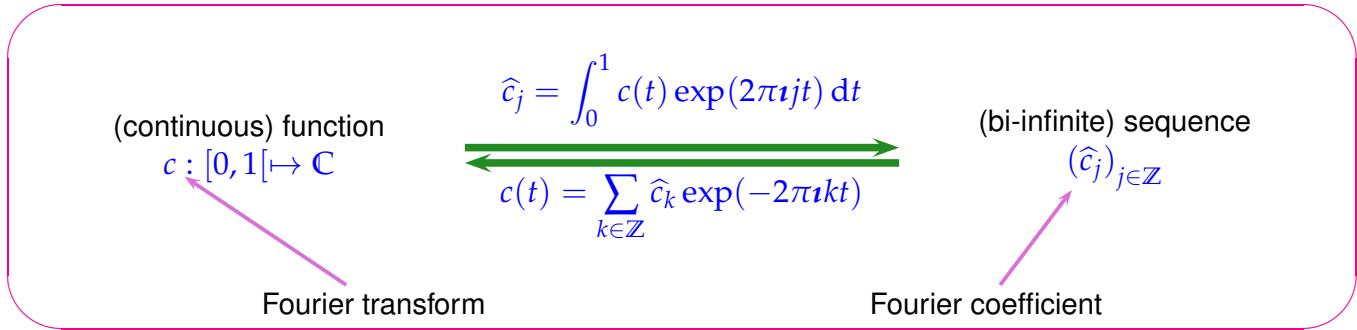
Terminology:  $y_j$  from (4.2.5.18) is called the ***j*-th Fourier coefficient** of the function  $c$ .

Notation:  $\hat{c}_j := y_j$  with  $y_j$  defined by (4.2.5.18)  $\hat{c}_j$  ***j*-th Fourier coefficient** of  $c : [0, 1] \rightarrow \mathbb{C}$

Note: Both, the mapping  $(\hat{c}_k)_{k \in \mathbb{Z}} \mapsto c$  from (4.2.5.6), and the mapping  $c \mapsto (\hat{c}_k)_{k \in \mathbb{Z}}$  from (4.2.5.18) are **linear**!

(Recall the concept of a linear mapping as explained in [NS02, Ch. 6].)

Summary of a fundamental correspondence:



**Remark 4.2.5.19 (Filtering in Fourier domain)** What happens to the Fourier transform of a bi-infinite signal, if it passes through a channel?

Consider a (bi-)infinite signal  $(x_k)_{k \in \mathbb{Z}}$  sent through a **finite** ( $\rightarrow$  Def. 4.1.0.4), **linear** ( $\rightarrow$  Def. 4.1.0.9) **time-invariant** ( $\rightarrow$  Def. 4.1.0.7) causal ( $\rightarrow$  Def. 4.1.0.11) channel with impulse response ( $\rightarrow$  4.1.0.3)  $(\dots, 0, h_0, \dots, h_{n-1}, 0, \dots)$  ( $\rightarrow$  § 4.1.0.1).

➤ this results in the output signal, see (4.1.0.14),

$$y_k = \sum_{j=0}^{n-1} h_j x_{k-j}, \quad k \in \mathbb{Z}. \quad (4.2.5.20)$$

Fourier transforms of signals:

$(y_k)_{k \in \mathbb{Z}} \leftrightarrow c(t) \quad , \quad (x_j)_{j \in \mathbb{Z}} \leftrightarrow b(t)$

(Assume that  $(x_k)_{k \in \mathbb{Z}}$  satisfies (4.2.5.10))

$$\begin{aligned} c(t) &= \sum_{k \in \mathbb{Z}} y_k \exp(-2\pi i kt) = \sum_{k \in \mathbb{Z}} \sum_{j=0}^{n-1} h_j x_{k-j} \exp(-2\pi i kt) \\ &\text{[shift summation index } k \text{]} = \sum_{j=0}^{n-1} \sum_{k \in \mathbb{Z}} h_j x_k \exp(-2\pi i jt) \exp(-2\pi i kt) \\ &= \underbrace{\left( \sum_{j=0}^{n-1} h_j \exp(-2\pi i jt) \right)}_{\text{trigonometric polynomial of degree } n-1} b(t). \end{aligned} \quad (4.2.5.21)$$

### Definition 4.2.5.22. Trigonometric polynomial

A **trigonometric polynomial** is a function  $\mathbb{R} \rightarrow \mathbb{C}$  that is a weighted sum of finitely many terms  $t \mapsto \exp(-2\pi i kt)$ ,  $k \in \mathbb{Z}$ .

We summarize the insight gleaned from (4.2.5.21):

#### Discrete convolution in Fourier domain

The discrete convolution of a signal with finite impulse response amounts to a multiplication of its Fourier transform with a trigonometric polynomial whose coefficients are given by the impulse response.

§4.2.5.24 (Isometry property of Fourier transform) We find a **conservation of power** through Fourier transform:

Lemma 4.2.0.14  $\Rightarrow$  for Fourier matrix  $\mathbf{F}_n$ , see (4.2.0.13),  $\frac{1}{\sqrt{n}}\mathbf{F}_n$  is **unitary** ( $\rightarrow$  Def. 6.2.1.2)

$$\xrightarrow{\text{Thm. 3.3.2.2}} \left\| \frac{1}{\sqrt{n}} \mathbf{F}_n \mathbf{y} \right\|_2 = \|\mathbf{y}\|_2. \quad (4.2.5.25)$$

Since DFT boils down to multiplication with  $\mathbf{F}_n$  ( $\rightarrow$  Def. 4.2.0.18), we conclude from (4.2.5.25)

$$c_k \text{ from (4.2.5.1)} \Rightarrow \frac{1}{n} \sum_{k=0}^{n-1} |c_k|^2 = \sum_{j=-m}^m |y_j|^2. \quad (4.2.5.26)$$

Now we adopt the function perspective again and associated  $c_k \leftrightarrow c(t_k)$ . Then we pass to the limit  $m \rightarrow \infty$ , appeal to Riemann summation (see above), and conclude

$$(4.2.5.26) \xrightarrow{m \rightarrow \infty} \int_0^1 |c(t)|^2 dt = \sum_{j \in \mathbb{Z}} |y_j|^2. \quad (4.2.5.27)$$

### Theorem 4.2.5.28. Isometry property of Fourier transform

If  $\sum_{k \in \mathbb{Z}} |\hat{c}_j|^2 < \infty$ , then the Fourier series

$$c(t) = \sum_{k \in \mathbb{Z}} \hat{c}_k \exp(-2\pi i kt)$$

yields a function  $c \in L^2([0, 1])$  that satisfies

$$\|c\|_{L^2([0,1])}^2 := \int_0^1 |c(t)|^2 dt = \sum_{k \in \mathbb{Z}} |\hat{c}_j|^2.$$

Recalling the concept of the  $L^2$ -norm of a function, see (5.2.4.5), the theorem can be stated as follows:

Thm. 4.2.5.28  $\leftrightarrow$  The  $L^2$ -norm of a Fourier transform agrees with the Euclidean norm of the corresponding signal.

Note: Euclidean norm of a sequence  $\|(y_k)_{k \in \mathbb{Z}}\|_2^2 := \sum_{k \in \mathbb{Z}} |y_j|^2$ .

## 4.3 Fast Fourier Transform (FFT)

You might have been wondering why the reduction to DFTs received so much attention in Section 4.2.1. An explanation is given now.



**Supplementary literature.** [DR08, Sect. 8.7.3], [Han02, Sect. 53], [QSS00, Sect. 10.9.2]

At first glance (at (4.2.0.19)): DFT in  $\mathbb{C}^n$  seems to require asymptotic computational effort of  $O(n^2)$  (matrix  $\times$  vector multiplication with dense matrix).

### C++-code 4.3.0.1: Naive DFT-implementation

```

2 // DFT (4.2.0.19) of vector y returned in c
3 void naivedft(const VectorXcd& y, VectorXcd& c) {
4     using idx_t = VectorXcd::Index;
5     const idx_t n = y.size();
6     const std::complex<double> i(0,1);
7     c.resize(n);
8     // root of unity  $\omega_n$ , w holds its powers
9     std::complex<double> w = std::exp(-2*M_PI/n*i), s = w;
10    c(0) = y.sum();
11    for (long j = 1; j < n; ++j) {
12        c(j) = y(n-1);
13        for (long k = n-2; k >= 0; --k) c(j) = c(j)*s + y(k);
14        s *= w;
15    }
16 }
```

### EXAMPLE 4.3.0.2 (Efficiency of fft)

#### C++-code 4.3.0.3: timing of different implementations of DFT

```

1 # include <complex>
2 # include <Eigen/Dense>
3 # include <unsupported/Eigen/FFT>
4 # include <figure/figure.hpp>
5 # include "timer.h"
6 # include "meshgrid.hpp"
7 # include "eigen_fft_backend_name.hpp"
8
9 using namespace Eigen;
10
11 //! Benchmarking discrete fourier transformations
12 // N = maximal length of the transformed vector
13 // nruns = no. of runs on which to take the minimal timing
14 // If the size of the transformed array is a prime number the kissfft
15 // implementation of EIGEN performs extremely bad, use FFTW or MKL
16 // in these cases
17 void benchmark(const int N, const int nruns=3) {
18     std::complex<double> i(0,1); // imaginary unit
19
20     /* Version which computes for every vector size (extremely slow)
21     MatrixXd res(N-1, 4); // save timing results and vector size
22     for (int n = 2; n <= N; ++n)
23     /* Version which computes for powers of two only
24     MatrixXd res(int(std::log2(N)), 4); // save timing results and vector size
25     for (int n = 2, j = 0; n <= N; n*=2, ++j) {
26         VectorXd y = VectorXd::Random(n);
27         VectorXcd c = VectorXcd::Zero(n);
28         // compute fourier transformed with a loop implementation
29         Timer t1; t1.start();
30         for (int k = 0; k < nruns; ++k) {
31             std::complex<double> omega = std::exp(-2*M_PI/n*i),
32                                         s = omega;
33             c(0) = y.sum();
34             for (int j = 1; j < n; ++j) {
35                 c(j) = y(n-1);
36                 for (int l = n-1; l > 0; --l){ c(j) = c(j)*s + y(l); }
37                 s *= omega;
38             }
39     }
```

```

39     t1.lap();
40 }
41 // compute fourier transformed through matrix multiplication
42 MatrixXd I, J; VectorXd lin = VectorXd::LinSpaced(n, 0, n-1);
43 meshgrid(lin, lin, I, J);
44 MatrixXcd F = (-2*M_PI/n*i*I.cwiseProduct(J)).array().exp().matrix();
45 Timer t2; t2.start();
46 for (int k = 0; k < nruns; ++k) {
47     c = F*y; t2.lap();
48 }
49 // use Eigen's FFT to compute fourier transformation
50 // Note: slow for large primes!
51 Eigen::FFT<double> fft;
52 Timer t3; t3.start();
53 for (int k = 0; k < nruns; ++k) {
54     fft.fwd(c,y); t3.lap();
55 }
56 // save timings
57 res(j,0) = n; res(j,1) = t1.min();
58 res(j,2) = t2.min(); res(j,3) = t3.min();
59 }

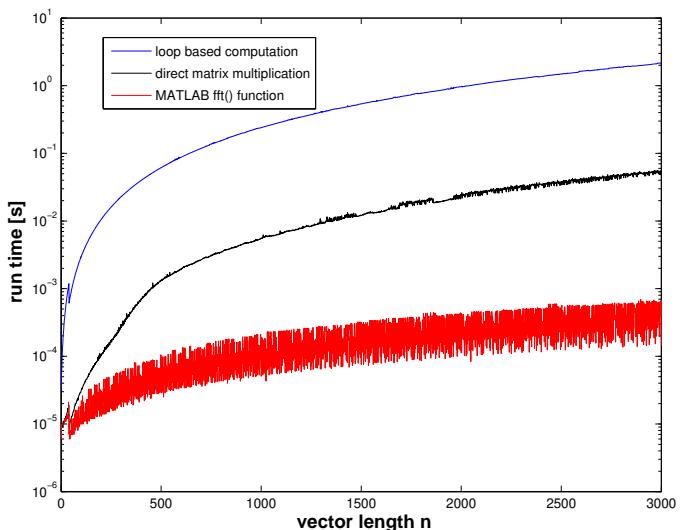
60 mgl::Figure fig;
61 fig.title("FFT timing");
62 fig.setlog(true, true);
63 fig.plot(res.col(0), res.col(1), "b*").label("Loop based computation");
64 fig.plot(res.col(0), res.col(2), "m^").label("Direct matrix multiplication");
65 fig.plot(res.col(0), res.col(3), "r<").label("Eigen's FFT module");
66 fig.fplot("x^2*1e-8", "k:").label("O(n^2)");
67 fig.fplot("x*lg(x)*1e-8", "k;").label("O(n log(n))");
68 fig.legend(0,1);
69 fig.save("ffttime");
70 std::cout << res << std::endl;
71 }

72 int main() {
73     // print the backend used by eigen for the fft
74     std::cout << "FFT Backend: " << eigen_fft_backend_name() << std::endl;
75
76     // WARNING: this code may take very long to run
77     benchmark(1024*8);
78     return 0;
79 }
80
81 }

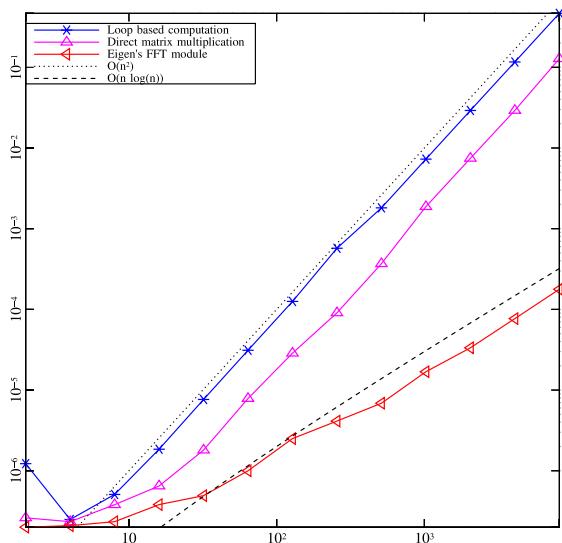
```

## Timings in MATLAB

- ▷:
1. Straightforward implementation involving MATLAB loops
  2. Multiplication with Fourier matrix (4.2.0.13)
  3. MATLAB's built-in function `fft()`



## FFT timing



## Timings in EIGEN

▷: (only for  $n = 2^L$ !)

1. Straightforward implementation involving C++ loops, see Code 4.3.0.1
2. Multiplication with Fourier matrix (4.2.0.13)
3. EIGEN's built-in **FFT** class, method `fwd()`

Note: Eigen uses KISS FFT as default backend in its FFT module, which falls back loop-based slow DFT when used on data sizes, which are large primes. An superior alternative is FFTW, see § 4.3.0.11.

The secret of EIGEN's `fft()`:

the **Fast Fourier Transform (FFT) algorithm** [DV90]

(discovered by C.F. Gauss in 1805, rediscovered by Cooley & Tuckey in 1965, one of the “top ten algorithms of the century”).

An elementary manipulation of (4.2.0.19) for  $n = 2m$ ,  $m \in \mathbb{N}$ :

$$\begin{aligned}
 c_k &= \sum_{j=0}^{n-1} y_j e^{-\frac{2\pi i}{n} jk} \\
 &= \sum_{j=0}^{m-1} y_{2j} e^{-\frac{2\pi i}{n} 2jk} + \sum_{j=0}^{m-1} y_{2j+1} e^{-\frac{2\pi i}{n} (2j+1)k} \\
 &= \underbrace{\sum_{j=0}^{m-1} y_{2j} e^{-\frac{2\pi i}{m} jk}}_{=: \tilde{c}_k^{\text{even}}} + e^{-\frac{2\pi i}{n} k} \cdot \underbrace{\sum_{j=0}^{m-1} y_{2j+1} e^{-\frac{2\pi i}{m} jk}}_{=: \tilde{c}_k^{\text{odd}}}.
 \end{aligned} \tag{4.3.0.4}$$

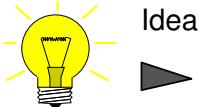
Note  $m$ -periodicity:  $\tilde{c}_k^{\text{even}} = \tilde{c}_{k+m}^{\text{even}}$ ,  $\tilde{c}_k^{\text{odd}} = \tilde{c}_{k+m}^{\text{odd}}$ .

Note:  $\tilde{c}_k^{\text{even}}, \tilde{c}_k^{\text{odd}}$  from DFTs of length  $m$ !

with  $\mathbf{y}_{\text{even}} := (y_0, y_2, \dots, y_{n-2})^\top \in \mathbb{C}^m$ :  $(\tilde{c}_k^{\text{even}})_{k=0}^{m-1} = \mathbf{F}_m \mathbf{y}_{\text{even}}$ ,  
 with  $\mathbf{y}_{\text{odd}} := (y_1, y_3, \dots, y_{n-1})^\top \in \mathbb{C}^m$ :  $(\tilde{c}_k^{\text{odd}})_{k=0}^{m-1} = \mathbf{F}_m \mathbf{y}_{\text{odd}}$ .

(4.3.0.4):

DFT of length  $2m = 2 \times$  DFT of length  $m + 2m$  additions & multiplications



Idea:

divide & conquer recursion

FFT-algorithm

Recursive demonstration code for DFT of length  $n = 2^L$ :

#### C++ code 4.3.0.5: Recursive FFT → GITLAB

```

2 // Recursive DFT for vectors of length n = 2L
3 VectorXcd fftrec(const VectorXcd& y) {
4     using idx_t = VectorXcd::Index;
5     using comp = std::complex<double>;
6     // Nothing to do for DFT of length 1
7     const idx_t n = y.size();
8     if (n == 1) return y;
9     if (n % 2 != 0) throw std::runtime_error("size(y) must be even!");
10    const Eigen::Map<const
11        Eigen::Matrix<comp, Eigen::Dynamic, Eigen::Dynamic, Eigen::RowMajor>>
12        Y(y.data(), n/2, 2); // for selecting even and odd components
13    const VectorXcd c1 = fftrec(Y.col(0)), c2 = fftrec(Y.col(1));
14    const comp omega = std::exp(-2*M_PI/n*comp(0,1)); // root of unity ωn
15    comp s(1.0, 0.0);
16    VectorXcd c(n);
17    // Scaling of DFT of odd components plus periodic copying
18    for (long k = 0; k < n; ++k) {
19        c(k) = c1(k%(n/2)) + c2(k%(n/2))*s;
20        s *= omega;
21    }
22    return c;
}
```

Computational cost of `fftrec`:

$1 \times$  DFT of length  $2^L$

$2 \times$  DFT of length  $2^{L-1}$

$4 \times$  DFT of length  $2^{L-2}$



$2^L \times$  DFT of length 1

Code 4.3.0.5: each level of the recursion requires  $O(2^L)$  elementary operations.

#### Asymptotic complexity

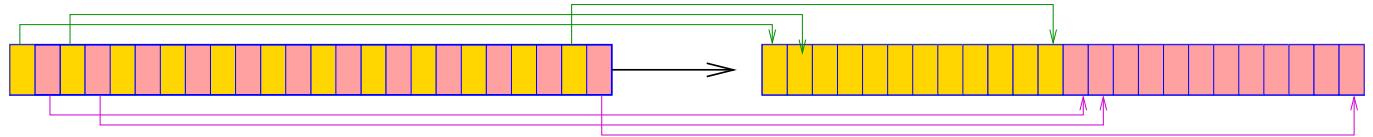
Asymptotic complexity of FFT algorithm,  $n = 2^L$ :  $O(L2^L) = O(n \log_2 n)$

(`fft.fwd()`/`fft.inv()`-function calls: computational cost is  $\approx 5n \log_2 n$ ).

**Remark 4.3.0.7 (FFT algorithm by matrix factorization)**

For  $n = 2m, m \in \mathbb{N}$ ,

consider even-odd sorting  $P_m^{\text{OE}}(1, \dots, n) = (1, 3, \dots, n-1, 2, 4, \dots, n)$ .



As  $\omega_n^{2j} = \omega_m^j$ :

$$\text{permutation of rows } P_m^{\text{OE}} \mathbf{F}_n = \left[ \begin{array}{c|c|c} & \mathbf{F}_m & \mathbf{F}_m \\ \hline \mathbf{F}_m & \begin{bmatrix} \omega_n^0 & \omega_n^1 & \ddots & \omega_n^{n/2-1} \end{bmatrix} & \mathbf{F}_m & \begin{bmatrix} \omega_n^{n/2} & \omega_n^{n/2+1} & \ddots & \omega_n^{n-1} \end{bmatrix} \\ \hline & \mathbf{F}_m & \mathbf{I} & \mathbf{I} \\ \hline & \mathbf{F}_m & \begin{bmatrix} \omega_n^0 & \omega_n^1 & \ddots & \omega_n^{n/2-1} \end{bmatrix} & \begin{bmatrix} -\omega_n^0 & -\omega_n^1 & \ddots & -\omega_n^{n/2-1} \end{bmatrix} \end{array} \right] =$$

This reveals how to apply a divide-and-conquer idea when evaluating  $\mathbf{F}_n \mathbf{x}$ .

Example: partitioning of Fourier matrix for  $n = 10$

$$P_5^{\text{OE}} \mathbf{F}_{10} = \left[ \begin{array}{cc|cc|cc} \omega^0 & \omega^0 \\ \omega^0 & \omega^2 & \omega^4 & \omega^6 & \omega^8 & \omega^0 & \omega^2 & \omega^4 & \omega^6 & \omega^8 \\ \omega^0 & \omega^4 & \omega^8 & \omega^2 & \omega^6 & \omega^0 & \omega^4 & \omega^8 & \omega^2 & \omega^6 \\ \omega^0 & \omega^6 & \omega^2 & \omega^8 & \omega^4 & \omega^0 & \omega^6 & \omega^2 & \omega^8 & \omega^4 \\ \omega^0 & \omega^8 & \omega^6 & \omega^4 & \omega^2 & \omega^0 & \omega^8 & \omega^6 & \omega^4 & \omega^2 \\ \hline \omega^0 & \omega^1 & \omega^2 & \omega^3 & \omega^4 & \omega^5 & \omega^6 & \omega^7 & \omega^8 & \omega^9 \\ \omega^0 & \omega^3 & \omega^6 & \omega^9 & \omega^2 & \omega^5 & \omega^8 & \omega^1 & \omega^4 & \omega^7 \\ \omega^0 & \omega^5 & \omega^0 & \omega^5 & \omega^0 & \omega^5 & \omega^0 & \omega^5 & \omega^0 & \omega^5 \\ \omega^0 & \omega^7 & \omega^4 & \omega^1 & \omega^8 & \omega^5 & \omega^2 & \omega^9 & \omega^6 & \omega^3 \\ \omega^0 & \omega^9 & \omega^8 & \omega^7 & \omega^6 & \omega^5 & \omega^4 & \omega^3 & \omega^2 & \omega^1 \end{array} \right], \quad \omega := \omega_{10}.$$

What if  $n \neq 2^L$ ? Quoted from MATLAB manual:

To compute an  $n$ -point DFT when  $n$  is composite (that is, when  $n = pq$ ), the FFTW library decomposes the problem using the Cooley-Tukey algorithm, which first computes  $p$  transforms of size  $q$ , and then computes  $q$  transforms of size  $p$ . The decomposition is applied recursively to both the  $p$ - and  $q$ -point DFTs until the problem can be solved using one of several machine-generated fixed-size "codelets." The codelets in turn use several algorithms in combination, including a variation of Cooley-Tukey, a prime factor algorithm, and a split-radix algorithm. The particular factorization of  $n$  is chosen heuristically.

The execution time for fft depends on the length of the transform. It is fastest for powers of two. It is almost as fast for lengths that have only small prime factors. It is typically several times slower for lengths that are prime or which have large prime factors → Ex. 4.3.0.12.

**Remark 4.3.0.8 (FFT based on general factorization)** Fast Fourier transform algorithm for DFT of length  $n = pq$ ,  $p, q \in \mathbb{N}$  (Cooley-Tukey-Algorithm)

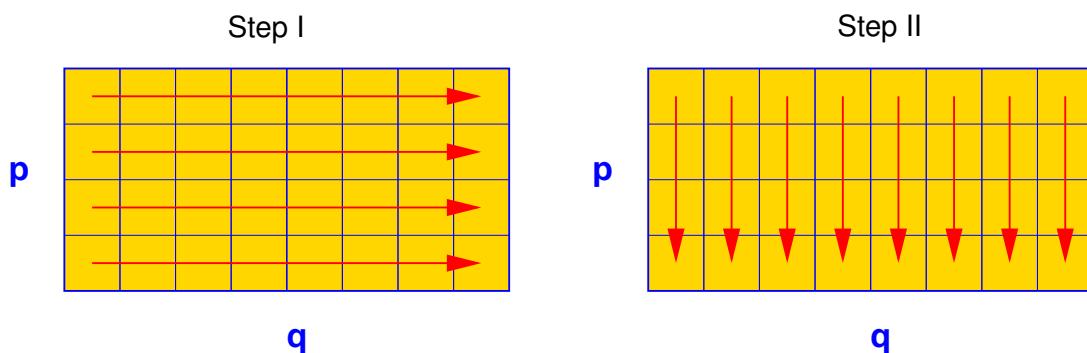
$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{jk} \stackrel{j=lp+m}{=} \sum_{m=0}^{p-1} \sum_{l=0}^{q-1} y_{lp+m} e^{-\frac{2\pi i}{pq}(lp+m)k} = \sum_{m=0}^{p-1} \omega_n^{mk} \sum_{l=0}^{q-1} y_{lp+m} \omega_q^{l(k \bmod q)}. \quad (4.3.0.9)$$

Step I: perform  $p$  DFTs of length  $q$        $z_{m,k} := \sum_{l=0}^{q-1} y_{lp+m} \omega_q^{lk}, \quad 0 \leq m < p, 0 \leq k < q$ .

Step II: for  $k =: rq + s, \quad 0 \leq r < p, 0 \leq s < q$

$$c_{rq+s} = \sum_{m=0}^{p-1} e^{-\frac{2\pi i}{pq}(rq+s)m} z_{m,s} = \sum_{m=0}^{p-1} (\omega_n^{ms} z_{m,s}) \omega_p^{mr}$$

and hence  $q$  DFTs of length  $p$  give all  $c_k$ .



**Remark 4.3.0.10 (FFT for prime  $n$ )** When  $n \neq 2^L$ , even the Cooley-Tukey algorithm of Rem. 4.3.0.8 will eventually lead to a DFT for a vector with prime length.

Quoted from the MATLAB manual:

When  $n$  is a prime number, the FFTW library first decomposes an  $n$ -point problem into three  $(n - 1)$ -point problems using Rader's algorithm [Rad68]. It then uses the Cooley-Tukey decomposition described above to compute the  $(n - 1)$ -point DFTs.

Details of Rader's algorithm: starting point is a theorem from number theory:

$$\forall p \in \mathbb{N} \text{ prime } \exists g \in \{1, \dots, p-1\}: \{g^k \bmod p : k = 1, \dots, p-1\} = \{1, \dots, p-1\},$$

- permutation  $P_{p,g} : \{1, \dots, p-1\} \mapsto \{1, \dots, p-1\}$ ,  $P_{p,g}(k) = g^k \bmod p$ ,  
 reversing permutation  $P_k : \{1, \dots, p-1\} \mapsto \{1, \dots, p-1\}$ ,  $P_k(i) = k - i + 1$ .

With these two permutations we can achieve something amazing:

For the Fourier matrix  $\mathbf{F}_p = (f_{ij})_{i,j=1}^p$  the permuted block  $P_{p-1} P_{p,g} (f_{ij})_{i,j=2}^{p-1} P_{p,g}^\top$  is circulant.

Example for  $p = 13$ :  $g = 2$ , permutation (2 4 8 3 6 12 11 9 5 10 7 1)

$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$	$\omega^0$
$\omega^0$	$\omega^2$	$\omega^4$	$\omega^8$	$\omega^3$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^7$	$\omega^1$
$\omega^0$	$\omega^1$	$\omega^2$	$\omega^4$	$\omega^8$	$\omega^3$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^7$
$\omega^0$	$\omega^7$	$\omega^1$	$\omega^2$	$\omega^4$	$\omega^8$	$\omega^3$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$
$\omega^0$	$\omega^{10}$	$\omega^7$	$\omega^1$	$\omega^2$	$\omega^4$	$\omega^8$	$\omega^3$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$
$\omega^0$	$\omega^5$	$\omega^{10}$	$\omega^7$	$\omega^1$	$\omega^2$	$\omega^4$	$\omega^8$	$\omega^3$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$
$\mathbf{F}_{13}$	→	$\omega^0$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^7$	$\omega^1$	$\omega^2$	$\omega^4$	$\omega^8$	$\omega^3$	$\omega^6$
		$\omega^0$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^7$	$\omega^1$	$\omega^2$	$\omega^4$	$\omega^8$	$\omega^3$
		$\omega^0$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^7$	$\omega^1$	$\omega^2$	$\omega^4$	$\omega^8$
		$\omega^0$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^7$	$\omega^1$	$\omega^2$	$\omega^4$
		$\omega^0$	$\omega^3$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^7$	$\omega^1$	$\omega^4$
		$\omega^0$	$\omega^8$	$\omega^3$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^7$	$\omega^1$
		$\omega^0$	$\omega^4$	$\omega^8$	$\omega^3$	$\omega^6$	$\omega^{12}$	$\omega^{11}$	$\omega^9$	$\omega^5$	$\omega^{10}$	$\omega^2$

Then apply fast algorithms for multiplication with circulant matrices (= discrete periodic convolution, see § 4.1.0.37) to right lower  $(n-1) \times (n-1)$  block of permuted Fourier matrix. These fast algorithms rely on DTFs of length  $n-1$ , see Code 4.2.1.3. □

Asymptotic complexity of `fft.fwd() / fft.inv()` for  $\mathbf{y} \in \mathbb{C}^n = O(n \log n)$ .

← Section 4.2.1

Asymptotic complexity of discrete periodic convolution, see Code 4.2.1.3:

$\text{Cost}(\text{pconvfft}(\mathbf{u}, \mathbf{x}), \mathbf{u}, \mathbf{x} \in \mathbb{C}^n) = O(n \log n)$ .

Asymptotic complexity of discrete convolution, see Code 4.2.1.4:

$\text{Cost}(\text{myconv}(\mathbf{h}, \mathbf{x}), \mathbf{h}, \mathbf{x} \in \mathbb{C}^n) = O(n \log n)$ .

#### §4.3.0.11 (FFTW - A highly-performing self-tuning library for FFT)



From [FFTW homepage](#): **FFTW** is a C subroutine library for computing the discrete Fourier transform (DFT) in one or more dimensions, of arbitrary input size, and of both real and complex data.

**FFTW** will perform well on most architectures without modification. Hence the name, "FFTW," which stands for the somewhat whimsical title of "Fastest Fourier Transform in the West."



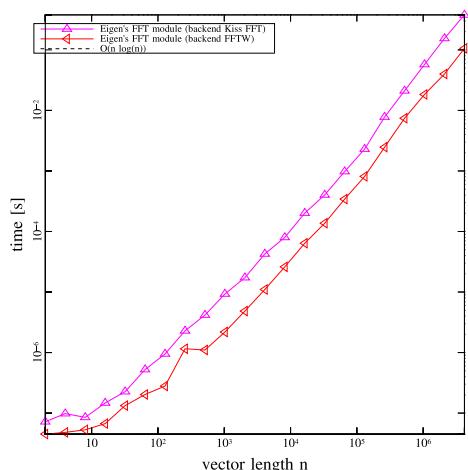
**Supplementary literature.** [FJ05] offers a comprehensive presentation of the design and implementation of the FFTW library (version 3.x). This paper also conveys the many tricks it takes to achieve satisfactory performance for DFTs of arbitrary length.

**FFTW** can be installed from source following the instructions from the [installation page](#) after downloading the source code of FFTW 3.3.5 from the [download page](#). Precompiled binaries for various linux distributions are available in their main package repositories:

- Ubuntu/Debian: `apt-get install fftw3 fftw3-dev`
  - Fedora: `dnf install fftw fftw-devel`

EIGEN's FFT module can use different backend implementations, one of which is the **FFTW** library. The backend may be enabled by defining the preprocessor directive **Eigen\_FFTW\_DEFAULT** (prior to inclusion of unsupported/Eigen/FFT) and linking with the FFTW library (-lfftw3). This setup procedure may be handled automatically by a build system like CMake (see `set_eigen_fft_backend` macro on → [GITLAB](#)).

#### **EXAMPLE 4.3.0.12 (Efficiency of `fft` for different backend implementations)**



## Platform:

- ◆ Linux (Ubuntu 16.04 64bit)
  - ◆ Intel(R) Core(TM) i7-4600U CPU @ 2.10GHz
  - ◆ L2 256KB, L3 4 MB, 8 GB DDR3 @ 1.60GHz
  - ◆ Clang 3.8.0, -O3

For reasonably high input sizes the FFTW backend gives, compared to EIGEN's default backend (Kiss FFT), a speedup of 2-4x.

## 4.4 Trigonometric transformations



*Supplementary literature.* [Han02, Sect. 55], see also [Str99] for an excellent presentation of various variants of the cosine transform.

Keeping in mind  $\exp(2\pi i x) = \cos(2\pi x) + i \sin(2\pi x)$  we may also consider the real/imaginary parts of the Fourier basis vectors  $(\mathbf{F}_n)_{:,j}$  as bases of  $\mathbb{R}^n$  and define the corresponding basis transformation. They can all be realized by means of `fft` with an asymptotic computational effort of  $\mathcal{O}(n \log n)$ . These transformations avoid the use of complex numbers.

Details are given in the sequel.

#### 4.4.1 Sine transform

Another trigonometric basis transform in  $\mathbb{R}^{n-1}$ ,  $n \in \mathbb{N}$ :

$$\text{Standard basis of } \mathbb{R}^{n-1} \quad \left\{ \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \right\} \quad \longleftarrow \quad \text{“Sine basis”} \quad \left\{ \begin{bmatrix} \sin(\frac{\pi}{n}) \\ \sin(\frac{2\pi}{n}) \\ \vdots \\ \sin(\frac{(n-1)\pi}{n}) \end{bmatrix}, \begin{bmatrix} \sin(\frac{2\pi}{n}) \\ \sin(\frac{4\pi}{n}) \\ \vdots \\ \sin(\frac{2(n-1)\pi}{n}) \end{bmatrix}, \dots, \begin{bmatrix} \sin(\frac{(n-1)\pi}{n}) \\ \sin(\frac{2(n-1)\pi}{n}) \\ \vdots \\ \sin(\frac{(n-1)^2\pi}{n}) \end{bmatrix} \right\}$$

Basis transform matrix (sine basis  $\rightarrow$  standard basis):  $\mathbf{S}_n := (\sin(jk\pi/n))_{j,k=1}^{n-1} \in \mathbb{R}^{n-1, n-1}$ .

##### Lemma 4.4.1.1. Properties of the sine matrix

$\sqrt{2/n} \mathbf{S}_n \in \mathbb{R}^{n,n}$  is real, symmetric and orthogonal ( $\rightarrow$  Def. 6.2.1.2)

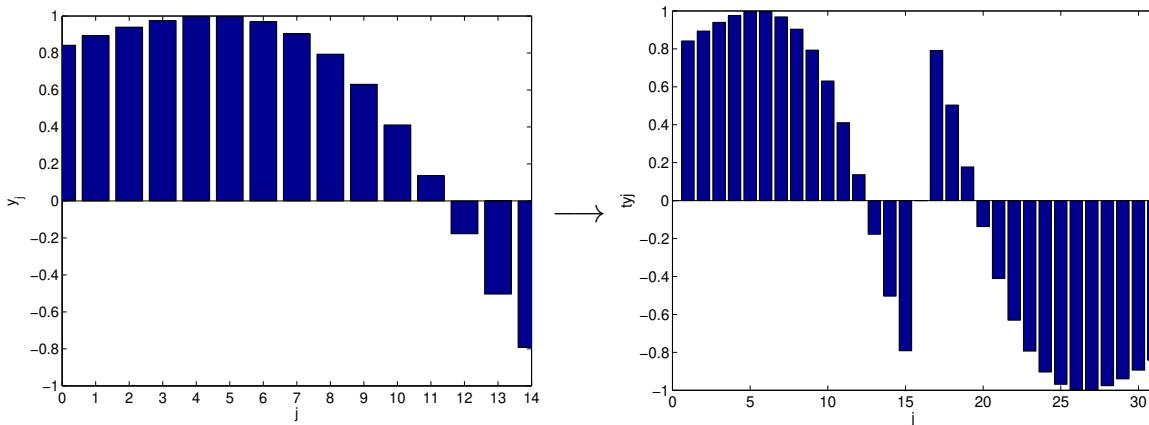
Sine transform of  $\mathbf{y} = [y_1, \dots, y_{n-1}]^\top \in \mathbb{R}^{n-1}$ : 
$$s_k = \sum_{j=1}^{n-1} y_j \sin(\pi j k / n) \quad , \quad k = 1, \dots, n-1 .$$

(4.4.1.2)

By elementary consideration we can devise a DFT-based algorithm for the sine transform ( $\triangleq \mathbf{S}_n \times$ vector):

tool: “wrap around”:  $\tilde{\mathbf{y}} \in \mathbb{R}^{2n}$ :  $\tilde{y}_j = \begin{cases} y_j & , \text{if } j = 1, \dots, n-1 , \\ 0 & , \text{if } j = 0, n , \\ -y_{2n-j} & , \text{if } j = n+1, \dots, 2n-1 . \end{cases}$  ( $\tilde{\mathbf{y}}$  “odd”)

This “wrap around” transformation can be visualized as follows:



Next we use  $\sin(x) = \frac{1}{2i}(\exp(ix) - \exp(-ix))$  to identify the DFT of a wrapped around vector as a sine transform:

$$(\mathbf{F}_{2n} \tilde{\mathbf{y}})_k \stackrel{(4.2.0.19)}{=} \sum_{j=1}^{2n-1} \tilde{y}_j e^{-\frac{2\pi i}{2n} kj} = \sum_{j=1}^{n-1} y_j e^{-\frac{\pi i}{n} kj} - \sum_{j=n+1}^{2n-1} y_{2n-j} e^{-\frac{\pi i}{n} kj}$$

$$= \sum_{j=1}^{n-1} y_j (e^{-\frac{\pi i}{n} k j} - e^{\frac{\pi i}{n} k j}) = -2i (\mathbf{S}_n \mathbf{y})_k , k = 1, \dots, n-1 .$$

**C++ code 4.4.1.3: Wrap-around implementation of sine transform → GITLAB**

```

2 // Simple sine transform of  $\mathbf{y} \in \mathbb{R}^n - 1$  into  $\mathbf{c} \in \mathbb{R}^{n-1}$  by (4.4.1.2)
3 void sinetrfwrap(const VectorXd &y, VectorXd& c)
4 {
5     VectorXd::Index n = y.size() + 1;
6     // Create wrapped vector  $\tilde{\mathbf{y}}$ 
7     VectorXd yt(2*n); yt << 0, y, 0, -y.reverse();
8
9     Eigen::VectorXcd ct;
10    Eigen::FFT<double> fft; // DFT helper class
11    fft.SetFlag(Eigen::FFT<double>::Flag::Unscaled);
12    fft.fwd(ct, yt);
13
14    const std::complex<double> v(0, 2); // factor  $2i$ 
15    c = (-ct.middleRows(1, n-1) / v).real();
16 }
```

**Remark 4.4.1.4 (Sine transform via DFT of half length)** The simple Code 4.4.1.3 relies on a DFT for vectors of length  $2n$ , which may be a waste of computational resources in some applications. A DFT of length  $n$  is sufficient as demonstrated by the following manipulations.

Step ①: transform of the coefficients

$$\tilde{y}_j = \sin(j\pi/n)(y_j + y_{n-j}) + \frac{1}{2}(y_j - y_{n-j}), \quad j = 1, \dots, n-1, \quad \tilde{y}_0 = 0 .$$

Step ②: real DFT (→ Section 4.2.3) of  $(\tilde{y}_0, \dots, \tilde{y}_{n-1}) \in \mathbb{R}^n$ :  $c_k := \sum_{j=0}^{n-1} \tilde{y}_j e^{-\frac{2\pi i}{n} j k}$

$$\begin{aligned} \text{Hence } \operatorname{Re}\{c_k\} &= \sum_{j=0}^{n-1} \tilde{y}_j \cos(-\frac{2\pi i}{n} j k) = \sum_{j=1}^{n-1} (y_j + y_{n-j}) \sin(\frac{\pi j}{n}) \cos(\frac{2\pi i}{n} j k) \\ &= \sum_{j=0}^{n-1} 2y_j \sin(\frac{\pi j}{n}) \cos(\frac{2\pi i}{n} j k) = \sum_{j=0}^{n-1} y_j \left( \sin(\frac{2k+1}{n}\pi j) - \sin(\frac{2k-1}{n}\pi j) \right) \\ &= s_{2k+1} - s_{2k-1} . \\ \operatorname{Im}\{c_k\} &= \sum_{j=0}^{n-1} \tilde{y}_j \sin(-\frac{2\pi i}{n} j k) = - \sum_{j=1}^{n-1} \frac{1}{2}(y_j - y_{n-j}) \sin(\frac{2\pi i}{n} j k) = - \sum_{j=1}^{n-1} y_j \sin(\frac{2\pi i}{n} j k) \\ &= -s_{2k} . \end{aligned}$$

Step ③: extraction of  $s_k$

$$\begin{aligned} s_{2k+1}, \quad k = 0, \dots, \frac{n}{2} - 1 &\quad \blacktriangleright \quad \text{from recursion } s_{2k+1} - s_{2k-1} = \operatorname{Re}\{c_k\}, \quad s_1 = \sum_{j=1}^{n-1} y_j \sin(\pi j/n), \\ s_{2k}, \quad k = 1, \dots, \frac{n}{2} - 2 &\quad \blacktriangleright \quad s_{2k} = -\operatorname{Im}\{c_k\} . \end{aligned}$$

Implementation (via a `fft` of length  $n/2$ ):

**C++11-code 4.4.1.5: Sine transform → GITLAB**

```

2 void sinetransform(const Eigen::VectorXd &y, Eigen::VectorXd& s)
3 {
4     int n = y.rows() + 1;
5     std::complex<double> i(0,1);
6
7     // Prepare sine terms
8     Eigen::VectorXd x = Eigen::VectorXd::LinSpaced(n-1, 1, n-1);
9     Eigen::VectorXd sinevals = x.unaryExpr([&](double z){ return
10         imag(std::pow(std::exp(i*M_PI/(double)n), z)); });
11
12     // Transform coefficients
13     Eigen::VectorXd yt(n);
14     yt(0) = 0;
15     yt.tail(n-1) = sinevals.array() * (y + y.reverse().array() +
16         0.5*(y-y.reverse()).array());
17
18     // FFT
19     Eigen::VectorXcd c;
20     Eigen::FFT<double> fft;
21     fft.fwd(c, yt);
22
23     s.resize(n);
24     s(0) = sinevals.dot(y);
25
26     for (int k=2; k<=n-1; ++k)
27     {
28         int j = k-1; // Shift index to consider indices starting from 0
29         if (k%2==0)
30             s(j) = -c(k/2).imag();
31         else
32             s(j) = s(j-2) + c((k-1)/2).real();
33     }
34 }
```

**EXAMPLE 4.4.1.6 (Diagonalization of local translation invariant linear grid operators)** We consider a so-called 5-points-stencil-operator on  $\mathbb{R}^{n,n}$ ,  $n \in \mathbb{N}$ , defined as follows

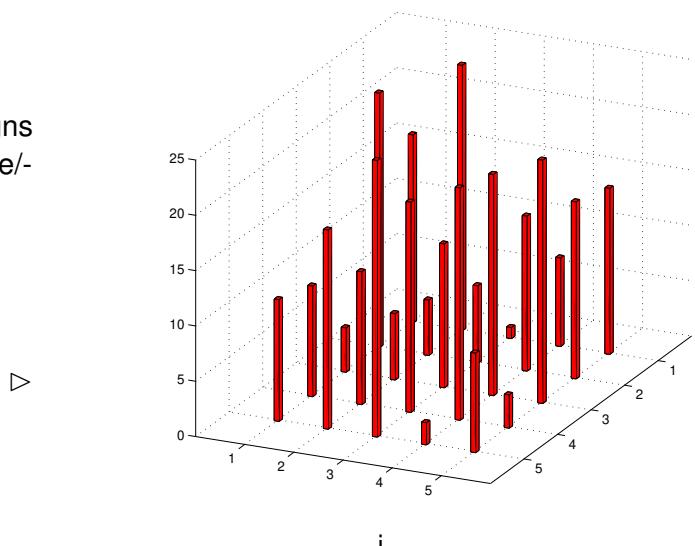
$$\mathcal{T} : \begin{cases} \mathbb{R}^{n,n} & \rightarrow \mathbb{R}^{n,n}, \\ \mathbf{x} & \mapsto \mathcal{T}(\mathbf{x}), \end{cases} \quad (\mathcal{T}(\mathbf{x}))_{ij} := c_x x_{ij} + c_y x_{i,j+1} + c_y x_{i,j-1} + c_x x_{i+1,j} + c_x x_{i-1,j} \quad (4.4.1.7)$$

with coefficients  $c, c_y, c_x \in \mathbb{R}$ , convention:  $x_{ij} := 0$  for  $(i, j) \notin \{1, \dots, n\}^2$ .

A matrix can be regarded as a function that assigns values (= matrix entries) to the points of a 2D lattice/-grid:

$$\begin{array}{c} \text{Matrix } \mathbf{X} \in \mathbb{R}^{n,n} \\ \Updownarrow \\ \text{grid function } \in \{1, \dots, n\}^2 \mapsto \mathbb{R} \end{array}$$

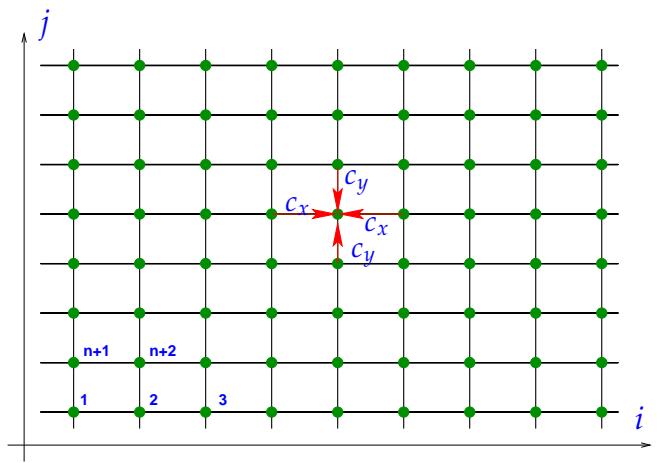
Visualization of a grid function



Identification  $\mathbb{R}^{n,n} \cong \mathbb{R}^{n^2}$ ,  $x_{ij} \sim \tilde{x}_{(j-1)n+i}$  (row-wise numbering) gives a matrix representation  $\mathbf{T} \in \mathbb{R}^{n^2, n^2}$  of  $\mathbf{T}$ :

$$\mathbf{T} = \begin{bmatrix} \mathbf{C} & c_y \mathbf{I} & 0 & \cdots & \cdots & 0 \\ c_y \mathbf{I} & \mathbf{C} & c_y \mathbf{I} & & & \vdots \\ 0 & \ddots & \ddots & \ddots & & \\ \vdots & & c_y \mathbf{I} & \mathbf{C} & c_y \mathbf{I} & \\ 0 & \cdots & \cdots & 0 & c_y \mathbf{I} & \mathbf{C} \end{bmatrix} \in \mathbb{R}^{n^2, n^2},$$

$$\mathbf{C} = \begin{bmatrix} c & c_x & 0 & \cdots & \cdots & 0 \\ c_x & c & c_x & & & \vdots \\ 0 & \ddots & \ddots & \ddots & & \\ \vdots & & c_x & c & c_x & \\ 0 & \cdots & \cdots & 0 & c_x & c \end{bmatrix} \in \mathbb{R}^{n,n}.$$

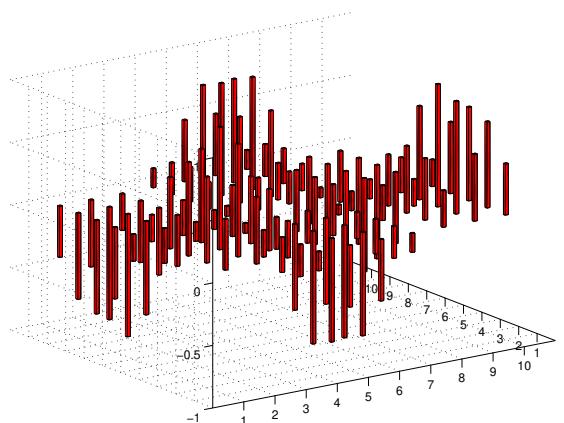


We already know the sine basis of  $\mathbb{R}^{n,n}$ :

$$\mathbf{B}^{kl} = (\sin(\frac{\pi}{n+1}ki) \sin(\frac{\pi}{n+1}lj))_{i,j=1}^n. \quad (4.4.1.8)$$

These matrices will also provide a basis of the vector space of grid functions  $\{1, \dots, n\}^2 \mapsto \mathbb{R}$ .

$n = 10$ : grid function  $\mathbf{B}^{2,3}$



The key observation is that elements of the sine basis are eigenvectors of  $\mathbf{T}$ :

$$\begin{aligned} (T(\mathbf{B}^{kl}))_{ij} &= c \sin(\frac{\pi}{n+1}ki) \sin(\frac{\pi}{n+1}lj) + c_y \sin(\frac{\pi}{n+1}ki) (\sin(\frac{\pi}{n+1}l(j-1)) + \sin(\frac{\pi}{n+1}l(j+1))) + \\ &\quad c_x \sin(\frac{\pi}{n+1}lj) (\sin(\frac{\pi}{n+1}k(i-1)) + \sin(\frac{\pi}{n+1}k(i+1))) \\ &= \sin(\frac{\pi}{n+1}ki) \sin(\frac{\pi}{n+1}lj) (c + 2c_y \cos(\frac{\pi}{n+1}l) + 2c_x \cos(\frac{\pi}{n+1}k)) \end{aligned}$$

Hence  $\mathbf{B}^{kl}$  is **eigenvector** of  $\mathbf{T}$  (or  $\mathbf{T}$  after row-wise numbering) and the corresponding eigenvalue is given by  $c + 2c_y \cos(\frac{\pi}{n+1}l) + 2c_x \cos(\frac{\pi}{n+1}k)$ . Recall very similar considerations for discrete (periodic) convolutions in 1D ( $\rightarrow$  § 4.2.0.6) and 2D ( $\rightarrow$  § 4.2.4.9)

The basis transform can be implemented efficiently based on the 1D sine transform:

$$\mathbf{X} = \sum_{k=1}^n \sum_{l=1}^n y_{kl} \mathbf{B}^{kl} \Rightarrow x_{ij} = \sum_{k=1}^n \sin\left(\frac{\pi}{n+1}ki\right) \sum_{l=1}^n y_{kl} \sin\left(\frac{\pi}{n+1}lj\right).$$

Hence nested sine transforms ( $\rightarrow$  Section 4.2.4) for rows/columns of  $\mathbf{Y} = (y_{kl})_{k,l=1}^n$ .

Here: implementation of sine transform (4.4.1.2) with “wrapping”-technique.

#### C++11-code 4.4.1.9: 2D sine transform $\rightarrow$ GITLAB

```

2 void sinetransform2d(const Eigen::MatrixXd& Y, Eigen::MatrixXd& S)
3 {
4     int m = Y.rows();
5     int n = Y.cols();
6
7     Eigen::VectorXcd c;
8     Eigen::FFT<double> fft;
9     std::complex<double> i(0,1);
10
11    Eigen::MatrixXcd C(2*m+2,n);
12    C.row(0) = Eigen::VectorXcd::Zero(n);
13    C.middleRows(1, m) = Y.cast<std::complex<double>>();
14    C.row(m+1) = Eigen::VectorXcd::Zero(n);
15    C.middleRows(m+2, m) = -Y.colwise().reverse().cast<std::complex<double>>();
16
17    // FFT on each column of C - Eigen::fft only operates on vectors
18    for (int i=0; i<n; ++i)
19    {
20        fft.fwd(c,C.col(i));
21        C.col(i) = c;
22    }
23
24    C.middleRows(1,m) = i*C.middleRows(1,m)/2.;
25
26    Eigen::MatrixXcd C2(2*n+2,m);
27    C2.row(0) = Eigen::VectorXcd::Zero(m);
28    C2.middleRows(1,n) = C.middleRows(1,m).transpose();
29    C2.row(n+1) = Eigen::VectorXcd::Zero(m);
30    C2.middleRows(n+2, n) = -C.middleRows(1,m).transpose().colwise().reverse();
31
32    // FFT on each column of C2 - Eigen::fft only operates on vectors
33    for (int i=0; i<m; ++i)
34    {
35        fft.fwd(c,C2.col(i));
36        C2.col(i) = c;
37    }
38
39    S = (i*C2.middleRows(1,n).transpose() / 2.).real();
40 }
```

#### C++11-code 4.4.1.10: FFT-based solution of local translation invariant linear operators $\rightarrow$ GITLAB

```

2 void fftbasedsolutionlocal(const Eigen::MatrixXd& B,
```

```

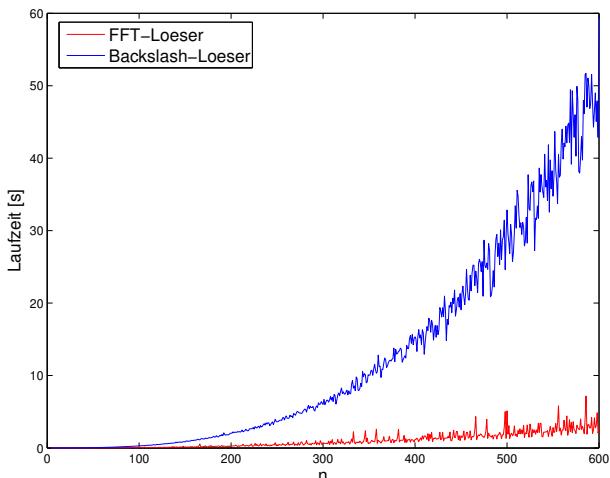
3 | double c, double cx, double cy, Eigen::MatrixXd& X)
4 | {
5 |     size_t m = B.rows();
6 |     size_t n = B.cols();
7 |
8 |     // Eigen's meshgrid
9 |     Eigen::MatrixXd I = Eigen::RowVectorXd::LinSpaced(n, 1, n).replicate(m, 1);
10 |    Eigen::MatrixXd J = Eigen::VectorXd::LinSpaced(m, 1, m).replicate(1, n);
11 |
12 |    // FFT
13 |    Eigen::MatrixXd X_;
14 |    sinetransform2d(B, X_);
15 |
16 |    // Translation
17 |    Eigen::MatrixXd T;
18 |    T = c + 2*cx*(M_PI/(n+1)*I).array().cos() +
19 |         2*cy*(M_PI/(m+1)*J).array().cos();
20 |    X_ = X_.cwiseQuotient(T);
21 |
22 |    sinetransform2d(X_, X);
23 |    X = 4*X/((m+1)*(n+1));
24 |

```

Thus the diagonalization of  $\mathbf{T}$  via 2D sine transform yields an efficient algorithm for solving linear system of equations  $\mathbf{T}(\mathbf{X}) = \mathbf{B}$ : computational cost  $O(n^2 \log n)$ . □

**EXPERIMENT 4.4.1.11 (Efficiency of FFT-based solver)** In the experiment we test the gain in runtime obtained by using DFT-based algorithms for solving linear systems of equations with coefficient matrix  $\mathbf{T}$  induced by the operator  $\mathbf{T}$  from (4.4.1.7)

**tic-toc-timing** (MATLAB V7, Linux, Intel Pentium 4 Mobile CPU 1.80GHz)



## 4.4.2 Cosine transform

Another trigonometric basis transform in  $\mathbb{R}^n$ ,  $n \in \mathbb{N}$ :

$$\text{standard basis of } \mathbb{R}^n \quad \left\{ \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \right\} \leftarrow \text{“cosine basis”} \quad \left\{ \begin{bmatrix} 2^{-1/2} \\ \cos(\frac{\pi}{2n}) \\ \cos(\frac{2\pi}{2n}) \\ \vdots \\ \cos(\frac{(n-1)\pi}{2n}) \end{bmatrix}, \begin{bmatrix} 2^{-1/2} \\ \cos(\frac{3\pi}{2n}) \\ \cos(\frac{6\pi}{2n}) \\ \vdots \\ \cos(\frac{3(n-1)\pi}{2n}) \end{bmatrix}, \dots, \begin{bmatrix} 2^{-1/2} \\ \cos(\frac{(2n-1)\pi}{2n}) \\ \cos(\frac{2(2n-1)\pi}{2n}) \\ \vdots \\ \cos(\frac{(n-1)(2n-1)\pi}{2n}) \end{bmatrix} \right\}.$$

Basis transform matrix (cosine basis  $\rightarrow$  standard basis):

$$\mathbf{C}_n = [c_{ij}]_{i,j=0}^{n-1} \in \mathbb{R}^{n,n} \quad \text{with} \quad c_{ij} = \begin{cases} 2^{-1/2} & , \text{if } i = 0 , \\ \cos(i \frac{2j+1}{2n} \pi) & , \text{if } i > 0 . \end{cases}$$

### Lemma 4.4.2.1. Properties of cosine matrix

The matrix  $\sqrt{2/n} \mathbf{C}_n \in \mathbb{R}^{n,n}$  is real and orthogonal ( $\rightarrow$  Def. 6.2.1.2).

Note:  $\mathbf{C}_n$  is not symmetric.

cosine transform of  $\mathbf{y} = [y_0, \dots, y_{n-1}]^\top$  :  $c_k = \sum_{j=0}^{n-1} y_j \cos(k \frac{2j+1}{2n} \pi)$  ,  $k = 1, \dots, n-1$  ,

(4.4.2.2)

$$c_0 = \frac{1}{\sqrt{2}} \sum_{j=0}^{n-1} y_j .$$

Implementation of  $\mathbf{Cy}$  using the "wrapping"-technique as in Code 4.4.1.3:

### C++11-code 4.4.2.3: Cosine transform $\rightarrow$ GITLAB

```

2 void cosinetransform (const Eigen::VectorXd& y, Eigen::VectorXd& c)
3 {
4     int n = y.size();
5
6     Eigen::VectorXd y_(2*n);
7     y_.head(n) = y;
8     y_.tail(n) = y.reverse();
9
10    // FFT
11    Eigen::VectorXcd z;
12    Eigen::FFT<double> fft;
13    fft.fwd(z,y_);
14
15    std::complex<double> i(0,1);
16    c.resize(n);
17    c(0) = z(0).real()/(2*sqrt(2));
18    for(size_t j=1; j<n; ++j) {

```

```

19     }
20   }
21 }
```

Implementation of  $\mathbf{C}_n^{-1}\mathbf{y}$  (“Wrapping”-technique):

#### C++11-code 4.4.2.4: Inverse cosine transform → GITLAB

```

2 void icosinetransform (const Eigen::VectorXd& c, Eigen::VectorXd& y)
3 {
4     size_t n = c.size();
5
6     std::complex<double> i(0,1);
7     Eigen::VectorXcd c_1(n);
8     c_1(0) = sqrt(2)*c(0);
9     for(size_t j=1; j<n; ++j) {
10         c_1(j) = pow(exp(-i*M_PI/(2*(double)n)), j) * c(j);
11     }
12
13     Eigen::VectorXcd c_2(2*n);
14     c_2.head(n) = c_1;
15     c_2(n) = 0;
16     c_2.tail(n-1) = c_1.tail(n-1).reverse().conjugate();
17
18 // FFT
19 Eigen::VectorXd z;
20 Eigen::FFT<double> fft;
21 fft.inv(z,c_2);
22
23 // To obtain the same result of Matlab,
24 // shift the inverse FFT result by 1.
25 Eigen::VectorXd y_(2*n);
26 y_.head(2*n-1) = z.tail(2*n-1);
27 y_(2*n-1) = z(0);
28
29 y = 2*y_.head(n);
30 }
```

The cosine transforms discussed above are named DCT-II and DCT-III.

Various cosine transforms arise by imposing various boundary conditions:

#### Remark 4.4.2.5 (Cosine transforms for compression)

- DCT-II: even around  $-1/2$  and  $N - 1/2$
  - DCT-III: even around  $0$  and odd around  $N$
- DCT-II is used in JPEG-compression while a slightly modified DCT-IV makes the main component of MP3, AAC and WMA formats.

## 4.5 Toeplitz Matrix Techniques

**EXAMPLE 4.5.0.1 (Parameter identification for linear time-invariant filters)** •  $(x_k)_{k \in \mathbb{Z}}$   $m$ -periodic

discrete signal = *known* input

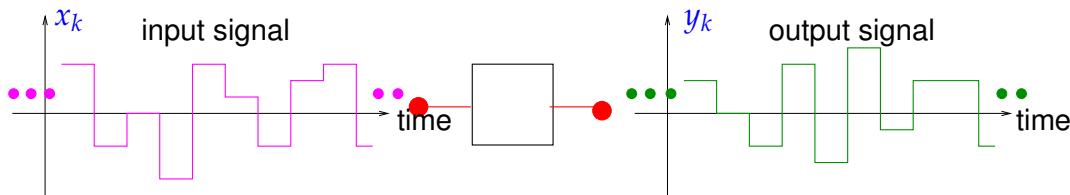
- $(y_k)_{k \in \mathbb{Z}}$  *m*-periodic measured<sup>(\*)</sup> output signal of a linear time-invariant filter, see § 4.1.0.1.  
<sup>(\*)</sup> → inevitably affected by measurement errors!
- Sought: Estimate for the impulse response ( $\rightarrow$  Def. 4.1.0.3) of the filter

This task reminds us of the parameter estimation problem from Ex. 3.0.1.4, which we tackled with least squares techniques. We employ similar ideas for the current problem

- Known: impulse response of filter has maximal duration  $n\Delta t$ ,  $n \in \mathbb{N}$ ,  $n \leq m$

cf. (4.1.0.14)

$$\exists \mathbf{h} = (h_0, \dots, h_{n-1})^\top \in \mathbb{R}^n, \quad n \leq m : \quad y_k = \sum_{j=0}^{n-1} h_j x_{k-j}. \quad (4.5.0.2)$$



If the  $y_k$  were exact, we could retrieve  $h_0, \dots, h_{n-1}$  by examining only  $y_0, \dots, y_{n-1}$  and inverting the discrete periodic convolution ( $\rightarrow$  Def. 4.1.0.33) using (4.2.0.17).

However, in case the  $y_k$  are affected by measurement errors it is advisable to use all available  $y_k$  for a least squares estimate of the impulse response.

We can now formulate the least squares parameter identification problem: seek  $\mathbf{h} = (h_0, \dots, h_{n-1})^\top \in \mathbb{R}^n$  with

$$\|\mathbf{A}\mathbf{h} - \mathbf{y}\|_2 = \left\| \begin{bmatrix} x_0 & x_{-1} & \cdots & \cdots & x_{1-n} \\ x_1 & x_0 & x_{-1} & & \vdots \\ \vdots & x_1 & x_0 & \ddots & \\ & & \ddots & \ddots & \vdots \\ & & & \ddots & x_{-1} \\ x_{n-1} & x_n & x_{n-1} & & x_1 & x_0 & x_1 \\ \vdots & & & & & \vdots & \\ x_{m-1} & \cdots & & \cdots & x_{m-n} \end{bmatrix} \begin{bmatrix} h_0 \\ \vdots \\ h_{n-1} \end{bmatrix} - \begin{bmatrix} y_0 \\ \vdots \\ y_{m-1} \end{bmatrix} \right\|_2 \rightarrow \min.$$

➤ Linear least squares problem,  $\rightarrow$  Chapter 3 with a coefficient matrix  $\mathbf{A}$  that enjoys the property that  $(\mathbf{A})_{ij} = x_{i-j}$  (constant entries of diagonals).

The coefficient matrix for the normal equations ( $\rightarrow$  Section 3.1.2, Thm. 3.1.2.1) corresponding to the above linear least squares problem is

$$\mathbf{M} := \mathbf{A}^\top \mathbf{A}, \quad (\mathbf{M})_{ij} = \sum_{k=1}^m x_{k-i} x_{k-j} = z_{i-j} \quad \text{due to periodicity of } (x_k)_{k \in \mathbb{Z}}.$$

➤

Again,  $\mathbf{M} \in \mathbb{R}^{n,n}$  is a matrix with constant diagonals & s.p.d.  
("constant diagonals"  $\Leftrightarrow$   $(\mathbf{M})_{i,j}$  depends only on  $i - j$ )

**EXAMPLE 4.5.0.3 (Linear regression for stationary Markov chains)** We consider a sequence of scalar random variables:  $(Y_k)_{k \in \mathbb{Z}}$ , a so-called **Markov chain**. These can be thought of as values for a random quantity sampled at equidistant points in time.

Assume: **stationary** (time-independent) **correlation**, that is, with  $(\mathcal{A}, \Omega, dP)$  denoting the underlying probability space,

$$\mathcal{E}(Y_{i-j} Y_{i-k}) = \int_{\Omega} Y_{i-j}(\omega) Y_{i-k}(\omega) dP(\omega) = u_{k-j} \quad \forall i, j, k \in \mathbb{Z}, \quad u_i = u_{-i}.$$

Here  $\mathcal{E}$  stands for the **expectation** of a random variable.

Model: We assume a finite linear dependency of the form

$$\exists \mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{R}^n: \quad Y_k = \sum_{j=1}^n x_j Y_{k-j} \quad \forall k \in \mathbb{Z}.$$

with *unknown* parameters  $x_j, j = 1, \dots, n$ : for fixed  $i \in \mathbb{Z}$

$$\text{Estimator} \quad \mathbf{x} = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \mathcal{E} \left| Y_i - \sum_{j=1}^n x_j Y_{i-j} \right|^2 \quad (4.5.0.4)$$

Next, we use the linearity of the expectation:

$$\begin{aligned} & \Rightarrow E|Y_i|^2 - 2 \sum_{j=1}^n x_j u_k + \sum_{k,j=1}^n x_k x_j u_{k-j} \rightarrow \min. \\ & \Rightarrow \mathbf{x}^T \mathbf{A} \mathbf{x} - 2 \mathbf{b}^T \mathbf{x} \rightarrow \min \quad \text{with } \mathbf{b} = (u_k)_{k=1}^n, \quad \mathbf{A} = (u_{i-j})_{i,j=1}^n. \end{aligned} \quad (4.5.0.5)$$

By definition  $\mathbf{A}$  is a so-called **covariance matrix** and, as such, has to be symmetric and positive definite ( $\rightarrow$  Def. 1.1.2.6). Also note that

$$\mathbf{x}^T \mathbf{A} \mathbf{x} - 2 \mathbf{b}^T \mathbf{x} = (\mathbf{x} - \mathbf{x}^*)^T \mathbf{A} (\mathbf{x} - \mathbf{x}^*) - \mathbf{x}^* \mathbf{A} \mathbf{x}^*, \quad (4.5.0.6)$$

with  $\mathbf{x}^* = \mathbf{A}^{-1} \mathbf{b}$ . Therefore  $\mathbf{x}^*$  is the unique minimizer of  $\mathbf{x}^T \mathbf{A} \mathbf{x} - 2 \mathbf{b}^T \mathbf{x}$ . The problem is reduced to solving the linear system of equations  $\mathbf{A} \mathbf{x} = \mathbf{b}$  (**Yule-Walker-equation**, see below).

Note that the covariance matrix  $\mathbf{A}$  has constant diagonals, too. □

**§4.5.0.7 (Matrices with constant diagonals)** Matrices with constant diagonals occur frequently in mathematical models, see Ex. 4.5.0.1, ???. They generalize circulant matrices ( $\rightarrow$  Def. 4.1.0.38).

Note: “Information content” of a matrix  $\mathbf{M} \in \mathbb{K}^{m,n}$  with constant diagonals, that is,  $(\mathbf{M})_{i,j} = m_{i-j}$ , is  $m+n-1$  numbers  $\in \mathbb{K}$ .

#### Definition 4.5.0.8. Toeplitz matrix

$\mathbf{T} = (t_{ij})_{i,j=1}^n \in \mathbb{K}^{m,n}$  is a **Toeplitz matrix**, if there is a vector  $\mathbf{u} = [u_{-m+1}, \dots, u_{n-1}] \in \mathbb{K}^{m+n-1}$  such that  $t_{ij} = u_{j-i}$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq n$ .

$$\mathbf{T} = \begin{bmatrix} u_0 & u_1 & \cdots & & \cdots & u_{n-1} \\ u_{-1} & u_0 & u_1 & & & \vdots \\ \vdots & \ddots & \ddots & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & u_1 \\ u_{1-m} & \cdots & & \cdots & u_{-1} & u_0 \end{bmatrix} \quad \square$$

### 4.5.1 Toeplitz Matrix Arithmetic

Given:  $\mathbf{T} = (u_{j-i}) \in \mathbb{K}^{m,n}$ , a Toeplitz matrix with generating vector  $\mathbf{u} = [u_{-m+1}, \dots, u_{n-1}]^\top \in \mathbb{K}^{m+n-1}$ , see Def. 4.5.0.8.

Task: Efficient evaluation of matrix  $\times$  vector product  $\mathbf{T}\mathbf{x}$ ,  $\mathbf{x} \in \mathbb{K}^n$

To motivate the approach we realize that we have already encountered Toeplitz matrices in the convolution of finite signals discussed in Rem. 4.1.0.17, see (4.1.0.18). The trick introduced in Rem. 4.1.0.40 was to extend the matrix to a circulant matrix by zero padding, compare (4.1.0.43).

Idea: Extend  $\mathbf{T} \in \mathbb{K}^{m,n}$  to a circulant matrix ( $\rightarrow$  Def. 4.1.0.38)  $\mathbf{C} \in \mathbb{K}^{m+n,m+n}$  generated by the  $m+n$ -periodic sequence  $(c_j)_{j \in \mathbb{Z}}$  given by



$$c_j = \begin{cases} u_j & \text{for } j = -m+1, \dots, n-1, \\ 0 & \text{for } j = n, \end{cases} \quad + \text{ periodic extension.}$$

The upper left  $m \times n$  block of  $\mathbf{C}$  contains  $\mathbf{T}$ :

$$(\mathbf{C})_{ij} = c_{i-j}, \quad 1 \leq i, j \leq m+n \Rightarrow (\mathbf{C})_{1:m, 1:n} = \mathbf{T}. \quad (4.5.1.1)$$

The following formula demonstrates the structure of  $\mathbf{C}$  in the case  $m = n$ .

$$\mathbf{C} = \left[ \begin{array}{cccccc|cccccc} u_0 & u_1 & \cdots & \cdots & u_{n-1} & 0 & u_{1-n} & \cdots & \cdots & u_{-1} \\ u_{-1} & u_0 & u_1 & & \vdots & u_{n-1} & 0 & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots & & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & u_1 & \vdots & & \ddots & \ddots & u_{1-n} \\ u_{1-n} & \cdots & \cdots & u_{-1} & u_0 & u_1 & & u_{n-1} & & 0 \\ \hline 0 & u_{1-n} & \cdots & \cdots & u_{-1} & u_0 & u_1 & \cdots & \cdots & u_{n-1} \\ u_{n-1} & 0 & \ddots & & \vdots & u_{-1} & u_0 & u_1 & & \vdots \\ \vdots & \ddots & \ddots & & \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & u_{1-n} & \vdots & & \ddots & \ddots & u_1 \\ u_1 & & u_{n-1} & 0 & u_{1-n} & u_{1-n} & \cdots & \cdots & u_{-1} & u_0 \end{array} \right]$$

Recall from 4.3 that the multiplication with a circulant  $(m+n) \times (m+n)$ -matrix (= discrete periodic convolution  $\rightarrow$  Def. 4.1.0.33) can be carried out by means of DFT/FFT with an asymptotic computational effort of  $\mathcal{O}(m+n) \log(m+n)$  for  $m, n \rightarrow \infty$ , see Code 4.2.1.3.

From (4.5.1.1) it is clear how to implement matrix  $\times$  vector for the Toeplitz matrix  $\mathbf{T}$

$$\xrightarrow{\text{zero padding}} \mathbf{C} \begin{bmatrix} \mathbf{x} \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{T}\mathbf{x} \\ * \end{bmatrix}$$

► Computational effort for computing  $\mathbf{T}\mathbf{x}$ :  $\mathcal{O}((n+m) \log(m+n))$  (FFT based, Section 4.3)  
Almost optimal in light of the data complexity  $\mathcal{O}(m+n)$  of a Toeplitz matrix.

## 4.5.2 The Levinson Algorithm

Given: Symmetric positive definite (s.p.d.) ( $\rightarrow$  Def. 1.1.2.6) Toeplitz matrix  $\mathbf{T} = (u_{j-i})_{i,j=1}^n \in \mathbb{K}^{n,n}$  with generating vector  $\mathbf{u} = [u_{-n+1}, \dots, u_{n-1}] \in \mathbb{C}^{2n-1}$

Note that the symmetry of a Toeplitz matrix is induced by the property  $u_{-k} = u_k$  of its generating vector.

Task: Find an efficient solution algorithm for the LSE  $\mathbf{T}\mathbf{x} = \mathbf{b}$ ,  $\mathbf{b} \in \mathbb{C}^n$ , the Yule-Walker problem from 4.5.0.3.

Without loss of generality we assume that  $\mathbf{T}$  has unit diagonal,  $u_0 = 1$ .

Recursive (inductive) solution strategy:

Define:

- ◆  $\mathbf{T}_k := [u_{j-i}]_{i,j=1}^k \in \mathbb{K}^{k,k}$  (left upper block of  $\mathbf{T}$ )  $\Rightarrow$   $\mathbf{T}_k$  is s.p.d. Toeplitz matrix,
- ◆  $\mathbf{x}^k \in \mathbb{K}^k$ :  $\mathbf{T}_k \mathbf{x}^k = [b_1, \dots, b_k]^\top \Leftrightarrow \mathbf{x}^k = \mathbf{T}_k^{-1} \mathbf{b}^k$ ,
- ◆  $\mathbf{u}^k := (u_1, \dots, u_k)^\top$

Thus we can block partition the LSE

$$\mathbf{T}_{k+1} \mathbf{x}^{k+1} = \begin{bmatrix} \mathbf{T}_k & \begin{matrix} u_k \\ \vdots \\ u_1 \end{matrix} \\ \hline u_k & \cdots & u_1 & 1 \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{x}}^{k+1} \\ \hline x_{k+1}^{k+1} \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_k \\ \hline b_{k+1} \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{b}}^{k+1} \\ b_{k+1} \end{bmatrix} \quad (4.5.2.1)$$

Now recall block Gaussian elimination/block-LU decomposition from Rem. 2.3.1.14, Rem. 2.3.2.19. They teach us how to eliminate  $\tilde{\mathbf{x}}^{k+1}$  and obtain an expression for  $x_{k+1}^{k+1}$ .

To state the formulas concisely, we introduce the reversing permutation:

$$P_k : \{1, \dots, k\} \mapsto \{1, \dots, k\} \quad , \quad P_k(i) := k - i + 1 . \quad (4.5.2.2)$$

Then we carry out block elimination:

$$\begin{aligned} \blacktriangleright \quad \tilde{\mathbf{x}}_{k+1} &= \mathbf{T}_k^{-1} (\tilde{\mathbf{b}}^{k+1} - x_{k+1}^{k+1} P_k \mathbf{u}^k) = \mathbf{x}^k - x_{k+1}^{k+1} \mathbf{T}_k^{-1} P_k \mathbf{u}^k , \\ x_{k+1}^{k+1} &= b_{k+1} - P_k \mathbf{u}^k \cdot \tilde{\mathbf{x}}^{k+1} = b_{k+1} - P_k \cdot \mathbf{x}^k + x_{k+1}^{k+1} P_k \cdot \mathbf{T}_k^{-1} P_k \mathbf{u}^k . \end{aligned} \quad (4.5.2.3)$$

The recursive idea is clear after introducing the *auxiliary vectors*:  $\mathbf{y}^k := \mathbf{T}_k^{-1} P_k \mathbf{u}^k$

$$\mathbf{x}^{k+1} = \begin{bmatrix} \tilde{\mathbf{x}}^{k+1} \\ x_{k+1}^{k+1} \end{bmatrix} \quad \text{with} \quad \begin{aligned} x_{k+1}^{k+1} &= (b_{k+1} - P_k \mathbf{u}^k) / \sigma_k \\ \tilde{\mathbf{x}}^{k+1} &= \mathbf{x}^k - x_{k+1}^{k+1} \mathbf{y}^k \end{aligned} \quad , \quad \sigma_k := 1 - P_k \mathbf{u}^k \cdot \mathbf{y}^k . \quad (4.5.2.4)$$

Below: **Levinson algorithm** for the solution of the Yule-Walker problem  $\mathbf{T}\mathbf{x} = \mathbf{b}$  with an s.p.d. Toeplitz matrix described by its generating vector  $\mathbf{u}$  (recursive,  $u_{n+1}$  not used!)

Linear recursion: Computational cost  $\sim (n - k)$  on level  $k$ ,  $k = 0, \dots, n - 1$

$\geq$  Asymptotic complexity  $O(n^2)$

### C++ code 4.5.2.5: Levinson algorithm → [GITLAB](#)

```

2 void levinson(const Eigen::VectorXd& u, const Eigen::VectorXd& b,
3 Eigen::VectorXd& x, Eigen::VectorXd& y)
4 {
5     size_t k = u.size() - 1;
6     if(k == 0) {
7         x.resize(1); y.resize(1);
8         x(0) = b(0); y(0) = u(0);
9         return;
10    }
11
12    Eigen::VectorXd xk, yk;
13    levinson(u.head(k), b.head(k), xk, yk);
14
15    double sigma = 1 - u.head(k).dot(yk);
16
17    double t = (b(k) - u.head(k).reverse().dot(xk)) / sigma;
18    x = xk - t * yk.head(k).reverse();
19    x.conservativeResize(x.size() + 1);
20    x(x.size() - 1) = t;
21
22    double s = (u(k) - u.head(k).reverse().dot(yk)) / sigma;
23    y = yk - s * yk.head(k).reverse();
24    y.conservativeResize(y.size() + 1);
25    y(y.size() - 1) = s;
26 }
```

### Remark 4.5.2.6 (Fast Toeplitz solvers)

FFT-based algorithms for solving  $\mathbf{T}\mathbf{x} = \mathbf{b}$  with asymptotic complexity  $O(n \log^3 n)$  [Ste03] !



*Supplementary literature.* [DR08, Sect. 8.5]: Very detailed and elementary presentation, but

the discrete Fourier transform through trigonometric interpolation, which is not covered in this chapter. Hardly addresses discrete convolution.

[Han02, Ch. IX] presents the topic from a mathematical point of few stressing approximation and trigonometric interpolation. Good reference for algorithms for circulant and Toeplitz matrices.

[Sau06, Ch. 10] also discusses the discrete Fourier transform with emphasis on interpolation and (least squares) approximation. The presentation of signal processing differs from that of the course.

There is a vast number of books and survey papers dedicated to discrete Fourier transforms, see, for instance, [Bri88; DV90]. Issues and technical details way beyond the scope of the course are discussed in these monographs.

# Bibliography

- [Bri88] E.O. Brigham. *The Fast Fourier Transform and Its Applications*. Englewood Cliffs, NJ: Prentice-Hall, 1988 (cit. on p. 317).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 296, 317).
- [DV90] P. Duhamel and M. Vetterli. “Fast Fourier transforms: a tutorial review and a state of the art”. In: *Signal Processing* 19 (1990), pp. 259–299 (cit. on pp. 299, 317).
- [FJ05] M. Frigo and S. G. Johnson. “The Design and Implementation of FFTW3”. In: *Proceedings of the IEEE* 93.2 (Feb. 2005), pp. 216–231. DOI: [10.1109/JPROC.2004.840301](https://doi.org/10.1109/JPROC.2004.840301) (cit. on p. 304).
- [Gut09] M.H. Gutknecht. *Lineare Algebra*. Lecture Notes. SAM, ETH Zürich, 2009 (cit. on p. 268).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 266, 296, 304, 317).
- [HR11] Georg Heinig and Karla Rost. “Fast algorithms for Toeplitz and Hankel matrices”. In: *Linear Algebra Appl.* 435.1 (2011), pp. 1–59.
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on pp. 268, 295).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 288, 296).
- [Rad68] C.M. Rader. “Discrete Fourier Transforms when the Number of Data Samples Is Prime”. In: *Proceedings of the IEEE* 56 (1968), pp. 1107–1108 (cit. on p. 302).
- [Sau06] T. Sauer. *Numerical analysis*. Boston: Addison Wesley, 2006 (cit. on p. 317).
- [Ste03] M. Stewart. “A Superfast Toeplitz Solver with Improved Numerical Stability”. In: *SIAM J. Matrix Analysis Appl.* 25.3 (2003), pp. 669–693 (cit. on p. 317).
- [Str99] Gilbert Strang. “The Discrete Cosine Transform”. In: *SIAM Review* 41.1 (1999), pp. 135–147. DOI: [10.1137/S0036144598336745](https://doi.org/10.1137/S0036144598336745) (cit. on p. 304).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on pp. 292, 294).

# Chapter 5

## Data Interpolation and Data Fitting in 1D

### Contents

---

5.1	Abstract Interpolation . . . . .	319
5.2	Global Polynomial Interpolation . . . . .	326
5.2.1	Polynomials . . . . .	326
5.2.2	Polynomial Interpolation: Theory . . . . .	327
5.2.3	Polynomial Interpolation: Algorithms . . . . .	331
5.2.4	Polynomial Interpolation: Sensitivity . . . . .	345
5.3	Shape-Preserving Interpolation . . . . .	348
5.3.1	Shape Properties of Functions and Data . . . . .	349
5.3.2	Piecewise Linear Interpolation . . . . .	351
5.4	Cubic Hermite Interpolation . . . . .	352
5.4.1	Definition and Algorithms . . . . .	353
5.4.2	Local Monotonicity-Preserving Hermite Interpolation . . . . .	356
5.5	Splines . . . . .	359
5.5.1	Cubic Spline Interpolation . . . . .	360
5.5.2	Structural Properties of Cubic Spline Interpolants . . . . .	363
5.5.3	Shape Preserving Spline Interpolation . . . . .	366
5.6	Trigonometric Interpolation . . . . .	370
5.6.1	Trigonometric Polynomials . . . . .	371
5.6.2	Reduction to Lagrange Interpolation . . . . .	372
5.6.3	Equidistant Trigonometric Interpolation . . . . .	374
5.7	Least Squares Data Fitting . . . . .	378

---

### 5.1 Abstract Interpolation

The task of (one-dimensional, scalar) **data interpolation** (point interpolation) can be described as follows:

#### One-dimensional interpolation

Given: data points  $(t_i, y_i)$ ,  $i = 0, \dots, n$ ,  $n \in \mathbb{N}$ ,  $t_i \in I \subset \mathbb{R}$ ,  $y_i \in \mathbb{R}$

Objective: *Reconstruction of a function  $f : I \mapsto \mathbb{R}$*   
 • satisfying the  $n + 1$  **interpolation conditions**

$$f(t_i) = y_i, \quad i = 0, \dots, n. \quad (5.1.0.2)$$

- and belonging to a set  $V$  of eligible functions.

The function  $f$  we find is called the **interpolant** of the given data set  $\{(t_i, y_i)\}_{i=0}^n$ .

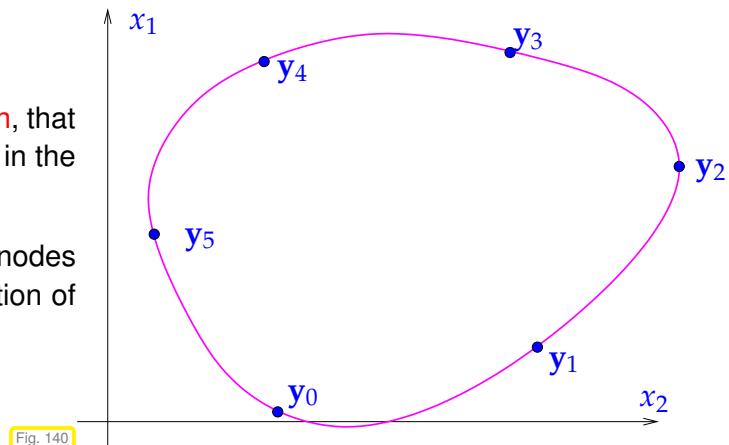
A minimal requirement on the data is that the  $t_i$  pairwise distinct:  $t_i \neq t_j$ , if  $i \neq j$ ,  $i, j \in \{0, \dots, n\}$ .

Parlance: The numbers  $t_i \in \mathbb{R}$  are called **nodes**, the  $y_i \in \mathbb{R}$  are the (data) **values**.

For ease of presentation we will usually assume that the nodes are ordered:  $t_0 < t_1 < \dots < t_n$  and  $[t_0, t_n] \subset I$ . However, algorithms often must not take for granted sorted nodes.

**Remark 5.1.0.3 (Interpolation of vector-valued data)** A natural generalization is data interpolation with **vector-valued data values**, seeking a function  $f : I \rightarrow \mathbb{R}^d$ ,  $d \in \mathbb{N}$ , such that, for given data points  $(t_i, y_i)$ ,  $t_i \in I$  mutually different,  $y_i \in \mathbb{R}^d$ , it satisfies the interpolation conditions  $f(t_i) = y_i$ ,  $i = 0, \dots, n$ .

In this case all methods available for scalar data can be applied component-wise.



An important application is **curve reconstruction**, that is the interpolation of points  $y_0, \dots, y_n \in \mathbb{R}^2$  in the plane.

A particular aspect of this problem is that the nodes  $t_i$  also have to be found, usually from the location of the  $y_i$  in a preprocessing step.

**Remark 5.1.0.4 (Multi-dimensional data interpolation)** In many applications (computer graphics, computer vision, numerical method for partial differential equations, remote sensing, geodesy, etc.) one has to reconstruct functions of several variables.

This leads to task of **multi-dimensional data interpolation**:

Given: data points  $(x_i, y_i)$ ,  $i = 0, \dots, n$ ,  $n \in \mathbb{N}$ ,  $x_i \in D \subset \mathbb{R}^m$ ,  $m > 1$ ,  $y_i \in \mathbb{R}^d$   
 Objective: *reconstruction of a (continuous) function  $f : D \mapsto \mathbb{R}^d$  satisfying the  $n + 1$  interpolation conditions*  $f(x_i) = y_i$ ,  $i = 0, \dots, n$ .

Significant additional challenges arise in a genuine multidimensional setting. A treatment is beyond the scope of this course. However, the one-dimensional techniques presented in this chapter are relevant even for multi-dimensional data interpolation, if the points  $x_i \in \mathbb{R}^m$  are points of a **finite lattice** also called **tensor product grid**.

For instance, for  $m = 2$  this is the case, if

$$\{\mathbf{x}_i\}_i = \left\{ [t_k, s_l]^\top \in \mathbb{R}^2 : k \in \{0, \dots, K\}, l \in \{0, \dots, L\} \right\}, \quad (5.1.0.5)$$

where  $t_k \in \mathbb{R}$ ,  $k = 0, \dots, K$ , and  $s_l$ ,  $l = 0, \dots, L$ ,  $K, L \in \mathbb{N}$ , are pairwise distinct nodes.  $\square$

**§5.1.0.6 (Interpolation schemes)** When we talk about “**interpolation schemes**” in 1D, we mean a mapping

$$I : \begin{cases} \mathbb{R}^{n+1} \times \mathbb{R}^{n+1} & \rightarrow \{f : I \rightarrow \mathbb{R}\} \\ ([t_i]_{i=0}^n, [y_i]_{i=0}^n) & \mapsto \text{interpolant} \end{cases}.$$

Once the function space to which the interpolant belongs is specified, then an interpolation scheme defines an “**interpolation problem**” in the sense of § 1.5.5.1. Sometimes, only the data values  $y_i$  are considered input data, whereas the dependence of the interpolant on the nodes  $t_i$  is suppressed, see Section 5.2.4.  $\square$

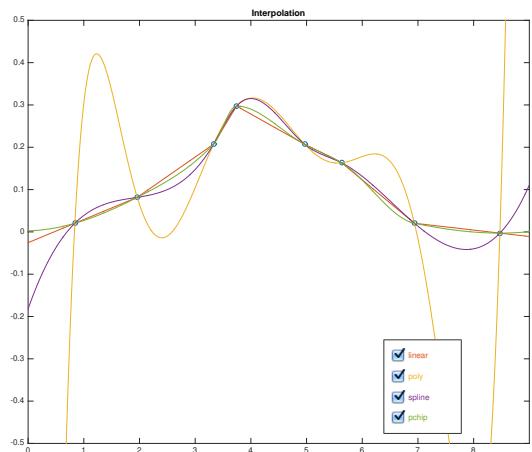


Fig. 141

▷ There are infinitely many ways to fix an interpolant for given data points.

Interpolants can have vastly different properties.

In this chapter we will discuss a few widely used methods to build interpolants and their different properties will become apparent.

► We may (have to!) impose additional requirements on the interpolant:

- minimal **smoothness** of  $f$ , e.g.  $f \in C^1$ , etc.
- special **shape** of  $f$  (positivity, monotonicity, convexity  $\rightarrow$  Section 5.3 below)

**EXAMPLE 5.1.0.7 (Constitutive relations from measurements)** This example addresses an important application of data interpolation in 1D.

In this context:  $t$ ,  $y$   $\hat{=}$  two state variables of a physical system, where  $t$  determines  $y$ : a functional dependence  $y = y(t)$  is assumed.

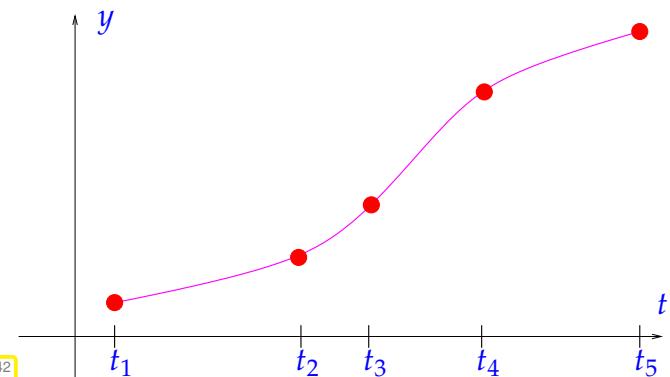
Examples:  $t$  and  $y$  could be

$t$	$y$
voltage $U$	current $I$
pressure $p$	density $\rho$
magnetic field $H$	magnetic flux $B$
...	...

Known: several *accurate* (\*) measurements

$$(t_i, y_i), \quad i = 1, \dots, m$$

Fig. 142



Imagine that  $t$ ,  $y$  correspond to the voltage  $U$  and current  $I$  measured for a 2-port non-linear circuit element (like a diode). This element will be part of a circuit, which we want to simulate based on nodal analysis as in Ex. 8.0.0.1. In order to solve the resulting non-linear system of equations  $F(\mathbf{u}) = 0$  for the nodal potentials (collected in the vector  $\mathbf{u}$ ) by means of Newton's method (→ Section 8.4) we need the voltage-current relationship for the circuit element as a continuously differentiable function  $I = f(U)$ .

(\*) Meaning of attribute “accurate”: justification for interpolation. If measured values  $y_i$  were affected by considerable errors, one would not impose the interpolation conditions (5.1.0.2), but opt for **data fitting** (→ Section 5.7). □

We can distinguish two aspects of the interpolation problem:

- ① Find interpolant  $f : I \subset \mathbb{R} \rightarrow \mathbb{R}$  and **store/represent** it (internally).
- ② Evaluate  $f$  at a few or many evaluation points  $x \in I$

**Remark 5.1.0.8 (Mathematical functions in a numerical code)** What does it mean to “represent” or “make available” a function  $f : I \subset \mathbb{R} \mapsto \mathbb{R}$  in a computer code?

! A general “mathematical” function  $f : I \subset \mathbb{R} \mapsto \mathbb{R}^d$ ,  $I$  an interval, contains an “infinite amount of information”.

Rather, in the context of numerical methods, “function” should be read as “subroutine”, a piece of code that can, for any  $x \in I$ , compute  $f(x)$  in finite time. Even this has to be qualified, because we can only pass machine numbers  $x \in I \cap \mathbb{M}$  (→ § 1.5.2.1) and, of course, in most cases,  $f(x)$  will be an approximation. In a C++ code a simple real valued function can be incarnated through a function object of a type as given in Code 5.1.0.9, see also Section 0.3.3.

#### C++-code 5.1.0.9: C++ data type representing a real-valued function

```

1 class Function {
2     private:
3         // various internal data describing f
4     public:
5         // Constructor: expects information for specifying the function
6         Function(/* ... */);
7         // Evaluation operator
8         double operator() (double t) const;
9     };

```

### §5.1.0.10 (Internal representation of classes of mathematical functions)

→ Idea: **parametrization**, a finite number of parameters  $c_0, \dots, c_m, m \in \mathbb{N}$ , characterizes  $f$ .

Special case: Representation with *finite linear combination* of basis functions  
 $b_j : I \subset \mathbb{R} \mapsto \mathbb{R}, j = 0, \dots, m$ :

$$f = \sum_{j=0}^m c_j b_j, \quad c_j \in \mathbb{R}^d. \quad (5.1.0.11)$$

→  $f$  ∈ finite dimensional **function space**  $V_m := \text{Span}\{b_0, \dots, b_m\}$ ,  $\dim V_m = m + 1$ .

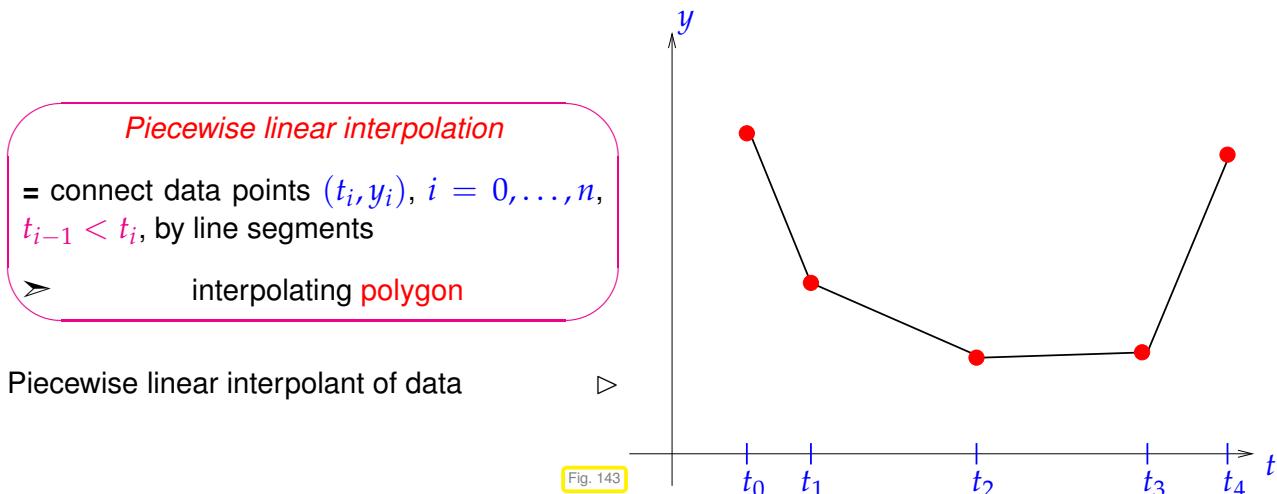
Of course, the basis functions  $b_j$  should be “simple” in the sense that  $b_j(x)$  can be computed efficiently for every  $x \in I$  and every  $j = 0, \dots, m$ .

Note that the basis functions *may depend on the nodes  $t_i$* , but they *must not depend on the values  $y_i$* .

→ The internal representation of  $f$  (in the data member section of the class **Function** from Code 5.1.0.9) will then boil down to storing the **coefficients/parameters**  $c_j, j = 0, \dots, m$ .

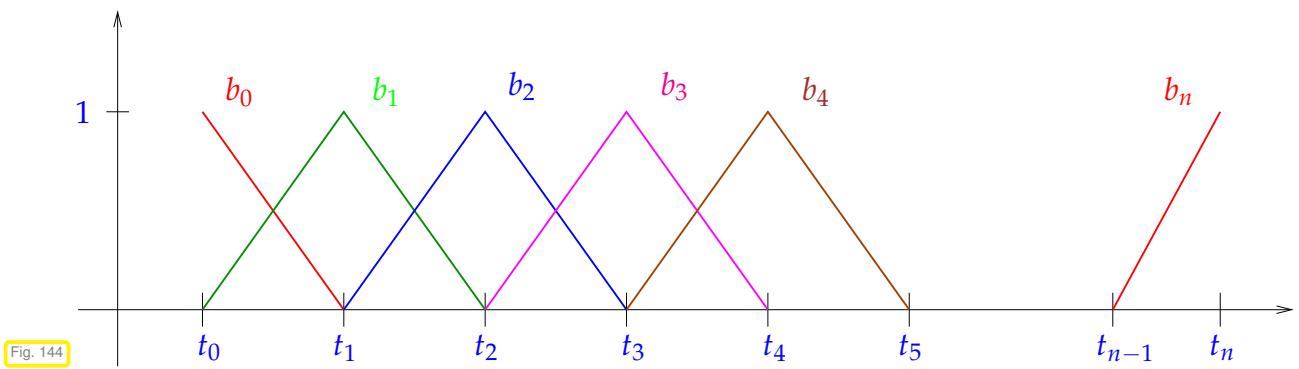
Note: The focus in this chapter will be on the special case that the data interpolants belong to a finite-dimensional space of functions spanned by “simple” basis functions.

**EXAMPLE 5.1.0.12 (Piecewise linear interpolation, see also Section 5.3.2)** Recall: A linear function in 1D is a function of the form  $x \mapsto a + bx, a, b \in \mathbb{R}$  (polynomial of degree 1).



What could be a convenient set of basis functions  $\{b_j\}_{j=0}^n$  for representing the piecewise linear interpolant through  $n + 1$  data points?

“Tent function” (“hat function”) basis:



Note: in Fig. 144 the basis functions have to be extended by zero outside the  $t$ -range where they are drawn.

Explicit formulas for these basis functions can be given and bear out that they are really “simple”:

$$\begin{aligned} b_0(t) &= \begin{cases} 1 - \frac{t-t_0}{t_1-t_0} & \text{for } t_0 \leq t < t_1, \\ 0 & \text{for } t \geq t_1. \end{cases} \\ b_j(t) &= \begin{cases} 1 - \frac{t_j-t}{t_{j+1}-t_{j-1}} & \text{for } t_{j-1} \leq t < t_j, \\ 1 - \frac{t-t_j}{t_{j+1}-t_j} & \text{for } t_j \leq t < t_{j+1}, \quad , \quad j = 1, \dots, n-1, \\ 0 & \text{elsewhere in } [t_0, t_n]. \end{cases} \\ b_n(t) &= \begin{cases} 1 - \frac{t_n-t}{t_n-t_{n-1}} & \text{for } t_{n-1} \leq t < t_n, \\ 0 & \text{for } t < t_{n-1}. \end{cases} \end{aligned} \tag{5.1.0.13}$$

Moreover, these basis functions are *uniquely* determined by the conditions

- $b_j$  is continuous on  $[t_0, t_n]$ ,
- $b_j$  is linear on each interval  $[t_{i-1}, t_i]$ ,  $i = 1, \dots, n$ ,
- $b_j(t_i) = \delta_{ij} := \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{else.} \end{cases}$   $\Rightarrow$  a so-called **cardinal basis** for the node set  $\{t_i\}_{i=0}^n$ .

This last condition implies a simple basis representation of a (the ?) piecewise linear interpolant of the data points  $(t_i, y_i)$ ,  $i = 0, \dots, n$ :

$$f(t) = \sum_{j=0}^n y_j b_j(t), \quad t_0 \leq t \leq t_n, \tag{5.1.0.14}$$

where the  $b_j$  are given by (5.1.0.13). □

**§5.1.0.15 (Interpolation as a linear mapping)** We consider the setting for interpolation that the interpolant belongs to a finite-dimension space  $V_m$  of functions spanned by basis functions  $b_0, \dots, b_m$ , see Rem. 5.1.0.8. Then the interpolation conditions imply that the basis expansion coefficients satisfy a linear system of equations:

$$(5.1.0.2) \& (5.1.0.11) \Rightarrow f(t_i) = \sum_{j=0}^m c_j b_j(t_i) = y_i, \quad i = 0, \dots, n, \tag{5.1.0.16}$$

$\Updownarrow$

$$\mathbf{Ac} := \begin{bmatrix} b_0(t_0) & \dots & b_m(t_0) \\ \vdots & & \vdots \\ b_0(t_n) & \dots & b_m(t_n) \end{bmatrix} \begin{bmatrix} c_0 \\ \vdots \\ c_m \end{bmatrix} = \begin{bmatrix} y_0 \\ \vdots \\ y_n \end{bmatrix} =: \mathbf{y}. \tag{5.1.0.17}$$

This is an  $(m+1) \times (n+1)$  linear system of equations !

The interpolation problem in  $V_m$  and the linear system (5.1.0.17) are really equivalent in the sense that (unique) solvability of one implies (unique) solvability of the other.



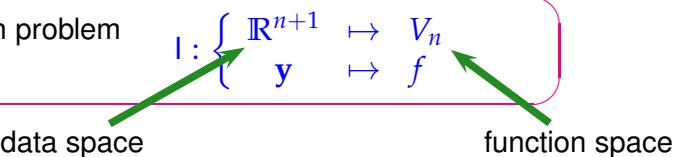
Necessary condition for unique solvability :  $m = n$   
of interpolation problem (5.1.0.16)

If  $m = n$  and  $\mathbf{A}$  from (5.1.0.17) regular ( $\rightarrow$  Def. 2.2.1.1), then for any values  $y_j, j = 0, \dots, n$  we can find coefficients  $c_j, j = 0, \dots, n$ , and, from them build the interpolant according to (5.1.0.11):

$$f = \sum_{j=0}^n (\mathbf{A}^{-1}\mathbf{y})_j b_j . \quad (5.1.0.18)$$



For fixed nodes  $t_i$  the interpolation problem  
(5.1.0.16) defines linear mapping



Beware, “linear” in the statement above has nothing to do with a linear function or piecewise linear interpolation discussed in Ex. 5.1.0.12!

### Definition 5.1.0.19. Linear interpolation operator

An interpolation operator  $I : \mathbb{R}^{n+1} \mapsto C^0([t_0, t_m])$  for the given nodes  $t_0 < t_1 < \dots < t_n$  is called linear, if

$$I(\alpha\mathbf{y} + \beta\mathbf{z}) = \alpha I(\mathbf{y}) + \beta I(\mathbf{z}) \quad \forall \mathbf{y}, \mathbf{z} \in \mathbb{R}^{n+1}, \alpha, \beta \in \mathbb{R} . \quad (5.1.0.20)$$

☞ Notation:  $C^0([t_0, t_m]) \doteq$  vector space of continuous functions on  $[t_0, t_m]$

Remark 5.1.0.21 (A data type designed for of interpolation problem) If a constitutive relationship for a circuit element is needed in a C++ simulation code ( $\rightarrow$  Ex. 5.1.0.7), the following data type could be used to represent it:

### C++-code 5.1.0.22: C++ class representing an interpolant in 1D

```

1
2 class Interpolant {
3     private:
4         // Various internal data describing f
5         // Can be the coefficients of a basis representation (5.1.0.11)
6     public:
7         // Constructor: computation of coefficients c_j of representation
8         // (5.1.0.11)
9         Interpolant(const vector<double>& t, const vector<double>& y);
10        // Evaluation operator for interpolant f

```

```

10     double operator() (double t) const;
11 };

```

Practical object oriented implementation of interpolation operator:

- ◆ Constructor: “setup phase”, e.g. building and solving linear system of equations (5.1.0.17)
- ◆ Evaluation operator, e.g., implemented as evaluation of linear combination (5.1.0.11)

Crucial issue:

computational effort for evaluation of interpolant at single point:  $O(1)$  or  $O(n)$  (or in between)?

□

## 5.2 Global Polynomial Interpolation

(Global) polynomial interpolation, that is, interpolation into spaces of functions spanned by polynomials up to a certain degree, is the simplest interpolation scheme and of great importance as building block for more complex algorithms.

### 5.2.1 Polynomials

Notation: Vector space of the polynomials of degree  $\leq k$ ,  $k \in \mathbb{N}$ :

$$\mathcal{P}_k := \{t \mapsto \alpha_k t^k + \alpha_{k-1} t^{k-1} + \cdots + \alpha_1 t + \alpha_0 \cdot 1, \alpha_j \in \mathbb{R}\}. \quad (5.2.1.1)$$

leading coefficient

Terminology: the functions  $t \mapsto t^k$ ,  $k \in \mathbb{N}_0$ , are called **monomials**

$t \mapsto \alpha_k t^k + \alpha_{k-1} t^{k-1} + \cdots + \alpha_0$  = **monomial representation** of a polynomial.

Monomial representation is a linear combination of basis functions  $t \mapsto t^k$ , see Rem. 5.1.0.8.

Obvious:  $\mathcal{P}_k$  is a vector space, see [NS02, Sect. 4.2, Bsp. 4]. What is its dimension?

#### Theorem 5.2.1.2. Dimension of space of polynomials

$$\dim \mathcal{P}_k = k + 1 \quad \text{and} \quad \mathcal{P}_k \subset C^\infty(\mathbb{R}).$$

*Proof.* Dimension formula by linear independence of monomials. □

As a consequence of Thm. 5.2.1.2 the monomial representation of a polynomial is unique.

**§5.2.1.3 (The charms of polynomials)** Why are polynomials important in computational mathematics?

- Easy to compute (only elementary operations required), integrate and differentiate
- Vector space & algebra
- Analysis: Taylor polynomials & power series

□

**Remark 5.2.1.4 (Monomial representation)** Polynomials (of degree  $k$ ) in monomial representation are stored as a vector of their coefficients  $\alpha_j$ ,  $j = 0, \dots, k$ . A convention for the ordering has to be fixed. For instance, MATLAB functions expect a monomial representation through a vector of their monomial coefficients in descending order:

MATLAB:  $\alpha_k t^k + \alpha_{k-1} t^{k-1} + \dots + \alpha_0 \rightarrow$  vector  $[\alpha_k, \alpha_{k-1}, \dots, \alpha_0]^\top$  (ordered!).

□

**Remark 5.2.1.5 (Horner scheme → [DR08, Bem. 8.11])** Efficient evaluation of a polynomial in monomial representation through Horner scheme as indicated by the following representation:

$$p(t) = t(\dots t(t(\alpha_n t + \alpha_{n-1}) + \alpha_{n-2}) + \dots + \alpha_1) + \alpha_0. \quad (5.2.1.6)$$

The following code gives an implementation based on vector data types of EIGEN. The function is vectorized in the sense that many evaluation points are processed in parallel.

#### C++-code 5.2.1.7: Horner scheme (vectorized version)

```

2 // Efficient evaluation of a polynomial in monomial representation
3 // using the Horner scheme (5.2.1.6)
4 // IN: p = vector of monomial coefficients, length = degree + 1
5 // (leading coefficient in p[0], MATLAB convention Rem. 5.2.1.4)
6 // t = vector of evaluation points t_i
7 // OUT: vector of values: polynomial evaluated at t_i
8 Eigen::VectorXd horner(const Eigen::VectorXd &p, const Eigen::VectorXd &t) {
9     const VectorXd::Index n = t.size();
10    Eigen::VectorXd y{p[0] * VectorXd::Ones(n)};
11    for (unsigned i = 1; i < p.size(); ++i)
12        y = t.cwiseProduct(y) + p[i] * VectorXd::Ones(n);
13    return y;
14 }
```

Optimal asymptotic complexity:  $O(n)$

Realized in PYTHON “built-in”-function `numpy.polyval(p,x)`. (The argument `x` can be a matrix or a vector. In this case the function evaluates the polynomial described by `p` for each entry/component.) □

## 5.2.2 Polynomial Interpolation: Theory



*Supplementary literature.* This topic is also presented in [DR08, Sect. 8.2.1], [QSS00, Sect. 8.1], [AG11, Ch. 10].

Now we consider the interpolation problem introduced in Section 5.1 for the special case that the sought interpolant belongs to the polynomial space  $\mathcal{P}_k$  (with suitable degree  $k$ ).

### Lagrange polynomial interpolation problem

Given the simple nodes  $t_0, \dots, t_n$ ,  $n \in \mathbb{N}$ ,  $-\infty < t_0 < t_1 < \dots < t_n < \infty$  and the values  $y_0, \dots, y_n \in \mathbb{R}$  compute  $p \in \mathcal{P}_n$  such that

$$p(t_j) = y_j \quad \text{for } j = 0, \dots, n. \quad (5.2.2.2)$$

Is this a well-defined problem? Obviously, it fits the framework developed in Rem. 5.1.0.8 and § 5.1.0.15, because  $\mathcal{P}_n$  is a finite-dimensional space of functions, for which we already know a basis, the monomials. Thus, in principle, we could examine the matrix  $\mathbf{A}$  from (5.1.0.17) to decide, whether the polynomial interpolant exists and is unique. However, there is a shorter way.

**§5.2.2.3 (Lagrange polynomials)** For nodes  $t_0 < t_1 < \dots < t_n$  ( $\rightarrow$  Lagrange interpolation) consider the

$$\text{Lagrange polynomials} \quad L_i(t) := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j}, \quad i = 0, \dots, n. \quad (5.2.2.4)$$

→ Evidently, the Lagrange polynomials satisfy  $L_i \in \mathcal{P}_n$  and

$$L_i(t_j) = \delta_{ij} := \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{else.} \end{cases}$$

From this relationship we infer that the Lagrange polynomials are linearly independent. Since there are  $n+1 = \dim \mathcal{P}_n$  different Lagrange polynomials, we conclude that they form a basis of  $\mathcal{P}_n$ , which is a cardinal basis for the node set  $\{t_i\}_{i=0}^n$ .  $\dashv$

### EXAMPLE 5.2.2.5 (Lagrange polynomials for uniformly spaced nodes)

Consider the equidistant nodes in  $[-1, 1]$ :

$$\mathcal{T} := \{t_j = -1 + \frac{2}{n}j\}, \quad j = 0, \dots, n.$$

The plot shows the Lagrange polynomials for this set of nodes that do not vanish in the nodes  $t_0$ ,  $t_2$ , and  $t_5$ , respectively.

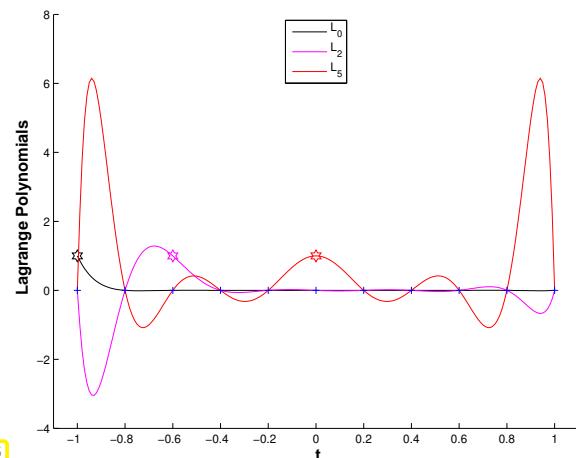


Fig. 145

The Lagrange polynomial interpolant  $p$  for data points  $(t_i, y_i)_{i=0}^n$  allows a straightforward representation with respect to the basis of Lagrange polynomials for the node set  $\{t_i\}_{i=0}^n$ :

$$p(t) = \sum_{i=0}^n y_i L_i(t) \quad \Leftrightarrow \quad p \in \mathcal{P}_n \quad \text{and} \quad p(t_i) = y_i. \quad (5.2.2.6)$$

**Theorem 5.2.2.7. Existence & uniqueness of Lagrange interpolation polynomial** → [QSS00, Thm. 8.1], [DR08, Satz 8.3]

The general Lagrange polynomial interpolation problem admits a **unique** solution  $p \in \mathcal{P}_n$  for any  $(t_i, y_i)$ .

*Proof.* Consider the *linear* evaluation operator

$$\text{eval}_{\mathcal{T}} : \begin{cases} \mathcal{P}_n &\mapsto \mathbb{R}^{n+1}, \\ p &\mapsto (p(t_i))_{i=0}^n, \end{cases}$$

which maps between finite-dimensional vector spaces of the same dimension, see Thm. 5.2.1.2.

Representation (5.2.2.6)  $\Rightarrow$  existence of interpolating polynomial  
 $\Rightarrow$   $\text{eval}_{\mathcal{T}}$  is **surjective** (“onto”)

Known from linear algebra: for a linear mapping  $T : V \mapsto W$  between finite-dimensional vector spaces with  $\dim V = \dim W$  holds the equivalence

$$T \text{ surjective} \Leftrightarrow T \text{ bijective} \Leftrightarrow T \text{ injective.}$$

Applying this equivalence to  $\text{eval}_{\mathcal{T}}$  yields the assertion of the theorem □

**Corollary 5.2.2.8. Lagrange interpolation as linear mapping** → § 5.1.0.15

The polynomial interpolation in the nodes  $\mathcal{T} := \{t_j\}_{j=0}^n$  defines a linear operator

$$I_{\mathcal{T}} : \begin{cases} \mathbb{R}^{n+1} &\rightarrow \mathcal{P}_n, \\ (y_0, \dots, y_n)^T &\mapsto \text{interpolating polynomial } p. \end{cases} \quad (5.2.2.9)$$

**Remark 5.2.2.10 (Vandermonde matrix)** Lagrangian polynomial interpolation leads to linear systems of equations also for the representation coefficients of the polynomial interpolant in **monomial basis**, see § 5.1.0.15:

$$\begin{aligned} p(t_j) = y_j &\iff \sum_{i=0}^n a_i t_j^i = y_j, \quad j = 0, \dots, n \\ &\iff \text{solution of } (n+1) \times (n+1) \text{ linear system } \mathbf{Va} = \mathbf{y} \text{ with matrix} \end{aligned}$$

$$\mathbf{V} = \begin{bmatrix} 1 & t_0 & t_0^2 & \cdots & t_0^n \\ 1 & t_1 & t_1^2 & \cdots & t_1^n \\ 1 & t_2 & t_2^2 & \cdots & t_2^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_n & t_n^2 & \cdots & t_n^n \end{bmatrix} \in \mathbb{R}^{n+1, n+1}. \quad (5.2.2.11)$$

A matrix in the form of  $\mathbf{V}$  is called **Vandermonde matrix**.

The following code initializes a Vandermonde matrix in EIGEN:

**C++ code 5.2.2.12: Initialization of Vandermonde matrix**

```
2 // Initialization of a Vandermonde matrix (5.2.2.11)
```

```

3 // from interpolation points  $t_i$ .
4 MatrixXd vander(const VectorXd &t) {
5     const VectorXd::Index n = t.size();
6     MatrixXd V(n,n); V.col(0) = VectorXd::Ones(n); V.col(1) = t;
7     // Store componentwise integer powers of point coordinate vector
8     // into the columns of the Vandermonde matrix
9     for (int j=2;j<n;j++) V.col(j) = (t.array().pow(j)).matrix();
10    return V;
11 }

```

**Remark 5.2.2.13 (Matrix representation of interpolation operator)** In the case of Lagrange interpolation:

- if Lagrange polynomials are chosen as basis for  $\mathcal{P}_n$ , then  $I_{\mathcal{T}}$  is represented by the identity matrix;
- if monomials are chosen as basis for  $\mathcal{P}_n$ , then  $I_{\mathcal{T}}$  is represented by the inverse of the Vandermonde matrix  $V$ , see Eq. (5.2.2.11).

**Remark 5.2.2.14 (Generalized polynomial interpolation** → [DR08, Sect. 8.2.7], [QSS00, Sect. 8.4])

The following generalization of Lagrange interpolation is possible: We still seek a polynomial interpolant, but beside function values also prescribe **derivatives** up to a certain order for interpolating polynomial at given nodes.

Convention: indicate occurrence of derivatives as interpolation conditions by *multiple nodes*.

#### Generalized polynomial interpolation problem

Given the (possibly multiple) nodes  $t_0, \dots, t_n$ ,  $n \in \mathbb{N}$ ,  $-\infty < t_0 \leq t_1 \leq \dots \leq t_n < \infty$  and the values  $y_0, \dots, y_n \in \mathbb{R}$  compute  $p \in \mathcal{P}_n$  such that

$$\frac{d^k}{dt^k} p(t_j) = y_j \quad \text{for } k = 0, \dots, l_j \quad \text{and } j = 0, \dots, n, \quad (5.2.2.15)$$

where  $l_j := \max\{i - i': t_j = t_i = t_{i'}, i, i' = 0, \dots, n\}$  is the **multiplicity** of the nodes  $t_j$ .

Admittedly, the statement of the generalized polynomial interpolation problem is hard to decipher. Let us look at a simple special case, which is also the most important case of generalized Lagrange interpolation. It is the case when all the multiplicities are equal to 2. It is called **Hermite interpolation** (or osculatory interpolation) and the generalized interpolation conditions read for nodes  $t_0 = t_1 < t_2 = t_3 < \dots < t_{n-1} = t_n$  (note the **double nodes!**) [QSS00, Ex. 8.6]:

$$p(t_{2j}) = y_{2j}, \quad p'(t_{2j}) = y_{2j+1}, \quad j = 0, \dots, n/2.$$

**Theorem 5.2.2.16. Existence & uniqueness of generalized Lagrange interpolation polynomials**

*The generalized polynomial interpolation problem Eq. (5.2.2.15) admits a unique solution  $p \in \mathcal{P}_n$  for any  $(t_i, y_i)$*

### Definition 5.2.2.17. Generalized Lagrange polynomials

The **generalized Lagrange polynomials** for the nodes  $\mathcal{T} = \{t_j\}_{j=0}^n \subset \mathbb{R}$  (multiple nodes allowed) are defined as  $L_i := |_{\mathcal{T}}(\mathbf{e}_{i+1})$ ,  $i = 0, \dots, n$ , where  $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^T \in \mathbb{R}^{n+1}$  are the unit vectors.

Note: The linear interpolation operator  $|_{\mathcal{T}}$  in this definition refers to generalized Lagrangian interpolation. Its existence is guaranteed by Thm. 5.2.2.16.

### EXAMPLE 5.2.2.18 (Generalized Lagrange polynomials for Hermite Interpolation)

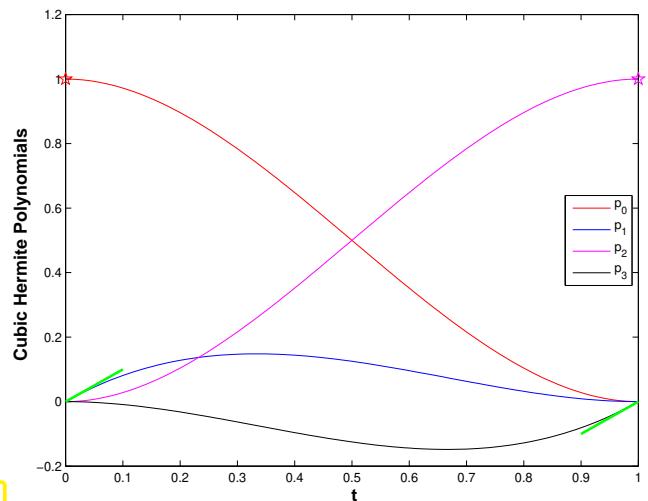
Consider the node set

$$\mathcal{T} = \{t_0 = 0, t_1 = 0, t_2 = 1, t_3 = 1\}.$$

The plot shows the four unique generalized Lagrange polynomials of degree  $n = 3$  for these nodes. They satisfy

$$\begin{aligned} p_0(0) &= 1, p_0(1) = p'_0(0) = p'_0(1) = 0, \\ p_1(1) &= 1, p_1(0) = p'_1(0) = p'_1(1) = 0, \\ p'_2(0) &= 1, p_2(1) = p_2(0) = p'_2(1) = 0, \\ p'_3(1) &= 1, p_3(1) = p_3(0) = p'_3(0) = 0. \end{aligned}$$

Fig. 146



More details are given in Section 5.4. For explicit formulas for the polynomials see (5.4.1.5).

## 5.2.3 Polynomial Interpolation: Algorithms

Now we consider the algorithmic realization of Lagrange interpolation as introduced in Section 5.2.2. The setting is as follows:

Given: nodes  $\mathcal{T} := \{-\infty < t_0 < t_1 < \dots < t_n < \infty\}$ ,  
values  $\mathbf{y} := \{y_0, y_1, \dots, y_n\}$ ,

we write  $p := |_{\mathcal{T}}(\mathbf{y})$  for the unique Lagrange interpolation polynomial whose existence is asserted by Thm. 5.2.2.7.

When used in a numerical code, different demands can be made for a routine that implements Lagrange interpolation. These demands determine, which algorithm is most suitable.

### 5.2.3.1 Multiple evaluations

Task: For ① a **fixed** set  $\{t_0, \dots, t_n\}$  of nodes,  
and ② **many** different given data values  $y_i$ ,  $i = 0, \dots, n$   
and ③ **many** arguments  $x_k$ ,  $k = 1, \dots, N$ ,  $N \gg 1$ ,  
**efficiently** compute all  $p(x_k)$  for  $p \in \mathcal{P}_n$  interpolating in  $(t_i, y_i)$ ,  $i = 0, \dots, n$ .

The definition of a possible interpolator data type could be as follows:

### C++-code 5.2.3.1: Polynomial Interpolation

```

1 class PolyInterp {
2     private:
3         // various internal data describing p
4         Eigen::VectorXd t;
5     public:
6         // Constructors taking node vector ( $t_0, \dots, t_n$ ) as argument
7         PolyInterp(const Eigen::VectorXd &_t);
8         template <typename SeqContainer>
9             PolyInterp(const SeqContainer &v);
10            // Evaluation operator for data ( $y_0, \dots, y_n$ ); computes
11            //  $p(x_k)$  for  $x_k$ 's passed in x
12            Eigen::VectorXd eval(const Eigen::VectorXd &y,
13                                const Eigen::VectorXd &x) const;
14 };

```

The member function `eval(y, x)` expects  $n$  data values in  $y$  and (any number of) evaluation points in  $x$  ( $\leftrightarrow [x_1, \dots, x_N]^\top$ ) and returns the vector  $[p(x_1), \dots, p(x_N)]^\top$ , where  $p$  is the Lagrange polynomial interpolant.

An implementation directly based on the evaluation of Lagrange polynomials (5.2.2.4) and (5.2.2.6) would incur an asymptotic computational effort of  $O(n^2N)$  for every single invocation of `eval` and large  $n, N$ .

### §5.2.3.2 (Barycentric interpolation formula)



By means of pre-computing parts of the Lagrange polynomials  $L_i$  the asymptotic effort for `eval` can be reduced substantially.

Simple manipulations starting from (5.2.2.6) give an alternative representation of  $p$ :

$$p(t) = \sum_{i=0}^n L_i(t) y_i = \sum_{i=0}^n \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j} y_i = \sum_{i=0}^n \lambda_i \prod_{\substack{j=0 \\ j \neq i}}^n (t - t_j) y_i = \prod_{j=0}^n (t - t_j) \cdot \sum_{i=0}^n \frac{\lambda_i}{t - t_i} y_i,$$

where  $\lambda_i = \frac{1}{(t_i - t_0) \cdots (t_i - t_{i-1})(t_i - t_{i+1}) \cdots (t_i - t_n)}$ ,  $i = 0, \dots, n$ : independent of  $y_i$ !

From the above formula, with  $p(t) \equiv 1, y_i = 1$ :

$$1 = \prod_{j=0}^n (t - t_j) \sum_{i=0}^n \frac{\lambda_i}{t - t_i} \quad \Rightarrow \quad \prod_{j=0}^n (t - t_j) = \frac{1}{\sum_{i=0}^n \frac{\lambda_i}{t - t_i}}$$

► Barycentric interpolation formula 
$$p(t) = \frac{\sum_{i=0}^n \frac{\lambda_i}{t - t_i} y_i}{\sum_{i=0}^n \frac{\lambda_i}{t - t_i}}. \quad (5.2.3.3)$$

with  $\lambda_i = \frac{1}{(t_i - t_0) \cdots (t_i - t_{i-1})(t_i - t_{i+1}) \cdots (t_i - t_n)}$ ,  $i = 0, \dots, n$ , independent of  $t$  and  $y_i$ . Hence, the values  $\lambda_i$  can be precomputed!

The use of (5.2.3.3) involves

- ◆ computation of weights  $\lambda_i, i = 0, \dots, n$ : cost  $O(n^2)$  (only once!),

- ◆ cost  $O(n)$  for every subsequent evaluation of  $p$ .
- ⇒ total asymptotic complexity  $O(Nn) + O(n^2)$

The following C++ class demonstrated the use of the barycentric interpolation formula for efficient multiple point evaluation of a Lagrange interpolation polynomial:

#### C++-code 5.2.3.4: Class for multiple data/multiple point evaluations

```

2 template <typename NODESCALAR = double> class BarycPolyInterp {
3     private:
4         using nodeVec_t = Eigen::Matrix<NODESCALAR, Eigen::Dynamic, 1>;
5         using idx_t = typename nodeVec_t::Index;
6         // Number n of interpolation points, deg polynomial +1
7         const idx_t n;
8         // Locations of n interpolation points
9         nodeVec_t t;
10        // Precomputed values λi, i = 0, ..., n - 1
11        nodeVec_t lambda;
12
13    public:
14        // Constructors taking node vector [t0, ..., tn]T as
15        // argument
16        BarycPolyInterp(const nodeVec_t &t);
17        // The interpolation points may also be passed in an STL container
18        template <typename SeqContainer> BarycPolyInterp(const SeqContainer &v);
19        // Computation of p(xk) for data values
20        // (y0, ..., yn) and evaluation points xk
21        template <typename RESVEC, typename DATAVEC>
22        RESVEC eval(const DATAVEC &y, const nodeVec_t &x) const;
23
24    private:
25        void init_lambda(void);
26    };

```

#### C++-code 5.2.3.5: Interpolation class: constructors

```

2 template <typename NODESCALAR>
3 BarycPolyInterp<NODESCALAR>::BarycPolyInterp(const nodeVec_t &t)
4     : n(t.size()), t(t), lambda(n) {
5     init_lambda();
6 }
7
8 template <typename NODESCALAR>
9 template <typename SeqContainer>
10 BarycPolyInterp<NODESCALAR>::BarycPolyInterp(const SeqContainer &v)
11     : n(v.size()), t(n), lambda(n) {
12     idx_t ti = 0;
13     for (auto tp : v)
14         t(ti++) = tp;
15     init_lambda();
16 }

```

#### C++-code 5.2.3.6: Interpolation class: precomputations

```

2 template <typename NODESCALAR>
3 void BarycPolyInterp<NODESCALAR>::init_lambda(void) {
4     // Precompute the weights λi with effort O(n2)

```

```

5   for (unsigned k = 0; k < n; ++k) {
6     // little workaround: in EIGEN cannot subtract a vector
7     // from a scalar; multiply scalar by vector of ones
8     lambda(k) =
9       1. / ((t(k) * nodeVec_t::Ones(k) - t.head(k)).prod() *
10      (t(k) * nodeVec_t::Ones(n - k - 1) - t.tail(n - k - 1)).prod());
11   }
12 }
```

### C++-code 5.2.3.7: Interpolation class: multiple point evaluations

```

2 template <typename NODESCALAR>
3 template <typename RESVEC, typename DATAVEC>
4 RESVEC BarycPolyInterp<NODESCALAR>::eval(const DATAVEC &y,
5                                         const nodeVec_t &x) const {
6   const idx_t N = x.size(); // No. of evaluation points
7   RESVEC p(N);           // Output vector
8   // Compute quotient of weighted sums of  $\frac{\lambda_i}{t-t_i}$ ,
9   // effort  $O(n)$ 
10  for (int i = 0; i < N; ++i) {
11    nodeVec_t z = (x[i] * nodeVec_t::Ones(n) - t);
12
13    // Check if we want to evaluate close to a node
14    const double tref{z.cwiseAbs().maxCoeff()}; // reference size
15    idx_t k;
16    if (z.cwiseAbs().minCoeff(&k) <
17        tref * std::abs(std::numeric_limits<NODESCALAR>::epsilon())) {
18      // evaluation at node  $t_k$ 
19      p[i] = y[k];
20    } else {
21      const nodeVec_t mu = lambda.cwiseQuotient(z);
22      p[i] = (mu.cwiseProduct(y)).sum() / mu.sum();
23    }
24  } // end for
25  return p;
26 }
```

Runtime measurements of direct evaluation of a polynomial in monomial representation vs. barycentric formula are reported in Exp. 5.2.3.13.

#### 5.2.3.2 Single evaluation



*Supplementary literature.* This topic is also discussed in [DR08, Sect. 8.2.2].

Task: Given a set of interpolation points  $(t_j, y_j)$ ,  $j = 0, \dots, n$ , with pairwise different interpolation nodes  $t_j$ , perform *a single* point evaluation of the Lagrange polynomial interpolant  $p$  at  $x \in \mathbb{R}$ .

We discuss the efficient implementation of the following function for  $n \gg 1$ . It is meant for a single evaluation of a Lagrange interpolant.

```
double eval(const Eigen::VectorXd &t, const Eigen::VectorXd &y,
            double x);
```

**§5.2.3.8 (Aitken-Neville scheme)** The starting point is a recursion formula for partial Lagrange interpolants: For  $0 \leq k \leq \ell \leq n$  define

$p_{k,\ell}$  := unique interpolating polynomial of degree  $\ell - k$  through  $(t_k, y_k), \dots, (t_\ell, y_\ell)$ ,

From the uniqueness of polynomial interpolants ( $\rightarrow$  Thm. 5.2.2.7) we find

$$\begin{aligned} p_{k,k}(x) &\equiv y_k \quad (\text{"constant polynomial"}) , \quad k = 0, \dots, n , \\ p_{k,\ell}(x) &= \frac{(x - t_k)p_{k+1,\ell}(x) - (x - t_\ell)p_{k,\ell-1}(x)}{t_\ell - t_k} \\ &= p_{k+1,\ell}(x) + \frac{x - t_\ell}{t_\ell - t_k} (p_{k+1,\ell}(x) - p_{k,\ell-1}(x)) , \quad 0 \leq k \leq \ell \leq n , \end{aligned} \quad (5.2.3.9)$$

because the left and right hand sides represent polynomials of degree  $\ell - k$  through the points  $(t_j, y_j)$ ,  $j = k, \dots, \ell$ .

Thus the values of the partial Lagrange interpolants can be computed sequentially and their dependencies can be expressed by the following so-called **Aitken-Neville scheme**:

$n =$	0	1	2	3
$t_0$	$y_0 =: p_{0,0}(x)$	$\xrightarrow{} p_{0,1}(x)$	$\xrightarrow{} p_{0,2}(x)$	$\xrightarrow{} p_{0,3}(x)$
$t_1$	$y_1 =: p_{1,1}(x)$	$\xrightarrow{} p_{1,2}(x)$	$\xrightarrow{} p_{1,3}(x)$	
$t_2$	$y_2 =: p_{2,2}(x)$	$\xrightarrow{} p_{2,3}(x)$		
$t_3$	$y_3 =: p_{3,3}(x)$			

(ANS)

Here, the arrows indicate contributions to the convex linear combinations of (5.2.3.9). The computation can advance from left to right, which is done in following C++ code.

#### C++-code 5.2.3.10: Aitken-Neville algorithm

```

2 // Aitken-Neville algorithm for evaluation of interpolating polynomial
3 // IN: t, y: (vectors of) interpolation data points
4 // x: (single) evaluation point
5 // OUT: value of interpolant in x
6 double ANipoleval(const VectorXd& t, VectorXd y, const double x) {
7     for (int i = 0; i < y.size(); ++i) {
8         for (int k = i - 1; k >= 0; --k) {
9             // Recursion (5.2.3.9)
10            y[k] = y[k + 1] + (y[k + 1] - y[k])*(x - t[i])/(t[i] - t[k]);
11        }
12    }
13    return y[0];
14 }
```

The vector  $y$  contains the diagonals (from bottom left to top right) of the above triangular tableaux:

i	y [ 0 ]	y [ 1 ]	y [ 2 ]	y [ 3 ]
0	$y_0$	$y_1$	$y_2$	$y_3$
1	$p_{0,1}(x)$	$y_1$	$y_2$	$y_3$
2	$p_{0,2}(x)$	$p_{1,2}(x)$	$y_2$	$y_3$
3	$p_{0,3}(x)$	$p_{1,3}(x)$	$p_{2,3}(x)$	$y_3$

Asymptotic complexity of `ANipoleval` in terms of number of data points is  $O(n^2)$  (two nested loops). This

is the same as for evaluation based on the barycentric formula, but the Aitken-Neville has a key advantage discussed in the next §.

**§5.2.3.11 (Polynomial interpolation with data updates)** The Aitken-Neville algorithm has another interesting feature, when we run through the Aitken-Neville scheme from the top left corner:

$n =$	0	1	2	3
$t_0$	$y_0 =: p_{0,0}(x)$	$\rightarrow p_{0,1}(x)$	$\rightarrow p_{0,2}(x)$	$\rightarrow p_{0,3}(x)$
$t_1$	$y_1 =: p_{1,1}(x)$	$\rightarrow p_{1,2}(x)$	$\rightarrow p_{1,3}(x)$	
$t_2$	$y_2 =: p_{2,2}(x)$	$\rightarrow p_{2,3}(x)$		
$t_3$	$y_3 =: p_{3,3}(x)$			

Thus, the values of partial polynomial interpolants at  $x$  can be computed before all data points are even processed. This results in an “update-friendly” algorithm that can efficiently supply the point values  $p_{0,k}(x)$ ,  $k = 0, \dots, n$ , while being supplied with the data points  $(t_i, y_i)$ . It can be used for the efficient implementation of the following interpolator class:

#### C++-code 5.2.3.12: Single point evaluation with data updates

```

1  class PolyEval {
2  private:
3      // evaluation point and various internal data describing the
4      // polynomials
4 public:
5     // Constructor taking the evaluation point as argument
6     PolyEval(double x);
7     // Add another data point and update internal information
8     void addPoint(t,y);
9     // Value of current interpolating polynomial at x
10    double eval(void) const;
11 };

```

#### EXPERIMENT 5.2.3.13 (Timing polynomial evaluations)

##### Timing for single-point evaluation

Comparison of the computational time needed for polynomial interpolation of

$$\{t_i = i\}_{i=1,\dots,n}, \quad \{y_i = \sqrt{i}\}_{i=1,\dots,n}$$

$n = 3, \dots, 200$ , and evaluation in a single point  $x \in [0, n]$ .

Minimum computational time over 100 runs →

The measurements were carried out with the code `polytiming.cpp` ➔ [GITLAB](#), gcc with -O3.

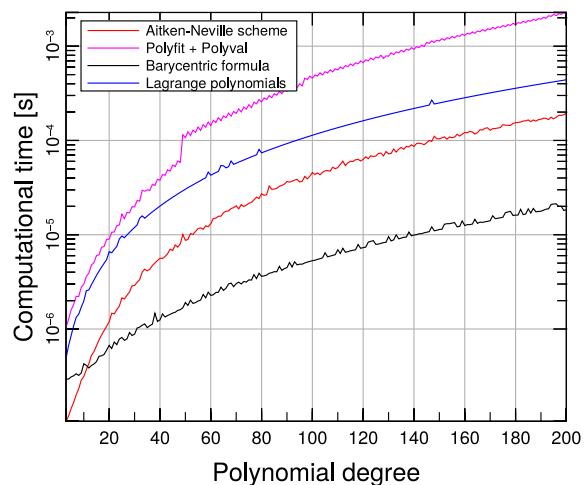


Fig. 147

This uses functions given in Code 5.2.3.7, Code 5.2.3.10 and the function `polyfit` (with a clearly greater

computational effort !)

`polyfit` is the equivalent to PYTHON's/MATLAB's built-in `polyfit`. The implementation can be found on [GitLab](#).

C++-code 5.2.3.14: Polynomial evaluation using `polyfit`

```

1 # include <Eigen/Dense>
2 # include "polyfit.hpp"
3 # include "polyval.hpp"

4

5 using Eigen::VectorXd;
6 /* Evaluation of the interpolation polynomials with polyfit+polyval
7 * IN: t = nodes
8 *      y = values in t
9 *      x = evaluation points
10 *      v will be used to save the values of interpolant in x */
11 void ipoleval(const VectorXd& t, const VectorXd& y, const VectorXd& x, VectorXd& v) {
12     // get coefficients using polyfit
13     VectorXd p = polyfit(t, y, y.size() - 1);
14     // evaluate using polyval (<-> horner scheme)
15     polyval(p, x, v);
16 }
```

C++-code 5.2.3.15: Lagrange polynomial interpolation and evaluation

```

2 // Evaluation of the interpolation polynomials with Lagrange polynomials
3 // IN: t = vector of interpolation nodes
4 // y = values in t
5 // x = evaluation points
6 // OUT p will be used to save the values of the interpolant evaluated in
7 //      x
8 void intpolyval_lag(const VectorXd& t, const VectorXd& y, const VectorXd& x,
9                     VectorXd& p) {
10    p = VectorXd::Zero(x.size());
11    for (unsigned k = 0; k < t.size(); ++k) {
12        // Compute values of k-th Lagrange polynomial in evaluation points
13        VectorXd L; lagrangepoly(x, k, t, L);
14        p += y(k)*L;
15    }
16 }
```

### 5.2.3.3 Extrapolation to zero

**Extrapolation** is interpolation with the evaluation point  $t$  outside the interval  $[\inf_{j=0,\dots,n} t_j, \sup_{j=0,\dots,n} t_j]$ . In the sequel we assume  $t = 0$ ,  $t_i > 0$ . Of course, Lagrangian polynomial interpolation can also be used for extrapolation. In this section we give a very important application of this “Lagrangian extrapolation”.

Task: compute the limit  $\lim_{h \rightarrow 0} \psi(h)$  with prescribed accuracy, though the evaluation of the function  $\psi = \psi(h)$  (maybe given in procedural form only) for very small arguments  $|h| \ll 1$  is difficult, usually because of numerically instability ( $\rightarrow$  Section 1.5.5).

The extrapolation technique introduced below works well, if

- ◆  $\psi$  is an even function of its argument:  $\psi(t) = \psi(-t)$ ,

- ◆  $\psi = \psi(h)$  behaves “nicely” around  $h = 0$ .

Theory: existence of an **asymptotic expansion** in  $h^2$

$$f(h) = f(0) + A_1 h^2 + A_2 h^4 + \cdots + A_n h^{2n} + R(h) , \quad A_k \in \mathbb{R} ,$$

with remainder estimate  $|R(h)| = O(h^{2n+2})$  for  $h \rightarrow 0$ .

### Idea: computing inaccessible limit by extrapolation to zero



- ① Pick  $h_0, \dots, h_n$  for which  $\psi$  can be evaluated “safely”.
- ② evaluation of  $\psi(h_i)$  for different  $h_i, i = 0, \dots, n, |t_i| > 0$ .
- ③  $\psi(0) \approx p(0)$  with interpolating polynomial  $p \in \mathcal{P}_n, p(h_i) = \psi(h_i)$ .

**§5.2.3.17 (Numerical differentiation through extrapolation)** In Ex. 1.5.4.7 we have already seen a situation, where we wanted to compute the limit of a function  $\psi(h)$  for  $h \rightarrow 0$ , but could not do it with sufficient accuracy. In this case  $\psi(h)$  was a one-sided difference quotient with span  $h$ , meant to approximate  $f'(x)$  for a differentiable function  $f$ . The cause of numerical difficulties was **cancellation** → § 1.5.4.5.

Now we will see how to dodge cancellation in difference quotients and how to use extrapolation to zero to computes derivatives with high accuracy:

Given: smooth function  $f : I \subset \mathbb{R} \mapsto \mathbb{R}$  in **procedural form**: function  $y = f(x)$

Sought: (approximation of)  $f'(x), x \in I$ .

Natural idea: approximation of derivative by (symmetric) **difference quotient**

$$\frac{df}{dx}(x) \approx \frac{f(x+h) - f(x-h)}{2h} . \quad (5.2.3.18)$$

► straightforward implementation fails due to cancellation in the numerator, see also Ex. 1.5.4.7.

### C++-code 5.2.3.19: Numeric differentiation through difference quotients

```

2 // numerical differentiation using difference quotient
3 // f'(x) = lim_{h \rightarrow 0} \frac{f(x+h)-f(x)}{h}
4 // IN: f (function object) = function to derive
5 // df = exact derivative (to compute error),
6 // name = string of function name (for plot filename)
7 // OUT: plot of error will be saved as "<name>numdiff.eps"
8 template <class Function, class Derivative>
9 void diff(const double x, Function& f, Derivative& df, const std::string name) {
10     std::vector<long double> error, h;
11     // build vector of widths of difference quotients
12     for (int i = -1; i >= -61; i -= 5) h.push_back(std::pow(2, i));
13     for (unsigned j = 0; j < h.size(); ++j) {
14         // compute approximate solution using difference quotient
15         double df_approx = (f(x + h[j]) - f(x))/h[j];
16         // compute relative error
17         double rel_err = std::abs((df_approx - df(x))/df(x));
18         error.push_back(rel_err);
19     }

```

$f(x) = \arctan(x)$		$f(x) = \sqrt{x}$		$f(x) = \exp(x)$	
$h$	Relative error	$h$	Relative error	$h$	Relative error
$2^{-1}$	0.20786640808609	$2^{-1}$	0.09340033543136	$2^{-1}$	0.29744254140026
$2^{-6}$	0.00773341103991	$2^{-6}$	0.00352613693103	$2^{-6}$	0.00785334954789
$2^{-11}$	0.00024299312415	$2^{-11}$	0.00011094838842	$2^{-11}$	0.00024418036620
$2^{-16}$	0.00000759482296	$2^{-16}$	0.00000346787667	$2^{-16}$	0.00000762943394
$2^{-21}$	0.00000023712637	$2^{-21}$	0.00000010812198	$2^{-21}$	0.00000023835113
$2^{-26}$	0.00000001020730	$2^{-26}$	0.00000001923506	$2^{-26}$	0.00000000429331
$2^{-31}$	0.00000005960464	$2^{-31}$	0.00000001202188	$2^{-31}$	0.00000012467100
$2^{-36}$	0.000000679016113	$2^{-36}$	0.000000198842224	$2^{-36}$	0.00000495453865

Recall the considerations elaborated in Ex. 1.5.4.7. Owing to the impact of **roundoff errors** amplified by cancellation,  $h \rightarrow 0$  does not achieve arbitrarily high accuracy. Rather, we observe fewer correct digits for very small  $h$ !

Extrapolation offers a numerically stable ( $\rightarrow$  Def. 1.5.5.19) alternative, because for a  $2(n+1)$ -times continuously differentiable function  $f : I \subset \mathbb{R} \mapsto \mathbb{R}$ ,  $x \in I$  we find that the symmetric difference quotient behaves like a polynomial in  $h^2$  in the vicinity of  $h = 0$ . Consider Taylor sum of  $f$  in  $x$  with Lagrange remainder term:

$$\psi(h) := \frac{f(x+h) - f(x-h)}{2h} \sim f'(x) + \sum_{k=1}^n \frac{1}{(2k)!} \frac{d^{2k}f}{dx^{2k}}(x) h^{2k} + \frac{1}{(2n+2)!} f^{(2n+2)}(\xi(x)).$$



Since  $\lim_{h \rightarrow 0} \psi(h) = f'(x)$

→

approximate  $f'(x)$  by interpolation of  $\psi$  in points  $h_i$ .

The following C++ function `diffex()` implements extrapolation to zero of symmetric difference quotients relying on the update-friendly version of the Aitken-Neville algorithm as presented in § 5.2.3.11, Code 5.2.3.12. Note that the extrapolated value taking into account all available difference quotients always resides in  $y[0]$ .

#### C++-code 5.2.3.20: Numerical differentiation by adaptive extrapolation to zero

```

2 // Extrapolation based numerical differentiation
3 // with a posteriori error control
4 // f: handle of a function defined in a neighbourhood of x ∈ ℝ
5 // x: point at which approximate derivative is desired
6 // h0: initial distance from x
7 // rtol: relative target tolerance, atol: absolute tolerance
8 template <class Function>
9 double diffex(Function& f, const double x, const double h0,
10               const double rtol, const double atol) {
11     const unsigned nit = 10; // Maximum depth of extrapolation
12     VectorXd h(nit); h[0] = h0; // Widths of difference quotients
13     VectorXd y(nit); // Approximations returned by difference quotients
14     y[0] = (f(x + h0) - f(x - h0))/(2*h0); // Widest difference quotients
15
16     // using Aitken-Neville scheme with x = 0, see Code 5.2.3.10
17     for (unsigned i = 1; i < nit; ++i) {
18         // create data points for extrapolation
19         h[i] = h[i-1]/2; // Next width half a big
20         y[i] = (f(x + h[i]) - f(x - h[i]))/h(i - 1);
21         // Aitken-Neville update
22         for (int k = i - 1; k >= 0; --k)

```

```

23     y[k] = y[k+1] - (y[k+1]-y[k])*h[i]/(h[i]-h[k]);
24     // termination of extrapolation when desired tolerance is reached
25     const double errest = std::abs(y[1]-y[0]); // error indicator
26     if (errest < rtol*std::abs(y[0]) || errest < atol) //
27         break;
28 }
29 return y[0]; // Return value extrapolated from largest number of
30 // difference quotients
}

```

While the extrapolation table ( $\rightarrow \S\ 5.2.3.11$ ) is computed, more and more accurate approximations of  $f'(x)$  become available. Thus, the difference between the two last approximations (stored in  $y[0]$  and  $y[1]$  in Code 5.2.3.20) can be used to gauge the error of the current approximation, it provides an **error indicator**, which can be used to decide when the level of extrapolation is sufficient, see Line 26.

	<code>auto f = [] (double x) {return std::atan(x);}</code>	<code>auto g = [] (double x) {return std::sqrt(x);}</code>	<code>auto h = [] (double x) {return std::exp(x);}</code>		
Degree	Relative error	Degree	Relative error	Degree	Relative error
0	0.04262829970946	0	0.02849215135713	0	0.04219061098749
1	0.02044767428982	1	0.01527790811946	1	0.02129207652215
2	0.00051308519253	2	0.00061205284652	2	0.00011487434095
3	0.00004087236665	3	0.00004936258481	3	0.00000825582406
4	0.00000048930018	4	0.00000067201034	4	0.00000000589624
5	0.00000000746031	5	0.00000001253250	5	0.00000000009546
6	0.00000000001224	6	0.000000000004816	6	0.00000000000002
		7	0.00000000000021		

Advantage: guaranteed accuracy  $\rightarrow$  efficiency

□

### 5.2.3.4 Newton basis and divided differences



*Supplementary literature.* We also refer to [DR08, Sect. 8.2.4], [QSS00, Sect. 8.2].

In § 5.2.3.8 we have seen a method to evaluate partial polynomial interpolants for a single or a few evaluation points efficiently. Now we want to do this for *many* evaluation points that may not be known when we receive information about the first interpolation points.

#### C++code 5.2.3.21: Polynomial evaluation

```

1 class PolyEval {
2     private:
3         // evaluation point and various internal data describing the
4         // polynomials
4     public:
5         // Idle Constructor
6         PolyEval();
7         // Add another data point and update internal information
8         void addPoint(t,y);
9         // Evaluation of current interpolating polynomial at x
10        Eigen::VectorXd operator () (const Eigen::VectorXd &x) const;
11    };

```

The challenge: Both `addPoint()` and the evaluation operator may be called many times and the implementation has to remain efficient under these circumstances.

Why not use the techniques from § 5.2.3.2? Drawback of the Lagrange basis or barycentric formula: adding another data point affects *all* basis polynomials/all precomputed values!

### §5.2.3.22 (Newton basis for $\mathcal{P}_n$ )

Our tool now: “update friendly” representation: **Newton basis** for  $\mathcal{P}_n$

$$N_0(t) := 1, \quad N_1(t) := (t - t_0), \quad \dots, \quad N_n(t) := \prod_{i=0}^{n-1} (t - t_i). \quad (5.2.3.23)$$

Note:  $N_n \in \mathcal{P}_n$  with *leading coefficient 1*  $\Rightarrow$  linear independence  $\Rightarrow$  basis property.

The abstract considerations of § 5.1.0.15 still apply and we get a linear system of equations for the coefficients  $a_j$  of the polynomial interpolant in Newton basis:

$$a_j \in \mathbb{R}: \quad a_0 N_0(t_j) + a_1 N_1(t_j) + \dots + a_n N_n(t_j) = y_j, \quad j = 0, \dots, n.$$

$\Leftrightarrow$  **triangular** linear system

$$\begin{array}{l} \left[ \begin{array}{cccc} N_0(t_0) & N_1(t_0) & \cdots & N_n(t_0) \\ N_0(t_1) & N_1(t_1) & \cdots & N_n(t_1) \\ \vdots & \vdots & & \vdots \\ N_0(t_n) & N_1(t_n) & \cdots & N_n(t_n) \end{array} \right] \left[ \begin{array}{c} a_0 \\ a_1 \\ \vdots \\ a_n \end{array} \right] = \left[ \begin{array}{c} y_0 \\ y_1 \\ \vdots \\ y_n \end{array} \right] \\ \Updownarrow \\ \left[ \begin{array}{ccccc} 1 & 0 & \cdots & 0 & \\ 1 & (t_1 - t_0) & \ddots & \vdots & \\ \vdots & \vdots & \ddots & 0 & \\ 1 & (t_n - t_0) & \cdots & \prod_{i=0}^{n-1} (t_n - t_i) & \end{array} \right] \left[ \begin{array}{c} a_0 \\ a_1 \\ \vdots \\ a_n \end{array} \right] = \left[ \begin{array}{c} y_0 \\ y_1 \\ \vdots \\ y_n \end{array} \right]. \end{array}$$

Solution of the system with forward substitution:

$$\begin{aligned} a_0 &= y_0, \\ a_1 &= \frac{y_1 - a_0}{t_1 - t_0} = \frac{y_1 - y_0}{t_1 - t_0}, \\ a_2 &= \frac{y_2 - a_0 - (t_2 - t_0)a_1}{(t_2 - t_0)(t_2 - t_1)} = \frac{y_2 - y_0 - (t_2 - t_0)\frac{y_1 - y_0}{t_1 - t_0}}{(t_2 - t_0)(t_2 - t_1)} = \frac{\frac{y_2 - y_0}{t_2 - t_0} - \frac{y_1 - y_0}{t_1 - t_0}}{t_2 - t_1}, \\ &\vdots \end{aligned}$$

Observation: same quantities computed again and again !

**§5.2.3.24 (Divided differences)** In order to reveal the pattern, we turn to a new interpretation of the coefficients  $a_j$  of the interpolating polynomials in Newton basis.

Newton basis polynomial  $N_j(t)$ : degree  $j$  and leading coefficient 1  
 $\Rightarrow a_j$  is the leading coefficient of the interpolating polynomial  $p_{0,j}$ .

(the notation  $p_{\ell,m}$  for partial polynomial interpolants through the data points  $(t_\ell, y_\ell), \dots, (t_m, y_m)$  was introduced in Section 5.2.3.2, see (5.2.3.9))

► The recursion

$$\begin{aligned} p_{k,k}(x) &\equiv y_k \quad (\text{"constant polynomial"}) , \quad k = 0, \dots, n , \\ p_{k,\ell}(x) &= \frac{(x - t_k)p_{k+1,\ell}(x) - (x - t_\ell)p_{k,\ell-1}(x)}{t_\ell - t_k} \\ &= p_{k+1,\ell}(x) + \frac{x - t_\ell}{t_\ell - t_k} (p_{k+1,\ell}(x) - p_{k,\ell-1}(x)) , \quad 0 \leq k \leq \ell \leq n , \end{aligned} \quad (5.2.3.9)$$

implies a recursion for the leading coefficients  $a_{\ell,m}$  of the interpolating polynomials  $p_{\ell,m}$ ,  $0 \leq \ell \leq m \leq n$ :

$$a_{\ell,m} = \frac{a_{\ell+1,m} - a_{\ell,m-1}}{t_m - t_\ell} . \quad (5.2.3.25)$$

Hence, instead of using elimination for a triangular linear system, we find a simpler and more efficient algorithm using the so-called **divided differences**:

$$\begin{aligned} y[t_i] &= y_i \\ y[t_i, \dots, t_{i+k}] &= \frac{y[t_{i+1}, \dots, t_{i+k}] - y[t_i, \dots, t_{i+k-1}]}{t_{i+k} - t_i} \quad (\text{recursion}) \end{aligned} \quad (5.2.3.26)$$

↓

**§5.2.3.27 (Efficient computation of divided differences)** Recursive calculation by **divided differences scheme**, cf. Aitken-Neville scheme, Code 5.2.3.10:

$$\begin{array}{c|c} t_0 & y[t_0] \\ t_1 & y[t_1] > y[t_0, t_1] \\ t_2 & y[t_2] > y[t_1, t_2] > y[t_0, t_1, t_2] \\ t_3 & y[t_3] > y[t_2, t_3] > y[t_1, t_2, t_3] > y[t_0, t_1, t_2, t_3], \end{array} \quad (5.2.3.28)$$

The elements can be computed from left to right, every “ $>$ ” indicates the evaluation of the recursion formula (5.2.3.26).

However, we can again resort to the idea of § 5.2.3.11 and traverse (5.2.3.28) along the diagonals from left bottom to right top: If a new datum  $(t_0, y_0)$  is added, it is enough to compute the  $n+1$  new divided differences

$$y[t_0] , \quad y[t_0, t_1] , \quad y[t_0, t_1, t_2] , \dots , y[t_0, \dots, t_n] .$$

The C++ function `difdiff()` listed in Code 5.2.3.29 computes divided differences for data points  $(t_i, y_i)$ ,  $i = 0, \dots, n$ , in this fashion. For  $n = 3$  the values of the outer loop variable  $\ell$  in the different combinations are as follows:

$$\begin{array}{c|c} t_0 & y[t_0] \\ t_1 & y[t_1] \quad \ell = 3 > y[t_0, t_1] \\ t_2 & y[t_2] \quad \ell = 2 > y[t_1, t_2] \quad \ell = 3 > y[t_0, t_1, t_2] \\ t_3 & y[t_3] \quad \ell = 1 > y[t_2, t_3] \quad \ell = 2 > y[t_1, t_2, t_3] \quad \ell = 3 > y[t_0, t_1, t_2, t_3], \end{array} \quad (5.2.3.28)$$

In `divdiff()` the divided differences  $y[t_0], y[t_0, t_1], \dots, y[t_0, \dots, t_n]$  overwrite the original data values  $y_j$  in the vector  $y$  (in-situ computation).

#### C++ code 5.2.3.29: In-situ computation of divided differences

```

2 // IN: t = node set (mutually different)
3 // y = nodal values
4 // OUT: y = coefficients of polynomial in Newton basis
5 void divdiff(const VectorXd &t, VectorXd &y) {
6     const int n = y.size() - 1;
7     // Follow scheme (5.2.3.28), recursion (5.2.3.25)
8     for (int l = 1; l < n; ++l)
9         for (int j = n - l; j < n; ++j)
10            y[j] = (y[j] - y[j - 1]) / (t[j] - t[n - 1 - l]);
11 }
```

By derivation the computed finite differences are the coefficients of interpolating polynomials in Newton basis:

$$p(t) = a_0 + a_1(t - t_0) + a_2(t - t_0)(t - t_1) + \dots + a_n \prod_{j=0}^{n-1} (t - t_j) \quad (5.2.3.30)$$

$$a_0 = y[t_0], a_1 = y[t_0, t_1], a_2 = y[t_0, t_1, t_2], \dots$$

Thus, `divdiff()` from Code 5.2.3.29 computes the coefficients  $a_j$ ,  $j = 0, \dots, n$ , of the polynomial interpolant with respect to the Newton basis. It uses only the first  $j+1$  data points to find  $a_j$ .  $\square$

**§5.2.3.31 (Efficient evaluation of polynomial in Newton form)** “Backward evaluation” of  $p(t)$  in the spirit of Horner’s scheme ( $\rightarrow$  Rem. 5.2.1.5, [DR08, Alg. 8.20]):

$$p \leftarrow a_n, \quad p \leftarrow (t - t_{n-1})p + a_{n-1}, \quad p \leftarrow (t - t_{n-2})p + a_{n-2}, \quad \dots$$

#### C++-code 5.2.3.32: Divided differences evaluation by modified Horner scheme

```

2 // Evaluation of polynomial in Newton basis (divided differences)
3 // IN: t = nodes (mutually different)
4 // y = values in t
5 // x = evaluation points (as Eigen::Vector)
6 // OUT: p = values in x */
7 void evaldivdiff(const VectorXd& t, const VectorXd& y, const VectorXd& x, VectorXd&
8     p) {
9     const unsigned n = y.size() - 1;
10
11    // get Newton coefficients of polynomial (non in-situ implementation!)
12    VectorXd coeffs; divdiff(t, y, coeffs);
13
14    // evaluate
15    VectorXd ones = VectorXd::Ones(x.size());
16    p = coeffs(n)*ones;
17    for (int j = n - 1; j >= 0; --j) {
18        p = (x - t(j)*ones).cwiseProduct(p) + coeffs(j)*ones;
19    }
}
```

Computational effort:

- ◆  $\mathcal{O}(n^2)$  for computation of divided differences (“setup phase”),
- ◆  $\mathcal{O}(n)$  for every single evaluation of  $p(t)$ .

(both operations can be interleaved, see Code 5.2.3.21) □

### EXAMPLE 5.2.3.33 (Class PolyEval)

Implementation of a C++ class supporting the efficient update and evaluation of an interpolating polynomial making use of

- presentation in Newton basis (5.2.3.23),
- computation of representation coefficients through divided difference scheme (5.2.3.28), see Code 5.2.3.29,
- evaluation by means of Horner scheme, see Code 5.2.3.32.

#### C++-code 5.2.3.34: Definition of a class for “update friendly” polynomial interpolant

```

1  class PolyEval {
2  private:
3      std::vector<double> t; // Interpolation nodes
4      std::vector<double> y; // Coefficients in Newton representation
5  public:
6      PolyEval(); // Idle constructor
7      void addPoint(double t, double y); // Add another data point
8      // evaluate value of current interpolating polynomial at x,
9      double operator() (double x) const;
10     private:
11     // Update internal representation, called by addPoint()
12     void divdiff();
13 }
```

#### C++-code 5.2.3.35: Implementation of class PolyEval

```

1 PolyEval::PolyEval() {}
2
3 void PolyEval::addPoint(double td, double yd) {
4     t.push_back(td); y.push_back(yd); divdiff();
5 }
6
7 void PolyEval::divdiff() {
8     int n = t.size();
9     for(int j=0; j<n-1; j++) y[n-1] = ((y[n-1]-y[j])/(t[n-1]-t[j]));
10 }
11
12 double PolyEval::operator() (double x) const {
13     double s = y.back();
14     for(int i = y.size()-2; i >= 0; --i) s = s*(x-t[i])+y[i];
15     return s;
16 }
```

**Remark 5.2.3.36 (Divided differences and derivatives)** If  $y_0, \dots, y_n$  are the values of a smooth function

$f$  in the points  $t_0, \dots, t_n$ , that is,  $y_j := f(t_j)$ , then

$$y[t_i, \dots, t_{i+k}] = \frac{f^{(k)}(\xi)}{k!}$$

for a certain  $\xi \in [t_i, t_{i+k}]$ , see [DR08, Thm. 8.21].  $\square$

## 5.2.4 Polynomial Interpolation: Sensitivity



*Supplementary literature.* For related discussions see [QSS00, Sect. 8.1.3].

This section addresses a *major shortcoming of polynomial interpolation* in case the interpolation knots  $t_i$  are imposed, which is usually the case when given data points have to be interpolated, cf. Ex. 5.1.0.7.

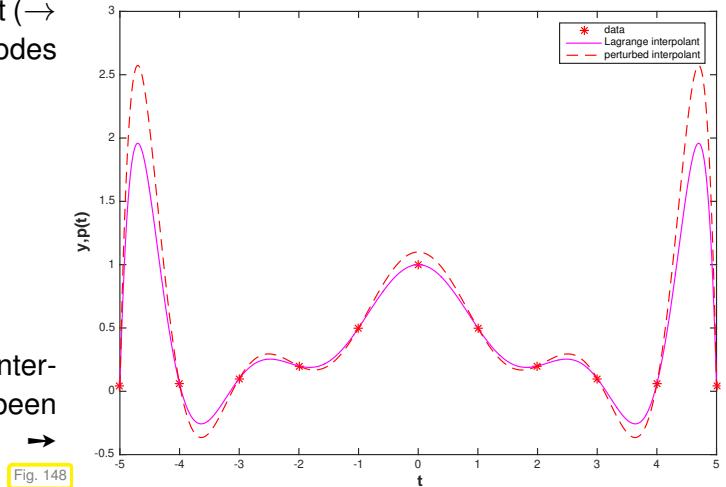
**EXAMPLE 5.2.4.1 (Oscillating polynomial interpolant (Runge's counterexample) → [DR08, Sect. 8.3], [QSS00, Ex. 8.1])** This example offers a glimpse of the problems haunting polynomial interpolation.

We examine the polynomial Lagrange interpolant ( $\rightarrow$  Section 5.2.2, (5.2.2.2)) for uniformly spaced nodes and the following data:

$$\begin{aligned} \mathcal{T} &:= \left\{ -5 + \frac{10}{n} j \right\}_{j=0}^n, \\ y_j &= \frac{1}{1+t_j^2}, \quad j = 0, \dots, n. \end{aligned}$$

Plotted is the interpolant for  $n = 10$  and the interpolant for which the data value at  $t = 0$  has been perturbed by 0.1

(See also Ex. 6.1.2.10 below.)



- ☞ possible strong oscillations of interpolating polynomials of *high degree* on *uniformly spaced nodes*!
- ☞ Slight perturbations of data values can engender strong variations of a *high-degree* Lagrange interpolant “far away”.

**§5.2.4.2 (Problem map and sensitivity)** In Section 2.2.2 we introduced the concept of **sensitivity** to describe how perturbations of the data affect the output for a problem map as defined in § 1.5.5.1. Concretely, in Section 2.2.2 we discussed the sensitivity for linear systems of equations. Motivated by Ex. 5.2.4.1, we now examine the sensitivity of Lagrange interpolation with respect to perturbations in the data values.

The problem: data  $\leftrightarrow$  values  $y_i$ ,  $i = 0, \dots, n \Rightarrow$  data space  $\mathbb{R}^{n+1}$ .  
problem map  $\leftrightarrow$  polynomial interpolation mapping  $I_{\mathcal{T}}$ , see Cor. 5.2.2.8, result space  $\mathcal{P}_n$ .

Thus, the (pointwise) sensitivity of polynomial interpolation will tell us to what extent perturbations in the  $y$ -data will affect the values of the interpolating function somewhere else. In the case of high sensitivity small perturbations in the data can cause big variations in some function values, which is clearly undesirable.  $\square$

**§5.2.4.3 (Norms on spaces of functions)** Necessary for studying sensitivity of polynomial interpolation in quantitative terms are norms ( $\rightarrow$  Def. 1.5.5.4) on the vector space of continuous functions  $C(I)$ ,  $I \subset \mathbb{R}$ . The following norms are the most relevant:

$$\text{supremum norm} \quad \|f\|_{L^\infty(I)} := \sup\{|f(t)| : t \in I\}, \quad (5.2.4.4)$$

$$L^2\text{-norm} \quad \|f\|_{L^2(I)}^2 := \int_I |f(t)|^2 dt, \quad (5.2.4.5)$$

$$L^1\text{-norm} \quad \|f\|_{L^1(I)} := \int_I |f(t)| dt. \quad (5.2.4.6)$$

Note the relationship with the vector norms introduced in § 1.5.5.3.  $\square$

**§5.2.4.7 (Sensitivity of linear problem maps)** In § 5.1.0.15 we have learned that (polynomial) interpolation gives rise to a **linear** problem map, see Def. 5.1.0.19. For this class of problem maps the investigation of sensitivity has to study **operator norms**, a generalization of matrix norms ( $\rightarrow$  Def. 1.5.5.10).

Let  $L : X \rightarrow Y$  be a *linear* problem map between two normed spaces, the data space  $X$  (with norm  $\|\cdot\|_X$ ) and the result space  $Y$  (with norm  $\|\cdot\|_Y$ ). Thanks to linearity, perturbations of the result  $y := L(x)$  for the input  $x \in X$  can be expressed as follows:

$$L(x + \delta x) = L(x) + L(\delta x) = y + L(\delta x).$$

Hence, the sensitivity (in terms of propagation of absolute errors) can be measured by the **operator norm**

$$\|L\|_{X \rightarrow Y} := \sup_{\delta x \in X \setminus \{0\}} \frac{\|L(\delta x)\|_Y}{\delta\|x\|_X}. \quad (5.2.4.8)$$

This can be read as the “matrix norm of  $L$ ”, *cf.* Def. 1.5.5.10.  $\square$

It seems challenging to compute the operator norm (5.2.4.8) for  $L = l_T$  ( $l_T$  the Lagrange interpolation operator for node set  $T \subset I$ ),  $X = \mathbb{R}^{n+1}$  (equipped with a vector norm), and  $Y = C(I)$  (endowed with a norm from § 5.2.4.3). The next lemma will provide surprisingly simple concrete formulas.

#### Lemma 5.2.4.9. Absolute conditioning of polynomial interpolation

Given a mesh  $T \subset \mathbb{R}$  with generalized Lagrange polynomials  $L_i$ ,  $i = 0, \dots, n$ , and fixed  $I \subset \mathbb{R}$ , the norm of the interpolation operator satisfies

$$\|l_T\|_{\infty \rightarrow \infty} := \sup_{y \in \mathbb{R}^{n+1} \setminus \{0\}} \frac{\|l_T(y)\|_{L^\infty(I)}}{\|y\|_\infty} = \left\| \sum_{i=0}^n |L_i| \right\|_{L^\infty(I)}, \quad (5.2.4.10)$$

$$\|l_T\|_{2 \rightarrow 2} := \sup_{y \in \mathbb{R}^{n+1} \setminus \{0\}} \frac{\|l_T(y)\|_{L^2(I)}}{\|y\|_2} \leq \left( \sum_{i=0}^n \|L_i\|_{L^2(I)}^2 \right)^{\frac{1}{2}}. \quad (5.2.4.11)$$

*Proof.* (for the  $L^\infty$ -Norm) By  $\triangle$ -inequality

$$\|l_T(y)\|_{L^\infty(I)} = \left\| \sum_{j=0}^n y_j L_j \right\|_{L^\infty(I)} \leq \sup_{t \in I} \sum_{j=0}^n |y_j| |L_j(t)| \leq \|y\|_\infty \left\| \sum_{j=0}^n |L_j| \right\|_{L^\infty(I)},$$

equality in (5.2.4.10) for  $\mathbf{y} := (\text{sgn}(L_j(t^*)))_{j=0}^n$ ,  $t^* := \text{argmax}_{t \in I} \sum_{i=0}^n |L_i(t)|$ .  $\square$

*Proof.* (for the  $L^2$ -Norm) By  $\triangle$ -inequality and Cauchy-Schwarz inequality

$$\|\mathcal{T}(\mathbf{y})\|_{L^2(I)} \leq \sum_{j=0}^n |y_j| \|L_j\|_{L^2(I)} \leq \left( \sum_{j=0}^n |y_j|^2 \right)^{\frac{1}{2}} \left( \sum_{j=0}^n \|L_j\|_{L^2(I)}^2 \right)^{\frac{1}{2}}.$$

 $\square$ 

Terminology:

$$\text{Lebesgue constant of } \mathcal{T}: \lambda_{\mathcal{T}} := \left\| \sum_{i=0}^n |L_i| \right\|_{L^\infty(I)} = \|\mathcal{T}\|_{\infty \rightarrow \infty}$$

**Remark 5.2.4.12 (Lebesgue constant for equidistant nodes)** We consider Lagrange interpolation for the special setting

$$I = [-1, 1], \quad \mathcal{T} = \{-1 + \frac{2k}{n}\}_{k=0}^n \quad (\text{uniformly spaced nodes}).$$

Asymptotic estimate (with (5.2.2.4) and Stirling formula): for  $n = 2m$

$$|L_m(1 - \frac{1}{n})| = \frac{\frac{1}{n} \cdot \frac{1}{n} \cdot \frac{3}{n} \cdots \frac{n-3}{n} \cdot \frac{n+1}{n} \cdots \frac{2n-1}{n}}{\left( \frac{2}{n} \cdot \frac{4}{n} \cdots \frac{n-2}{n} \cdot 1 \right)^2} = \frac{(2n)!}{(n-1)2^{2n}((n/2)!)^2 n!} \sim \frac{2^{n+3/2}}{\pi(n-1)n}$$

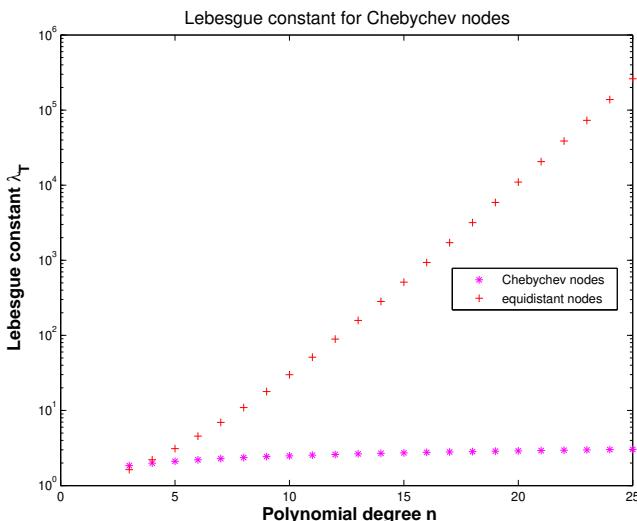


Fig. 149

Sophisticated theory [CR92] gives a lower bound for the Lebesgue constant for uniformly spaced nodes:

$$\lambda_{\mathcal{T}} \geq Ce^{n/2}$$

with  $C > 0$  independent of  $n$ .

We can also perform a numerical evaluation of the expression

$$\lambda_{\mathcal{T}} = \left\| \sum_{i=0}^n |L_i| \right\|_{L^\infty(I)},$$

for the Lebesgue constant of polynomial interpolation, see Lemma 5.2.4.9.

#### C++-code 5.2.4.13: C++ code for approximate computation of Lebesgue constants

```

2 // Computation of Lebesgue constant of polynomial interpolation
3 // with knots  $t_i$  passed in the vector  $t$  based on (5.2.4.10).
4 //  $N$  specifies the number of sampling points for the approximate
5 // computation of the maximum norm of the Lagrange polynimial
6 // on the interval  $[-1, 1]$ .
7 double lebesgue(const VectorXd& t, const unsigned& N) {
8     const unsigned n = t.size();
9     // compute denominators of normalized Lagrange polynomials relative to
10    // the nodes  $t$ 
11    VectorXd den(n);

```

```

11  for (unsigned i = 0; i < n; ++i) {
12      VectorXd tmp(n - 1);
13      // Note: comma initializer can't be used with vectors of length 0
14      if (i == 0) tmp = t.tail(n - 1);
15      else if (i == n - 1) tmp = t.head(n - 1);
16      else tmp << t.head(i), t.tail(n - (i + 1));
17      den(i) = (t(i) - tmp.array()).prod();
18  }
19
20  double l = 0; // return value
21  for (unsigned j = 0; j < N; ++j) {
22      const double x = -1 + j*(2./N); // sampling point for  $\|\cdot\|_{L^\infty([-1,1])}$ 
23      double s = 0;
24      for (unsigned k = 0; k < n; ++k) {
25          // v provides value of normalized Lagrange polynomials
26          VectorXd tmp(n - 1);
27          if (k == 0) tmp = t.tail(n - 1);
28          else if (k == n - 1) tmp = t.head(n - 1);
29          else tmp << t.head(k), t.tail(n - (k + 1));
30          double v = (x - tmp.array()).prod()/den(k);
31          s += std::abs(v); // sum over modulus of the polynomials
32      }
33      l = std::max(l, s); // maximum of sampled values
34  }
35  return l;
36 }
```

Note: In Code 5.2.4.13 the norm  $\|L_i\|_{L^\infty(I)}$  can be computed only approximate by taking the maximum modulus of function values in many sampling points.  $\square$

**§5.2.4.14 (Importance of knowing the sensitivity of polynomial interpolation)** In Ex. 5.1.0.7 we learned that interpolation is an important technique for obtaining a mathematical (and algorithmic) description of a constitutive relationship from measured data.

If the interpolation operator is poorly conditioned, tiny measurement errors will lead to big (local) deviations of the interpolant from its “true” form.

Since measurement errors are inevitable, poorly conditioned interpolation procedures are useless for determining constitutive relationships from measurements.  $\square$

## 5.3 Shape-Preserving Interpolation

When reconstructing a quantitative dependence of quantities from measurements, first principles from physics often stipulate qualitative constraints, which translate into *shape properties* of the function  $f$ , e.g., when modelling the material law for a gas:

$t_i$  pressure values,  $y_i$  densities  $\geq f$  positive & monotonic.

Notation: given data:  $(t_i, y_i) \in \mathbb{R}^2$ ,  $i = 0, \dots, n$ ,  $n \in \mathbb{N}$ ,  $t_0 < t_1 < \dots < t_n$ .

### EXAMPLE 5.3.0.1 (Magnetization curves)

For many materials physics stipulates properties of the functional dependence of magnetic flux  $B$  from magnetic field strength  $H$ :

- ◆  $H \mapsto B(H)$  smooth (at least  $C^1$ ),
- ◆  $H \mapsto B(H)$  monotonic (increasing),
- ◆  $H \mapsto B(H)$  concave

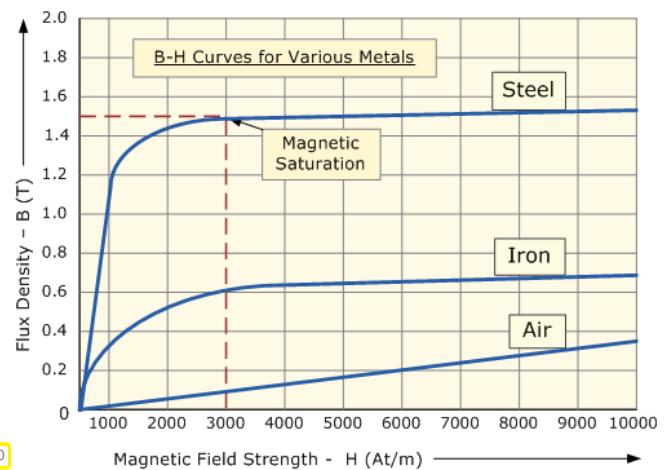


Fig. 150

### 5.3.1 Shape Properties of Functions and Data

**§5.3.1.1 (The “shape” of data)** The section is about “shape preservation”. In the previous example we have already seen a few properties that constitute the “shape” of a function: sign, monotonicity and curvature. Now we have to identify analogous properties of data sets in the form of sequences of interpolation points  $(t_j, y_j)$ ,  $j = 0, \dots, n$ ,  $t_j$  pairwise distinct.

#### Definition 5.3.1.2. monotonic data

The data  $(t_i, y_i)$  are called **monotonic** when  $y_i \geq y_{i-1}$  or  $y_i \leq y_{i-1}$ ,  $i = 1, \dots, n$ .

#### Definition 5.3.1.3. Convex/concave data

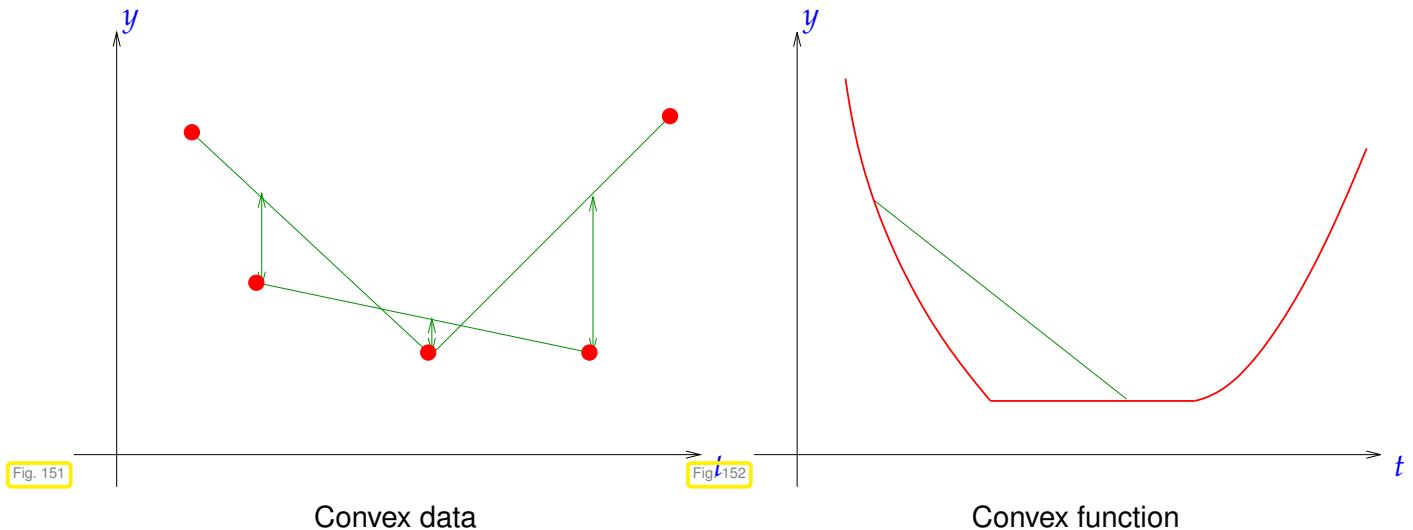
The data  $\{(t_i, y_i)\}_{i=0}^n$  are called **convex** (**concave**) if

$$\Delta_j \stackrel{(\geq)}{\leq} \Delta_{j+1}, \quad j = 1, \dots, n-1, \quad \Delta_j := \frac{y_j - y_{j-1}}{t_j - t_{j-1}}, \quad j = 1, \dots, n.$$

Mathematical characterization of convex data:

$$y_i \leq \frac{(t_{i+1} - t_i)y_{i-1} + (t_i - t_{i-1})y_{i+1}}{t_{i+1} - t_{i-1}} \quad \forall i = 1, \dots, n-1,$$

i.e., each data point lies below the line segment connecting the other data, *cf.* definition of convexity of a function [Str09, Def. 5.5.2].



**Definition 5.3.1.4. Convex/concave function** → [Str09, Def. 5.5.2]

$$f : I \subset \mathbb{R} \mapsto \mathbb{R} \quad \begin{array}{l} \text{convex} \\ \text{concave} \end{array} \quad :\Leftrightarrow \quad \begin{array}{ll} f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) & \forall 0 \leq \lambda \leq 1, \\ f(\lambda x + (1 - \lambda)y) \geq \lambda f(x) + (1 - \lambda)f(y) & \forall x, y \in I. \end{array}$$

**§5.3.1.5 ((Local) shape preservation)** Now we consider interpolation problem to build an interpolant  $f$  with special properties inherited from the given data  $(t_i, y_i), i = 0, \dots, n$ .

Goal: shape preserving interpolation:

positive data	→	positive interpolant $f$ ,
monotonic data	→	monotonic interpolant $f$ ,
convex data	→	convex interpolant $f$ .

More ambitious goal: local shape preserving interpolation: for each subinterval  $I = (t_i, t_{i+j})$

positive data in $I$	→	locally positive interpolant $f _I$ ,
monotonic data in $I$	→	locally monotonic interpolant $f _I$ ,
convex data in $I$	→	locally convex interpolant $f _I$ .

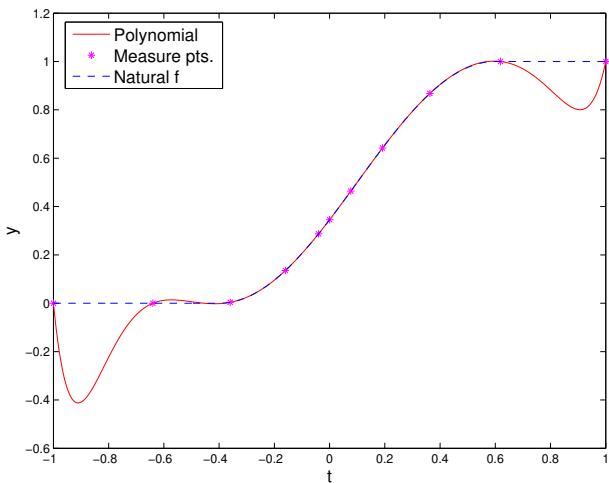
### EXPERIMENT 5.3.1.6 (Bad behavior of global polynomial interpolants)

We perform Lagrange interpolation for the following positive and monotonic data:

$t_i$	-1.0	-0.640	-0.3600	-0.1600	-0.0400	0.0000	0.0770	0.1918	0.3631	0.6187	1.0
$y_i$	0.0	0.000	0.0039	0.1355	0.2871	0.3455	0.4639	0.6422	0.8678	1.0000	1.0

created by taking points on the graph of

$$f(t) = \begin{cases} 0 & \text{if } t < -\frac{2}{5}, \\ \frac{1}{2}(1 + \cos(\pi(t - \frac{3}{5}))) & \text{if } -\frac{2}{5} < t < \frac{3}{5}, \\ 1 & \text{otherwise.} \end{cases}$$



← Interpolating polynomial, degree = 10

We observe oscillations at the endpoints of the interval (see Fig. 148), which means that we encounter

- no locality,
- no positivity,
- no monotonicity,
- no local conservation of the curvature, in the case of global polynomial interpolation.

### 5.3.2 Piecewise Linear Interpolation

There is a very simple method of achieving perfect shape preservation by means of a linear ( $\rightarrow$  § 5.1.0.15) interpolation operator into the space of continuous functions:

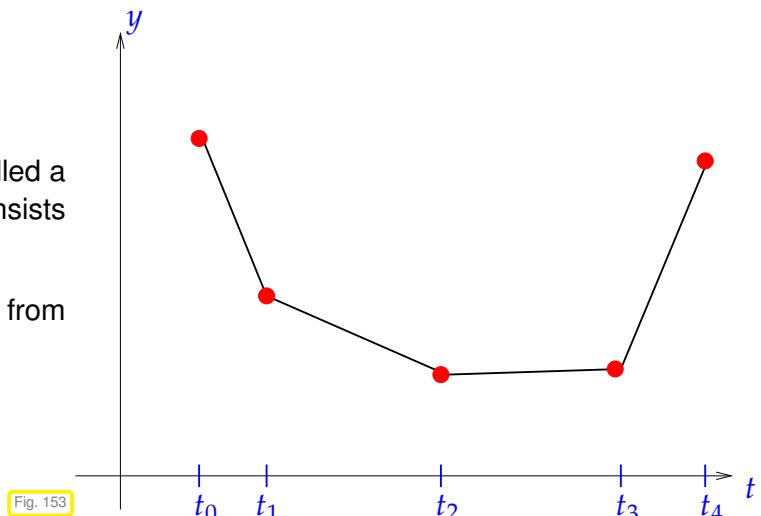
Data:  $(t_i, y_i) \in \mathbb{R}^2, i = 0, \dots, n, n \in \mathbb{N}, t_0 < t_1 < \dots < t_n$ .

Then the piecewise linear interpolant  $s : [t_0, t_n] \rightarrow \mathbb{R}$  is defined as, cf. Ex. 5.1.0.12:

$$s(t) = \frac{(t_{i+1} - t)y_i + (t - t_i)y_{i+1}}{t_{i+1} - t_i} \quad \text{for } t \in [t_i, t_{i+1}] . \quad (5.3.2.1)$$

The piecewise linear interpolant is also called a polygonal curve. It is continuous and consists of  $n$  line segments.

Piecewise linear interpolant of data from Fig. 151 ▷



Piecewise linear interpolation means simply “connect the data points in  $\mathbb{R}^2$  using straight lines”.

Obvious: linear interpolation is **linear** (as mapping  $y \mapsto s$ , see Def. 5.1.0.19) and **local** in the following sense:

$$y_j = \delta_{ij}, \quad i, j = 0, \dots, n \Rightarrow \text{supp}(s) \subset [t_{i-1}, t_{i+1}] . \quad (5.3.2.2)$$

As obvious are the properties asserted in the following theorem. The local preservation of curvature is a straightforward consequence of Def. 5.3.1.3.

**Theorem 5.3.2.3. Local shape preservation by piecewise linear interpolation**

Let  $s \in C([t_0, t_n])$  be the piecewise linear interpolant of  $(t_i, y_i) \in \mathbb{R}^2$ ,  $i = 0, \dots, n$ , for every subinterval  $I = [t_j, t_k] \subset [t_0, t_n]$ :

- |   |   |
|---|---|
| if $(t_i, y_i) _I$ are positive/negative<br>if $(t_i, y_i) _I$ are monotonic (increasing/decreasing)<br>if $(t_i, y_i) _I$ are convex/concave | $\Rightarrow s _I$ is positive/negative,<br>$\Rightarrow s _I$ is monotonic (increasing/decreasing),<br>$\Rightarrow s _I$ is convex/concave. |
|---|---|

Local shape preservation = perfect shape preservation!

Bad news: none of this properties carries over to local polynomial interpolation of higher polynomial degree  $d > 1$ .

**EXAMPLE 5.3.2.4 (Piecewise quadratic interpolation)** We consider the following generalization of piecewise linear interpolation of data points  $(t_j, y_j) \in \mathbb{R} \times \mathbb{R}$ ,  $j = 0, \dots, n$ .

From Thm. 5.2.2.7 we know that a parabola (polynomial of degree 2) is uniquely determined by 3 data points. Thus, the idea is to form groups of three adjacent data points and interpolate each of these triplets by a 2nd-degree polynomial (parabola).

Assume:  $n = 2m$  even

► piecewise quadratic interpolant  $q : [\min\{t_i\}, \max\{t_i\}] \mapsto \mathbb{R}$  is defined by

$$q_j := q|_{[t_{2j-2}, t_{2j}]} \in \mathcal{P}_2, \quad q_j(t_i) = y_i, \quad i = 2j-2, 2j-1, 2j, \quad j = 1, \dots, m. \quad (5.3.2.5)$$

Nodes as in Exp. 5.3.1.6

Piecewise linear (blue) and quadratic (red) interpolants



No shape preservation for piecewise quadratic interpolant

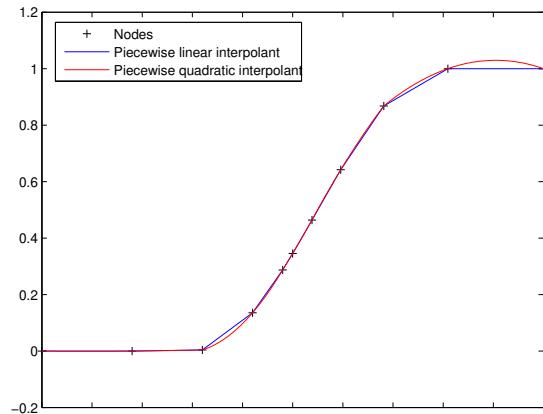


Fig. 154

The “only” drawback of piecewise linear interpolation:

interpolant is only  $C^0$  but not  $C^1$  (no continuous derivative).

However: Interpolant usually serves as input for other numerical methods like a Newton-method for solving non-linear systems of equations, see Section 8.4, which requires derivatives.

## 5.4 Cubic Hermite Interpolation

Aim: construct local shape-preserving ( $\rightarrow$  Section 5.3) (linear ?) interpolation operator that fixes shortcoming of piecewise linear interpolation by ensuring  $C^1$ -smoothness of the interpolant.

notation:  $C^1([a, b]) \hat{=} \text{space of continuously differentiable functions } [a, b] \mapsto \mathbb{R}$ .

### 5.4.1 Definition and Algorithms

Given: mesh points  $(t_i, y_i) \in \mathbb{R}^2, i = 0, \dots, n, t_0 < t_1 < \dots < t_n$ .

Goal: build function  $f \in C^1([t_0, t_n])$  satisfying the interpolation conditions  $f(t_i) = y_i, i = 0, \dots, n$ .

#### Definition 5.4.1.1. Cubic Hermite polynomial interpolant

Given data points  $(t_j, y_j) \in \mathbb{R} \times \mathbb{R}, j = 0, \dots, n$ , with pairwise distinct ordered nodes  $t_j$ , and slopes  $c_j \in \mathbb{R}$ , the piecewise cubic Hermite interpolant  $s : [t_0, t_n] \rightarrow \mathbb{R}$  is defined by the requirements

$$s|_{[t_{i-1}, t_i]} \in \mathcal{P}_3, \quad i = 1, \dots, n, \quad s(t_i) = y_i, \quad s'(t_i) = c_i, \quad i = 0, \dots, n.$$

#### Corollary 5.4.1.2. Smoothness of cubic Hermite polynomial interpolant

Piecewise cubic Hermite interpolants are continuously differentiable on their interval of definition.

*Proof.* The assertion of the corollary follows from the agreement of function values and first derivative values on nodes shared by two intervals, on each of which the piecewise cubic Hermite interpolant is a polynomial of degree 3.  $\square$

**§5.4.1.3 (Local representation of piecewise cubic Hermite interpolant)** Locally, we can write a piecewise cubic Hermite interpolant as a linear combination of generalized cardinal basis functions with coefficients supplied by the data values  $y_j$  and the slopes  $c_j$ :

$$\Rightarrow s(t) = y_{i-1}H_1(t) + y_iH_2(t) + c_{i-1}H_3(t) + c_iH_4(t), \quad t \in [t_{i-1}, t_i], \quad (5.4.1.4)$$

where the basic functions  $H_k, k = 1, 2, 3, 4$ , are as follows:

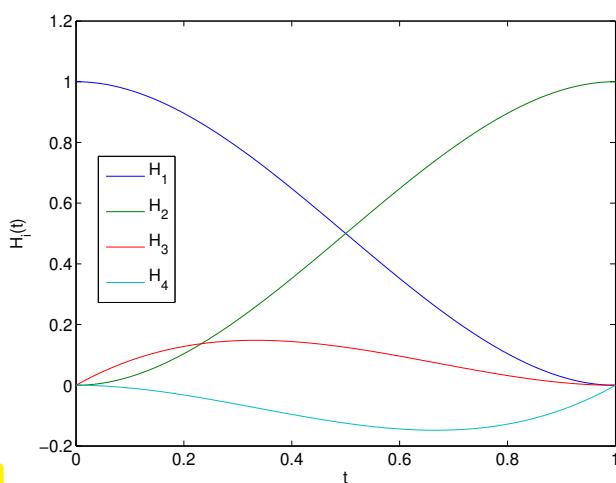


Fig. 155

Local basis polynomial on  $[0, 1]$

By tedious, but straightforward computations using the chain rule we find the following values for  $H_k$  and  $H'_k$  at the endpoints of the interval  $[t_{i-1}, t_i]$ .

$$H_1(t) := \phi\left(\frac{t_i-t}{h_i}\right), \quad (5.4.1.5a)$$

$$H_2(t) := \phi\left(\frac{t-t_{i-1}}{h_i}\right), \quad (5.4.1.5b)$$

$$H_3(t) := -h_i\psi\left(\frac{t_i-t}{h_i}\right), \quad (5.4.1.5c)$$

$$H_4(t) := h_i\psi\left(\frac{t-t_{i-1}}{h_i}\right), \quad (5.4.1.5d)$$

$$h_i := t_i - t_{i-1}, \quad (5.4.1.5e)$$

$$\phi(\tau) := 3\tau^2 - 2\tau^3, \quad (5.4.1.5f)$$

$$\psi(\tau) := \tau^3 - \tau^2. \quad (5.4.1.5g)$$

	$H(t_{i-1})$	$H(t_i)$	$H'(t_{i-1})$	$H'(t_i)$
$H_1$	1	0	0	0
$H_2$	0	1	0	0
$H_3$	0	0	1	0
$H_4$	0	0	0	1

This amounts to a proof for (5.4.1.4) (why?).

The formula (5.4.1.4) is handy for the local evaluation of piecewise cubic Hermite interpolants. The function `hermloceval` in Code 5.4.1.6 performs the efficient evaluation (in multiple points) of a piecewise cubic polynomial  $s$  on  $t_1, t_2$  uniquely defined by the constraints  $s(t_1) = y_1, s(t_2) = y_2, s'(t_1) = c_1, s'(t_2) = c_2$ :

#### C++ code 5.4.1.6: Local evaluation of cubic Hermite polynomial

```

2 // Multiple point evaluation of Hermite polynomial
3 // y1, y2: data values
4 // c1, c2: slopes
5 VectorXd hermloceval(VectorXd t, double t1, double t2,
6                     double y1, double y2,
7                     double c1, double c2) {
8     const double h = t2-t1, a1 = y2-y1, a2 = a1-h*c1, a3 = h*c2-a1-a2;
9     t = ((t.array()-t1)/h).matrix();
10    return (y1+(a1+(a2+a3*t.array())*(t.array()-1))*t.array()).matrix();
11 }
```

**§5.4.1.7 (Linear Hermite interpolation)** However, the data for an interpolation problem ( $\rightarrow$  Section 5.1) are merely the interpolation points  $(t_j, y_j), j = 0, \dots, n$ , but not the slopes of the interpolant at the nodes. Thus, in order to define an interpolation operator into the space of piecewise cubic Hermite functions, we have supply a mapping  $\mathbb{R}^{n+1} \times \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$  computing the slopes  $c_j$  from the data points.

Since this mapping should be local it is natural to rely on (weighted) averages of the local slopes  $\Delta_j$  ( $\rightarrow$  Def. 5.3.1.3) of the data, for instance

$$c_i = \begin{cases} \Delta_1 & , \text{ for } i = 0 , \\ \Delta_n & , \text{ for } i = n , \\ \frac{t_{i+1}-t_i}{t_{i+1}-t_{i-1}}\Delta_i + \frac{t_i-t_{i-1}}{t_{i+1}-t_{i-1}}\Delta_{i+1} & , \text{ if } 1 \leq i < n . \end{cases} \quad \Delta_j := \frac{y_j - y_{j-1}}{t_j - t_{j-1}}, j = 1, \dots, n . \quad (5.4.1.8)$$

► Leads to a **linear** ( $\rightarrow$  Def. 5.1.0.19), **local** Hermite interpolation operator

“Local” means, that, if the values  $y_j$  are non-zero for only a few adjacent data points with indices  $j = k, \dots, k+m$ ,  $m \in \mathbb{N}$  small, then the Hermite interpolant  $s$  is supported on  $[t_{k-\ell}, t_{k+m+\ell}]$  for small  $\ell \in \mathbb{N}$  independent of  $k$  and  $m$ .

#### EXAMPLE 5.4.1.9 (Piecewise cubic Hermite interpolation)

## Hermite interpolation

Data points:

- ◆ 11 equispaced nodes

$$t_j = -1 + 0.2 j, \quad j = 0, \dots, 10.$$

in the interval  $I = [-1, 1]$ ,

- ◆  $y_i = f(t_i)$  with

$$f(x) := \sin(5x) e^x.$$

Here we used weighted averages of slopes as in Eq. (5.4.1.8).

For details see code `hermintp1.hpp`  
[→ GITLAB](#).

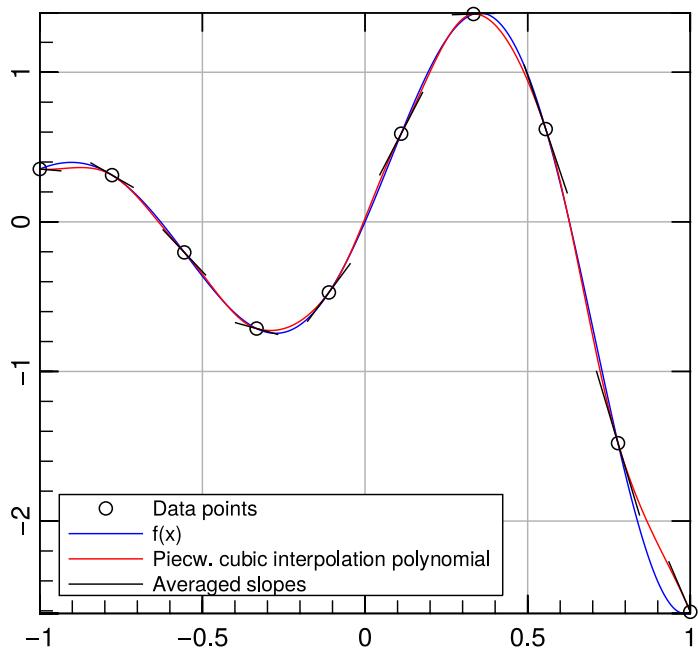


Fig. 156



No local/global preservation of monotonicity!

□

## 5.4.2 Local Monotonicity-Preserving Hermite Interpolation

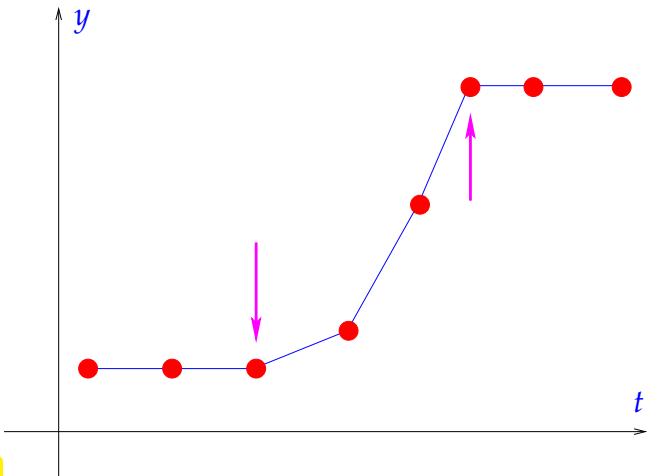
From Ex. 5.4.1.9 we learn that, if the slopes are chosen according to Eq. (5.4.1.8), then the resulting Hermite interpolation does not preserve monotonicity.

Consider the situation sketched on the right  $\Rightarrow$

The red circles ( $\bullet$ ) represent data points, the blue line (—) the piecewise linear interpolant  $\rightarrow$  Section 5.3.2.

In the nodes marked with  $\rightarrow$  the first derivative of a monotonicity preserving  $C^1$ -smooth interpolant must vanish! Otherwise an “overshoot” occurs, see also Fig. 154.

Of course, this will be violated, when a (weighted) arithmetic average is used for the computation of slopes for cubic Hermite interpolation.



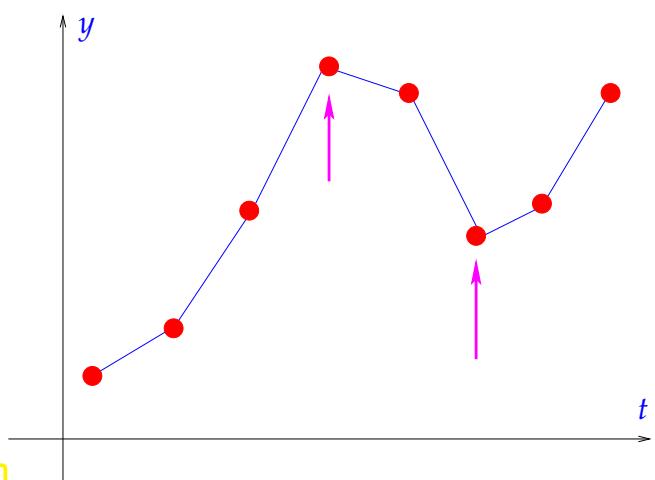
$\triangleleft$  Consider the situation sketched on the left.

The red circles ( $\bullet$ ) represent data points, the blue line (—) the piecewise linear interpolant  $\rightarrow$  Section 5.3.2.

A *local* monotonicity preserving  $C^1$ -smooth interpolant ( $\rightarrow$  § 5.3.1.5)  $s$  must be flat (= vanishing first derivative) in data points  $(t_j, y_j)$ , for which

$$\begin{aligned} y_{j-1} < y_j \quad \text{and} \quad y_{j+1} < y_j, \\ y_{j-1} > y_j \quad \text{and} \quad y_{j+1} > y_j, \end{aligned}$$

in “local extrema” of the data set.



Otherwise, overshoots or undershoots would destroy local monotonicity on one side of the extremum.

**§5.4.2.1 (Limiting of local slopes)** From the discussion of Fig. 157 and Fig. 158 it is clear that local monotonicity preservation entails that the local slopes  $c_i$  of a cubic Hermite interpolant ( $\rightarrow$  Def. 5.4.1.1) have to fulfill

$$c_i = \begin{cases} 0 & , \text{if } \operatorname{sgn}(\Delta_i) \neq \operatorname{sgn}(\Delta_{i+1}) , \\ \text{some “average” of } \Delta_i, \Delta_{i+1} & \text{otherwise} \end{cases}, \quad i = 1, \dots, n-1 . \quad (5.4.2.2)$$

☞ notation: sign function  $\operatorname{sgn}(\xi) = \begin{cases} 1 & , \text{if } \xi > 0 , \\ 0 & , \text{if } \xi = 0 , \\ -1 & , \text{if } \xi < 0 . \end{cases}$

A slope selection rule that enforces (5.4.2.2) is called a **limiter**.

Of course, testing for equality with zero does not make sense for data that may be affected by measurement or roundoff errors. Thus, the “average” in (5.4.2.2) must be close to zero already when either  $\Delta_i \approx 0$  or  $\Delta_{i+1} \approx 0$ . This is satisfied by the **weighted harmonic mean**

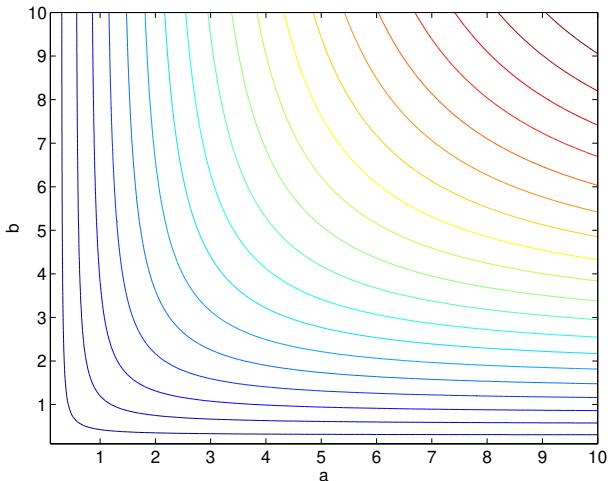
$$c_i = \frac{1}{\frac{w_a}{\Delta_i} + \frac{w_b}{\Delta_{i+1}}} , \quad (5.4.2.3)$$

with weights  $w_a > 0, w_b > 0, (w_a + w_b = 1)$ .

The harmonic mean = “smoothed  $\min(\cdot, \cdot)$ -function”.

Obviously  $\Delta_i \rightarrow 0$  or  $\Delta_{i+1} \rightarrow 0$  in (5.4.2.3), then  $c_i \rightarrow 0$ .

Contour plot of the harmonic mean of  $a$  and  $b \rightarrow (w_a = w_b = 1/2)$ .



A good choice of the weights is:

$$w_a = \frac{2h_{i+1} + h_i}{3(h_{i+1} + h_i)}, \quad w_b = \frac{h_{i+1} + 2h_i}{3(h_{i+1} + h_i)},$$

This yields the following local slopes, unless (5.4.2.2) enforces  $c_i = 0$ :

$$\begin{aligned} \text{sgn}(\Delta_1) = \text{sgn}(\Delta_2) \quad c_i = \begin{cases} \Delta_1 & , \text{ if } i = 0, \\ \frac{3(h_{i+1} + h_i)}{\frac{2h_{i+1} + h_i}{\Delta_i} + \frac{2h_i + h_{i+1}}{\Delta_{i+1}}} & , \text{ for } i \in \{1, \dots, n-1\}, \quad h_i := t_i - t_{i-1}. \\ \Delta_n & , \text{ if } i = n, \end{cases} \quad (5.4.2.4) \end{aligned}$$

Piecewise cubic Hermite interpolation with local slopes chosen according to (5.4.2.2) and (5.4.2.4) is available through the MATLAB function/PYTHON class  $v = \text{pchip}(t, y, x); / \text{scipy.interpolate.PchipInterpolator}$ . The argument  $t$  passes the interpolation nodes,  $y$  the corresponding data values, and  $x$  is a vector of evaluation points.  $\square$

### EXAMPLE 5.4.2.5 (Monotonicity preserving piecewise cubic polynomial interpolation)

Data from Exp. 5.3.1.6

Plot created with MATLAB-function call

$v = \text{pchip}(t, y, x);$

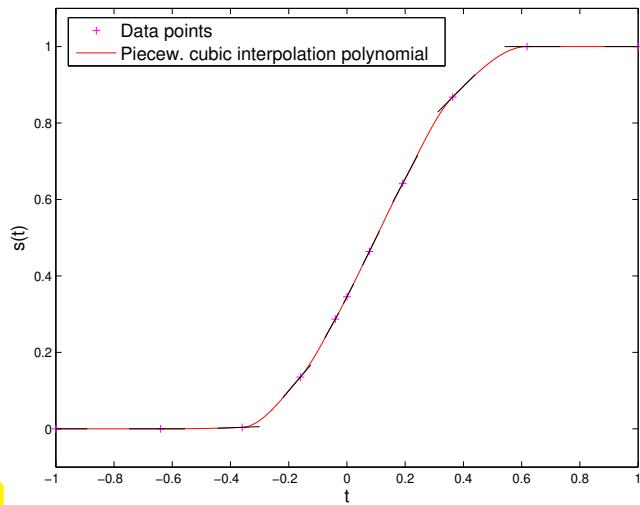
$t$ : Data nodes  $t_j$

$y$ : Data values  $y_j$

$x$ : Evaluation points  $x_i$

$v$ : Vector  $s(x_i)$

We observe perfect local monotonicity preservation, no under- or overshoots at extrema.  $\triangleright$



**Remark 5.4.2.6 (Non-linear cubic Hermite interpolation)** Note that the mapping  $\mathbf{y} := [y_0, \dots, y_n] \rightarrow c_i$  defined by (5.4.2.2) and (5.4.2.4) is **not linear**.

- The “pchip interpolator operator” does **not** provide a **linear** mapping from data space  $\mathbb{R}^{n+1}$  into  $C^1([t_0, t_n])$  (in the sense of Def. 5.1.0.19).

In fact, the non-linearity of the piecewise cubic Hermite interpolation operator is necessary for (only global) monotonicity preservation:

### Theorem 5.4.2.7. Property of linear, monotonicity preserving interpolation into $C^1$

If, for fixed node set  $\{t_j\}_{j=0}^n$ ,  $n \geq 2$ , an interpolation scheme  $I : \mathbb{R}^{n+1} \rightarrow C^1(I)$  is linear as a mapping from data values to continuous functions on the interval covered by the nodes ( $\rightarrow$  Def. 5.1.0.19), and monotonicity preserving, then  $I(\mathbf{y})'(t_j) = 0$  for all  $\mathbf{y} \in \mathbb{R}^{n+1}$  and  $j = 1, \dots, n - 1$ .

Of course, an interpolant that is flat in all data points, as stipulated by Thm. 5.4.2.7 for a linear, monotonicity preserving,  $C^1$ -smooth interpolation scheme, does not make much sense.

At least, the piecewise cubic Hermite interpolation operator is local (in the sense discussed in § 5.4.1.7).  $\square$

### Theorem 5.4.2.8. Monotonicity preservation of limited cubic Hermite interpolation

The cubic Hermite interpolation polynomial with slopes as in Eq. (5.4.2.4) provides a local monotonicity-preserving  $C^1$ -interpolant.

*Proof.* See F. FRITSCH UND R. CARLSON, *Monotone piecewise cubic interpolation*, SIAM J. Numer. Anal., 17 (1980), S. 238–246.  $\square$

The next code demonstrates the calculation of the slopes  $c_i$  in MATLAB's pchip (details in [FC80]):

### C++ code 5.4.2.9: Monotonicity preserving slopes in pchip

```

1 # include <Eigen/Dense>
2
3 using Eigen::VectorXd;
4
5 // using forward declaration of the function pchipend, implementation
6 // below
6 double pchipend(const double, const double, const double, const double);
7
8 void pchipslopes(const VectorXd& t, const VectorXd& y, VectorXd& c) {
9     // Calculation of local slopes  $c_i$  for shape preserving cubic Hermite
10    // interpolation, see (5.4.2.2), (5.4.2.4)
11    //  $t$ ,  $y$  are vectors passing the data points
12    const unsigned n = t.size();
13    const VectorXd h = t.tail(n - 1) - t.head(n - 1),
14        delta = (y.tail(n - 1) - y.head(n - 1)).cwiseQuotient(h); // linear
15        slopes
16    c = VectorXd::Zero(n);
17
18    // compute reconstruction slope according to (5.4.2.4)
19    for (unsigned i = 0; i < n - 2; ++i) {
20        if (delta(i)*delta(i + 1) > 0) {
21            const double w1 = 2*h(i + 1) + h(i),
22                w2 = h(i + 1) + 2*h(i);
23            c(i + 1) = (w1 + w2)/(w1/delta(i) + w2/delta(i + 1));
24        }
25    }

```

```

24 // Special slopes at endpoints, beyond (5.4.2.4)
25 c(0) = pchipend(h(0), h(1), delta(0), delta(1));
26 c(n - 1) = pchipend(h(n - 2), h(n - 3), delta(n - 2), delta(n - 3));
27 }
28
29 double pchipend(const double h1, const double h2, const double del1, const double
   del2) {
30 // Non-centered, shape-preserving, three-point formula
31 double d = ((2*h1 + h2)*del1 - h1*del2)/(h1 + h2);
32 if (d*del1 < 0) {
33     d = 0;
34 }
35 else if (del1*del2 < 0 && std::abs(d) > std::abs(3*del1)) {
36     d = 3*del1;
37 }
38 return d;
39 }
```

## 5.5 Splines



*Supplementary literature.* Splines are also presented in [DR08, Ch. 9].

Piecewise cubic Hermite Interpolation presented in Section 5.4 entailed determining reconstruction slopes  $c_i$ . Now we learn about a way how to do piecewise polynomial interpolation, which results in  $C^k$ -interpolants,  $k > 0$ , and dispenses with auxiliary slopes. The idea is to obtain the missing conditions implicitly from extra continuity conditions.

### Definition 5.5.0.1. Spline space → [QSS00, Def. 8.1]

Given an interval  $I := [a, b] \subset \mathbb{R}$  and a knot set/mesh  $\mathcal{M} := \{a = t_0 < t_1 < \dots < t_{n-1} < t_n = b\}$ , the vector space  $\mathcal{S}_{d,\mathcal{M}}$  of the **spline functions** of degree  $d$  (or order  $d+1$ ) is defined by

$$\mathcal{S}_{d,\mathcal{M}} := \{s \in C^{d-1}(I) : s_j := s|_{[t_{j-1}, t_j]} \in \mathcal{P}_d \quad \begin{matrix} \text{$d-1$-times continuously differentiable} \\ \text{locally polynomial of degree $d$} \end{matrix} \forall j = 1, \dots, n\}.$$

Obviously, spline spaces are mapped onto each other by differentiation & integration:

$$s \in \mathcal{S}_{d,\mathcal{M}} \Rightarrow s' \in \mathcal{S}_{d-1,\mathcal{M}} \wedge \int_a^t s(\tau) d\tau \in \mathcal{S}_{d+1,\mathcal{M}}.$$

Spline spaces of the lowest degrees:

- $d = 0$ :  $\mathcal{M}$ -piecewise constant *discontinuous* functions
- $d = 1$ :  $\mathcal{M}$ -piecewise linear *continuous* functions
- $d = 2$ : *continuously differentiable*  $\mathcal{M}$ -piecewise quadratic functions

The dimension of spline space can be found by a **counting argument** (heuristic): We count the number of “degrees of freedom” (d.o.f.s) possessed by a  $\mathcal{M}$ -piecewise polynomial of degree  $d$ , and subtract the

number of **linear constraints** implicit contained in Def. 5.5.0.1:

$$\dim \mathcal{S}_{d,\mathcal{M}} = n \cdot \dim \mathcal{P}_d - \#\{C^{d-1} \text{ continuity constraints}\} = n \cdot (d+1) - (n-1) \cdot d = n+d .$$

### Theorem 5.5.0.2. Dimension of spline space

*The space  $\mathcal{S}_{d,\mathcal{M}}$  from Def. 5.5.0.1 has dimension*

$$\dim \mathcal{S}_{d,\mathcal{M}} = n + d .$$

We already know the special case of interpolation in  $\mathcal{S}_{1,\mathcal{M}}$ , when the interpolation nodes are the knots of  $\mathcal{M}$ , because this boils down to simple piecewise linear interpolation, see Section 5.3.2.

## 5.5.1 Cubic Spline Interpolation



*Supplementary literature.* More details in [Han02, pp. XIII, 46], [QSS00, Sect. 8.6.1].

Cognitive psychology teaches us that the human eye perceives a  $C^2$ -functions as “smooth”, while it can still spot the abrupt change of curvature at the possible discontinuities of the second derivatives of a cubic Hermite interpolant (→ Def. 5.4.1.1).

For this reason the simplest spline functions featuring  $C^2$ -smoothness are of great importance in computer aided design (CAD). They are the **cubic splines**,  $\mathcal{M}$ -piecewise polynomials of degree 3 contained in  $\mathcal{S}_{3,\mathcal{M}}$  (→ Def. 5.5.0.1).

**§5.5.1.1 (Definition of cubic spline interpolant)** In this section we study **cubic spline interpolation** (related to cubic Hermite interpolation, Section 5.4)

Task: Given a mesh  $\mathcal{M} := \{t_0 < t_1 < \dots < t_n\}$ ,  $n \in \mathbb{N}$ , “find” a cubic spline  $s \in \mathcal{S}_{3,\mathcal{M}}$  that complies with the interpolation conditions

$$s(t_j) = y_j \quad , \quad j = 0, \dots, n . \tag{5.5.1.2}$$

$\hat{=}$  interpolation at knots !

From **dimensional considerations** it is clear that the interpolation conditions will fail to fix the interpolating cubic spline uniquely:

$$\dim \mathcal{S}_{3,\mathcal{M}} - \#\{\text{interpolation conditions}\} = (n+3) - (n+1) = 2 \text{ free d.o.f.}$$

“two conditions are missing”  $\Rightarrow$  interpolation problem is not yet well defined!

**§5.5.1.3 (Computing cubic spline interpolants)** We opt for a linear interpolation scheme (→ Def. 5.1.0.19) into the spline space  $\mathcal{S}_{3,\mathcal{M}}$ . As explained in § 5.1.0.15, this will lead to an equivalent linear system of equations for expansion coefficients with respect to a suitable basis.

We reuse the local representation of a cubic spline through cubic Hermite cardinal basis polynomials from

(5.4.1.5):

$$\begin{array}{l} \xrightarrow{(5.4.1.4)} \\ s|_{[t_{j-1}, t_j]}(t) = s(t_{j-1}) \cdot (1 - 3\tau^2 + 2\tau^3) + \\ s(t_j) \cdot (3\tau^2 - 2\tau^3) + \\ h_j s'(t_{j-1}) \cdot (\tau - 2\tau^2 + \tau^3) + \\ h_j s'(t_j) \cdot (-\tau^2 + \tau^3), \end{array} \quad (5.5.1.4)$$

with  $h_j := t_j - t_{j-1}$ ,  $\tau := (t - t_{j-1})/h_j$ .

➤ Task of cubic spline interpolation boils down to finding slopes  $s'(t_j)$  in the knots of the mesh  $\mathcal{M}$ .

Once these slopes are known, the efficient local evaluation of a cubic spline function can be done as for a cubic Hermite interpolant, see Section 5.4.1, Code 5.4.1.6.

Note: if  $s(t_j)$ ,  $s'(t_j)$ ,  $j = 0, \dots, n$ , are fixed, then the representation Eq. (5.5.1.4) already guarantees  $s \in C^1([t_0, t_n])$ , cf. the discussion for cubic Hermite interpolation, Section 5.4.

➤ only continuity of  $s''$       ♦ has to be enforced by choice of  $s'(t_j)$   
 $\Updownarrow$   
 ♦ will yield extra conditions to fix the  $s'(t_j)$

However, do the

- ♦ interpolation conditions Eq. (5.5.1.2)  $s(t_j) = y_j$ ,  $j = 0, \dots, n$ , and the
- ♦ smoothness constraint  $s \in C^2([t_0, t_n])$

uniquely determine the unknown slopes  $c_j := s'(t_j)$  ?

From  $s \in C^2([t_0, t_n])$  we obtain  $n - 1$  continuity constraints for  $s''(t)$  at the internal nodes

$$s''|_{[t_{j-1}, t_j]}(t_j) = s''|_{[t_j, t_{j+1}]}(t_j), \quad j = 1, \dots, n - 1. \quad (5.5.1.5)$$

Based on Eq. (5.5.1.4), we express Eq. (5.5.1.5) in concrete terms, using

$$\begin{aligned} s''|_{[t_{j-1}, t_j]}(t) &= s(t_{j-1})h_j^{-2}6(-1 + 2\tau) + s(t_j)h_j^{-2}6(1 - 2\tau) \\ &\quad + h_j^{-1}s'(t_{j-1})(-4 + 6\tau) + h_j^{-1}s'(t_j)(-2 + 6\tau), \quad \tau := (t - t_{j-1})/h_j, \end{aligned} \quad (5.5.1.6)$$

which can be obtained by the chain rule and from  $\frac{d\tau}{dt} = h_j^{-1}$ .

$$\begin{aligned} \xrightarrow{\text{Eq. (5.5.1.6)}} \quad s''|_{[t_{j-1}, t_j]}(t_{j-1}) &= -6 \cdot s(t_{j-1})h_j^{-2} + 6 \cdot s(t_j)h_j^{-2} - 4 \cdot h_j^{-1}s'(t_{j-1}) - 2 \cdot h_j^{-1}s'(t_j), \\ s''|_{[t_{j-1}, t_j]}(t_j) &= 6 \cdot s(t_{j-1})h_j^{-2} + -6 \cdot s(t_j)h_j^{-2} + 2 \cdot h_j^{-1}s'(t_{j-1}) + 4 \cdot h_j^{-1}s'(t_j). \end{aligned}$$

Eq. (5.5.1.5) →  $n - 1$  linear equations for  $n$  slopes  $c_j := s'(t_j)$ ; with  $s(t_i) = y_i$ ,

$$\frac{1}{h_j}c_{j-1} + \left(\frac{2}{h_j} + \frac{2}{h_{j+1}}\right)c_j + \frac{1}{h_{j+1}}c_{j+1} = 3\left(\frac{y_j - y_{j-1}}{h_j^2} + \frac{y_{j+1} - y_j}{h_{j+1}^2}\right), \quad (5.5.1.7)$$

for  $j = 1, \dots, n - 1$ .

Eq. (5.5.1.7)  $\Leftrightarrow$  underdetermined  $(n - 1) \times (n + 1)$  linear system of equations.

$n - 1$ : no. of interpolation conditions

$n + 1$ : dimension of cubic spline space on knot set  $\{t_0 < t_1 < \dots < t_n\}$

$$\begin{bmatrix} b_0 & a_1 & b_1 & 0 & \cdots & \cdots & 0 \\ 0 & b_1 & a_2 & b_2 & & & \\ 0 & \ddots & \ddots & \ddots & & \vdots & \\ \vdots & & \ddots & \ddots & \ddots & & \\ & & & \ddots & a_{n-2} & b_{n-2} & 0 \\ 0 & \cdots & \cdots & 0 & b_{n-2} & a_{n-1} & b_{n-1} \end{bmatrix} \begin{bmatrix} c_0 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} 3\left(\frac{y_1-y_0}{h_1^2} + \frac{y_2-y_1}{h_2^2}\right) \\ \vdots \\ 3\left(\frac{y_{n-1}-y_{n-2}}{h_{n-1}^2} + \frac{y_n-y_{n-1}}{h_n^2}\right) \end{bmatrix}. \quad (5.5.1.8)$$

with

$$b_i := \frac{1}{h_{i+1}}, \quad a_i := \frac{2}{h_i} + \frac{2}{h_{i+1}}, \quad i = 0, 1, \dots, n - 1, \quad [b_i, a_i > 0, \quad a_i = 2(b_i + b_{i-1})];.$$

→ two additional constraints are required, as already noted in § 5.5.1.1.  $\square$

**§5.5.1.9 (Types of cubic spline interpolants)** To saturate the remaining two degrees of freedom the following three approaches are popular:

① **Complete cubic spline interpolation:**  $s'(t_0) = c_0, s'(t_n) = c_n$  prescribed.

Then the first and last column can be removed from the system matrix of (5.5.1.8). Their products with  $c_0$  and  $c_n$ , respectively, have to be subtracted from the right hand side of (5.5.1.8).

② **Natural cubic spline interpolation:**  $s''(t_0) = s''(t_n) = 0$

$$\frac{2}{h_1}c_0 + \frac{1}{h_1}c_1 = 3\frac{y_1 - y_0}{h_1^2}, \quad \frac{1}{h_n}c_{n-1} + \frac{2}{h_n}c_n = 3\frac{y_n - y_{n-1}}{h_n^2}.$$

Combining these two extra equations with (5.5.1.8), we arrive at a linear system of equations with tridiagonal s.p.d. ( $\rightarrow$  Def. 1.1.2.6, Lemma 2.8.0.12) system matrix and unknowns  $c_0, \dots, c_n$ . Due to Thm. 2.7.5.4 it can be solved with an asymptotic computational effort of  $O(n)$ .

③ **Periodic cubic spline interpolation:**  $s'(t_0) = s'(t_n)$  ( $\Rightarrow c_0 = c_n$ ),  $s''(t_0) = s''(t_n)$

This removes one unknown and adds another equations so that we end up with an  $n \times n$ -linear system with s.p.d. ( $\rightarrow$  Def. 1.1.2.6) system matrix

$$\mathbf{A} := \begin{bmatrix} a_1 & b_1 & 0 & \cdots & 0 & b_0 \\ b_1 & a_2 & b_2 & & & 0 \\ 0 & \ddots & \ddots & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & 0 \\ 0 & & & \ddots & a_{n-1} & b_{n-1} \\ b_0 & 0 & \cdots & 0 & b_{n-1} & a_0 \end{bmatrix}, \quad b_i := \frac{1}{h_{i+1}}, \quad i = 0, 1, \dots, n - 1, \\ a_i := \frac{2}{h_i} + \frac{2}{h_{i+1}}, \quad i = 0, 1, \dots, n - 1.$$

This linear system can be solved with rank-1-modifications techniques (see § 2.6.0.12, Lemma 2.6.0.21) + tridiagonal elimination: asymptotic computational effort  $O(n)$ .  $\square$

### Remark 5.5.1.10 (Piecewise cubic interpolation schemes)

- ◆ Piecewise cubic local Lagrange interpolation
  - Extra degrees of freedom fixed by putting four nodes in one interval
  - ➡ yields merely  $C^0$ -interpolant; perfectly local.
- ◆ Cubic Hermite interpolation
  - Extra degrees of freedom fixed by reconstruction slopes
  - ➡ yields  $C^1$ -interpolant; still local.
- ◆ Cubic spline interpolation
  - Extra degrees of freedom fixed by  $C^2$ -smoothness, complete/natural/periodic constraint.
  - ➡ yields  $C^2$ -interpolant; non-local.

□

## 5.5.2 Structural Properties of Cubic Spline Interpolants

**§5.5.2.1 (Extremal properties of natural cubic spline interpolants → [QSS00, Sect. 8.6.1, Property 8.2])** For a function  $f : [a, b] \mapsto \mathbb{R}$ ,  $f \in C^2([a, b])$ , the term

$$E_{bd}(f) := \frac{1}{2} \int_a^b |f''(t)|^2 dt ,$$

models the elastic **bending energy** of a rod, whose shape is described by the graph of  $f$  (Soundness check: zero bending energy for straight rod). We will show that cubic spline interpolants have minimal bending energy among all  $C^2$ -smooth interpolating functions.

Given: mesh  $\mathcal{M} := \{a = t_0 < t_1 < \dots < t_n = b\}$  of  $[a, b]$  with knots  $t_j$

Set  $s \in \mathcal{S}_{3,\mathcal{M}} :=$  natural cubic spline interpolant of data points  $(t_i, y_i) \in \mathbb{R}^2$ ,  $i = 0, \dots, n$ .

### Theorem 5.5.2.2. Optimality of natural cubic spline interpolant

The natural cubic spline interpolant  $s$  minimizes the elastic curvature energy among all interpolating functions in  $C^2([a, b])$ , that is

$$E_{bd}(s) \leq E_{bd}(f) \quad \forall f \in C^2([a, b]), f(t_i) = y_i, i = 0, \dots, n .$$

Idea of proof: **variational calculus**

We show that any small perturbation of  $s$  such that the perturbed spline still satisfies the interpolation conditions leads to an increase in elastic energy.

Pick *perturbation direction*  $k \in C^2([t_0, t_n])$  satisfying  $k(t_i) = 0, i = 0, \dots, n$ :

$$\begin{aligned} E_{\text{bd}}(s + k) &= \frac{1}{2} \int_a^b |s'' + \lambda k''|^2 dt \\ &= E_{\text{bd}}(s) + \underbrace{\int_a^b s''(t)k''(t) dt}_{:=I} + \underbrace{\frac{1}{2} \int_a^b |k''|^2 dt}_{\geq 0}. \end{aligned} \quad (5.5.2.3)$$

Scrutiny of  $I$ : split in interval contributions, integrate by parts twice, and use  $s^{(4)} \equiv 0$ :

$$\begin{aligned} I &= \sum_{j=1}^n \int_{t_{j-1}}^{t_j} s''(t)k''(t) dt \\ &= - \sum_{j=1}^n \left( s'''(t_j^-) \underbrace{k(t_j)}_{=0} - s'''(t_{j-1}^+) \underbrace{k(t_{j-1})}_{=0} \right) + \underbrace{s''(t_n)}_{=0} k'(t_n) - \underbrace{s''(t_0)}_{=0} k'(t_0) = 0. \end{aligned}$$

In light of (5.5.2.3): non perturbation compatible with interpolation conditions can make the bending energy of  $s$  decrease!  $\square$

#### Remark 5.5.2.4 (Origin of the term “Spline”)

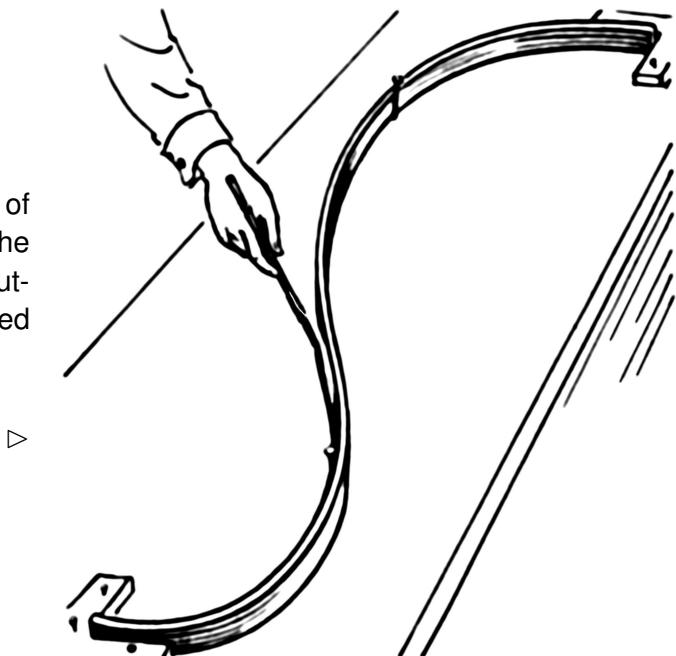
§ 5.5.2.1: (Natural) cubic spline interpolant provides  $C^2$ -curve of minimal elastic bending energy that travels through prescribed points.

$\Updownarrow$

Nature: A thin elastic rod fixed at certain points attains a shape that minimizes its potential bending energy (virtual work principle of statics).

- Cubic spline interpolation approximates shape of elastic rods. Such rods were in fact used in the manufacturing of ship hulls as “analog computers” for “interpolating points” that were specified by the designer of the ship.

Cubic spline interpolation  
before the age of computers



#### Remark 5.5.2.5 (Shape preservation of cubic spline interpolation)

Data  $s(t_j) = y_j$  from Exp. 5.3.1.6 and

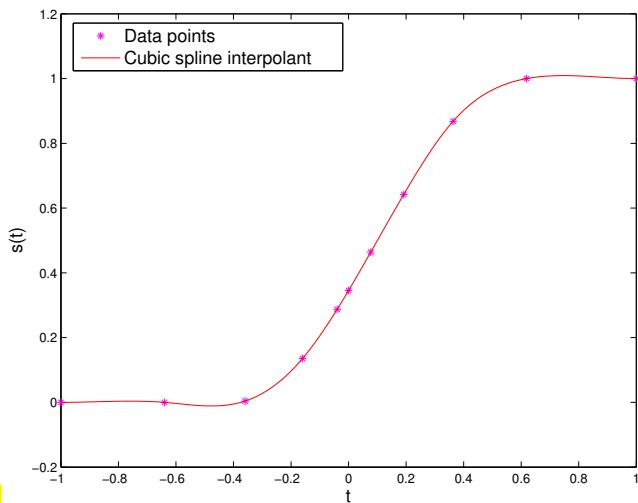
$$c_0 := \frac{y_1 - y_0}{t_1 - t_0},$$

$$c_n := \frac{y_n - y_{n-1}}{t_n - t_{n-1}}.$$

The cubic spline interpolant is **nor** monotonicity **nor** curvature preserving

This is not surprising in light of Thm. 5.4.2.7, because cubic spline interpolation is a linear interpolation scheme.

Fig. 161



**§5.5.2.6 (Weak locality of an interpolation scheme)** In terms of locality of interpolation schemes, in the sense of § 5.4.1.7, we have seen:

- Piecewise linear interpolation (→ Section 5.3.2) is strictly local: Changing a single data value  $y_j$  affects the interpolant only on the interval  $[t_{j-1}, t_{j+1}]$ .
- Monotonicity preserving piecewise cubic Hermite interpolation (→ Section 5.4.2) is still local, because changing  $y_j$  will lead to a change in the interpolant only in  $[t_{j-2}, t_{j+2}]$  (the remote intervals are affected through the averaging of local slopes).
- Polynomial Lagrange interpolation is highly non-local, see Ex. 5.2.4.1.

We can weaken the notion of **locality** of an interpolation scheme on an ordered node set  $\{t_i\}_{i=0}^n$ :

- (weak) locality measures the impact of a perturbation of a data value  $y_j$  at points  $t \in [t_0, t_n]$  as a function of  $|t - t_j|$ .
- an interpolation scheme is **weakly local**, if the impact of the perturbation of  $y_i$  displays a rapid (e.g. exponential) decay as  $|t - t_i|$  increases.

For a **linear** interpolation scheme (→ § 5.1.0.15) locality can be deduced from the decay of the **cardinal interpolants/cardinal basis functions** (→ Lagrange polynomials of § 5.2.2.3), that is, the functions  $b_j := I(\mathbf{e}_j)$ , where  $\mathbf{e}_j$  is the  $j$ -th unit vector, and  $I$  the interpolation operator. Then weak locality can be quantified as

$$\exists \lambda > 0: |b_j(t)| \leq \exp(-\lambda|t - t_j|), \quad t \in [t_0, t_n]. \quad (5.5.2.7)$$

Remember:

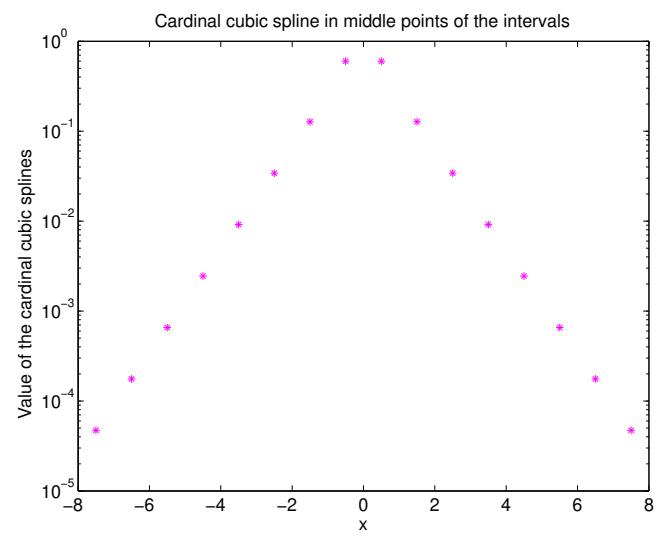
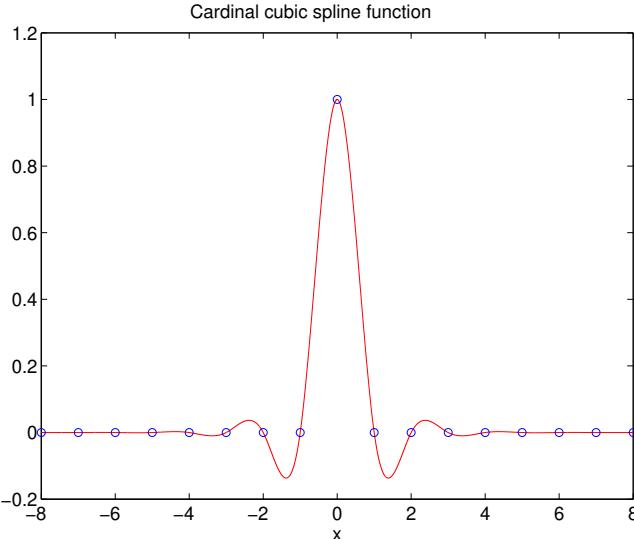
- Lagrange polynomials satisfying (5.2.2.4) provide cardinal interpolants for polynomial interpolation → § 5.2.2.3. As is clear from Fig. 145, they do not display any decay away from their “base node”. Rather, they grow strongly. Hence, there is no locality in global polynomial interpolation.
- Tent functions (→ Fig. 144) are the cardinal basis functions for piecewise linear interpolation, see Ex. 5.1.0.12. Hence, this scheme is perfectly local, see (5.3.2.2).

**EXPERIMENT 5.5.2.8 (Weak locality of the natural cubic spline interpolation)** Given a grid  $\mathcal{M} := \{t_0 < t_1 < \dots < t_n\}$  the  $i$ th natural **cardinal spline** is defined as

$$L_i \in \mathcal{S}_{3,\mathcal{M}}, \quad L_i(t_j) = \delta_{ij}, \quad L_i''(t_0) = L_i''(t_n) = 0. \quad (5.5.2.9)$$

► Natural spline interpolant:  $s(t) = \sum_{j=0}^n y_j L_j(t)$ .

Decay of  $L_i$   $\leftrightarrow$  weak locality of natural cubic spline interpolation.



Exponential decay of the cardinal splines  $\rightarrow$  cubic spline interpolation is weakly local. □

### 5.5.3 Shape Preserving Spline Interpolation

According to Rem. 5.5.2.5 is cubic spline interpolation neither monotonicity preserving nor curvature preserving. Necessarily so, because it is a *linear* interpolation scheme, see Thm. 5.4.2.7.

This section presents a non-linear *quadratic spline* ( $\rightarrow$  Def. 5.5.0.1,  $C^1$ -functions) based interpolation scheme that manages to preserve both monotonicity and curvature of data even in a local sense, cf. Section 5.3.

Given: data points  $(t_i, y_i) \in \mathbb{R}^2, i = 0, \dots, n$ , assume ordering  $t_0 < t_1 < \dots < t_n$ .

Sought:

- ◆ extended knot set  $\mathcal{M} \subset [t_0, t_n]$  ( $\rightarrow$  Def. 5.5.0.1),
- ◆ an interpolating quadratic spline function  $s \in \mathcal{S}_{2,\mathcal{M}}, s(t_i) = y_i, i = 0, \dots, n$  that preserves the “shape” of the data in the sense of § 5.3.1.1.

Notice that here  $\mathcal{M} \neq \{t_j\}_{j=0}^n$ :  $s$  interpolates the data in the points  $t_i$  but is piecewise polynomial with respect to  $\mathcal{M}$ ! The interpolation nodes will usually not belong to  $\mathcal{M}$ .

We proceed in four steps:

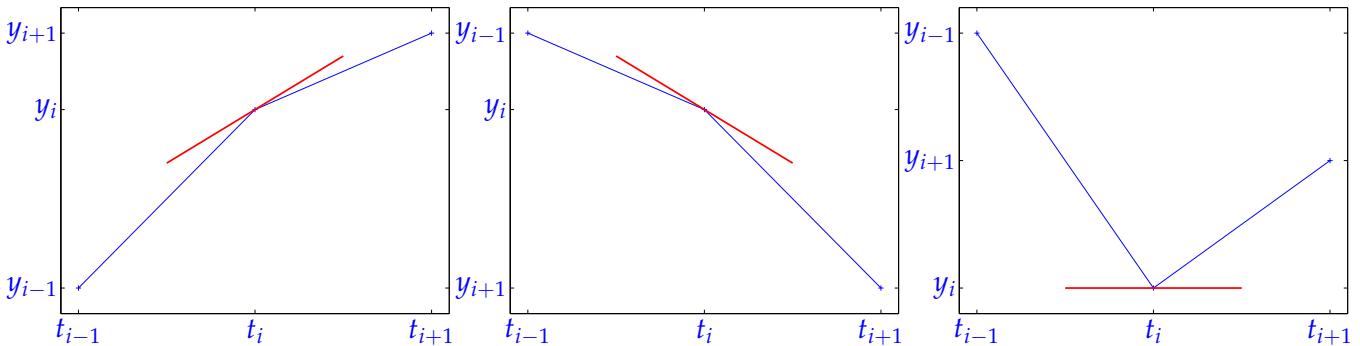
- ① Shape preserving choice of slopes  $c_i, i = 0, \dots, n$  [MR81; QQY01], analogous to Section 5.4.2

Recall Eq. (5.4.2.2) and Eq. (5.4.2.4): we fix the slopes  $c_i$  in the nodes using the harmonic mean of data slopes  $\Delta_j$ , the final interpolant will be tangent to these segments in the points  $(t_i, y_i)$ . If  $(t_i, y_i)$  is a local maximum or minimum of the data,  $c_i$  is set to zero ( $\rightarrow$  § 5.4.2.1)

Limiter  $c_i := \begin{cases} \frac{2}{\Delta_i^{-1} + \Delta_{i+1}^{-1}} & , \text{ if } \text{sign}(\Delta_i) = \text{sign}(\Delta_{i+1}), \\ 0 & \text{otherwise,} \end{cases} \quad i = 1, \dots, n-1 .$

$$c_0 := 2\Delta_1 - c_1, \quad c_n := 2\Delta_n - c_{n-1},$$

where  $\Delta_j = \frac{y_j - y_{j-1}}{t_j - t_{j-1}}$ .

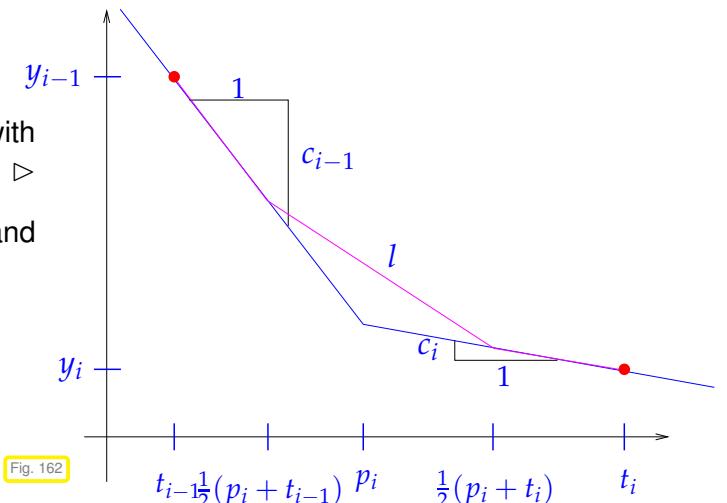


Figures: slopes according to limited harmonic mean formula

- ② Choice of “extra knots”  $p_i \in ]t_{i-1}, t_i]$ ,  $i = 1, \dots, n$ :

Rule:

Let  $T_i$  be the unique straight line through  $(t_i, y_i)$  with slope  $c_i$ ; — in figure



If the intersection of  $T_{i-1}$  and  $T_i$  is non-empty and has a  $t$ -coordinate in  $]t_{i-1}, t_i]$ ,

- ☞ then  $p_i := t$ -coordinate of  $T_{i-1} \cap T_i$ ,
- ☞ otherwise  $p_i = \frac{1}{2}(t_{i-1} + t_i)$ .

These points will be used to build the knot set for the final quadratic spline:

$$\mathcal{M} = \{t_0 < p_1 \leq t_1 < p_2 \leq \dots < p_n \leq t_n\}.$$

#### C++-code 5.5.3.1: Quadratic spline: selection of $p_i$

```

1 # include <Eigen/Dense>
2 using Eigen::VectorXd;
3
4 // choice of extra knots  $p_i \in ]t_{i-1}, t_i[$ 
5 // for given nodeset  $t$  and coefficients  $c$ 
6 void extra_knots(const VectorXd& t, const VectorXd& c, VectorXd& p) {
7     const unsigned n = t.size();
8     p = (t(n - 1) - 1)*VectorXd::Ones(n - 1);
9     for (unsigned j = 0; j < n - 1; ++j) {
10         if (c(j) != c(j + 1)) {
11             p(j) = (y(j+1) - y(j) + t(j)*c(j) - t(j+1)*c(j+1)) / (c(j+1) - c(j));
12         }
13         if (p(j) < t(j) || p(j) > t(j+1)) {
14             p(j) = 0.5*(t(j) + t(j+1));
15         }
16     }
17 }
```

- ③ Set  $\underline{l}$  = linear spline (polygon) on the knot set  $\mathcal{M}'$  (middle points of  $\mathcal{M}$ )

$$\mathcal{M}' = \{t_0 < \frac{1}{2}(t_0 + p_1) < \frac{1}{2}(p_1 + t_1) < \frac{1}{2}(t_1 + p_2) < \dots < \frac{1}{2}(t_{n-1} + p_n) < \frac{1}{2}(p_n + t_n) < t_n\},$$

with  $\underline{l}(t_i) = y_i$ ,  $\underline{l}'(t_i) = c_i$ .

- In each interval  $(\frac{1}{2}(p_j + t_j), \frac{1}{2}(t_j + p_{j+1}))$  the spline corresponds to the segment of slope  $c_j$  passing through the data point  $(t_j, y_j)$ .
- In each interval  $(\frac{1}{2}(t_j + p_{j+1}), \frac{1}{2}(p_{j+1} + t_{j+1}))$  the spline corresponds to the segment connecting the previous ones, see Fig. 162.

$\underline{l}$  "inherits" local monotonicity and curvature from the data.

#### EXAMPLE 5.5.3.2 (Auxiliary construction for shape preserving quadratic spline interpolation)

Data points:  $t = (0 : 12); y = \cos(t);$

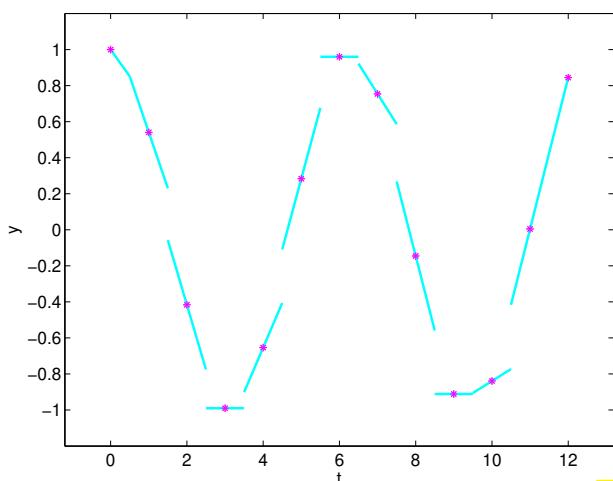


Fig. 163

Local slopes  $c_i, i = 0, \dots, n$

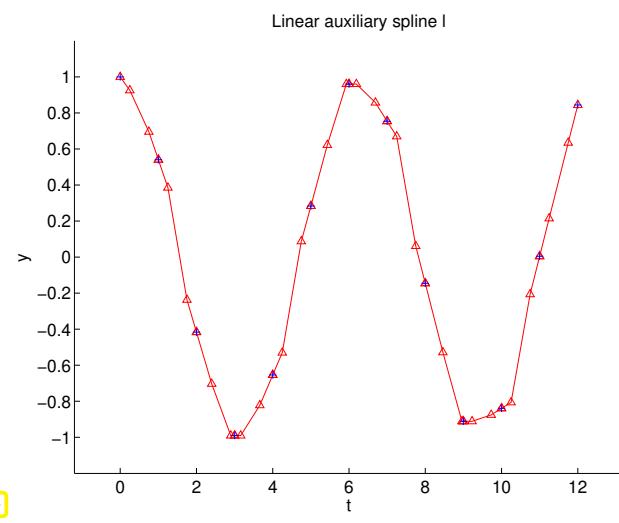


Fig. 164

Linear auxiliary spline  $\underline{l}$

- ④ Local quadratic approximation / interpolation of  $\underline{l}$ :

Tedious, but elementary calculus, confirms the following fact:

**Lemma 5.5.3.3.**

If  $g$  is a linear spline through the three points

$$(a, y_a), \left(\frac{1}{2}(a+b), w\right), (b, y_b) \quad \text{with } a < b, \quad y_a, y_b, w \in \mathbb{R},$$

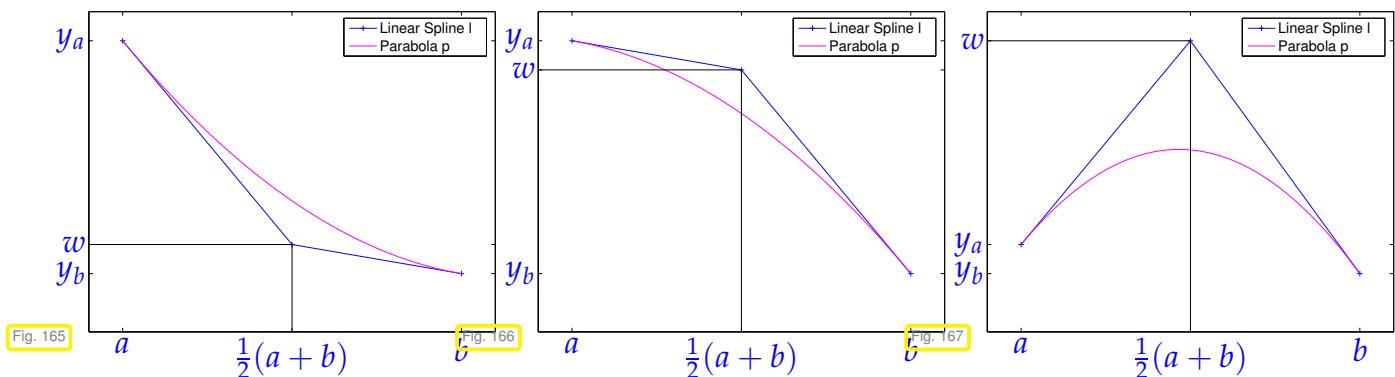
then the parabola

$$p(t) := (y_a(b-t)^2 + 2w(t-a)(b-t) + y_b(t-a)^2)/(b-a)^2, \quad a \leq t \leq b,$$

satisfies

1.  $p(a) = y_a, \quad p(b) = y_b, \quad p'(a) = g'(a), \quad p'(b) = g'(b),$
2.  $g$  monotonic increasing / decreasing  $\Rightarrow p$  monotonic increasing / decreasing,
3.  $g$  convex / concave  $\Rightarrow p$  convex / concave.

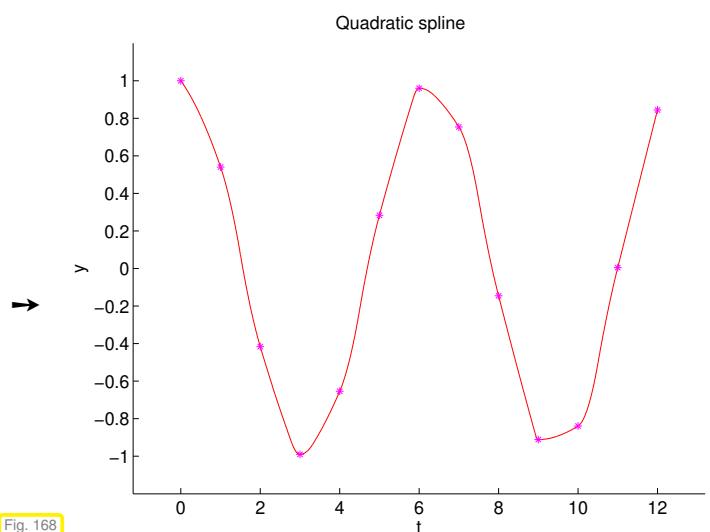
The proof boils down to discussing many cases as indicated in the following plots:



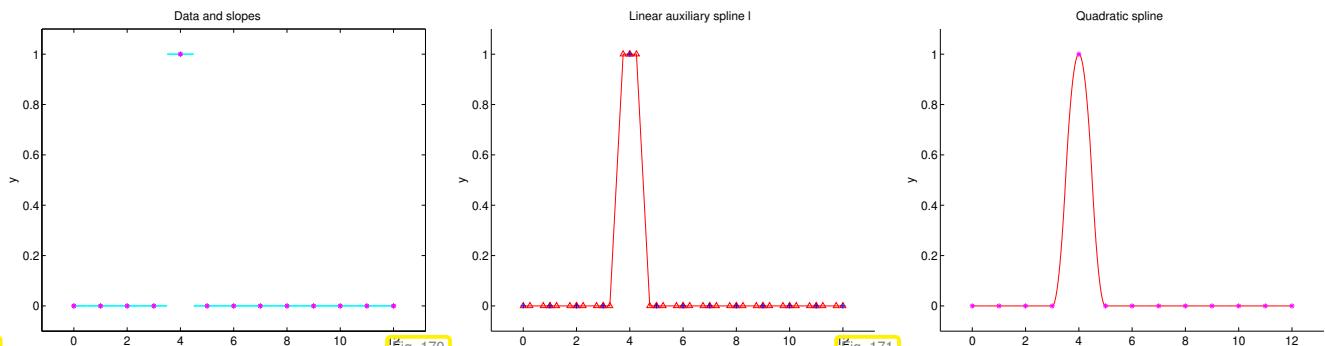
Lemma 5.5.3.3 implies that the final quadratic spline that passes through the points  $(t_j, y_j)$  with slopes  $c_j$  can be built locally as the parabola  $p$  using the linear spline  $l$  that plays the role of  $g$  in the lemma.

Continuation of Ex. 5.5.3.2:

Interpolating quadratic spline



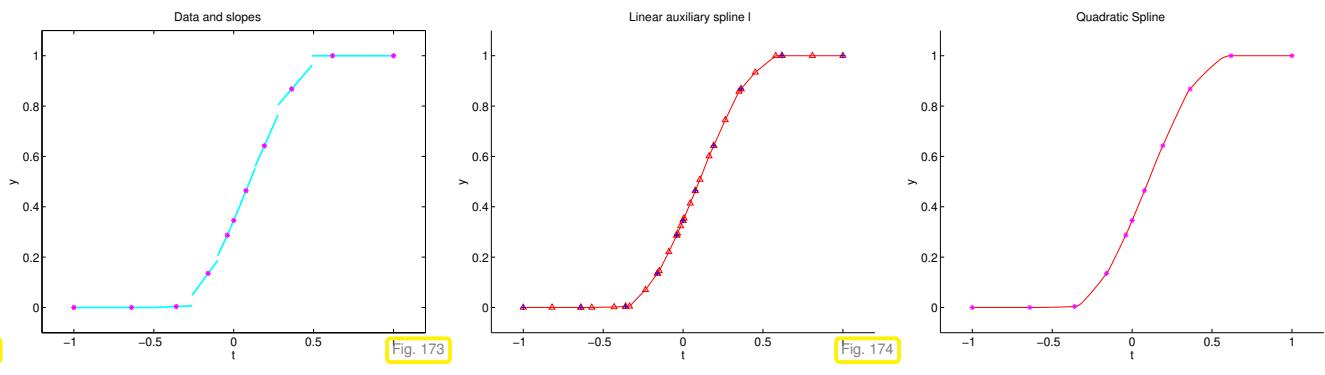
**EXAMPLE 5.5.3.4 (Cardinal shape preserving quadratic spline)** We examine the shape preserving quadratic spline that interpolates data values  $y_j = 0, j \neq i, y_i, i \in \{0, \dots, n\}$ , on an equidistant node set.



Shape preserving quadratic spline interpolation = local but not linear

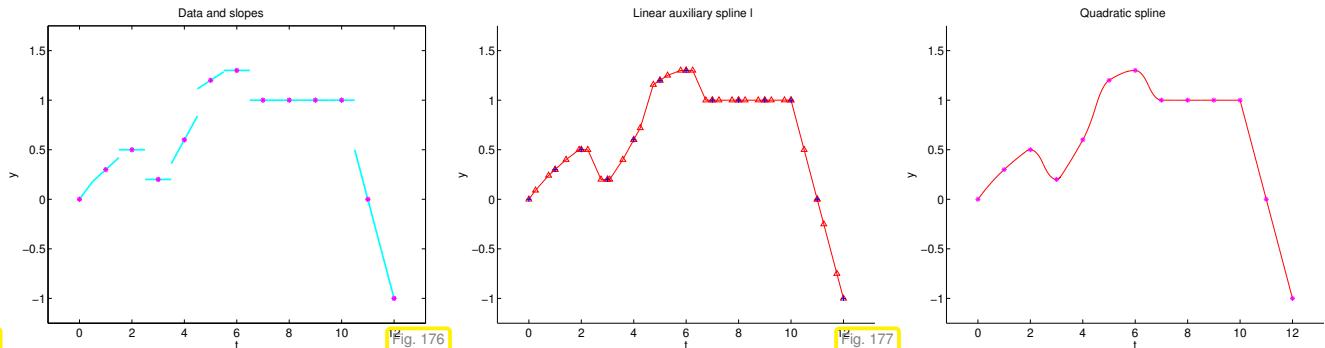
### EXAMPLE 5.5.3.5 (Shape preserving quadratic spline interpolation)

Data from Exp. 5.3.1.6:



$t_i$	0	1	2	3	4	5	6	7	8	9	10	11	12
$y_i$	0	0.3	0.5	0.2	0.6	1.2	1.3	1	1	1	0	-1	

Data from [MR81]:



In all cases we observe excellent local shape preservation. The implementation can be found in `shapepresintp.hpp` → [GITLAB](#).

## 5.6 Trigonometric Interpolation



*Supplementary literature.* This topic is also presented in [DR08, Sect. 8.5].

We consider time series data  $(t_i, y_i)$ ,  $i = 0, \dots, n$ ,  $t_i, y_i \in \mathbb{R}$ , obtained by sampling a time-dependent scalar physical quantity  $t \mapsto \varphi(t)$ . We know that  $\varphi$  is a  $T$ -periodic function with period  $T > 0$ , that is

$\varphi(t) = \varphi(t + T)$  for all  $t \in \mathbb{R}$ . In the spirit of shape preservation ( $\rightarrow$  Section 5.3) an interpolant  $f$  of the time series should also be  $T$ -periodic:  $f(t + T) = f(t)$  for all  $t \in \mathbb{R}$ .

#### Assumption 5.6.0.1. Sampling in a period

We assume the period  $T > 0$  to be known and  $t_i \in [0, T[$  for all interpolation nodes  $t_i, i = 0, \dots, n$ .

In the sequel, for the case of simplicity, we consider only  $T = 1$ .

Task: Given data points  $(t_i, y_i), y_i \in \mathbb{K}, t_i \in [0, T[$ , find a **1-periodic** function  $f : \mathbb{R} \rightarrow \mathbb{K}$  (the interpolant),  $f(t + T) = f(t) \forall t \in \mathbb{R}$ , that satisfies the **interpolation conditions**

$$f(t_i) = y_i, \quad i = 0, \dots, n. \quad (5.6.0.2)$$

### 5.6.1 Trigonometric Polynomials

The most fundamental periodic functions are derived from the trigonometric functions **sin** and **cos** and dilations of them (A **dilation** of a function  $t \mapsto \psi(t)$  is a function of the form  $t \mapsto \psi(ct)$  with some  $c > 0$ ).

#### Definition 5.6.1.1. Space of trigonometric polynomials

The vector space of 1-periodic **trigonometric polynomials** of degree  $2n, n \in \mathbb{N}$ , is given by

$$\mathcal{P}_{2n}^T := \text{Span}\{t \mapsto \cos(2\pi jt), t \mapsto \sin(2\pi jt)\}_{j=0}^n \subset C^\infty(\mathbb{R}).$$

The terminology is natural after recalling expressions for trigonometric functions via **complex** exponentials (“Euler’s formula”)

$$e^{it} = \cos t + i \sin t \Rightarrow \begin{cases} \cos t = \frac{1}{2}(e^{it} + e^{-it}) \\ \sin t = \frac{1}{2i}(e^{it} - e^{-it}) \end{cases}. \quad (5.6.1.2)$$

Thus we can rewrite  $q \in \mathcal{P}_{2n}^T$  given in the form

$$q(t) = \alpha_0 + \sum_{j=1}^n \alpha_j \cos(2\pi jt) + \beta_j \sin(2\pi jt), \quad \alpha_j, \beta_j \in \mathbb{R}, \quad (5.6.1.3)$$

by means of complex exponentials making use of Euler’s formula (5.6.1.2):

$$\begin{aligned} q(t) &= \alpha_0 + \sum_{j=1}^n \alpha_j \cos(2\pi jt) + \beta_j \sin(2\pi jt) \\ &= \alpha_0 + \frac{1}{2} \left\{ \sum_{j=1}^n (\alpha_j - i\beta_j) e^{2\pi ijt} + (\alpha_j + i\beta_j) e^{-2\pi ijt} \right\} \\ &= \alpha_0 + \frac{1}{2} \sum_{j=-n}^{-1} (\alpha_{-j} + i\beta_{-j}) e^{2\pi ijt} + \frac{1}{2} \sum_{j=1}^n (\alpha_j - i\beta_j) e^{2\pi ijt} \\ &= e^{-2\pi int} \sum_{j=0}^{2n} \gamma_j e^{2\pi ijt}, \quad \text{with} \quad \gamma_j = \begin{cases} \frac{1}{2}(\alpha_{n-j} + i\beta_{n-j}) & \text{for } j = 0, \dots, n-1, \\ \alpha_0 & \text{for } j = n, \\ \frac{1}{2}(\alpha_{j-n} - i\beta_{j-n}) & \text{for } j = n+1, \dots, 2n. \end{cases} \end{aligned} \quad (5.6.1.4)$$

Note:  $\gamma_j \in \mathbb{C}$  ➤ work in  $\mathbb{C}$  in the context of trigonometric interpolation! Admit  $y_k \in \mathbb{C}$ .

From the above manipulations we conclude

$$q \in \mathcal{P}_{2n+1}^T \Rightarrow q(t) = e^{-2\pi i t} \cdot p(e^{2\pi i t}) \quad \text{with} \quad p(z) = \sum_{j=0}^{2n} \gamma_j z^j \in \mathcal{P}_{2n}, \quad (5.6.1.5)$$

a polynomial!

and  $\gamma_j$  from (5.6.1.4).

(After scaling) a trigonometric polynomial of degree  $2n$  is a regular polynomial  $\in \mathcal{P}_{2n}$  (in  $\mathbb{C}$ ) restricted to the unit circle  $S^1 \subset \mathbb{C}$ .

☞ notation:  $S^1 := \{z \in \mathbb{C} : |z| = 1\}$  is the unit circle in the complex plane.

The relationship (5.6.1.5) justifies calling  $\mathcal{P}_{2n}^T$  a space of trigonometric *polynomials* and immediately reveals the dimension of  $\mathcal{P}_{2n}^T$ :

### Corollary 5.6.1.6. Dimension of $\mathcal{P}_{2n}^T$

The vector space  $\mathcal{P}_{2n}^T$  has dimension  $\dim \mathcal{P}_{2n}^T = 2n + 1$ .

## 5.6.2 Reduction to Lagrange Interpolation

We observed that trigonometric polynomials are standard (complex) polynomials in disguise. Next we can relate trigonometric interpolation to well-known standard Lagrangian interpolation discussed in Section 5.2.2. In fact, we slightly extend the method, because now we admit **complex interpolation nodes**. All results obtained earlier carry over to this setting.

The key tool is a smooth bijective mapping between  $I := [0, 1[$  and  $S^1$  defined as

$$\Phi_{S^1} : I \rightarrow S^1, \quad t \mapsto z := \exp(2\pi i t). \quad (5.6.2.1)$$

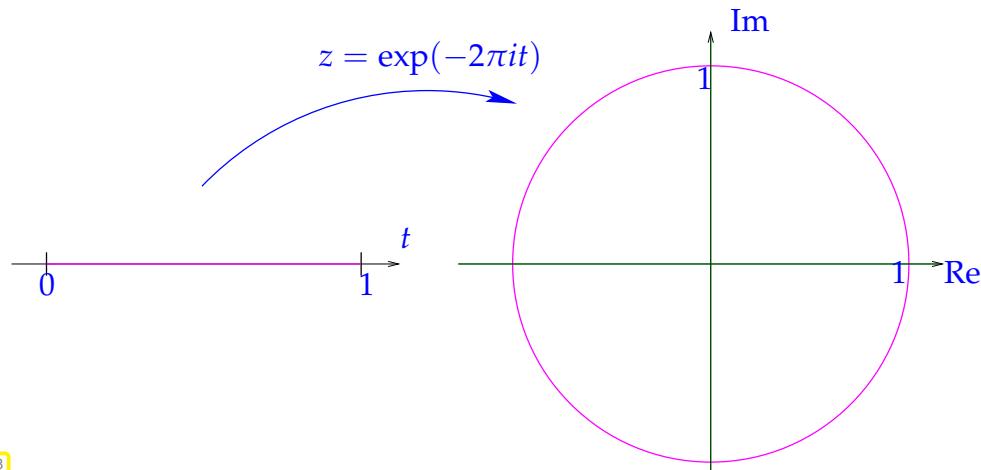


Fig. 178

Trigonometric interpolation through data points $(t_k, y_k)$	$\iff$	Polynomial interpolation through data points $(e^{2\pi i t_k}, y_k)$
Trigonometric interpolation of $f$ on $I$	$\iff$	Polynomial interpolation of pullback $(\Phi_{S^1}^{-1})^* f$ on $S^1$

Here we deal with a non-affine pullback, but the definition is the same as the one given in (6.1.1.16) for an

affine pullback:

$$(\Phi_{S^1}^{-1})^* : C^0([0,1]) \rightarrow C^0(S^1) \quad , \quad ((\Phi_{S^1}^{-1})^* f)(z) := f(\Phi_{S^1}^{-1}(z)) \quad , \quad z \in S^1. \quad (5.6.2.2)$$

Trigonometric interpolation  $\longleftrightarrow$  polynomial interpolation on  $S^1$

- All theoretical results and algorithms from polynomial interpolation carry over to trigonometric interpolation
  - ◆ Existence and uniqueness of trigonometric interpolation polynomial, see Thm. 5.2.2.7,
  - ◆ Concept of Lagrange polynomials, see (5.2.2.4),
  - ◆ the algorithms and representations discussed in Section 5.2.3.

The next code demonstrates the use of standard routines for polynomial interpolation provided by **BarycPolyInterp** ( $\rightarrow$  Code 5.2.3.7) for trigonometric interpolation.

#### C++-code 5.6.2.3: Evaluation of trigonometric interpolation polynomial in many points

```

2 // Evaluation of trigonometric interpolant at numerous points
3 // IN : t = vector of nodes  $t_0, \dots, t_n \in [0,1]$ 
4 // y = vector of data  $y_0, \dots, y_n$ 
5 // x = vector of evaluation points  $x_1, \dots, x_N$ 
6 // OUT : vector of values of the interpolant at points in x
7 VectorXd trigpolyval(const VectorXd& t, const VectorXd& y, const VectorXd& x) {
8     using idx_t = VectorXd::Index;
9     using comp = std::complex<double>;
10    const idx_t N = y.size(); // Number of data points
11    if (N % 2 == 0) throw std::runtime_error("Number of points must be odd!");
12    const idx_t n = (N - 1)/2;
13    const std::complex<double> M_I(0,1); // imaginary unit
14    // interpolation nodes and evalutation points on unit circle
15    VectorXcd tc = (2*M_PI*M_I*t).array().exp().matrix(),
16                xc = (2*M_PI*M_I*x).array().exp().matrix();
17    // Rescaled values, according to  $q(t) = e^{-2\pi i t} \cdot p(e^{2\pi i t})$ , see (5.6.1.4)
18    VectorXcd z = ((2*n*M_PI*M_I*t).array().exp() * y.array()).matrix();
19    // Evaluation of interpolating polynomial on unit circle using the
20    // barycentric interpolation formula in C, see Code 5.2.3.4
21    BarycPolyInterp<comp> Interpol(tc);
22    VectorXcd p = Interpol.eval<VectorXcd>(z, xc);
23    // Undo the scaling, see (5.6.1.5)
24    VectorXcd qc = ((-2*n*M_PI*M_I*x).array().exp() * p.array()).matrix();
25    return (qc.real()); // imaginary part is zero, cut it off
26 }
```

The computational effort is  $O(n^2 + Nn)$ ,  $N \doteq$  no. of evaluation points, as for barycentric polynomial interpolation studied in § 5.2.3.2.

The next code finds the coefficients  $\alpha_j, \beta_j \in \mathbb{R}$  of a trigonometric interpolation polynomial in the real-valued representation (5.6.1.3) for real-valued data  $y_j \in \mathbb{R}$  by simply solving the linear system of equations arising from the interpolation conditions (5.6.0.2).

#### C++-code 5.6.2.4: Computation of coefficients of trigonometric interpolation polynomial, general nodes

```

2 //Computes expansion coefficients of trigonometric polynomials (5.6.1.3)
3 // IN : t = vector of nodes  $t_0, \dots, t_n \in [0, 1]$ 
4 // y = vector of data  $y_0, \dots, y_n$ 
5 // OUT: pair of vectors storing the basis expansion coefficients  $\alpha_j, \beta_j$ ,
6 // see Def. 5.6.1.1
7 std::pair<VectorXd,VectorXd>
8 trigpolycoeff(const VectorXd& t, const VectorXd& y) {
9     const unsigned N = y.size(), n = (N-1)/2;
10    if (N % 2 == 0) throw "Number of points must be odd!\n";
11
12    // build system matrix M
13    MatrixXd M(N, N);
14    M.col(0) = VectorXd::Ones(N);
15    for (unsigned c = 1; c <= n; ++c) {
16        M.col(c) = (2*M_PI*c*t).array().cos().matrix();
17        M.col(n + c) = (2*M_PI*c*t).array().sin().matrix();
18    }
19    // solve LSE and extract coefficients  $\alpha_j$  and  $\beta_j$ 
20    VectorXd c = M.lu().solve(y);
21    return std::pair<VectorXd,VectorXd>(c.head(n + 1), c.tail(n));
}

```

The asymptotic computational effort of this implementation is dominated by the cost for Gaussian elimination applied to a fully populated (dense) matrix, see Thm. 2.5.0.2:  $O(n^3)$  for  $n \rightarrow \infty$ .

### 5.6.3 Equidistant Trigonometric Interpolation

Often timeseries data for a time-periodic quantity are measured with a constant rhythm over the entire (known) period of duration  $T > 0$ , that is,  $t_j = j\Delta t$ ,  $\Delta t = \frac{T}{n+1}$ ,  $j = 0, \dots, n$ . In this case, the formulas for computing the coefficients of the interpolating trigonometric polynomial ( $\rightarrow$  Def. 5.6.1.1) become special versions of the **discrete Fourier transform** (DFT, see Def. 4.2.0.18) studied in Section 4.2. Efficient implementation can thus harness the speed of FFT introduced in Section 4.3.

Now: 1-periodic setting, uniformly distributed interpolation nodes  $t_k = \frac{k}{2n+1}$ ,  $k = 0, \dots, 2n$

►  $(2n+1) \times (2n+1)$  linear system of equations:

$$\sum_{j=0}^{2n} \gamma_j \exp\left(2\pi i \frac{jk}{2n+1}\right) = (\mathbf{b})_k := \exp\left(2\pi i \frac{nk}{2n+1}\right) y_k, \quad k = 0, \dots, 2n. \quad (5.6.3.1)$$

$$\overleftarrow{\mathbf{F}}_{2n+1} \mathbf{c} = \mathbf{b}, \quad \mathbf{c} = [\gamma_0, \dots, \gamma_{2n}]^\top \quad \xrightarrow{\text{Lemma 4.2.0.14}} \quad \mathbf{c} = \frac{1}{2n+1} \mathbf{F}_{2n+1} \mathbf{b}. \quad (5.6.3.2)$$

$(2n+1) \times (2n+1)$  (conjugate) Fourier matrix, see (4.2.0.13)

► Fast solution by means of FFT:  $O(n \log n)$  asymptotic complexity, see Section 4.3

#### C++-code 5.6.3.3: Efficient computation of coefficient of trigonometric interpolation polynomial (equidistant nodes)

```

2 // Efficient FFT-based computation of coefficients in expansion (5.6.1.3)
3 // for a trigonometric interpolation polynomial in equidistant points

```

```

4 // $(\frac{j}{2n+1}, y_j)$ ,  $j = 0, \dots, 2n$ .
5 // IN : y has to be a row vector of odd length, return values are column
   // vectors
6 // OUT: vectors  $[\alpha_j]_j$ ,  $[\beta_j]_j$  of expansion coefficients
7 // with respect to trigonometric basis from Def. 5.6.1.1
8 std::pair<VectorXd, VectorXd> trigipequid(const VectorXd& y) {
9     using index_t = VectorXcd::Index;
10    const index_t N = y.size(), n = (N - 1)/2;
11    if (N % 2 != 1) throw "Number of points must be odd!";
12    // prepare data for fft
13    std::complex<double> M_I(0,1); // imaginary unit
14    // right hand side vector b from (5.6.3.2)
15    VectorXcd b(N);
16    for (index_t k = 0; k < N; ++k)
17        b(k) = y(k) * std::exp(2*M_PI*I*(double(n)/N*k));
18    Eigen::FFT<double> fft; // DFT helper class
19    VectorXcd c = fft.fwd(b); // means that "c = fft(b)"
20
21    // By (5.6.1.4) we can recover
22    //  $\alpha_j = \frac{1}{2}(\gamma_{n-j} + \gamma_{n+j})$  and  $\beta_j = \frac{1}{2i}(\gamma_{n-j} - \gamma_{n+j})$ ,  $j = 1, \dots, n$ ,  $\alpha_0 = \gamma_n$ .
23    VectorXcd alpha(n + 1), beta(n); alpha(0) = c(0)/((double)N);
24    for (index_t l = 1; l <= n; ++l) {
25        alpha(l) = (c(n-l)+c(n+l))/((double)N);
26        beta(l-1) = -M_I*(c(n-l)-c(n+l))/((double)N);
27    }
28    return std::pair<VectorXd, VectorXd>(alpha.real(), beta.real());
29 }
```

### EXAMPLE 5.6.3.4 (Runtime comparison for computation of coefficient of trigonometric interpolation polynomials)

Runtime measurements for MATLAB equivalents of codes Code 5.6.2.4 and Code 5.6.3.3

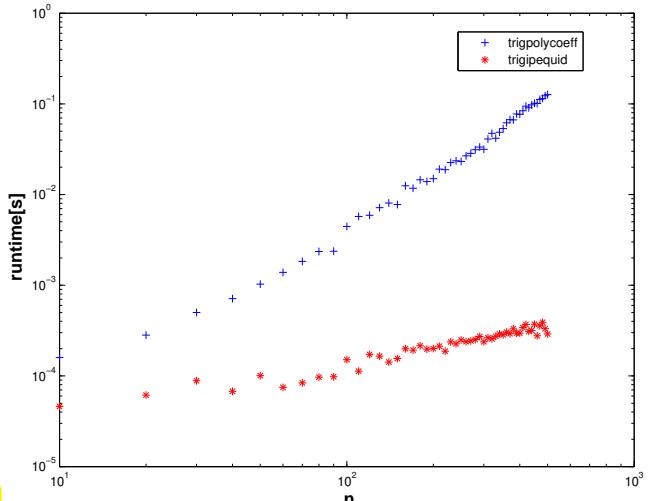
tic-toe-timings



MATLAB 7.10.0.499 (R2010a)

CPU: Intel Core i7, 2.66 GHz, 2 cores, L2 256 KB, L3 4 MB

OS: Mac OS X, Darwin Kernel Version 10.5.0



Same observation as in Ex. 4.3.0.12: massive gain in efficiency through relying on FFT.

### Remark 5.6.3.5 (Efficient evaluation of trigonometric interpolation polynomials)

Task: Find an efficient way to evaluate the trigonometric polynomial

$$q(t) = \alpha_0 + \sum_{j=1}^n \alpha_j \cos(2\pi jt) + \beta_j \sin(2\pi jt), \quad \alpha_j, \beta_j \in \mathbb{R}, \quad (5.6.1.3)$$

at *equidistant* points  $\frac{k}{N}, N > 2n, k = 0, \dots, N-1$ .

As in (5.6.1.4) we can rewrite

$$q(t) = e^{-2\pi int} \sum_{j=0}^{2n} \gamma_j e^{2\pi ijt}, \quad \text{with} \quad \gamma_j = \begin{cases} \frac{1}{2}(\alpha_{n-j} + i\beta_{n-j}) & \text{for } j = 0, \dots, n-1, \\ \alpha_0 & \text{for } j = n, \\ \frac{1}{2}(\alpha_{j-n} - i\beta_{j-n}) & \text{for } j = n+1, \dots, 2n. \end{cases}$$

$$(5.6.1.4) \quad \begin{aligned} \blacktriangleright \quad q(k/N) &= e^{-2\pi i nk/N} \sum_{j=0}^{2n} \gamma_j \exp(2\pi i \frac{kj}{N}), \quad k = 0, \dots, N-1. \\ \blacktriangleright \quad q(k/N) &= e^{-2\pi i kn/N} v_j \quad \text{with} \quad \mathbf{v} = \bar{\mathbf{F}}_N \tilde{\mathbf{c}} \end{aligned} \quad (5.6.3.6)$$

Fourier matrix, see (4.2.0.13).

where  $\tilde{\mathbf{c}} \in \mathbb{C}^N$  is obtained by *zero padding* of  $\mathbf{c} := (\gamma_0, \dots, \gamma_{2n})^T$ :

$$(\tilde{\mathbf{c}})_k = \begin{cases} \gamma_j & , \text{for } k = 0, \dots, 2n, \\ 0 & , \text{for } k = 2n+1, \dots, N-1. \end{cases}$$

The FFT-based implementation has been fine in the following code.

#### C++-code 5.6.3.7: Fast evaluation of trigonometric polynomial at *equidistant* points

```

2 void trigipequidcomp (const VectorXcd& a, const VectorXcd& b, const unsigned N,
3   VectorXcd& y) {
4   const unsigned n = a.size() - 1;
5   if (N < (2*n - 1)) {
6     std::cerr << "N is too small! Must be larger than 2*n";
7     return;
8   }
9   const std::complex<double> iu(0,1); // imaginary unit
10  // build vector γ
11  VectorXcd gamma(2*n + 1);
12  gamma(n) = a(0);
13  for (unsigned k = 0; k < n; ++k) {
14    gamma(k) = 0.5*( a(n - k) + iu*b(n - k - 1) );
15    gamma(n + k + 1) = 0.5*( a(k + 1) - iu*b(k) );
16  }
17  // zero padding to obtain c̃
18  VectorXcd ch(N); ch << gamma, VectorXcd::Zero(N - (2*n + 1));
19
20  // realize multiplication with conjugate fourier matrix
21  Eigen::FFT<double> fft;
22  VectorXcd chCon = ch.conjugate();
23  VectorXcd v = fft.fwd(chCon).conjugate();
24
25  // final scaling, implemented without efficiency considerations
26  y = VectorXcd(N);
27  for (unsigned k = 0; k < N; ++k)
    y(k) = v(k) * std::exp( -2.*k*n*M_PI/N*iu );

```

28 }

The next code merges the steps of computing the coefficient of the trigonometric interpolation polynomial in equidistant points and its evaluation in another set of equidistant points.

#### C++ code 5.6.3.8: *Equidistant points: fast on the fly evaluation of trigonometric interpolation polynomial*

```

2 // Evaluation of trigonometric interpolation polynomial through  $(\frac{j}{2n+1}, y_j)$ ,  

3 //  $j = 0, \dots, 2n$   

4 // in equidistant points  $\frac{k}{N}$ ,  $k = 0, N - 1$   

5 // IN : y = vector of values to be interpolated  

6 // q (COMPLEX!) will be used to save the return values  

7 void trigpolyvalequid(const VectorXd y, const int M, VectorXd& q) {  

8     const int N = y.size();  

9     if (N % 2 == 0) {  

10         std::cerr << "Number of points must be odd!\n";  

11         return;  

12     }  

13     const int n = (N - 1)/2;  

14     // computing coefficient  $\gamma_j$ , see (5.6.3.2)  

15     VectorXcd a, b;  

16     trigipequid(y, a, b);  

17  

18     std::complex<double> i(0,1);  

19     VectorXcd gamma(2*n + 1);  

20     gamma(n) = a(0);  

21     for (int k = 0; k < n; ++k) {  

22         gamma(k) = 0.5*( a(n - k) + i*b(n - k - 1) );  

23         gamma(n + k + 1) = 0.5*( a(k + 1) - i*b(k) );  

24     }  

25  

26     // zero padding  

27     VectorXcd ch(M); ch << gamma, VectorXcd::Zero(M - (2*n + 1));  

28  

29     // build conjugate fourier matrix  

30     Eigen::FFT<double> fft;  

31     const VectorXcd chCon = ch.conjugate();  

32     const VectorXcd v = fft.fwd(chCon).conjugate();  

33  

34     // multiplicate with conjugate fourier matrix  

35     VectorXcd q_complex = VectorXcd(M);  

36     for (int k = 0; k < M; ++k) {  

37         q_complex(k) = v(k) * std::exp( -2.*k*n*M_PI/M*i );  

38     }  

39     // complex part is zero up to machine precision, cut off!  

40     q = q_complex.real();

```

## 5.7 Least Squares Data Fitting

As remarked in Ex. 5.1.0.7 the basic assumption underlying the reconstructing of the functional dependence of two quantities by means of interpolation as that of accurate data. In case of data uncertainty or measurement errors the exact satisfaction of interpolation conditions ceases to make sense and we are better off reconstructing a fitting function that is merely “close to the data” in a sense to be made precise next.

The task of (multidimensional, vector-valued) **least squares data fitting** can be described as follows:

### Least square data fitting

Given: data points  $(\mathbf{t}_i, \mathbf{y}_i)$ ,  $i \in \{1, \dots, m\}$ ,  $m \in \mathbb{N}$ ,  $\mathbf{t}_i \in D \subset \mathbb{R}^k$ ,  $\mathbf{y}_i \in \mathbb{R}^d$

Objective: Find a (continuous) function  $\mathbf{f} : D \mapsto \mathbb{R}^d$  in some set  $S \subset C^0(D)$  of admissible functions satisfying

$$\mathbf{f} \in \underset{\mathbf{g} \in S}{\operatorname{argmin}} \sum_{i=1}^m \|\mathbf{g}(\mathbf{t}_i) - \mathbf{y}_i\|_2^2. \quad (5.7.0.2)$$

The function  $\mathbf{f}$  is called the (best) least squares fit for the data in  $S$ .

Focus in this section:  $k = 1$ ,  $d = 1$  (one-dimensional scalar setting)

$\Leftrightarrow$  fitting of scalar data depending on one parameter ( $t \in \mathbb{R}$ ),

$\Leftrightarrow$  Data points  $(t_i, y_i)$ ,  $t_i \in I \subset \mathbb{R}$ ,  $y_i \in \mathbb{R}$   $\Rightarrow$   $S \subset C^0(I)$

$\Leftrightarrow$  (best) least squares fit  $f : I \subset \mathbb{R} \rightarrow \mathbb{R}$  is a function of one real variable.

$$k = 1, d = 1: \quad (5.7.0.2) \Leftrightarrow f \in \underset{g \in S}{\operatorname{argmin}} \sum_{i=1}^m |g(t_i) - y_i|^2. \quad (5.7.0.3)$$

**§5.7.0.4 (Linear spaces of admissible functions)** Consider a special variant of the general least squares data fitting problem: The set  $S$  of admissible continuous functions is now chosen as a finite-dimensional vector space  $V_n \subset C^0(D)$ ,  $\dim V_n = n \in \mathbb{N}$ , cf. the discussion in § 5.1.0.15 for interpolation.

Choose basis of  $V_n$ :  $V_n = \operatorname{Span}\{\mathbf{b}_1, \dots, \mathbf{b}_n\}$ ,  $\mathbf{b}_j : D \rightarrow \mathbb{R}^d$  continuous

► The best least squares fit  $\mathbf{f} \in V_n$  can be represented by a finite linear combination of the basis functions  $\mathbf{b}_j$ :

$$\mathbf{f}(t) = \sum_{j=1}^n x_j \mathbf{b}_j(t) , \quad x_j \in \mathbb{R}. \quad (5.7.0.5)$$

Often, in the case  $d > 1$ ,  $V_n$  is chosen as a product space

$$V_n = \underbrace{W \times \cdots \times W}_{d \text{ factors}}, \quad (5.7.0.6)$$

of a space  $W \subset C^0(D)$  of  $\mathbb{R}$ -valued functions  $D \mapsto \mathbb{R}$ . In this case

$$\dim V_n = d \cdot \dim W.$$

If  $\ell := \dim W$  and  $q_1, \dots, q_\ell$  is a basis of  $W$ , then

$$\{\mathbf{e}_i q_j\}_{\substack{i=1, \dots, d \\ j=1, \dots, \ell}} = \{\mathbf{e}_1 q_1, \mathbf{e}_2 q_1, \dots, \mathbf{e}_d q_1, \mathbf{e}_1 q_2, \dots, \mathbf{e}_d q_2, \mathbf{e}_1 q_3, \dots, \mathbf{e}_1 q_\ell, \dots, \mathbf{e}_d q_\ell\} \quad (5.7.0.7)$$

is a basis of  $V_n$  ( $\mathbf{e}_i \hat{=} i$ -th unit vector).  $\square$

**§5.7.0.8 (Linear data fitting)  $\rightarrow$  [DR08, Sect. 4.1]** We adopt the setting of § 5.7.0.4 of an  $n$ -dimensional space  $V_n$  of admissible functions with basis  $\{\mathbf{b}_q, \dots, \mathbf{b}_n\}$ . Then the least squares data fitting problem can be recast as follows.

General Linear least squares fitting problem:

Given:

- ◆ data points  $(\mathbf{t}_i, \mathbf{y}_i) \in \mathbb{R}^k \times \mathbb{R}^d$ ,  $i = 1, \dots, m$
- ◆ basis functions  $\mathbf{b}_j : D \subset \mathbb{R}^k \mapsto \mathbb{R}$ ,  $j = 1, \dots, n$ ,  $n < m$

Sought: coefficients  $x_j \in \mathbb{R}$ ,  $j = 1, \dots, n$ , such that

$$\mathbf{x} := [x_1, \dots, x_n]^\top := \underset{z_j \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^m \left\| \sum_{j=1}^n z_j \mathbf{b}_j(\mathbf{t}_i) - \mathbf{y}_i \right\|_2^2. \quad (5.7.0.9)$$

Special cases:

- For  $k = 1$ ,  $d = 1$ , data points  $(t_i, y_i) \in \mathbb{R} \times \mathbb{R}$  (scalar, one-dimensional setting), and  $V_n = \operatorname{Span}\{\mathbf{b}_1, \dots, \mathbf{b}_n\}$ , we seek coefficients  $x_j \in \mathbb{R}$ ,  $j = 1, \dots, n$ , as the components of a vector  $\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$  satisfying

$$\mathbf{x} = \underset{\mathbf{z} \in \mathbb{R}^n}{\operatorname{argmin}} \sum_{i=1}^m \left| \sum_{j=1}^n (\mathbf{z})_j b_j(t_i) - y_i \right|^2. \quad (5.7.0.10)$$

- If  $V_n$  is a product space according to (5.7.0.6) with basis (5.7.0.7), then (5.7.0.9) amounts to finding vectors  $\mathbf{x}_j \in \mathbb{R}^d$ ,  $j = 1, \dots, \ell$  with

$$(\mathbf{x}_1, \dots, \mathbf{x}_\ell) = \underset{\mathbf{z}_j \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^m \left\| \sum_{j=1}^n \mathbf{z}_j q_j(\mathbf{t}_i) - \mathbf{y}_i \right\|_2^2. \quad (5.7.0.11)$$

**EXAMPLE 5.7.0.12 (Linear parameter estimation = linear data fitting)  $\rightarrow$  Ex. 3.0.1.4, Ex. 3.1.1.3)** The linear parameter estimation/**linear regression** problem presented in Ex. 3.0.1.4 can be recast as a linear data fitting problem with

- $k = n$ ,  $d = 1$ , data points  $(\mathbf{x}_i, y_i) \in \mathbb{R}^k \times \mathbb{R}$ ,
- an  $k+1$ -dimensional space  $V_n = \{\mathbf{x} \mapsto \mathbf{a}^\top \mathbf{x} + \beta, \mathbf{a} \in \mathbb{R}^k, \beta \in \mathbb{R}\}$  of **affine linear** admissible functions,
- the choice of basis  $\{\mathbf{x} \mapsto (\mathbf{x})_1, \dots, \mathbf{x} \mapsto (\mathbf{x})_k, \mathbf{x} \mapsto 1\}$ .

**§5.7.0.13 (Linear data fitting as a lineare least squares problem)** Linear (least squares) data fitting leads to an overdetermined linear system of equations for which we seek a least squares solution ( $\rightarrow$  Def. 3.1.1.1) as in Section 3.1.1. To see this rewrite

$$\sum_{i=1}^m \left\| \sum_{j=1}^n z_j \mathbf{b}_j(\mathbf{t}_i) - \mathbf{y}_i \right\|_2^2 = \sum_{i=1}^m \left( \sum_{j=1}^n \sum_{r=1}^d (\mathbf{b}_j(\mathbf{t}_i))_r z_j - (\mathbf{y}_i)_r \right)^2.$$

#### Theorem 5.7.0.14. Least squares solution of data fitting problem

The solution  $\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$  of the linear least squares fitting problem (5.7.0.9) is the least squares solution of the overdetermined linear system of equations

$$\begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_d \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_d \end{bmatrix}, \quad (5.7.0.15)$$

with  $\mathbf{A}_r := \begin{bmatrix} (\mathbf{b}_1(\mathbf{t}_1))_r & \dots & (\mathbf{b}_n(\mathbf{t}_1))_r \\ \vdots & & \vdots \\ (\mathbf{b}_q(\mathbf{t}_m))_r & \dots & (\mathbf{b}_n(\mathbf{t}_m))_r \end{bmatrix} \in \mathbb{R}^{m,n}$ ,  $\mathbf{b}_r := \begin{bmatrix} (\mathbf{y}_1)_r \\ \vdots \\ (\mathbf{y}_m)_r \end{bmatrix} \in \mathbb{R}^m$ ,  $r = 1, \dots, d$ .

In the one-dimensional, scalar case ( $= 1, d = 1$ ) of (5.7.0.10) the related overdetermined linear system of equations is

$$\begin{bmatrix} b_1(t_1) & \dots & b_n(t_1) \\ \vdots & & \vdots \\ b_1(t_m) & \dots & b_n(t_m) \end{bmatrix} \mathbf{x} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}. \quad (5.7.0.16)$$

Obviously, for  $m = n$  we recover the 1D data interpolation problem from § 5.1.0.15.

Having reduced the linear least squares data fitting problem to finding the least squares solution of an overdetermined linear system of equations, we can now apply theoretical results about least squares solutions, for instance, Cor. 3.1.2.13. The key issue is, whether the coefficient matrix of (5.7.0.16) has full rank  $n$ . Of course, this will depend on the location of the  $t_i$ .

#### Lemma 5.7.0.17. Unique solvability of linear least squares fitting problem

The scalar one-dimensional linear least squares fitting problem (5.7.0.10) with  $\dim V_n = n$ ,  $V_n$  the vector space of admissible functions, has a unique solution, if and only if there are  $t_{i_1}, \dots, t_{i_n}$  such that

$$\begin{bmatrix} b_1(t_{i_1}) & \dots & b_n(t_{i_1}) \\ \vdots & & \vdots \\ b_1(t_{i_n}) & \dots & b_n(t_{i_n}) \end{bmatrix} \in \mathbb{R}^{n,n} \quad \text{is invertible,} \quad (5.7.0.18)$$

which is independent of the choice of basis of  $V_n$ .

Equivalent to (5.7.0.18) is the requirement, that there is an  $n$ -subset of  $\{t_1, \dots, t_n\}$  such that the corresponding interpolation problem for  $V_n$  has a unique solution for any data values  $y_i$ .  $\square$

#### EXAMPLE 5.7.0.19 (Polynomial fitting)

Special variant of scalar ( $d = 1$ ), one-dimensional  $k = 1$  linear data fitting ( $\rightarrow$  § 5.7.0.8): we choose the space of admissible functions as **polynomials** of degree  $n - 1$ ,

$$V_n = \mathcal{P}_{n-1} \quad , \quad \text{e.g. with basis } b_j(t) = t^{j-1} \quad (\text{monomial basis, Section 5.2.1}) .$$

The corresponding overdetermined linear system of equations (5.7.0.16) now reads:

$$\begin{bmatrix} 1 & t_1 & t_1^2 & \dots & t_1^{n-1} \\ 1 & t_2 & t_2^2 & \dots & t_2^{n-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & t_m & t_m^2 & \dots & t_m^{n-1} \end{bmatrix} \mathbf{x} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} , \quad (5.7.0.20)$$

which, for  $M \geq n$ , has full rank, because it contains invertible Vandermonde matrices (5.2.2.11), Rem. 5.2.2.10.

The next code demonstrates the computation of the fitting polynomial with respect to the monomial basis of  $\mathcal{P}_{n-1}$ :

#### C++11-code 5.7.0.21: Polynomial fitting → GITLAB

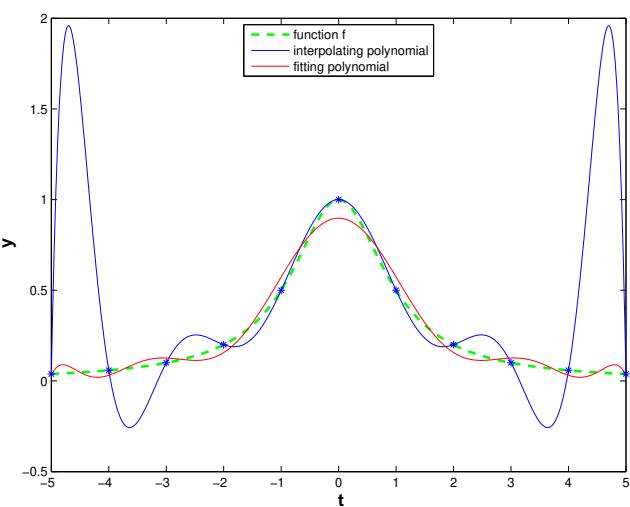
```

2 // Solver for polynomial linear least squares data fitting problem
3 // data points passed in t and y, 'order' = degree + 1
4 VectorXd polyfit(const VectorXd& t, const VectorXd& y, const unsigned& order) {
5     // Initialize the coefficient matrix of (5.7.0.20)
6     Eigen::MatrixXd A = Eigen::MatrixXd::Ones(t.size(), order + 1);
7     for (unsigned j = 1; j < order + 1; ++j)
8         A.col(j) = A.col(j - 1).cwiseProduct(t);
9     // Use EIGEN's built-in least squares solver, see Code 3.3.4.2
10    Eigen::VectorXd coeffs = A.householderQr().solve(y);
11    // leading coefficients have low indices.
12    return coeffs.reverse();
13 }
```

The function `polyfit` returns a vector  $[x_1, x_2, \dots, x_n]^T$  describing the fitting polynomial according to the convention

$$p(t) = x_1 t^{n-1} + x_2 t^{n-2} + \dots + x_{n-1} t + x_n . \quad (5.7.0.22)$$

#### EXPERIMENT 5.7.0.23 (Polybomial interpolation vs. polynomial fitting)



Data from function  $f(t) = \frac{1}{1+t^2}$ ,

- ◆ polynomial degree  $d = 10$ ,
- ◆ interpolation through data points  $(t_j, f(t_j))$ ,  $j = 0, \dots, d$ ,  $t_j = -5 + j$ , see Ex. 5.2.4.1,
- ◆ fitting to data points  $(t_j, f(t_j))$ ,  $j = 0, \dots, 3d$ ,  $t_j = -5 + j/3$ .

Experiment carried out with the code `interpfit.cpp` → GITLAB.

Fitting helps curb oscillations that haunt polynomial interpolation: in terms of “shape preservation” the fitted polynomial is clearly superior to the interpolating polynomial.

**Remark 5.7.0.24 (“Overfitting”)** In the previous example we saw that the degree-10 polynomial fitted to the data in 31 points rather well reflects the “shape of the data”, see Fig. 180. If we had fit a polynomial of degree 30 to the data instead, the result would have been a function with humongous oscillations, because we would have recovered the interpolating polynomial at equidistant nodes, *cf.* Ex. 5.2.4.1.

Thus, in Exp. 5.7.0.23 we see a manifestation of a widely observed phenomenon, which plays a big role in machine learning:

### Overfitting

Fitting data with functions from a large space often produces poorer results (in terms of “shape preservation” and “generalization error”) than the use of a smaller subspace for fitting.

## Learning outcomes

After you have studied this chapter you should

- understand the use of basis functions for representing functions on a computer,
- know the concept of a interpolation operator and what its linearity means,
- know the connection between linear interpolation operators and linear systems of equations,
- be familiar with efficient algorithms for polynomial interpolation in different settings,
- know the meaning and significance of “sensitivity” in the context of interpolation,
- Be familiar with the notions of “shape preservation” for an interpolation scheme and its different aspects (monotonicity, curvature),
- know the details of cubic Hermite interpolation and how to ensure that it is monotonicity preserving.
- know what splines are and how cubic spline interpolation with different endpoint constraints works.

# Bibliography

- [Aki70] H. Akima. “A New Method of Interpolation and Smooth Curve Fitting Based on Local Procedures”. In: *J. ACM* 17.4 (1970), pp. 589–602.
- [AG11] Uri M. Ascher and Chen Greif. *A first course in numerical methods*. Vol. 7. Computational Science & Engineering. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2011, pp. xxii+552. DOI: [10.1137/1.9780898719987](https://doi.org/10.1137/1.9780898719987) (cit. on p. 327).
- [BT04] Jean-Paul Berrut and Lloyd N. Trefethen. “Barycentric Lagrange interpolation”. In: *SIAM Rev.* 46.3 (2004), 501–517 (electronic). DOI: [10.1137/S0036144502417715](https://doi.org/10.1137/S0036144502417715).
- [CR92] D. Coppersmith and T.J. Rivlin. “The growth of polynomials bounded at equally spaced points”. In: *SIAM J. Math. Anal.* 23.4 (1992), pp. 970–983 (cit. on p. 347).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 327, 329, 330, 334, 340, 343, 345, 359, 370, 379).
- [FC80] F.N. Fritsch and R.E. Carlson. “Monotone Piecewise Cubic Interpolation”. In: *SIAM J. Numer. Anal.* 17.2 (1980), pp. 238–246 (cit. on p. 358).
- [Gou05] E. Gourgoulhon. *An introduction to polynomial interpolation*. Slides for School on spectral methods. 2005.
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on p. 360).
- [Kva14] Boris I. Kvasov. “Monotone and convex interpolation by weighted quadratic splines”. In: *Adv. Comput. Math.* 40.1 (2014), pp. 91–116. DOI: [10.1007/s10444-013-9300-9](https://doi.org/10.1007/s10444-013-9300-9).
- [MR81] D. McAllister and J. Roulier. “An algorithm for computing a shape-preserving osculatory quadratic spline”. In: *ACM Trans. Math. Software* 7.3 (1981), pp. 331–347 (cit. on pp. 366, 370).
- [Moh15] A. Mohsen. “CORRECTING THE FUNCTION DERIVATIVE ESTIMATION USING LAGRANGIAN INTERPOLATION”. In: *ZAMP* (2015).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on p. 326).
- [QQY01] A.L. Dontchev and H.-D. Qi, L.-Q. Qi, and H.-X. Yin. “Convergence of Newton’s method for convex best interpolation”. In: *Numer. Math.* 87.3 (2001), pp. 435–456 (cit. on p. 366).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 327, 329, 330, 340, 345, 359, 360, 363).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on pp. 349, 350).
- [Tre14] N. Trefethen. *SIX MYTHS OF POLYNOMIAL INTERPOLATION AND QUADRATURE*. Slides. 2014.

# Chapter 6

## Approximation of Functions in 1D

In Chapter 5 we aimed to fill the gaps between given data points by constructing a function connecting them. In this chapter is available already, at least in principle, but its evaluation is too costly, which forces us to replace it with a “cheaper” function.

### §6.0.0.1 (General function approximation problem)

#### Approximation of functions: Generic view

Given: a **function**  $f : D \subset \mathbb{R}^n \mapsto \mathbb{R}^d$ ,  $n, d \in \mathbb{N}$ , often in **procedural form**, e.g. for  $n = d = 1$ , as  
`double f(double)` ( $\rightarrow$  Rem. 5.1.0.8).

Goal: Find a “*simple*” function  $\tilde{f} : D \mapsto \mathbb{R}^d$  such that the **approximation error**  $f - \tilde{f}$  is “*small*”.

Let us clarify the meaning of some terms:

Made more precise: “*simple*”

The function  $\tilde{f}$  can be encoded by a **small amount of information** and is easy to evaluate.  
For instance, this is the case for polynomial or piecewise polynomial  $\tilde{f}$ .

Made more precise: “*small*” approximation error

$\|f - \tilde{f}\|$  is small for some **norm**  $\|\cdot\|$  on the space  $C^0(\overline{D})$  of (piecewise) continuous functions.

The most commonly used norms are

◆ the **supremum norm**  $\|g\|_\infty := \|g\|_{L^\infty(D)} := \max_{x \in D} |g(x)|$ , see (5.2.4.4).

If the approximation error is small with respect to the supremum norm,  $\tilde{f}$  is also called a good **uniform** approximant of  $f$ .

◆ the  **$L^2$ -norm**  $\|g\|_2^2 := \|g\|_{L^2(D)}^2 = \int_D |g(x)|^2 dx$ , see (5.2.4.5).

Below we consider only the case  $n = d = 1$ : approximation of scalar valued functions defined on an interval. The techniques can be applied componentwise in order to cope with the case of vector valued function ( $d > 1$ ).  $\square$

### EXAMPLE 6.0.0.3 (Model reduction in circuit simulation)

The non-linear circuit ( $\triangleright$ ) stands for a **diode**, a circuit element, whose resistance strongly depends on the voltage) sketched beside has two ports.

For the sake of circuit simulation it should be replaced by a non-linear lumped circuit element characterized by a single voltage-current constitutive relationship  $I = I(U)$ . For any value of the voltage  $U$  the current  $I$  can be computed by solving a (large) non-linear system of equations, see Ex. 8.0.0.1.

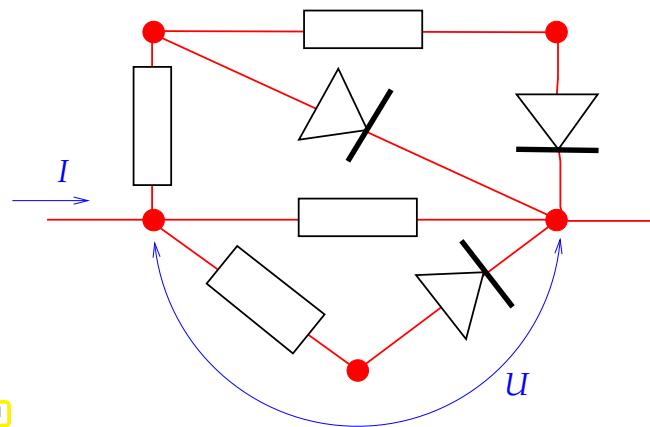


Fig. 181

A faster alternative is the advance approximation of the function  $U \mapsto I(U)$  based on a few computed values  $I(U_i)$ ,  $i = 0, \dots, n$ , followed by the fast evaluation of the approximant  $U \mapsto \tilde{I}(U)$  during actual circuit simulations. This is an example of **model reduction** by approximation of functions: a complex subsystem in a mathematical model is replaced by a **surrogate function**.

In this example we also encounter a typical situation: we have nothing at our disposal but, possibly expensive, point evaluations of the function  $U \mapsto I(U)$  ( $U \mapsto I(U)$  in “procedural form”, see Rem. 5.1.0.8). The number of evaluations of  $I(U)$  will largely determine the cost of building  $\tilde{I}$ .

This application displays a fundamental difference compared to the reconstruction of constitutive relationships from a priori measurements → Ex. 5.1.0.7: Now we are free to choose the number and location of the data points, because we can simply evaluate the function  $U \mapsto I(U)$  for any  $U$  and as often as needed.

#### C++ code 6.0.0.4: Class describing a 2-port circuit element for circuit simulation

```

1 class CircuitElement {
2     private:
3         // internal data describing  $U \mapsto \tilde{I}(U)$ .
4     public:
5         // Constructor taking some parameters and building  $\tilde{I}$ 
6         CircuitElement(const Parameters &P);
7         // Point evaluation operators for  $\tilde{I}$  and  $\frac{d}{dU}\tilde{I}$ 
8         double I(double U) const;
9         double dIdU(double U) const;
10    };

```

**§6.0.0.5 (Approximation schemes)** We define an abstract concept for the sake of clarity: When in this chapter we talk about an “**approximation scheme**” (in 1D) we refer to a mapping  $A : X \rightarrow V$ , where  $X$  and  $V$  are spaces of functions  $I \mapsto \mathbb{K}$ ,  $I \subset \mathbb{R}$  an interval.

Examples are

- $X = C^k(I)$ , the spaces of functions  $I \mapsto \mathbb{K}$  that are  $k$  times continuously differentiable,  $k \in \mathbb{N}$ .
- $V = \mathcal{P}_m(I)$ , the space of polynomials of degree  $\leq k$ , see Section 5.2.1
- $V = \mathcal{S}_{d,M}$ , the space of splines of degree  $d$  on the knot set  $M \subset I$ , see Def. 5.5.0.1.
- $V = \mathcal{P}_{2n}^T$ , the space of trigonometric polynomials of degree  $2n$ , see Def. 5.6.1.1.

**§6.0.0.6 (Approximation by interpolation)** In Chapter 5 we discussed ways to construct functions whose graph runs through given data points, see 5.1. We can hope that the interpolant will approximate the function, if the data points are also located on the graph of that function. Thus every interpolation scheme, see § 5.1.0.6, spawns a corresponding approximation scheme.

**Interpolation scheme + sampling → approximation scheme**

$$f : I \subset \mathbb{R} \rightarrow \mathbb{K} \xrightarrow{\text{sampling}} (t_i, y_i := f(t_i))_{i=0}^m \xrightarrow{\text{interpolation}} \tilde{f} := \mathbf{I}_{\mathcal{T}} \mathbf{y} \quad (\tilde{f}(t_i) = y_i).$$

↑  
free choice of nodes  $t_i \in I$

In this chapter we will mainly study approximation by interpolation relying on the interpolation schemes (→ § 5.1.0.6) introduced in Section 5.2, Section 5.4, and Section 5.5.

There is *additional freedom* compared to data interpolation: we can choose the interpolation nodes in a smart way in order to obtain an accurate interpolant  $\tilde{f}$ .

**Remark 6.0.0.8 (Interpolation and approximation: enabling technologies)** Approximation and interpolation (→ Chapter 5) are key components of many numerical methods, like for integration, differentiation and computation of the solutions of differential equations, as well as for computer graphics and generation of smooth curves and surfaces.



This chapter is a “foundations” part of the course

## Contents

6.1	Approximation by Global Polynomials . . . . .	387
6.1.1	Polynomial Approximation: Theory . . . . .	387
6.1.2	Error Estimates for Polynomial Interpolation . . . . .	393
6.1.3	Chebychev Interpolation . . . . .	402
6.2	Mean Square Best Approximation . . . . .	416
6.2.1	Abstract Theory . . . . .	416
6.2.2	Polynomial Mean Square Best Approximation . . . . .	420
6.3	Uniform Best Approximation . . . . .	426
6.4	Approximation by Trigonometric Polynomials . . . . .	430
6.4.1	Approximation by Trigonometric Interpolation . . . . .	431
6.4.2	Trigonometric Interpolation Error Estimates . . . . .	432
6.4.3	Trigonometric Interpolation of Analytic Periodic Functions . . . . .	436
6.5	Approximation by Piecewise Polynomials . . . . .	439
6.5.1	Piecewise Polynomial Lagrange Interpolation . . . . .	441
6.5.2	Cubic Hermite Interpolation: Error Estimates . . . . .	444
6.5.3	Cubic Spline Interpolation: Error Estimates [Han02, Ch. 47] . . . . .	446

## 6.1 Approximation by Global Polynomials

The space  $\mathcal{P}_k$  of polynomials of degree  $\leq k$  has been introduced in Section 5.2.1. For reasons listed in § 5.2.1.3 polynomials are the most important theoretical and practical tool for the approximation of functions. The next example presents an important case of approximation by polynomials.

**EXAMPLE 6.1.0.1 (Taylor approximation) → [Str09, Sect. 5.5]** The local approximation of sufficiently smooth functions by polynomials is a key idea in calculus, which manifests itself in the importance of approximation by **Taylor polynomials**: For  $f \in C^k(I)$ ,  $k \in \mathbb{N}$ ,  $I \subset \mathbb{R}$  an interval, we approximate

$$f(t) \approx f(t_0) + f'(t_0)(t - t_0) + \underbrace{\frac{f''(t_0)}{2}(t - t_0)^2 + \cdots + \frac{f^{(k)}(t_0)}{k!}(t - t_0)^k}_{=:T_k(t) \in \mathcal{P}_k}, \quad \text{for some } t_0 \in I.$$

☞ Notation:  $f^{(k)} \doteq k$ -th derivative of function  $f : I \subset \mathbb{R} \rightarrow \mathbb{K}$

The Taylor polynomial  $T_k$  of degree  $k$  approximates  $f$  in a neighbourhood  $J \subset I$  of  $t_0$  ( $J$  can be small!). This can be quantified by the **remainder formulas** [Str09, Bem. 5.5.1]

$$f(t) - T_k(t) = \int_{t_0}^t f^{(k+1)}(\tau) \frac{(t - \tau)^k}{k!} d\tau \tag{6.1.0.2a}$$

$$= f^{(k+1)}(\xi) \frac{(t - t_0)^{k+1}}{(k+1)!}, \quad \xi = \xi(t, t_0) \in ]\min(t, t_0), \max(t, t_0)[, \tag{6.1.0.2b}$$

which shows that for  $f \in C^{k+1}(I)$  the Taylor polynomial  $T_k$  is pointwise close to  $f \in C^{k+1}(I)$ , if the interval  $I$  is small and  $f^{(k+1)}$  is bounded pointwise.

Approximation by Taylor polynomials is easy and direct but inefficient: a polynomial of lower degree often gives the same accuracy. Moreover, when  $f$  is available only in procedural form as **double**  $f$  (**double**), (approximations of) higher order derivatives are difficult to obtain. ↴

**§6.1.0.3 (Nested approximation spaces of polynomials)** Obviously, for every interval  $I \subset \mathbb{R}$ , the spaces of polynomials are **nested** in the following sense:

$$\mathcal{P}_0 \subset \mathcal{P}_1 \subset \cdots \subset \mathcal{P}_m \subset \mathcal{P}_{m+1} \subset \cdots \subset C^\infty(I), \tag{6.1.0.4}$$

with finite, but increasing dimensions  $\dim \mathcal{P}_m = m + 1$  according to Thm. 5.2.1.2.

With this family of nested spaces of polynomials at our disposal, it is natural to study *associated families of approximation schemes*, one for each degree, mapping into  $\mathcal{P}_m$ ,  $m \in \mathbb{N}_0$ . ↴

### 6.1.1 Polynomial Approximation: Theory

**§6.1.1.1 (Scope of polynomial approximation)** Sloppily speaking, according to (6.1.0.2b) the Taylor polynomials from Ex. 6.1.0.1 provide **uniform** (→ § 6.0.0.1) approximation of a *smooth* function  $f$  in (small)

intervals, provided that its derivatives do not blow up “too fast” (We do not want to make this precise here).

The question is, whether polynomials still offer uniform approximation on arbitrary bounded closed intervals and for functions that are merely continuous, but not any smoother. The answer is YES and this profound result is known as the **Weierstrass Approximation Theorem**. Here we give an extended version with a concrete formula due to Bernstein, see [Dav75, Section 6.2].

### Theorem 6.1.1.2. Uniform approximation by polynomials

For  $f \in C^0([0, 1])$ , define the ***n*-th Bernstein approximant** as

$$p_n(t) = \sum_{j=0}^n f(j/n) \binom{n}{j} t^j (1-t)^{n-j}, \quad p_n \in \mathcal{P}_n. \quad (6.1.1.3)$$

It satisfies  $\|f - p_n\|_\infty \rightarrow 0$  for  $n \rightarrow \infty$ . If  $f \in C^m([0, 1])$ , then even  $\|f^{(k)} - p_n^{(k)}\|_\infty \rightarrow 0$  for  $n \rightarrow \infty$  and all  $0 \leq k \leq m$ .

☞ Notation:  $g^{(k)} \triangleq k$ -th derivative of a function  $g : I \subset \mathbb{R} \rightarrow \mathbb{K}$ .

In (6.1.1.3) the function  $f$  is approximated by a linear combination of **Bernstein polynomials** of degree  $n$

$$B_j^n(t) := \binom{n}{j} t^j (1-t)^{n-j}, \quad B_j^n \in \mathcal{P}_n. \quad (6.1.1.4)$$

Plots of Bernstein polynomials of degree  $n = 7$  ▷

Bernstein polynomials provide a positive **partition of unity** over  $[0, 1]$ :

$$\sum_{j=0}^n B_j^n(t) = 1 \quad \forall t \in \mathbb{R}, \quad (6.1.1.5)$$

$$0 \leq B_j^n(t) \leq 1 \quad \forall 0 \leq t \leq 1. \quad (6.1.1.6)$$

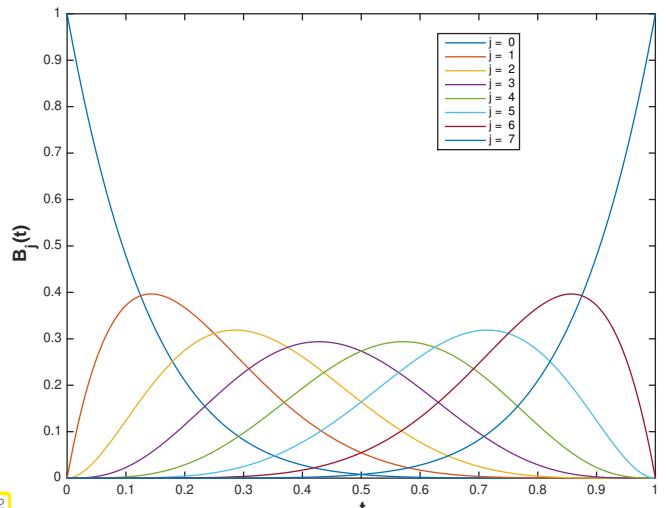


Fig. 182

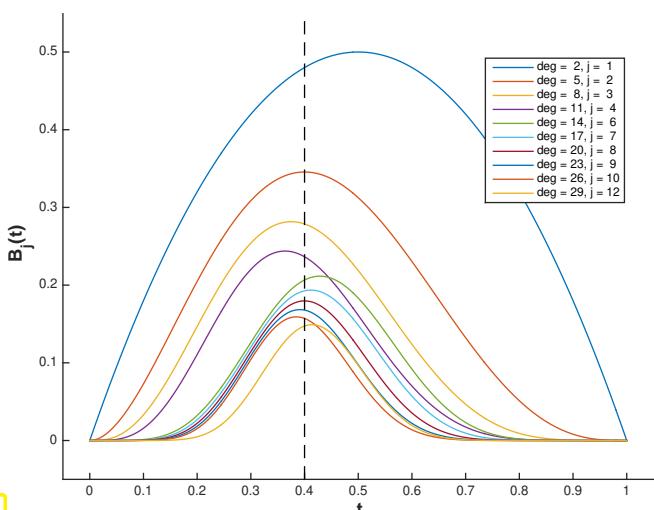


Fig. 183

▷ Since  $\frac{d}{dt} B_j^n(t) = B_j^n(t) \left( \frac{j}{t} - \frac{n-j}{1-t} \right)$ ,  $B_j^n$  has its unique local maximum in  $[0, 1]$  at the site  $t_{\max} := \frac{j}{n}$ . As  $n \rightarrow \infty$  the Bernstein polynomials become more and more concentrated around the maximum.

*Proof. (of Thm. 6.1.1.2, first part)* Fix  $t \in [0, 1]$ . Using the notations from (6.1.1.3) and the identity (6.1.1.5) we find

$$f(t) - p_n(t) = \sum_{k=0}^n (f(t) - f(j/n)) B_j^n(t). \quad (6.1.1.7)$$

As we see from Fig. 183, for large  $n$  the bulk of sum will be contributed by Bernstein polynomial with index  $j/n \approx t$ , because for every  $\delta > 0$

$$\sum_{|j/n-t|>\delta} B_j^n(t) \leq \frac{1}{\delta^2} \sum_{|j/n-t|>\delta} (j/n - t)^2 B_j^n(t) \leq \frac{1}{\delta^2} \sum_{j=0}^n (j/n - t)^2 B_j^n(t) \stackrel{(*)}{=} \frac{nt(1-t)}{\delta^2 n^2} \leq \frac{1}{4n\delta^2}.$$

$\sum_{|j/n-t|>\delta}$  means summation over  $j \in \mathbb{N}_0$  with summation indices confined to the set  $\{j : |j/n - t| > \delta\}$ .

The identity  $(*)$  can be established by direct but tedious computations.

Combining this estimate with (6.1.1.6) and (6.1.1.7) we arrive at

$$|f(t) - p_n(t)| \leq \sum_{|j/n-t|>\delta} \frac{1}{4n\delta^2} |f(t) - f(j/n)| + \sum_{|j/n-t|\leq\delta} |f(t) - f(j/n)|.$$

Since,  $f$  is uniformly continuous on  $[0, 1]$ , given  $\epsilon > 0$  we can choose  $\delta > 0$  independently of  $t$  such that  $|f(s) - f(t)| < \epsilon$ , if  $|s - t| < \delta$ . Then, if we choose  $n > (\epsilon\delta^2)^{-1}$ , we can bound

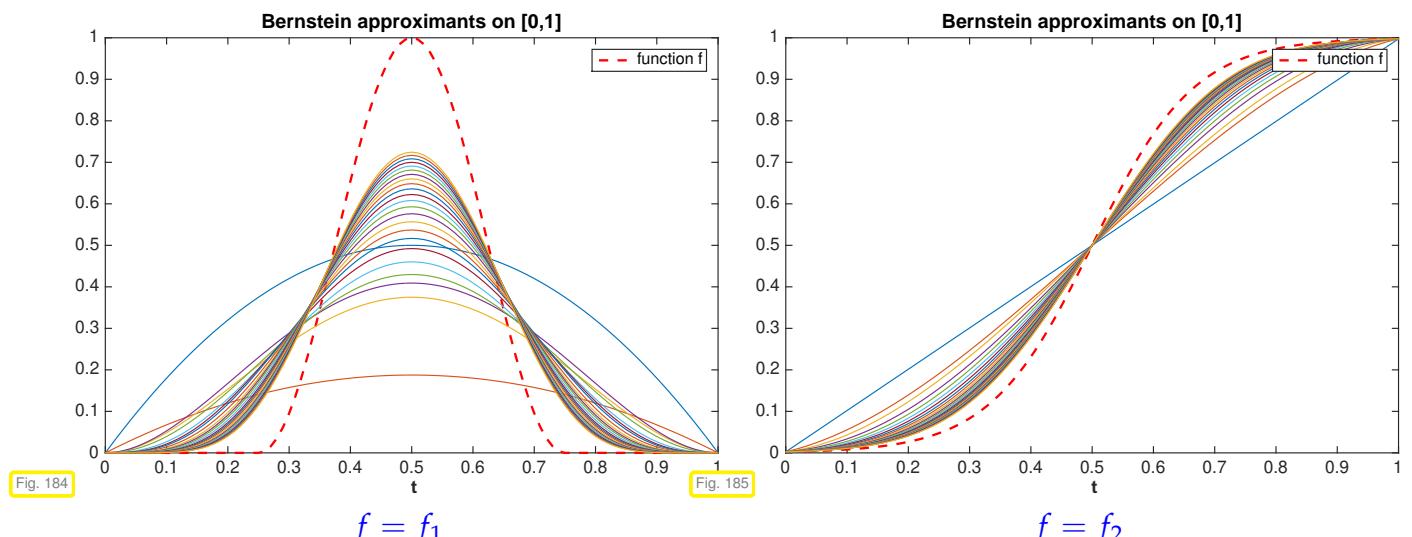
$$|f(t) - p_n(t)| \leq (\|f\|_\infty + 1)\epsilon \quad \forall t \in [0, 1].$$

This means that  $p_n$  is arbitrarily close to  $f$  for sufficiently large  $n$ .  $\square$

**EXPERIMENT 6.1.1.8 (Bernstein approximants)** We compute and plot  $p_n$ ,  $n = 1, \dots, 25$ , for two functions

$$f_1(t) := \begin{cases} 0 & , \text{if } |2t - 1| > \frac{1}{2}, \\ \frac{1}{2}(1 + \cos(2\pi(2t - 1))) & \text{else,} \end{cases} \quad , \quad f_2(t) := \frac{1}{1 + e^{-12(x-1/2)}}.$$

The following plots display the sequences of the polynomials  $p_n$  for  $n = 2, \dots, 25$ .



We see that the Bernstein approximants “slowly” edge closer and closer to  $f$ . Apparently it takes a very large degree to get really close to  $f$ .  $\square$

**§6.1.1.9 (Best approximation)** Now we introduce a concept needed to gauge how close an approximation scheme gets to the best possible performance.

#### Definition 6.1.1.10. (Size of) best approximatont error

Let  $\|\cdot\|$  be a (semi-)norm on a space  $X$  of functions  $I \mapsto \mathbb{K}$ ,  $I \subset \mathbb{R}$  an interval. The (size of the) **best approximation error** of  $f \in X$  in the space  $\mathcal{P}_k$  of polynomials of degree  $\leq k$  with respect to  $\|\cdot\|$  is

$$\text{dist}_{\|\cdot\|}(f, \mathcal{P}_k) := \inf_{p \in \mathcal{P}_k} \|f - p\|.$$

The notation  $\text{dist}_{\|\cdot\|}$  is motivated by the notation of “distance” as distance to the nearest point in a set.

For the  $L^2$ -norm  $\|\cdot\|_2$  and the supremum norm  $\|\cdot\|_\infty$  the best approximation error is well defined for  $C = C^0(I)$ .

The polynomial realizing best approximation w.r.t.  $\|\cdot\|$  may neither be unique nor computable with reasonable effort. Often one is content with rather sharp upper bounds like those asserted in the next theorem, due to Jackson [Dav75, Thm. 13.3.7].

#### Theorem 6.1.1.11. $L^\infty$ polynomial best approximation estimate

If  $f \in C^r([-1, 1])$  ( $r$  times continuously differentiable),  $r \in \mathbb{N}$ , then, for any polynomial degree  $n \geq r$ ,

$$\inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty([-1, 1])} \leq (1 + \pi^2/2)^r \frac{(n-r)!}{n!} \|f^{(r)}\|_{L^\infty([-1, 1])},$$

where

$$\|f^{(r)}\|_{L^\infty([-1, 1])} := \max_{x \in [-1, 1]} |f^{(r)}(x)|.$$

As above,  $f^{(r)}$  stands for the  $r$ -th derivative of  $f$ . Using Stirling's formula

$$\sqrt{2\pi} n^{n+1/2} e^{-n} \leq n! \leq e n^{n+1/2} e^{-n} \quad \forall n \in \mathbb{N}, \quad (6.1.1.12)$$

we can get a looser bound of the form

$$\inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty([-1, 1])} \leq C(r) n^{-r} \|f^{(r)}\|_{L^\infty([-1, 1])}, \quad (6.1.1.13)$$

with  $C(r)$  dependent on  $r$ , but *independent of  $f$*  and, in particular, the polynomial *degree  $n$* . Using the Landau symbol from Def. 1.4.1.2 we can rewrite the statement of (6.1.1.13) in asymptotic form

$$(6.1.1.13) \Rightarrow \inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty([-1, 1])} = O(n^{-r}) \quad \text{for } n \rightarrow \infty.$$

□

**Remark 6.1.1.14 (Transformation of polynomial approximation schemes)** What if a polynomial approximation scheme is defined only on a special interval, say  $[-1, 1]$ . Then by the following trick it can be transferred to any interval  $[a, b] \subset \mathbb{R}$ .

Assume that an interval  $[a, b] \subset \mathbb{R}$ ,  $a < b$ , and a polynomial approximation scheme  $\hat{A} : C^0([-1, 1]) \rightarrow$

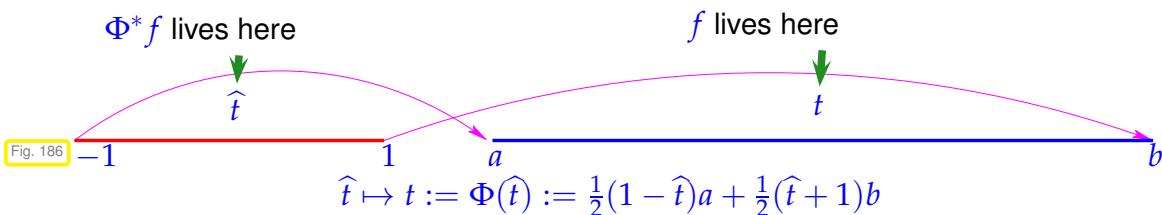
$\mathcal{P}_n$  are given. Based on the **affine linear** mapping

$$\Phi : [-1, 1] \rightarrow [a, b] , \quad \Phi(\hat{t}) := a + \frac{1}{2}(\hat{t} + 1)(b - a) , \quad -1 \leq \hat{t} \leq 1 , \quad (6.1.1.15)$$

we can introduce the **affine pullback** of functions:

$$\Phi^* : C^0([a, b]) \rightarrow C^0([-1, 1]) , \quad \Phi^*(f)(\hat{t}) := (f \circ \Phi)(\hat{t}) := f(\Phi(\hat{t})) , \quad -1 \leq \hat{t} \leq 1 .$$

(6.1.1.16)



We add the important observations that affine pullbacks are linear and bijective, they are isomorphisms of the involved vector spaces of functions (what is the inverse?).

### Lemma 6.1.1.17. Affine pullbacks preserve polynomials

If  $\Phi^* : C^0([a, b]) \rightarrow C^0([-1, 1])$  is an affine pullback according to (6.1.1.15) and (6.1.1.16), then  $\Phi^* : \mathcal{P}_n \rightarrow \mathcal{P}_n$  is a **bijection** linear mapping for any  $n \in \mathbb{N}_0$ .

*Proof.* This is a consequence of the fact that translations and dilations take polynomials to polynomials of the *same degree*: for monomials we find

$$\Phi^*\{t \rightarrow t^n\} = \{\hat{t} \rightarrow (a + \frac{1}{2}(\hat{t} + 1)(b - a))^n\} \in \mathcal{P}_n .$$

□

The lemma tells us that the spaces of polynomials of some maximal degree are *invariant under affine pullback*. Thus, we can *define* a polynomial approximation scheme  $\mathbf{A}$  on  $C^0([a, b])$  by

$$\mathbf{A} : C^0([a, b]) \rightarrow \mathcal{P}_n , \quad \mathbf{A} := (\Phi^*)^{-1} \circ \widehat{\mathbf{A}} \circ \Phi^* , \quad (6.1.1.18)$$

whenever  $\widehat{\mathbf{A}}$  is a polynomial approximation scheme on  $[-1, 1]$ . □

**Remark 6.1.1.19 (Transforming approximation error estimates)** Thm. 6.1.1.11 targets only the special interval  $[-1, 1]$ . What does it imply for polynomial best approximation on a general interval  $[a, b]$ ? To answer this question we apply techniques from Rem. 6.1.1.14, in particular the pullback (6.1.1.16).

We first have to study the change of norms of functions under the action of affine pullbacks:

### Lemma 6.1.1.20. Transformation of norms under affine pullbacks

For every  $f \in C^0([a, b])$  we have

$$\|f\|_{L^\infty([a,b])} = \|\Phi^*f\|_{L^\infty([-1,1])} , \quad \|f\|_{L^2([a,b])} = \sqrt{|b-a|} \|\Phi^*f\|_{L^2([-1,1])} . \quad (6.1.1.21)$$

*Proof.* The first estimate should be evident, and the second is a consequence of the transformation formula for integrals [Str09, Satz 6.1.5]

$$\int_{-1}^1 \Phi^* t(\hat{t}) d\hat{t} = \frac{b-a}{2} \int_a^b f(t) dt , \quad (6.1.1.22)$$

and the definition of the  $L^2$ -norm from (5.2.4.5).  $\square$

Thus, for norms of the approximation errors of polynomial approximation schemes defined by affine transformation (6.1.1.18) we get

$$\begin{aligned} \|f - Af\|_{L^\infty([a,b])} &= \left\| \Phi^* f - \widehat{\mathbf{A}}(\Phi^* f) \right\|_{L^\infty([-1,1])} , \\ \|f - Af\|_{L^2([a,b])} &= \sqrt{|b-a|} \left\| \Phi^* f - \widehat{\mathbf{A}}(\Phi^* f) \right\|_{L^2([-1,1])} , \quad \forall f \in C^0([a,b]) . \end{aligned} \quad (6.1.1.23)$$

Equipped with approximation error estimates for  $\mathbf{A}$ , we can infer corresponding estimates for  $\widehat{\mathbf{A}}$ .

The bounds for approximation errors often involve norms of derivatives as in Thm. 6.1.1.11. Hence, it is important to understand the interplay of pullback and differentiation: By the 1D chain rule

$$\frac{d}{d\hat{t}}(\Phi^* f)(\hat{t}) = \frac{df}{dt}(\Phi(\hat{t})) \frac{d\Phi}{d\hat{t}} = \frac{df}{dt}(\Phi(\hat{t})) \cdot \frac{1}{2}(b-a) ,$$

which implies a simple scaling rule for derivatives of arbitrary order  $r \in \mathbb{N}_0$ :

$$(\Phi^* f)^{(r)} = \left( \frac{b-a}{2} \right)^r \Phi^*(f^{(r)}) . \quad (6.1.1.24)$$

Lemma 6.1.1.20

$$\left\| (\Phi^* f)^{(r)} \right\|_{L^\infty([-1,1])} = \left( \frac{b-a}{2} \right)^r \left\| f^{(r)} \right\|_{L^\infty([a,b])} , \quad f \in C^r([a,b]), r \in \mathbb{N}_0 . \quad (6.1.1.25)$$

**Remark 6.1.1.26 (Polynomial best approximation on general intervals)** The estimate (6.1.1.24) together with Thm. 6.1.1.11 paves the way for bounding the polynomial best approximation error on arbitrary intervals  $[a, b]$ ,  $a, b \in \mathbb{R}$ . Based on the affine mapping  $\Phi : [-1, 1] \rightarrow [a, b]$  from (6.1.1.15) and writing  $\Phi^*$  for the pullback according to (6.1.1.16) we can chain estimates. If  $f \in C^r([a, b])$  and  $n \geq r$ , then

$$\begin{aligned} \inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty([a,b])} &\stackrel{(*)}{=} \inf_{p \in \mathcal{P}_n} \|\Phi^* f - p\|_{L^\infty([-1,1])} \\ &\stackrel{\text{Thm. 6.1.1.11}}{\leq} (1 + \pi^2/2)^r \frac{(n-r)!}{n!} \left\| (\Phi^* f)^{(r)} \right\|_{L^\infty([-1,1])} \\ &\stackrel{(6.1.1.24)}{=} (1 + \pi^2/2)^r \frac{(n-r)!}{n!} \left( \frac{b-a}{2} \right)^r \left\| f^{(r)} \right\|_{L^\infty([a,b])} . \end{aligned}$$

In step (\*) we used the result of Lemma 6.1.1.17 that  $\Phi^* p \in \mathcal{P}_n$  for all  $p \in \mathcal{P}_n$ . Invoking the arguments that gave us (6.1.1.13), we end up with the simpler bound

$$\inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty([a,b])} \leq C(r) \left( \frac{b-a}{n} \right)^r \left\| f^{(r)} \right\|_{L^\infty([a,b])} . \quad (6.1.1.27)$$

Observe that the length of the interval enters the bound in  $r$ -th power.  $\square$

### 6.1.2 Error Estimates for Polynomial Interpolation

In Section 5.2.2, Cor. 5.2.2.8, we introduced the Lagrangian polynomial interpolation operator  $\mathbf{l}_{\mathcal{T}} : \mathbb{K}^{n+1} \rightarrow \mathcal{P}_n$  belonging to a node set  $\mathcal{T} = \{t_j\}_{j=0}^n$ . In the spirit of § 6.0.0.6 it induces an approximation scheme on  $C^0(I)$ ,  $I \subset \mathbb{R}$  an interval, if  $\mathcal{T} \subset I$ .

#### Definition 6.1.2.1. Lagrangian (interpolation polynomial) approximation scheme

Given an interval  $I \subset \mathbb{R}$ ,  $n \in \mathbb{N}$ , a node set  $\mathcal{T} = \{t_0, \dots, t_n\} \subset I$ , the **Lagrangian (interpolation polynomial) approximation scheme**  $\mathbf{L}_{\mathcal{T}} : C^0(I) \rightarrow \mathcal{P}_n$  is defined by

$$\begin{aligned}\mathbf{L}_{\mathcal{T}}(f) &:= \mathbf{l}_{\mathcal{T}}(\mathbf{y}) \in \mathcal{P}_n \quad \text{with} \quad \mathbf{y} := (f(t_0), \dots, f(t_n))^T \in \mathbb{K}^{n+1}, \\ \mathbf{l}_{\mathcal{T}}(\mathbf{y})(t_j) &= (\mathbf{y})_j, \quad j = 0, \dots, n.\end{aligned}$$

Our goal in this section will be

to estimate the norm of the **interpolation error**  $\|f - \mathbf{l}_{\mathcal{T}}f\|$  (for relevant norm on  $C(I)$ ).

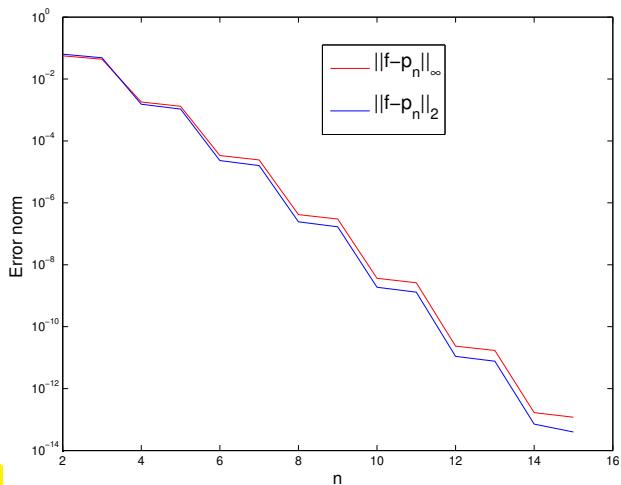
**§6.1.2.2 (Families of Lagrangian interpolation polynomial approximation schemes)** Already Thm. 6.1.1.11 considered the size of the best approximation error in  $\mathcal{P}_n$  as a function of the polynomial degree  $n$ . In the same vein, we may study a *family* of Lagrange interpolation schemes  $\{\mathbf{L}_{\mathcal{T}_n}\}_{n \in \mathbb{N}_0}$  on  $I \subset \mathbb{R}$  induced by a *family of node sets*  $\{\mathcal{T}_n\}_{n \in \mathbb{N}_0}$ ,  $\mathcal{T}_n \subset I$ , according to Def. 6.1.2.1.

An example for such a family of node sets on  $I := [a, b]$  are the **equidistant** or **equispaced** nodes

$$\mathcal{T}_n := \{t_j^{(n)} := a + (b - a) \frac{j}{n} : j = 0, \dots, n\} \subset I. \quad (6.1.2.3)$$

For families of Lagrange interpolation schemes  $\{\mathbf{L}_{\mathcal{T}_n}\}_{n \in \mathbb{N}_0}$  we can shift the focus onto estimating the **asymptotic** behavior of the norm of the interpolation error for  $n \rightarrow \infty$ .  $\dashv$

**EXPERIMENT 6.1.2.4 (Asymptotic behavior of Lagrange interpolation error)** We perform polynomial interpolation of  $f(t) = \sin t$  on equispaced nodes in  $I = [0, \pi]$ :  $\mathcal{T}_n = \{j\pi/n\}_{j=0}^n$ . Write  $p$  for the polynomial interpolant:  $p := \mathbf{L}_{\mathcal{T}_n}f \in \mathcal{P}_n$ .



In the numerical experiment the norms of the interpolation errors can be computed only approximately as follows.

- $L^\infty$ -norm: approximated by sampling on a grid of meshsize  $\pi/1000$ .
- $L^2$ -norm: numerical quadrature ( $\rightarrow$  Chapter 7) with trapezoidal rule (7.4.0.4) on a grid of meshsize  $\pi/1000$ .

$\triangleleft$  approximate norms  $\|f - \mathbf{L}_{\mathcal{T}_n}f\|_*, * = 2, \infty$ .

Fig. 187

**§6.1.2.5 (Classification of asymptotic behavior of norms of the interpolation error)** In the previous experiment we observed a clearly visible regular behavior of  $\|f - L_{T_n}f\|$  as we increased the polynomial degree  $n$ . The prediction of the decay law for  $\|f - L_{T_n}f\|$  for  $n \rightarrow \infty$  is one goal in the study of interpolation errors.

Often this goal can be achieved, even if a rigorous quantitative bound for a norm of the interpolation error remains elusive. In other words, in many cases

no bound for  $\|f - L_{T_n}f\|$  can be given, but its decay for increasing  $n$  can be described precisely.

Now we introduce some important terminology for the qualitative description of the behavior of  $\|f - L_{T_n}f\|$  as a function of the polynomial degree  $n$ . We assume that

$$\exists C \neq C(n) > 0: \|f - L_{T_n}f\| \leq C T(n) \quad \text{for } n \rightarrow \infty. \quad (6.1.2.6)$$

#### Definition 6.1.2.7. Types of asymptotic convergence of approximation schemes

Writing  $T(n)$  for the bound of the norm of the interpolation error according to (6.1.2.6) we distinguish the following *types of asymptotic behavior*:

$$\begin{aligned} \exists p > 0: \quad T(n) \leq n^{-p} & : \text{algebraic convergence, with rate } p > 0, \quad \forall n \in \mathbb{N}. \\ \exists 0 < q < 1: \quad T(n) \leq q^n & : \text{exponential convergence,} \end{aligned}$$

The bounds are assumed to be sharp in the sense, that no bounds with larger rate  $p$  (for algebraic convergence) or smaller  $q$  (for exponential convergence) can be found.

Convergence behavior of norms of the interpolation error is often expressed by means of the Landau-O-notation, cf. Def. 1.4.1.2:

Algebraic convergence:	$\ f - I_T f\  = O(n^{-p})$	
Exponential convergence:	$\ f - I_T f\  = O(q^n)$	for $n \rightarrow \infty$ (“asymptotic!”)

□

**Remark 6.1.2.8 (Different meanings of “convergence”)** Unfortunately, as in many other fields of mathematics and beyond, also in numerical analysis the meaning of terms is context-dependent:

Beware: same concept  $\leftrightarrow$  different meanings:

- **convergence** of a sequence (e.g. of iterates  $x^{(k)}$ ) → Section 8.1)
- **convergence** of an approximation (dependent on an approximation parameter, e.g.  $n$ )

□

#### Remark 6.1.2.9 (Determining the type of convergence in numerical experiments → § 1.4.1.6)

Given pairs  $(n_i, \epsilon_i)$ ,  $i = 1, 2, 3, \dots$ ,  $n_i \hat{=}$  polynomial degrees,  $\epsilon_i \hat{=}$  (measured) norms of interpolation errors, how can we tease out the likely type of convergence according to Def. 6.1.2.7? A similar task was already encountered in § 1.4.1.6, where we had to extract information about asymptotic complexity from runtime measurements.

① Conjectured: **algebraic convergence**:  $\epsilon_i \approx C n^{-p}$

$$\log(\epsilon_i) \approx \log(C) - p \log n_i \quad (\text{affine linear in log-log scale}).$$

Slope of line approximating  $(\log n_i, \log(\epsilon_i))$  predicts rate of algebraic convergence: Apply linear regression as explained in Ex. 3.1.1.3 for data points  $(\log n_i, \log \epsilon_i) \geq$  least squares estimate for rate  $p$ .

① Conjectured: exponential convergence:  $\epsilon_i \approx C \exp(-\beta n_i)$

$$\log \epsilon_i \approx \log(C) - \beta n_i \quad (\text{affine linear in lin-log scale}) .$$

Apply linear regression ( $\rightarrow$  Ex. 3.1.1.3) to points  $(n_i, \log \epsilon_i) \geq$  estimate for  $q := \exp(-\beta)$ .

☞ Fig. 187: we suspect exponential convergence in Exp. 6.1.2.4. ↗

**EXAMPLE 6.1.2.10 (Runge's example → Ex. 5.2.4.1)** We examine the polynomial interpolant of  $f(t) = \frac{1}{1+t^2}$  for equispaced nodes:

$$\mathcal{T}_n := \left\{ t_j := -5 + \frac{10}{n} j \right\}_{j=0}^n, \quad j = 0, \dots, n \quad \geq \quad y_j = \frac{1}{1+t_j^2} .$$

We rely on an approximate computation of the supremum norm of the interpolation error by means of sampling as in Exp. 6.1.2.4; here we used 1000 equidistant sampling points, see Code 6.1.2.11.

#### C++ code 6.1.2.11: Computing the interpolation error for Runge's example

```

2 // Note: "quick & dirty" implementation !
3 // Lambda function representing  $x \mapsto (1+x^2)^{-1}$ 
4 auto f = [] (double x) { return 1. / (1 + x * x); };
5 // 1000 sampling points for approximate maximum norm
6 const VectorXd x = VectorXd::LinSpaced(1000, -5, 5);
7 // Sample function
8 const VectorXd fx = x.unaryExpr(f); // evaluate f at x
9
10 std::vector<double> err; // Accumulate error norms here
11 for (int d = 1; d <= 20; ++d) {
12     // Interpolation nodes
13     const VectorXd t = Eigen::VectorXd::LinSpaced(d + 1, -5, 5);
14     // Interpolation data values
15     const VectorXd ft = feval(f, t);
16     // Compute interpolating polynomial in monomial representation
17     const VectorXd p = polyfit(t, ft, d);
18     // Evaluate polynomial interpolant in sampling points
19     const VectorXd y = polyval(p, x);
20     // Approximate supremum norm of interpolation error
21     err.push_back((y - fx).cwiseAbs().maxCoeff());
22 }
```

Here, `polyfit()` computes the monomial coefficients of a polynomial interpolant, while `polyval()` uses the vectorized Horner scheme of Code 5.2.1.7 to evaluate the polynomial in given points. The names of the functions are borrowed from PYTHON, see `numpy.poly1d`.

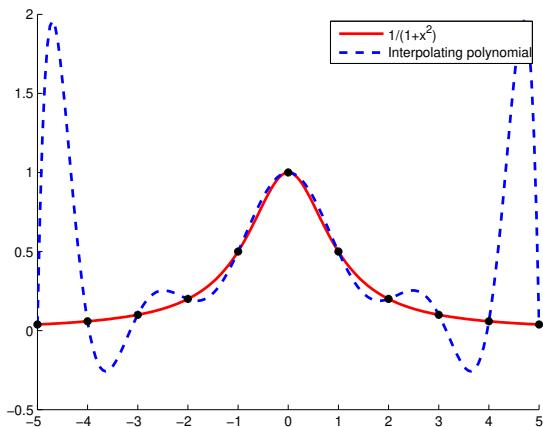


Fig. 188

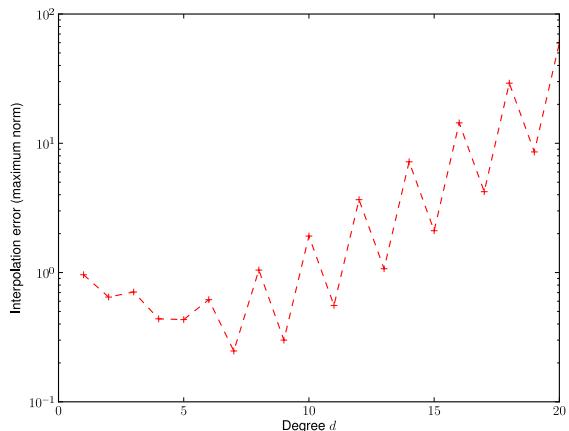
Interpolating polynomial,  $n = 10$ 

Fig. 189

Approximate  $\|f - L_{T_n}f\|_\infty$  on  $[-5, 5]$ Note: approximation of  $\|f - L_{T_n}f\|_\infty$  by *sampling* in 1000 equidistant points.Observation: Strong oscillations of  $L_{T_n}f$  near the endpoints of the interval, which seem to cause

$$\|f - L_{T_n}f\|_{L^\infty([-5, 5])} \xrightarrow{n \rightarrow \infty} \infty.$$

Though polynomials possess great power to approximate functions, see Thm. 6.1.1.11 and Thm. 6.1.1.2, here polynomial interpolants fail completely. Approximation theorists even discovered the following “negative result”:

### Theorem 6.1.2.12. Divergent polynomial interpolants

Given a sequence of meshes of increasing size  $\{\mathcal{T}_n\}_{n=1}^\infty$ ,  $\mathcal{T}_j = \{t_0^{(n)}, \dots, t_n^{(n)}\} \subset [a, b]$ ,  $a \leq t_0^{(n)} < t_2^{(j)} < \dots < t_n^{(n)} \leq b$ , there exists a continuous function  $f$  such that the sequence of interpolating polynomials  $(L_{\mathcal{T}_n}f)_{n=1}^\infty$  does not converge to  $f$  uniformly as  $n \rightarrow \infty$ .

Now we aim to establish bounds for the supremum norm of the interpolation error of Lagrangian interpolation similar to the result of Thm. 6.1.1.11.

### Theorem 6.1.2.13. Representation of interpolation error [DR08, Thm. 8.22], [Han02, Thm. 37.4]

We consider  $f \in C^{n+1}(I)$  and the Lagrangian interpolation approximation scheme ( $\rightarrow$  Def. 6.1.2.1) for a node set  $\mathcal{T} := \{t_0, \dots, t_n\} \subset I$ . Then,

for every  $t \in I$  there exists a  $\tau_t \in [\min\{t, t_0, \dots, t_n\}, \max\{t, t_0, \dots, t_n\}]$  such that

$$f(t) - L_{\mathcal{T}}(f)(t) = \frac{f^{(n+1)}(\tau_t)}{(n+1)!} \cdot \prod_{j=0}^n (t - t_j). \quad (6.1.2.14)$$

*Proof.* Write  $w_{\mathcal{T}}(t) := \prod_{j=0}^n (t - t_j) \in \mathcal{P}_{n+1}$  and fix  $t \in I \setminus \mathcal{T}$ .

$$t \neq t_j \Rightarrow w_{\mathcal{T}}(t) \neq 0 \Rightarrow \exists c = c(t) \in \mathbb{R}: f(t) - L_{\mathcal{T}}(f)(t) = c(t)w_{\mathcal{T}}(t) \quad (6.1.2.15)$$

Consider the *auxiliary function*  $\varphi(x) := f(x) - \mathcal{L}_T(f)(x) - cw_T(x)$  that has  $n+2$  distinct zeros  $t_0, \dots, t_n, t$ . By iterated application of the **mean value theorem** [Str09, Thm .5.2.1]/**Rolle's theorem**

$$f \in C^1([a, b]), f(a) = f(b) = 0 \Rightarrow \exists \xi \in ]a, b[ : f'(\xi) = 0,$$

to higher and higher derivatives, we conclude that

$$\varphi^{(m)} \text{ has } n+2-m \text{ distinct zeros in } I.$$

$$\stackrel{m:=n+1}{\Rightarrow} \exists \tau_t \in I : \varphi^{(n+1)}(\tau_t) = f^{(n+1)}(\tau_t) - c(n+1)! = 0.$$

This fixes the value of  $c = \frac{f^{(n+1)}(\tau_t)}{(n+1)!}$  and by (6.1.2.15) this amounts to the assertion of the theorem.  $\square$

**Remark 6.1.2.16 (Explicit representation of error of polynomial interpolation)** The previous theorem can be refined:

**Lemma 6.1.2.17. Error representation for polynomial Lagrange interpolation**

For  $f \in C^{n+1}(I)$  let  $\mathbf{l}_T \in \mathcal{P}_n$  stand for the unique Lagrange interpolant ( $\rightarrow$  Thm. 5.2.2.7) of  $f$  in the node set  $T := \{t_0, \dots, t_n\} \subset I$ . Then for all  $t \in I$  the interpolation error is

$$\begin{aligned} f(t) - \mathbf{l}_T(f)(t) &= \int_0^1 \int_0^{\tau_1} \cdots \int_0^{\tau_{n-1}} \int_0^{\tau_n} f^{(n+1)}(t_0 + \tau_1(t_1 - t_0) + \cdots \\ &\quad + \tau_n(t_n - t_{n-1}) + \tau(t - t_n)) d\tau d\tau_n \cdots d\tau_1 \cdot \prod_{j=0}^n (t - t_j). \end{aligned}$$

*Proof.* By induction on  $n$ , use (5.2.3.9) and the fundamental theorem of calculus [Ran00, Sect. 3.1]:  $\square$   $\square$

**Remark 6.1.2.18 (Error representation for generalized Lagrangian interpolation)** A result analogous to Lemma 6.1.2.17 holds also for general polynomial interpolation with multiple nodes as defined in (5.2.2.15).  $\square$

Lemma 6.1.2.17 provides an *exact formula* (6.1.2.14) for the interpolation error. From it we can derive *estimates* for the supremum norm of the interpolation error on the interval  $I$  as follows:

- ① first bound the right hand side via  $f^{(n+1)}(\tau_t) \leq \|f^{(n+1)}\|_{L^\infty(I)}$ ,
- ② then increase the right hand side further by switching to the maximum (in modulus) w.r.t.  $t$  (the resulting bound does no longer depend on  $t!$ ),
- ③ and, finally, take the maximum w.r.t.  $t$  on the left of  $\leq$ .

This yields the following **interpolation error estimate** for degree- $n$  Lagrange interpolation on the node set  $\{t_0, \dots, t_n\}$ :

$$\text{??} \Rightarrow \|f - \mathcal{L}_T f\|_{L^\infty(I)} \leq \frac{\|f^{(n+1)}\|_{L^\infty(I)}}{(n+1)!} \max_{t \in I} |(t - t_0) \cdots (t - t_n)|. \quad (6.1.2.19)$$

**Remark 6.1.2.20 (Significance of smoothness of interpoland)** The estimate (6.1.2.19) hinges on bounds for (higher) derivatives of the interpoland  $f$ , which, essentially, should belong to  $C^{n+1}(I)$ . The same can be said about the estimate of Thm. 6.1.1.11.

This reflects a general truth about estimates of norms of the interpolation error:

Quantitative interpolation error estimates rely on smoothness!

↓

**EXAMPLE 6.1.2.21 (Error of polynomial interpolation)** **Ex. 6.1.2.4 cnt'd** Now we are in a position to give a theoretical explanation for exponential convergence observed for polynomial interpolation of  $f(t) = \sin(t)$  on equidistant nodes: by Lemma 6.1.2.17 and (6.1.2.19)

$$\begin{aligned} \|f^{(k)}\|_{L^\infty(I)} &\leq 1, \quad \Rightarrow \quad \|f - p\|_{L^\infty(I)} \leq \frac{1}{(1+n)!} \max_{t \in I} |(t-0)(t-\frac{\pi}{n})(t-\frac{2\pi}{n}) \cdots (t-\pi)| \\ \forall k \in \mathbb{N}_0 \quad &\leq \frac{1}{n+1} \left(\frac{\pi}{n}\right)^{n+1}. \end{aligned}$$

- Uniform asymptotic (even more than) exponential convergence of the interpolation polynomials (independently of the set of nodes  $\mathcal{T}$ ). In fact,  $\|f - p\|_{L^\infty(I)}$  decays even faster than exponential!

↓

**EXAMPLE 6.1.2.22 (Runge's example)** **Ex. 6.1.2.10 cnt'd** How can the blow-up of the interpolation error observed in Ex. 6.1.2.10 be reconciled with Lemma 6.1.2.17?

Here  $f(t) = \frac{1}{1+t^2}$  allows only to conclude  $|f^{(n)}(t)| = 2^n n! \cdot O(|t|^{-2-n})$  for  $n \rightarrow \infty$ .

- Possible blow-up of error bound from ??  $\rightarrow \infty$  for  $n \rightarrow \infty$ .

↓

**Remark 6.1.2.23 ( $L^2$ -error estimates for polynomial interpolation)** ?? gives error estimates for the  $L^\infty$ -norm. What about other norms?

From Lemma 6.1.2.17 using the Cauchy-Schwarz inequality for integrals

$$\left| \int_a^b f(t)g(t) dt \right|^2 \leq \int_a^b |f(t)|^2 dt \cdot \int_a^b |g(t)|^2 dt, \quad \forall f, g \in C^0([a, b]), \quad (6.1.2.24)$$

repeatedly, we can estimate

$$\begin{aligned} \|f - L_T(f)\|_{L^2(I)}^2 &= \int_I \left| \int_0^1 \int_0^{\tau_1} \cdots \int_0^{\tau_{n-1}} \int_0^{\tau_n} f^{(n+1)}(\dots) d\tau d\tau_n \cdots d\tau_1 \cdot \underbrace{\prod_{j=0}^n (t-t_j)}_{|t-t_j| \leq |I|} \right|^2 dt \\ &\leq \int_I |I|^{2n+2} \underbrace{\text{vol}_{(n+1)}(S_{n+1})}_{=1/(n+1)!} \int_{S_{n+1}} |f^{(n+1)}(\dots)|^2 d\tau dt \\ &= \int_I \frac{|I|^{2n+2}}{(n+1)!} \underbrace{\int_I \text{vol}_n(C_{t,\tau})}_{\leq 2^{(n-1)/2}/n!} |f^{(n+1)}(\tau)|^2 d\tau dt, \end{aligned}$$

where

$$\begin{aligned} S_{n+1} &:= \{\mathbf{x} \in \mathbb{R}^{n+1}: 0 \leq x_n \leq x_{n-1} \leq \dots \leq x_1 \leq 1\} \quad (\text{unit simplex}), \\ C_{t,\tau} &:= \{\mathbf{x} \in S_{n+1}: t_0 + x_1(t_1 - t_0) + \dots + x_n(t_n - t_{n-1}) + x_{n+1}(t - t_n) = \tau\}. \end{aligned}$$

This gives the bound for the  $L^2$ -norm of the error:

$$\Rightarrow \|f - L_T(f)\|_{L^2(I)} \leq \frac{2^{(n-1)/4} |I|^{n+1}}{\sqrt{(n+1)! n!}} \left( \int_I |f^{(n+1)}(\tau)|^2 d\tau \right)^{1/2}. \quad (6.1.2.25)$$

Notice:

$f \mapsto \|f^{(n)}\|_{L^2(I)}$  defines a **seminorm** on  $C^{n+1}(I)$   
(**Sobolev-seminorm**, measure of the smoothness of a function).

Estimates like (6.1.2.25) play a key role in the analysis of numerical methods for solving partial differential equations (→ course “Numerical methods for partial differential equations”).

**Remark 6.1.2.26 (Interpolation error estimates and the Lebesgue constant)** The sensitivity of a (polynomial) interpolation scheme  $I_T : \mathbb{K}^{n+1} \rightarrow C^0(I)$ ,  $T \subset I$  a node set, as introduced in Section 5.2.4 and expressed by the **Lebesgue constant** (→ Lemma 5.2.4.9)

$$\lambda_T := \|I_T\|_{\infty \rightarrow \infty} := \sup_{\mathbf{y} \in \mathbb{R}^{n+1} \setminus \{0\}} \frac{\|I_T(\mathbf{y})\|_{L^\infty(I)}}{\|\mathbf{y}\|_\infty},$$

establishes an important connection between the norms of the interpolation error and of the best approximation error.

We first observe that the polynomial approximation scheme  $L_T$  induced by  $I_T$  preserves polynomials of degree  $\leq n$ :

$$L_T p = I_T([p(t)]_{t \in T}) = p \quad \forall p \in \mathcal{P}_n. \quad (6.1.2.27)$$

Thus, by the triangle inequality, for a generic norm on  $C^0(I)$  and  $\|L_T\|$  designating the associated operator norm of the linear mapping  $L_T$ , cf. (5.2.4.8),

$$\begin{aligned} \|f - L_T f\| &\stackrel{(6.1.2.27)}{=} \|(f - p) - L_T(f - p)\| \leq (1 + \|L_T\|) \|f - p\| \quad \forall p \in \mathcal{P}_n, \\ \Rightarrow \|f - L_T f\| &\leq (1 + \|L_T\|) \underbrace{\inf_{p \in \mathcal{P}_n} \|f - p\|}_{\text{best approximation error}}. \end{aligned} \quad (6.1.2.28)$$

Note that for  $\|\cdot\| = \|\cdot\|_{L^\infty(I)}$ , since  $\|[f(t)]_{t \in T}\|_\infty \leq \|f\|_{L^\infty(I)}$ , we can estimate the operator norm, cf. (5.2.4.8),

$$\|L_T\|_{L^\infty(I) \rightarrow L^\infty(I)} \leq \|I_T\|_{\mathbb{R}^{n+1} \rightarrow L^\infty(I)} = \lambda_T, \quad (6.1.2.29)$$

$$\Rightarrow \|f - L_T f\|_{L^\infty(I)} \leq (1 + \lambda_T) \inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty(I)} \quad \forall f \in C^0(I). \quad (6.1.2.30)$$

Hence, if a bound for  $\lambda_T$  is available, the best approximation error estimate of Thm. 6.1.1.11 immediately yields interpolation error estimates.

**§6.1.2.31 (Interpolation error estimates for analytic interpolants)** Exponential convergence can often be observed for families of Lagrangian approximation schemes, when they are applied to an analytic interpolant.

### Definition 6.1.2.32. Analyticity of a complex valued function

let  $D \subset \mathbb{C}$  be an *open* set in the complex plane. A function  $f : D \rightarrow \mathbb{C}$  is called **analytic/holomorphic** in  $D$ , if  $f \in C^\infty(D)$  and it possesses a convergent Taylor series at every point  $z \in D$ :

$$\forall z \in D: \exists \rho = \rho(z) > 0: f(w) = \sum_{n=0}^{\infty} \frac{f^{(n)}(z)}{n!} (w - z)^n \quad \forall w \in D: |z - w| < \rho(z).$$

We may say that an analytic function locally agrees with a “polynomial of degree  $\infty$ ”, because this is exactly what a convergent **power series**

$$s(w) = \sum_{n=0}^{\infty} a_n (w - z)^n, \quad a_n \in \mathbb{C},$$

represents.

The mathematical area of complex analysis ( $\rightarrow$  course in the BSc program CSE) studies analytic functions. Analyticity gives access to powerful tools provided by complex analysis. One of these tools is the residue theorem.

### Theorem 6.1.2.33. Residue theorem [Rem84, Ch. 13]

Let  $D \subset \mathbb{C}$  be an open set,  $G \subset D$  a closed set contained in  $D$ ,  $\gamma := \partial G$  its (piecewise smooth and oriented) boundary, and  $\Pi$  a finite set contained in the interior of  $G$ .

Then for each function  $f$  that is analytic in  $D \setminus \Pi$  holds

$$\frac{1}{2\pi i} \int_{\gamma} f(z) dz = \sum_{p \in \Pi} \text{res}_p f,$$

where  $\text{res}_p f$  is the **residual** of  $f$  in  $p \in \mathbb{C}$ .

- Note that the integral  $\int_{\gamma}$  in Thm. 6.1.2.33 is a path integral in the complex plane (“contour integral”): If the path of integration  $\gamma$  is described by a parameterization  $\tau \in J \mapsto \gamma(\tau) \in \mathbb{C}$ ,  $J \subset \mathbb{R}$ , then

$$\int_{\gamma} f(z) dz := \int_J f(\gamma(\tau)) \cdot \dot{\gamma}(\tau) d\tau, \tag{6.1.2.34}$$

where  $\dot{\gamma}$  designates the derivative of  $\gamma$  with respect to the parameter, and  $\cdot$  indicates multiplication in  $\mathbb{C}$ . For contour integrals we have the estimate

$$\left| \int_{\gamma} f(z) dz \right| \leq |\gamma| \max_{z \in \gamma} |f(z)|. \tag{6.1.2.35}$$

- $\Pi$  often stands for the set of **poles** of  $f$ , that is, points where “ $f$  attains the value  $\infty$ ”.

The residue theorem is very useful, because there are simple formulas for  $\text{res}_p f$ :

**Lemma 6.1.2.36. Residual formula for quotients**

let  $g$  and  $h$  be complex valued functions that are both analytic in a neighborhood of  $p \in \mathbb{C}$ , and satisfy  $h(p) = 0, h'(p) \neq 0$ . Then

$$\text{res}_p \frac{g}{h} = \frac{g(p)}{h'(p)}.$$

Now we consider a polynomial Lagrangian approximation scheme on the interval  $I := [a, b] \subset \mathbb{R}$ , based on the node set  $\mathcal{T} := \{t_0, \dots, t_n\} \subset I$ .

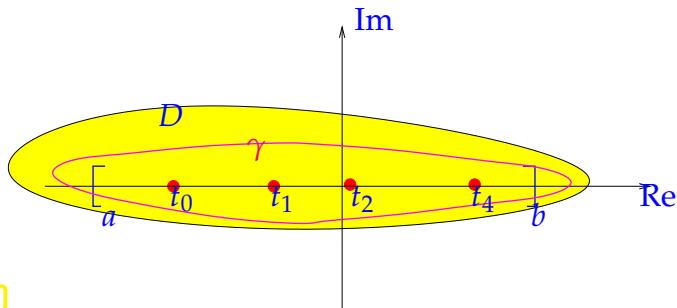


Fig. 190

**Assumption 6.1.2.37. Analyticity of interpoland**

We assume that the interpoland  $f : I \rightarrow \mathbb{C}$  can be extended to a function  $f : D \subset \mathbb{C} \rightarrow \mathbb{C}$ , which is **analytic** ( $\rightarrow$  Def. 6.1.2.32) on the open set  $D \subset \mathbb{C}$  with  $[a, b] \subset D$ .

Key is the following representation of the Lagrange polynomials (5.2.2.4) for node set  $\mathcal{T} = \{t_0, \dots, t_n\}$ :

$$L_j(t) = \prod_{k=0, k \neq j}^n \frac{t - t_k}{t_j - t_k} = \frac{w(t)}{(t - t_j) \prod_{k=0, k \neq j}^n (t_j - t_k)} = \frac{w(t)}{(t - t_j) w'(t_j)}, \quad (6.1.2.38)$$

where  $w(t) = (t - t_0) \cdots (t - t_n) \in \mathcal{P}_{n+1}$ .

Consider the following parameter dependent function  $g_t$ , whose set of poles in  $D$  is  $\Pi = \{t, t_0, \dots, t_n\}$

$$g_t(z) := \frac{f(z)}{(z - t) w(z)}, \quad z \in \mathbb{C} \setminus \Pi, \quad t \in [a, b] \setminus \{t_0, \dots, t_n\}.$$

- $g_t$  is analytic on  $D \setminus \{t, t_0, \dots, t_n\}$  ( $t$  must be regarded as parameter!)
- Apply residue theorem Thm. 6.1.2.33 to  $g_t$  and a closed path of integration  $\gamma \subset D$  winding once around  $[a, b]$ , such that its interior is simply connected, see the magenta curve in Fig. 190:

$$\frac{1}{2\pi i} \int_{\gamma} g_t(z) dz = \text{res}_t g_t + \sum_{j=0}^n \text{res}_{t_j} g_t \stackrel{\text{Lemma 6.1.2.36}}{=} \frac{f(t)}{w(t)} + \sum_{j=0}^n \frac{f(t_j)}{(t_j - t) w'(t_j)}$$

Possible, because all zeros of  $w$  are single zeros!

$$\Rightarrow f(t) = - \sum_{j=1}^n f(t_j) \underbrace{\frac{w(t)}{(t_j - t) w'(t_j)}}_{\substack{-\text{Lagrange polynomial !} \\ \text{polynomial interpolant !}}} + \underbrace{\frac{w(t)}{2\pi i} \int_{\gamma} g_t(z) dz}_{\text{interpolation error !}}. \quad (6.1.2.39)$$

This is another representation formula for the interpolation error, an alternative to that of ?? and Lemma 6.1.2.17. We conclude that for all  $t \in [a, b]$

$$|f(t) - \mathcal{L}_T f(t)| \leq \left| \frac{w(t)}{2\pi i} \int_{\gamma} \frac{f(z)}{(z-t)w(z)} dz \right| \leq \frac{|\gamma|}{2\pi} \frac{\max_{a \leq \tau \leq b} |w(\tau)|}{\min_{z \in \gamma} |w(z)|} \cdot \frac{\max_{z \in \gamma} |f(z)|}{\text{dist}([a, b], \gamma)}. \quad (6.1.2.40)$$

In concrete setting, in order to exploit the estimate (6.1.2.40) to study the  $n$ -dependence of the supremum norm of the interpolation error, we need to know

- an upper bound for  $|w(t)|$  for  $a \leq t \leq b$ ,
- an a lower bound for  $|w(z)|$ ,  $z \in \gamma$ , for a suitable path of integration  $\gamma \subset D$ ,
- a lower bound for the distance of the path  $\gamma$  and the interval  $[a, b]$  in the complex plane.

□

**Remark 6.1.2.41 (Determining the domain of analyticity)** The subset of  $\mathbb{C}$ , where a function  $f$  given by a formula, is analytic can often be determined without computing derivatives using the following consequence of the chain rule:

#### Theorem 6.1.2.42. Composition of analytic functions

If  $f : D \subset \mathbb{C} \rightarrow \mathbb{C}$  and  $g : U \subset \mathbb{C} \rightarrow \mathbb{C}$  are analytic in the open sets  $D$  and  $U$ , respectively, then their composition  $f \circ g$  is analytic in  $\{z \in U : g(z) \in D\}$ .

This can be combined with the following facts:

- Polynomials,  $\exp(z)$ ,  $\sin(z)$ ,  $\cos(z)$ ,  $\sinh(z)$ ,  $\cosh(z)$  are analytic on  $\mathbb{C}$  (entire functions).
- Rational functions (quotients of polynomials) are analytic everywhere except in the zeros of their denominator.
- The square root  $z \rightarrow \sqrt{z}$  is analytic in  $\mathbb{C} \setminus [-\infty, 0]$ .

For example, according to these rules the function  $f(t) = (1+t^2)^{-1}$  can be extended to an analytic function on  $\mathbb{C} \setminus \{-i, i\}$ . If  $\mathbf{A} \in \mathbb{C}^{n,n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ ,  $n \in \mathbb{N}$ , then  $z \mapsto (\mathbf{A} + z\mathbf{I})^{-1}\mathbf{b} \in \mathbb{C}^n$  is (componentwise) analytic in  $\mathbb{C} \setminus \sigma(\mathbf{A})$ , where  $\sigma(\mathbf{A})$  denotes the set of eigenvalues of  $\mathbf{A}$  (the spectrum). □

### 6.1.3 Chebychev Interpolation

As pointed out in § 6.0.0.6, when we build approximation schemes from interpolation schemes, we have the extra freedom to choose the sampling points (= interpolation nodes). Now, based on the insight into the structure of the interpolation error gained from ??, we seek to choose “optimal” sampling points. They will give rise to the so-called Chebychev polynomial approximation schemes, also known as Chebychev interpolation.

### 6.1.3.1 Motivation and Definition

Setting:

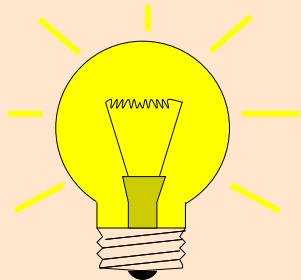
- ◆ Without loss of generality ( $\rightarrow$  Rem. 6.1.1.14):  $I = [-1, 1]$ ,
- ◆ interpoland  $f : I \rightarrow \mathbb{R}$  at least continuous,  $f \in C^0(I)$ ,
- ◆ set of interpolation nodes  $\mathcal{T} := \{-1 \leq t_0 < t_1 < \dots < t_{n-1} < t_n \leq 1\}, n \in \mathbb{N}$ .

Recall ??:

$$\|f - L_{\mathcal{T}} f\|_{L^\infty(I)} \leq \frac{1}{(n+1)!} \|f^{(n+1)}\|_{L^\infty(I)} \|w\|_{L^\infty(I)},$$

with **nodal polynomial**  $w(t) := (t - t_0) \cdots (t - t_n)$ .

#### Optimal choice of interpolation nodes independent of interpoland



Idea: choose nodes  $t_0, \dots, t_n$  such that  $\|w\|_{L^\infty(I)}$  is minimal!

This is equivalent to finding a polynomial  $q \in \mathcal{P}_{n+1}$

- ◆ with leading coefficient = 1,
- ◆ such that it minimizes the norm  $\|q\|_{L^\infty(I)}$ .

Then choose nodes  $t_0, \dots, t_n$  as **zeros** of  $q$ .  
(caution: all  $t_j$  must lie in  $I$ !)

**Remark 6.1.3.2 (A priori and a posteriori choice of optimal interpolation nodes)** We stress that we aim for an “optimal” *a priori* choice of interpolation nodes, a choice that is made **before** any information about the interpoland becomes available.

Of course, an *a posteriori* choice based on information gleaned from evaluations of the interpoland  $f$  may yield much better interpolants (in the sense of smaller norm of the interpolation error). Many modern algorithms employ this **a posteriori adaptive approximation policy**, but this chapter will not cover them.

However, see Section 7.5 for the discussion of an a posteriori adaptive approach for the numerical approximation of definite integrals.  $\square$

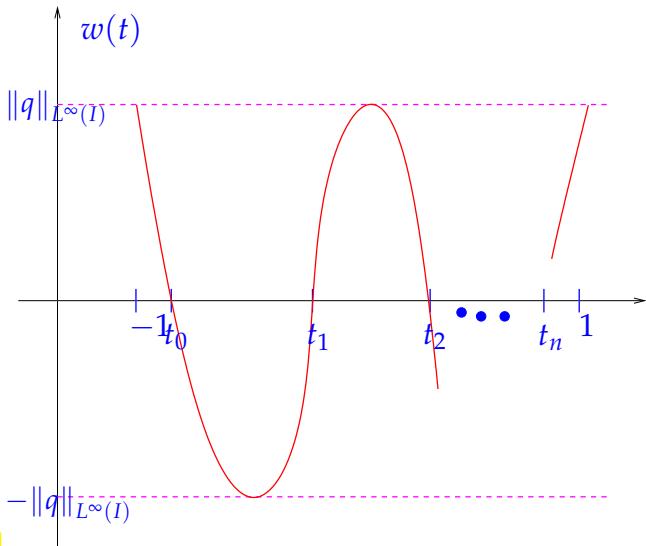


Fig. 191

Requirements on  $q$  (by heuristic reasoning)

Optimal polynomials  $q$  will exist, but they seem to be elusive. First, we develop some insights in how they must “look like”:

- If  $t^*$  is an extremal point of  $q$   
     $\rightarrow |q(t^*)| = \|q\|_{L^\infty(I)}$ ,
  - $q$  has  $n+1$  zeros in  $I$  (\*),
  - $|q(-1)| = |q(1)| = \|q\|_{L^\infty(I)}$ .
- $q$  has  $n+2$  extrema in  $[-1, 1]$

(\*) is motivated by an indirect argument:

If  $q(t) = (t - t_0) \cdots (t - t_{n+1})$  with  $t_0 < -1$ , then

$$|p(t)| := |(t+1)(t-t_1) \cdots (t-t_{n+1})| < |q(t)| \quad \forall t \in I \quad (\text{why?}),$$

which contradicts the minimality property of  $q$ . Same argument for  $t_0 > 1$ . The reasonings leading to the above heuristic demands will be elaborated in the proof of Thm. 6.1.3.8.

Are there polynomials satisfying these requirements? If so, do they allow a simple characterization?

### Definition 6.1.3.3. Chebychev polynomials → [Han02, Ch. 32]

The  $n^{\text{th}}$  Chebychev polynomial is  $T_n(t) := \cos(n \arccos t)$ ,  $-1 \leq t \leq 1, n \in \mathbb{N}$ .

The next result confirms that the  $T_n$  are polynomials, indeed.

### Theorem 6.1.3.4. 3-term recursion for Chebychev polynomials → [Han02, (32.2)]

The function  $T_n$  defined in Def. 6.1.3.3 satisfy the 3-term recursion

$$T_{n+1}(t) = 2t T_n(t) - T_{n-1}(t), \quad T_0 \equiv 1, \quad T_1(t) = t, \quad n \in \mathbb{N}. \quad (6.1.3.5)$$

*Proof.* Just use the trigonometric identity  $\cos(n+1)x = 2\cos nx \cos x - \cos(n-1)x$  with  $\cos x = t$ . □

The theorem implies:

- $T_n \in \mathcal{P}_n$ ,
- their leading coefficients are equal to  $2^{n-1}$ ,
- the  $T_n$  are linearly independent,
- $\{T_j\}_{j=0}^n$  is a basis of  $\mathcal{P}_n = \text{Span}\{T_0, \dots, T_n\}, n \in \mathbb{N}_0$ .

See Code 6.1.3.6 for algorithmic use of the 3-term recursion (6.1.3.5).

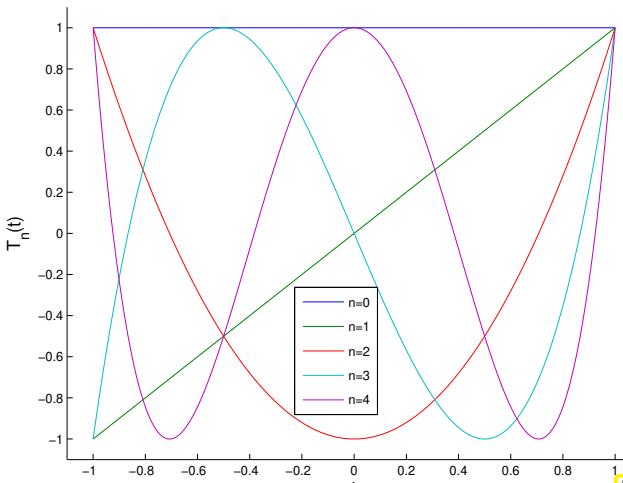
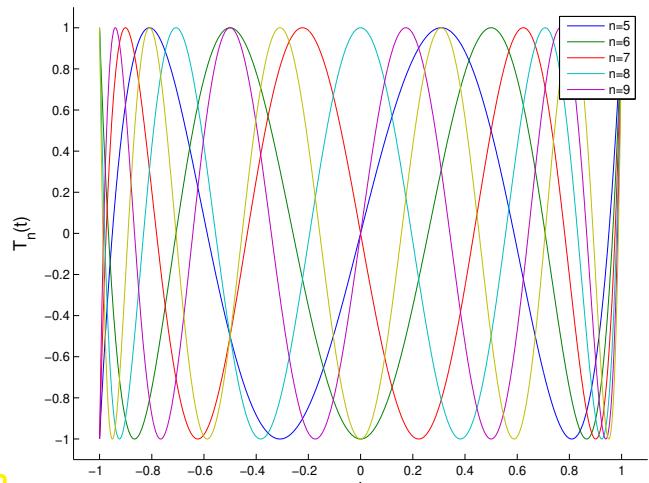


Fig. 192

Chebychev polynomials  $T_0, \dots, T_4$ Chebychev polynomials  $T_5, \dots, T_9$ 

### C++ code 6.1.3.6: Efficient evaluation of Chebychev polynomials up to a certain degree

```

2 // Computes the values of the Chebychev polynomials  $T_0, \dots, T_d$ ,  $d \geq 2$ 
3 // at points passed in  $x$  using the 3-term recursion (6.1.3.5).
4 // The values  $T_k(x_j)$ , are returned in row  $k+1$  of  $V$ .
5 void chebpoltmult(const int d, const RowVectorXd& x, MatrixXd& V) {
6     const unsigned n = x.size();
7     V = MatrixXd::Ones(d + 1, n); //  $T_0 \equiv 1$ 

```

```

8   V.block(1, 0, 1, n) = x; //  $T_1(x) = x$ 
9   for (int k = 1; k < d; ++k) {
10      const RowVectorXd p = V.block(k, 0, 1, n); //  $p = T_k$ 
11      const RowVectorXd q = V.block(k - 1, 0, 1, n); //  $q = T_{k-1}$ 
12      V.block(k + 1, 0, 1, n) = 2*x.cwiseProduct(p) - q; // 3-term recursion
13   }
14 }
```

From Def. 6.1.3.3 we conclude that  $T_n$  attains the values  $\pm 1$  in its extrema with alternating signs, thus matching our heuristic demands:

$$|T_n(\bar{t}_k)| = 1 \Leftrightarrow \exists k = 0, \dots, n: \bar{t}_k = \cos \frac{k\pi}{n}, \quad \|T_n\|_{L^\infty([-1,1])} = 1. \quad (6.1.3.7)$$

What is still open is the validity of the heuristics guiding the choice of the optimal nodes. The next fundamental theorem will demonstrate that, after scaling, the  $T_n$  really supply polynomials on  $[-1, 1]$  with fixed leading coefficient and minimal supremum norm.

**Theorem 6.1.3.8. Minimax property of the Chebychev polynomials** [DH03, Section 7.1.4.], [Han02, Thm. 32.2]

The polynomials  $T_n$  from Def. 6.1.3.3 minimize the supremum norm in the following sense:

$$\|T_n\|_{L^\infty([-1,1])} = \inf\{\|p\|_{L^\infty([-1,1])}: p \in \mathcal{P}_n, p(t) = 2^{n-1}t^n + \dots\}, \quad \forall n \in \mathbb{N}.$$

*Proof.* (indirect) Assume

$$\exists q \in \mathcal{P}_n, \text{ leading coefficient } = 2^{n-1}: \|q\|_{L^\infty([-1,1])} < \|T_n\|_{L^\infty([-1,1])}. \quad (6.1.3.9)$$

- $(T_n - q)(x) > 0$  in local maxima of  $T_n$   
 $(T_n - q)(x) < 0$  in local minima of  $T_n$

From knowledge of local extrema of  $T_n$ , see (6.1.3.7):

$T_n - q$  changes sign at least  $n + 1$  times

$\Rightarrow T_n - q$  has at least  $n$  zeros

$T_n - q \equiv 0$ , because  $T_n - q \in \mathcal{P}_{n-1}$  (same leading coefficient!)

This cannot be reconciled with the properties (6.1.3.9) of  $q$  and, thus, leads to a contradiction.  $\square$

The zeros of  $T_n$  are  $t_k = \cos\left(\frac{2k+1}{2n}\pi\right), \quad k = 0, \dots, n-1$ . (6.1.3.10)

To see this, notice

$$T_n(t) = 0 \stackrel{\text{zeros of } \cos}{\Leftrightarrow} n \arccos t \in (2\mathbb{Z} + 1)\frac{\pi}{2}$$

$$\arccos \in [0, \pi] \Leftrightarrow t \in \left\{\cos\left(\frac{2k+1}{n}\frac{\pi}{2}\right), k = 0, \dots, n-1\right\}.$$

Thus, we have identified the  $t_k$  from (6.1.3.10) as optimal interpolation nodes for a Lagrangian approximation scheme. The  $t_k$  are known as **Chebychev nodes**. Their distribution in  $[-1, 1]$  is plotted in Fig. 194, a geometric construction is indicated in Fig. 195.

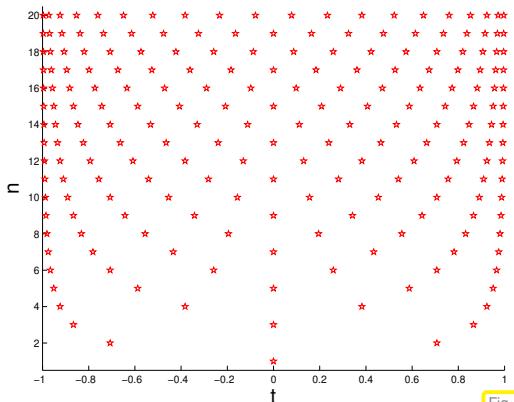


Fig. 194

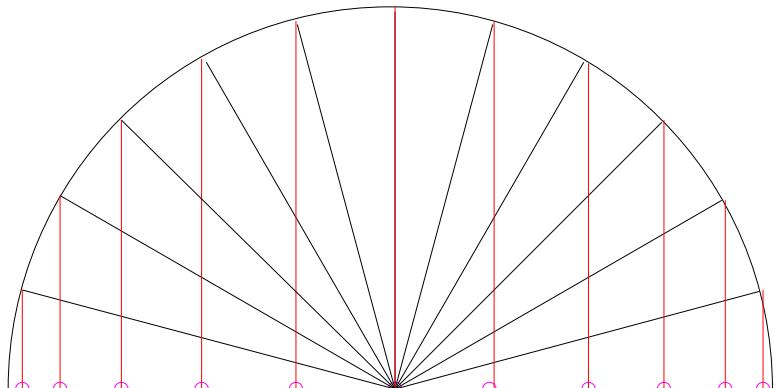


Fig. 195

When we use Chebychev nodes for polynomial interpolation we call the resulting Lagrangian approximation scheme **Chebychev interpolation**. On the interval  $[-1, 1]$  it is characterized by:

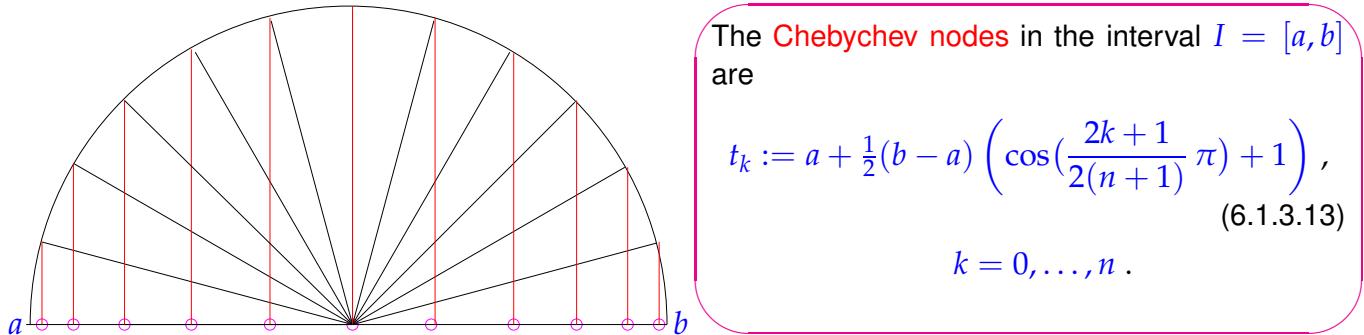
- “optimal” interpolation nodes  $\mathcal{T} = \left\{ \cos\left(\frac{2k+1}{2(n+1)}\pi\right), k = 0, \dots, n \right\}$ ,
- $w(t) = (t - t_0) \cdots (t - t_n) = 2^{-n} T_{n+1}(t)$ ,  $\|w\|_{L^\infty(I)} = 2^{-n}$ , with leading coefficient 1.

Then, by ??, we immediately get an interpolation error estimate for Chebychev interpolation of  $f \in C^{n+1}([-1, 1])$ :

$$\|f - I_T(f)\|_{L^\infty([-1, 1])} \leq \frac{2^{-n}}{(n+1)!} \|f^{(n+1)}\|_{L^\infty([-1, 1])}. \quad (6.1.3.11)$$

**Remark 6.1.3.12 (Chebychev polynomials on arbitrary interval)** Following the recipe of Rem. 6.1.1.14 Chebychev interpolation on an arbitrary interval  $[a, b]$  can immediately be defined. The same polynomial Lagrangian approximation scheme is obtained by transforming the Chebychev nodes (6.1.3.10) from  $[-1, 1]$  to  $[a, b]$  using the unique **affine transformation** (6.1.1.15):

$$\hat{t} \in [-1, 1] \mapsto t := a + \frac{1}{2}(\hat{t} + 1)(b - a) \in [a, b].$$



Estimates for the Chebychev interpolation error on  $[a, b]$  are easily derived from (6.1.3.11):

$$p \in \mathcal{P}_n \wedge p(t_j) = f(t_j) \Leftrightarrow \hat{p} \in \mathcal{P}_n \wedge \hat{p}(\hat{t}_j) = \hat{f}(\hat{t}_j).$$

With transformation formula for the integrals &  $\frac{d^n \hat{f}}{dt^n}(\hat{t}) = (\frac{1}{2}|I|)^n \frac{d^n f}{dt^n}(t)$ :

$$\|f - I_T(f)\|_{L^\infty(I)} = \|\hat{f} - I_{\hat{T}}(\hat{f})\|_{L^\infty([-1, 1])} \leq \frac{2^{-n}}{(n+1)!} \left\| \frac{d^{n+1} \hat{f}}{d\hat{t}^{n+1}} \right\|_{L^\infty([-1, 1])}$$

$$\leq \frac{2^{-2n-1}}{(n+1)!} |I|^{n+1} \|f^{(n+1)}\|_{L^\infty(I)}. \quad (6.1.3.14)$$

### 6.1.3.2 Chebychev Interpolation Error Estimates

**EXAMPLE 6.1.3.15 (Polynomial interpolation: Chebychev nodes versus equidistant nodes)** We consider Runge's function  $f(t) = \frac{1}{1+t^2}$ , see Ex. 6.1.2.10, and compare polynomial interpolation based on uniformly spaced nodes and Chebychev nodes in terms of behavior of interpolants.

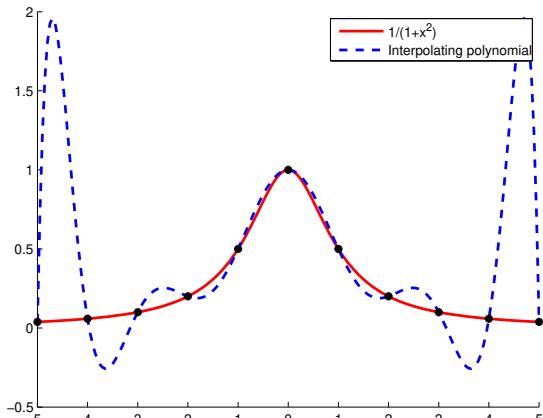


Fig. 196

Equidistant nodes

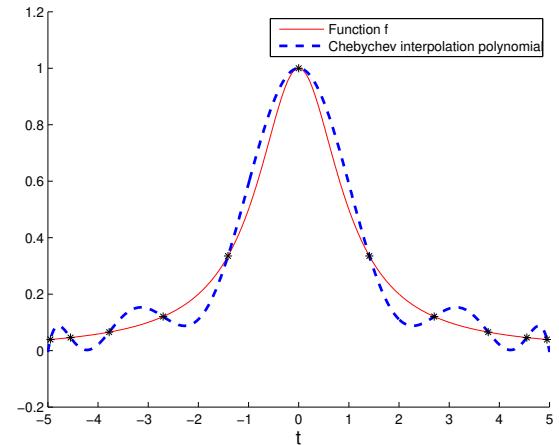


Fig. 197

Chebychev nodes

We observe that the Chebychev nodes cluster at the endpoints of the interval, which successfully suppresses the huge oscillations haunting equidistant interpolation there. □

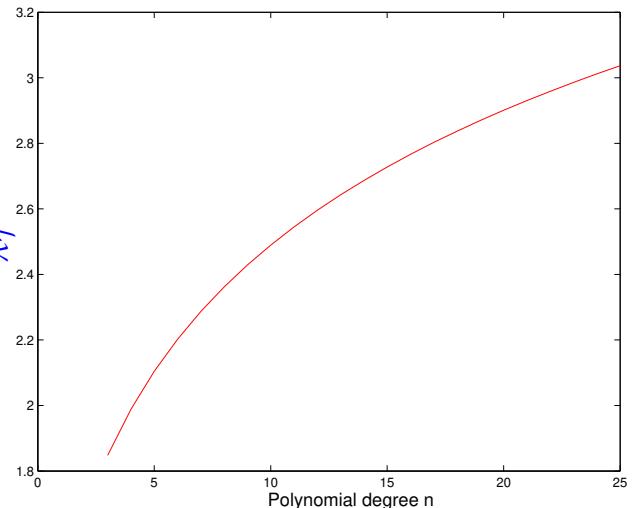
**Remark 6.1.3.16 (Lebesgue Constant for Chebychev nodes)** → [Section 5.2.4](#)) We saw in Rem. 5.2.4.12 that the Lebesgue constant  $\lambda_T$  that measures the sensitivity of a polynomial interpolation scheme, blows up exponentially with increasing number of *equispaced* interpolation nodes. In stark contrast  $\lambda_T$  grows only *logarithmically* in the number of Chebychev nodes.

More precisely, sophisticated theory [CB95; Ver86; Ver90] supplies the bound

$$\lambda_T \leq \frac{2}{\pi} \log(1+n) + 1.$$

$$(6.1.3.17)$$

Measured Lebesgue constant for Chebychev nodes based on approximate evaluation of (5.2.4.10) by sampling. ▷



Combining (6.1.3.17) with the general estimate from Rem. 6.1.2.26

$$\|f - L_T f\|_{L^\infty(I)} \leq (1 + \lambda_T) \inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty(I)} \quad \forall f \in C^0(I), \quad (6.1.2.30)$$

and the bound for the best approximation error for polynomials from Thm. 6.1.1.11

$$\inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty([-1,1])} \leq (1 + \pi^2/2)^r \frac{(n-r)!}{n!} \|f^{(r)}\|_{L^\infty([-1,1])},$$

we end up with a bound for the supremum norm of the interpolation error in the case of Chebychev interpolation on  $[-1, 1]$

$$\|f - L_T f\|_{L^\infty([-1,1])} \leq (2/\pi \log(1+n) + 2)(1 + \pi^2/2)^r \frac{(n-r)!}{n!} \|f^{(r)}\|_{L^\infty([-1,1])}. \quad (6.1.3.18)$$

□

**EXPERIMENT 6.1.3.19 (Chebychev interpolation errors)** Now we empirically investigate the behavior of norms of the interpolation error for Chebychev interpolation and functions with different (smoothness) properties as we increase the number of interpolation nodes.

In the experiments, for  $I = [a, b]$  we set  $x_l := a + \frac{b-a}{N}l$ ,  $l = 0, \dots, N$ ,  $N = 1000$ , and we approximate the norms of the interpolation error as follows ( $p \hat{=} \text{interpolating polynomial}$ ):

$$\|f - p\|_\infty \approx \max_{0 \leq l \leq N} |f(x_l) - p(x_l)| \quad (6.1.3.20)$$

$$\|f - p\|_2^2 \approx \frac{b-a}{2N} \sum_{0 \leq l < N} (|f(x_l) - p(x_l)|^2 + |f(x_{l+1}) - p(x_{l+1})|^2) \quad (6.1.3.21)$$

- ①  $f(t) = (1+t^2)^{-1}$ ,  $I = [-5, 5]$  (see Ex. 6.1.2.10): analytic in a neighborhood of  $I$ .  
Interpolation with  $n = 10$  Chebychev nodes (plot on the left).

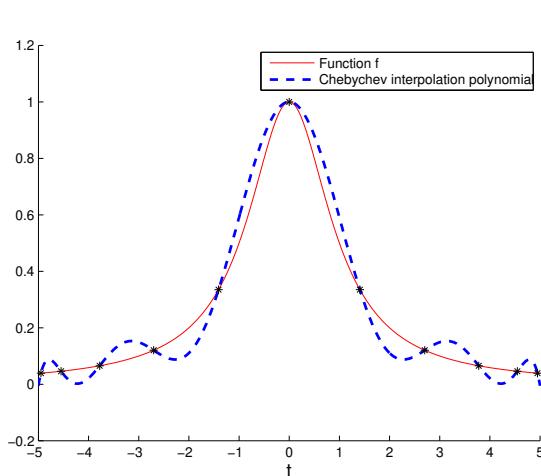
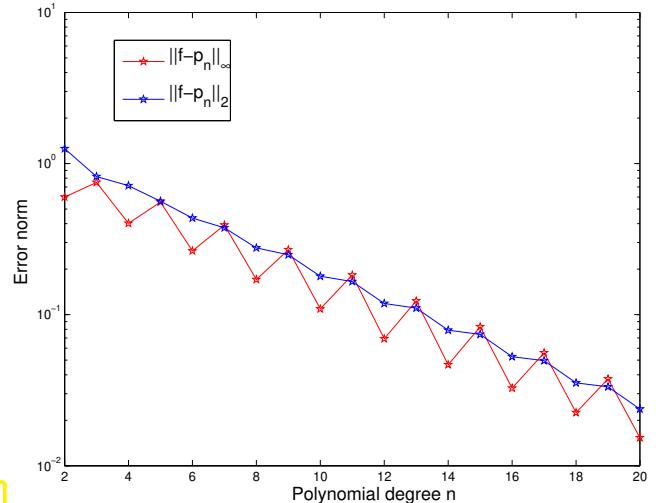


Fig. 198



Notice: exponential convergence (→ Def. 6.1.2.7) of the Chebychev interpolation:

$$p_n \rightarrow f, \quad \|f - L_n f\|_{L^\infty([-5,5])} \approx 0.8^n$$

- ②  $f(t) = \max\{1 - |t|, 0\}$ ,  $I = [-2, 2]$ ,  $n = 10$  nodes (plot on the left).

Now  $f \in C^0(I)$  but  $f \notin C^1(I)$ .

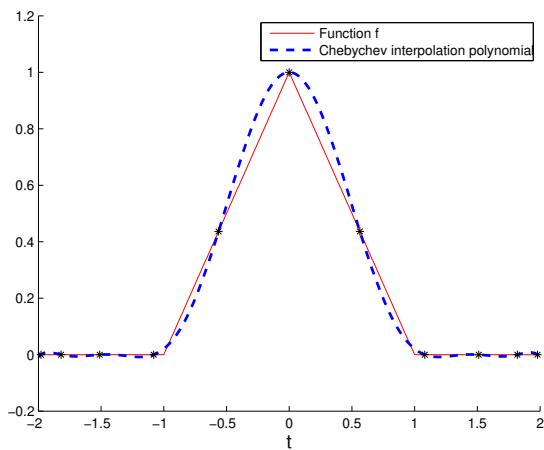


Fig. 199

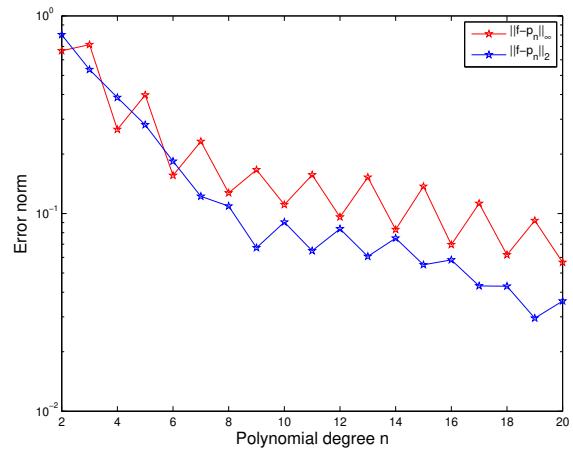


Fig. 200

From the doubly logarithmic plot we conclude →

- no exponential convergence
- algebraic convergence (?)

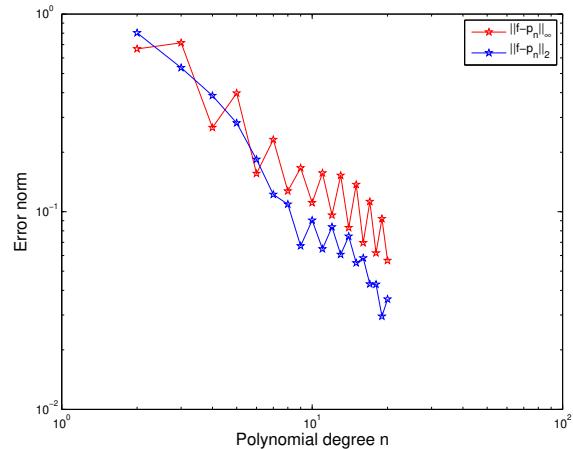


Fig. 201

$$\textcircled{3} \quad f(t) = \begin{cases} \frac{1}{2}(1 + \cos \pi t) & |t| < 1 \\ 0 & 1 \leq |t| \leq 2 \end{cases} \quad I = [-2, 2], \quad n = 10 \quad (\text{plot on the left}).$$

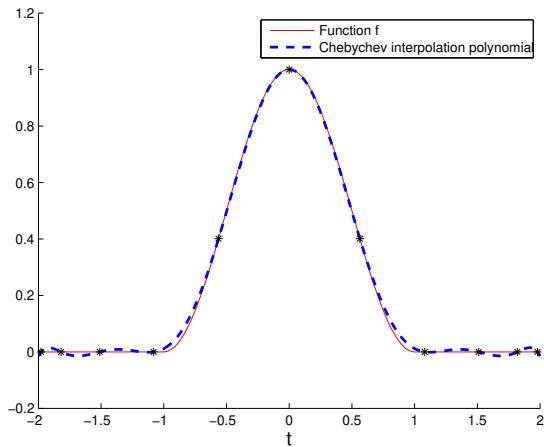


Fig. 202

Notice: only (vaguely) algebraic convergence.

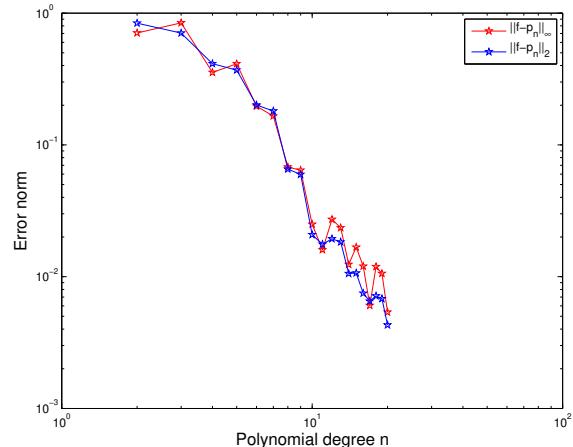


Fig. 203

Summary of observations, cf. Rem. 6.1.2.9:

- ◆ Essential role of smoothness of  $f$ : slow convergence of approximation error of the Cheychev interpolant if  $f$  enjoys little smoothness, cf. also (6.1.2.19),
- ◆ for analytic  $f \in C^\infty$  (→ Def. 6.1.2.32) the approximation error of the Cheychev interpolant seems to decay to zero exponentially in the polynomial degree  $n$ .

**Remark 6.1.3.22 (Chebychev interpolation of analytic functions)** Assuming that the interpoland  $f$  possesses an analytic extension to a complex neighborhood  $D$  of  $[-1, 1]$ , we now apply the theory of § 6.1.2.31 to bound the supremum norm of the Chebychev interpolation error of  $f$  on  $[-1, 1]$ .

To convert

$$|f(t) - \mathcal{L}_T f(t)| \leq \left| \frac{w(t)}{2\pi i} \int_{\gamma} \frac{f(z)}{(z-t)w(z)} dz \right| \leq \frac{|\gamma|}{2\pi} \frac{\max_{a \leq \tau \leq b} |w(\tau)|}{\min_{z \in \gamma} |w(z)|} \cdot \frac{\max_{z \in \gamma} |f(z)|}{\text{dist}([a, b], \gamma)}, \quad (6.1.2.40)$$

as obtained in § 6.1.2.31, into a more concrete estimate, we have to study the behavior of

$$w_n(t) = (t - t_0)(t - t_1) \cdots (t - t_n), \quad t_k = \cos\left(\frac{2k+1}{2n+2}\pi\right), \quad k = 0, \dots, n,$$

where the  $t_k$  are the Chebychev nodes according to (6.1.3.13). They are the zeros of the Chebychev polynomial ( $\rightarrow$  Def. 6.1.3.3) of degree  $n+1$ . Since  $w$  has leading coefficient 1, we conclude  $w = 2^{-n}T_{n+1}$ , and

$$\max_{-1 \leq t \leq 1} |w(t)| \leq 2^{-n}. \quad (6.1.3.23)$$

Next, we fix a suitable path  $\gamma \subset D$  for integration: For a constant  $\rho > 1$  we set

$$\begin{aligned} \gamma &:= \{z = \cos(\theta - i \log \rho), 0 \leq \theta \leq 2\pi\} \\ &= \left\{ z = \frac{1}{2} (\exp(i(\theta - i \log \rho)) + \exp(-i(\theta - i \log \rho))), 0 \leq \theta \leq 2\pi \right\} \\ &= \left\{ z = \frac{1}{2} (\rho e^{i\theta} + \rho^{-1} e^{-iu\theta}), 0 \leq \theta \leq 2\pi \right\} \\ &= \left\{ z = \frac{1}{2} (\rho + \rho^{-1}) \cos \theta + i \frac{1}{2} (\rho - \rho^{-1}) \sin \theta, 0 \leq \theta \leq 2\pi \right\} \end{aligned}$$

Thus, we see that  $\gamma$  is an ellipse with foci  $\pm 1$ , large axis  $\frac{1}{2}(\rho + \rho^{-1}) > 1$  and small axis  $\frac{1}{2}(\rho - \rho^{-1}) > 0$ .

Elliptical integration contours for different values of  $\rho$

▷

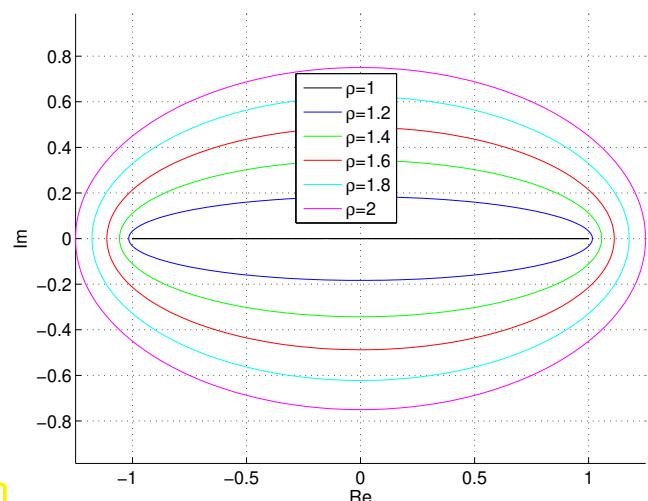


Fig. 204

Appealing to geometric evidence, we find  $\text{dist}(\gamma, [-1, 1]) = \frac{1}{2}(\rho + \rho^{-1}) - 1$ , which gives another term in (6.1.2.40).

The rationale for choosing this particular integration contour is that the  $\cos$  in its definition nicely cancels

the  $\arccos$  in the formula for the Chebychev polynomials. This lets us compute ( $s := n + 1$ )

$$\begin{aligned}
 |T_s(\cos(\theta - i \log \rho))|^2 &= |\cos(s(\theta - i \log \rho))|^2 \\
 &= \cos(s(\theta - i \log \rho)) \cdot \overline{\cos(s(\theta - i \log \rho))} \\
 &= \frac{1}{4}(\rho^s e^{is\theta} + \rho^{-s} e^{-is\theta})(\rho^s e^{-is\theta} + \rho^{-s} e^{is\theta}) \\
 &= \frac{1}{4}(\rho^{2s} + \rho^{-2s} + e^{2is\theta} + e^{-2is\theta}) \\
 &= \frac{1}{4}(\rho^s - \underbrace{\rho^{-s}}_{<1})^2 + \underbrace{\frac{1}{4}(e^{is\theta} + e^{-is\theta})^2}_{=\cos^2(s\theta)\geq 0} \geq \frac{1}{4}(\rho^s - 1)^2,
 \end{aligned}$$

for all  $0 \leq \theta \leq 2\pi$ , which provides a lower bound for  $|w_n|$  on  $\gamma$ . Plugging all these estimates into (6.1.2.40) we arrive at

$$\|f - L_T f\|_{L^\infty([-1,1])} \leq \frac{2|\gamma|}{\pi} \frac{1}{(\rho^{n+1} - 1)(\rho + \rho^{-1} - 2)} \cdot \max_{z \in \gamma} |f(z)|. \quad (6.1.3.24)$$

Note that instead on the nodal polynomial  $w$  we have inserted  $T_{n+1}$  into (6.1.2.40), which is a simple multiple. The factor will cancel.

► The supremum norm of the interpolation converges exponentially ( $\rho > 1$ ):

$$\|f - L_T f\|_{L^\infty([-1,1])} = O((2\rho)^{-n}) \quad \text{for } n \rightarrow \infty.$$

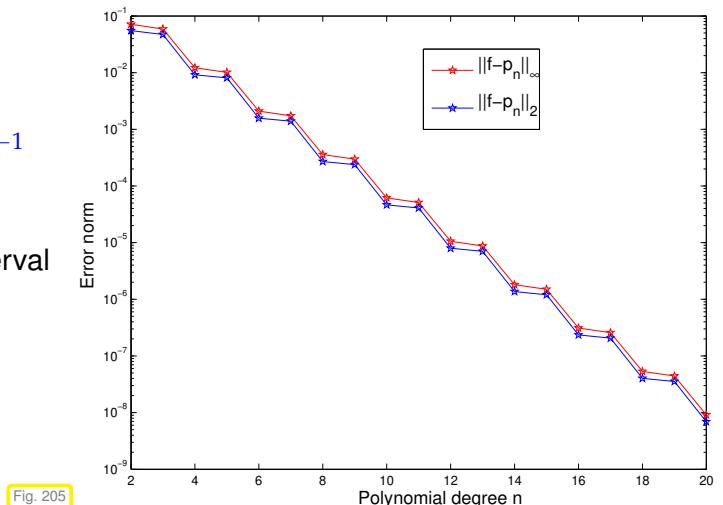
□

### EXPERIMENT 6.1.3.25 (Chebychev interpolation of analytic function → Exp. 6.1.3.19 cnt'd)

Modification: the same function  $f(t) = (1 + t^2)^{-1}$   
on a smaller interval  $I = [-1, 1]$ .

(Faster) exponential convergence than on the interval  $I = [-5, 5]$ :

$$\|f - I_n f\|_{L^2([-1,1])} \approx 0.42^n.$$



Explanation, cf. Rem. 6.1.3.22: for  $I = [-1, 1]$  the poles  $\pm i$  of  $f$  are farther away relative to the size of the interval than for  $I = [-5, 5]$ . □

### 6.1.3.3 Chebychev Interpolation: Computational Aspects

Task: Given: given degree  $n \in \mathbb{N}$ , continuous function  $f : [-1, 1] \mapsto \mathbb{R}$

Sought: efficient representation/evaluation of Chebychev interpolant  $p \in \mathcal{P}_n$  (= polynomial Lagrange interpolant of degree  $\leq n$  in Chebychev nodes (6.1.3.13) on  $[-1, 1]$ ).

More concretely, this boils down to a implementation of the following class:

#### C++ code 6.1.3.26: Definition of class for Chebychev interpolation

```

1  class ChebInterp {
2  private:
3      // various internal data describing Chebychev interpolating polynomial
4      p
5  public:
6      // Constructor taking function f and degree n as arguments
7      template <typename Function>
8      PolyInterp(const Function &f, unsigned int n);
9      // Evaluation operator: y_j = p(x_j), j = 1, ..., m (m "large")
10     void eval(const vector<double> &x, vector <double> &y) const;
11 };

```



Trick: internally represent  $p$  as a linear combination of Chebychev polynomials, a Chebychev expansion:

$$p = \sum_{j=0}^n \alpha_j T_j, \quad \alpha_j \in \mathbb{R}. \quad (6.1.3.27)$$

The representation (6.1.3.27) is always possible, because  $\{T_0, \dots, T_n\}$  is a basis of  $\mathcal{P}_n$ , owing to  $\deg T_n = n$ . The representation is amenable to efficient evaluation and computation by means of special algorithms.

**§6.1.3.28 (Fast evaluation of Chebychev expansion → [Han02, Alg. 32.1])** [Fast evaluation of Chebychev expansion] Let us assume that the Chebychev expansion coefficients  $\alpha_j$  in (6.1.3.27) are given and wonder, how we can efficiently compute  $p(x)$  for some  $x \in \mathbb{R}$ :

Task: Given  $n \in \mathbb{N}$ ,  $x \in \mathbb{R}$ , and the Chebychev expansion coefficients  $\alpha_j \in \mathbb{R}$ ,  $j = 0, \dots, n$ , compute  $p(x)$  with

$$p(x) = \sum_{j=0}^n \alpha_j T_j(x), \quad \alpha_j \in \mathbb{R}. \quad (6.1.3.27)$$



Idea: Use the 3-term recurrence (6.1.3.5)

$$T_j(x) = 2xT_{j-1}(x) - T_{j-2}(x), \quad j = 2, 3, \dots, \quad (6.1.3.5)$$

to design a recursive evaluation scheme.

By means of (6.1.3.5) rewrite (6.1.3.27) as

$$\begin{aligned} p(x) &= \sum_{j=0}^{n-1} \alpha_j T_j(x) + \alpha_n T_n(x) \stackrel{(6.1.3.5)}{=} \sum_{j=0}^{n-1} \alpha_j T_j(x) + \alpha_n (2xT_{n-1}(x) - T_{n-2}(x)) \\ &= \sum_{j=0}^{n-3} \alpha_j T_j(x) + (\alpha_{n-2} - \alpha_n) T_{n-2}(x) + (\alpha_{n-1} + 2x\alpha_n) T_{n-1}(x). \end{aligned}$$

We recover the point value  $p(x)$  as the point value of another polynomial of degree  $n - 1$  with known Chebychev expansion:

$$p(x) = \sum_{j=0}^{n-1} \tilde{\alpha}_j T_j(x) \quad \text{with} \quad \tilde{\alpha}_j = \begin{cases} \alpha_j + 2x\alpha_{j+1} & , \text{ if } j = n - 1, \\ \alpha_j - \alpha_{j+2} & , \text{ if } j = n - 2, \\ \alpha_j & \text{else.} \end{cases} \quad (6.1.3.29)$$

► recursive algorithm, see Code 6.1.3.30.

#### C++ code 6.1.3.30: Recursive evaluation of Chebychev expansion (6.1.3.27)

```

2 // Recursive evaluation of a polynomial  $p = \sum_{j=1}^{n+1} a_j T_{j-1}$  at point x
3 // based on (6.1.3.29)
4 // IN : Vector of coefficients a
5 // evaluation point x
6 // OUT: Value at point x
7 double recclenshaw(const VectorXd& a, const double x) {
8     const VectorXd::Index n = a.size() - 1;
9     if (n == 0) return a(0); // Constant polynomial
10    else if (n == 1) return (x*a(1) + a(0)); // Value  $\alpha_1 * x + \alpha_0$ 
11    else {
12        VectorXd new_a(n);
13        new_a << a.head(n - 2), a(n - 2) - a(n), a(n - 1) + 2*x*a(n);
14        return recclenshaw(new_a, x); // recursion
15    }
16}

```

Non-recursive version: [Clenshaw algorithm](#)

#### C++ code 6.1.3.31: Clenshaw algorithm for evalution of Chebychev expansion (6.1.3.27)

```

2 // Clenshaw algorithm for evaluating  $p = \sum_{j=1}^{n+1} a_j T_{j-1}$ 
3 // at points passed in vector x
4 // IN : a = [ $\alpha_j$ ], coefficients for  $p = \sum_{j=1}^{n+1} \alpha_j T_{j-1}$ 
5 // x = (many) evaluation points
6 // OUT: values  $p(x_j)$  for all j
7 VectorXd clenshaw(const VectorXd& a, const VectorXd& x) {
8     const int n = a.size() - 1; // degree of polynomial
9     MatrixXd d(n + 1, x.size()); // temporary storage for intermediate values
10    for (int c = 0; c < x.size(); ++c) d.col(c) = a;
11    for (int j = n - 1; j > 0; --j) {
12        d.row(j) += 2*x.transpose().cwiseProduct(d.row(j+1)); // see (6.1.3.29)
13        d.row(j-1) -= d.row(j + 1);
14    }
15    return d.row(0) + x.transpose().cwiseProduct(d.row(1));
16}

```



Computational effort :  $\mathcal{O}(nm)$  for evaluation at  $m$  points



**§6.1.3.32 (Computation of Chebychev expansions of interpolants)** Chebychev interpolation is a linear interpolation scheme, see § 5.1.0.15. Thus, the expansion  $\alpha_j$  in (6.1.3.27) can be computed by solving a linear system of equations of the form (5.1.0.17). However, for Chebychev interpolation this linear system can be cast into a very special form, which paves the way for its fast direct solution:

Task: Efficiently compute the Chebychev expansion coefficients  $\alpha_j$  in (6.1.3.27) from the interpolation conditions

$$p(t_k) = f(t_k), \quad k = 0, \dots, n, \quad \text{for } t_k := \cos\left(\frac{2k+1}{2(n+1)}\pi\right). \quad (6.1.3.33)$$

Chebychev nodes

Trick: transformation of  $p$  into a 1-periodic function, which turns out to be a **Fourier sum** (= finite Fourier series):

$$\begin{aligned} q(s) &:= p(\cos 2\pi s) = \sum_{j=0}^n \alpha_j T_j(\cos 2\pi s) \stackrel{\text{Def. 6.1.3.3}}{=} \sum_{j=0}^n \alpha_j \cos(2\pi js) \\ &= \sum_{j=0}^n \frac{1}{2}\alpha_j (\exp(2\pi js) + \exp(-2\pi js)) \quad [\text{by } \cos z = \frac{1}{2}(e^z + e^{-z})] \\ &= \sum_{j=-n}^{n+1} \beta_j \exp(-2\pi js), \quad \text{with } \beta_j := \begin{cases} 0 & , \text{for } j = n+1, \\ \frac{1}{2}\alpha_j & , \text{for } j = 1, \dots, n, \\ \alpha_0 & , \text{for } j = 0, \\ \frac{1}{2}\alpha_{n-j} & , \text{for } j = -n, \dots, -1. \end{cases} \end{aligned} \quad (6.1.3.34)$$

Transformed interpolation conditions (6.1.3.33) for  $q$ :

$$t = \cos(2\pi s) \stackrel{(6.1.3.33)}{\implies} q\left(\frac{2k+1}{4(n+1)}\right) = y_k := f(t_k), \quad k = 0, \dots, n. \quad (6.1.3.35)$$

This is an interpolation problem for *equidistant points on the unit circle* as we have seen them in Section 5.6.3.

Also observe the **symmetry**

$$\begin{aligned} q(s) &= q(1-s) \\ \Downarrow \leftarrow (6.1.3.35) \\ q\left(1 - \frac{2k+1}{4(n+1)}\right) &= y_k, \quad k = 0, \dots, n. \end{aligned}$$

It ensures that the coefficients  $\beta_j$  actually satisfy the constraints implied by their relationship with  $\alpha_j$ .

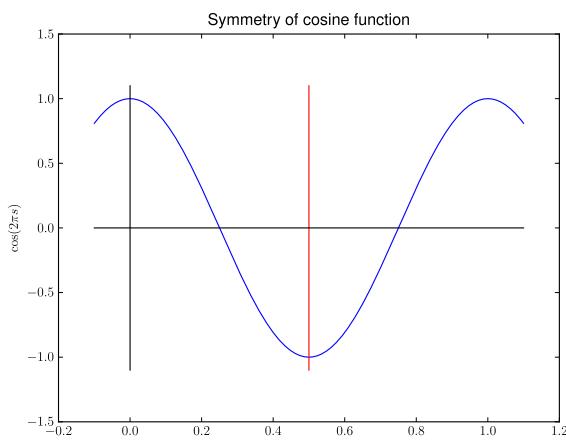


Fig. 206

Extend interpolation conditions (6.1.3.35) by symmetry, see Fig. 207

$$q\left(\frac{k}{2(n+1)} + \frac{1}{4(n+1)}\right) = z_k := \begin{cases} y_k & , \text{for } k = 0, \dots, n, \\ y_{2n+1-k} & , \text{for } k = n+1, \dots, 2n+1. \end{cases} \quad (6.1.3.36)$$

In a sense, we can mirror the interpolation conditions at  $x = \frac{1}{2}$ :

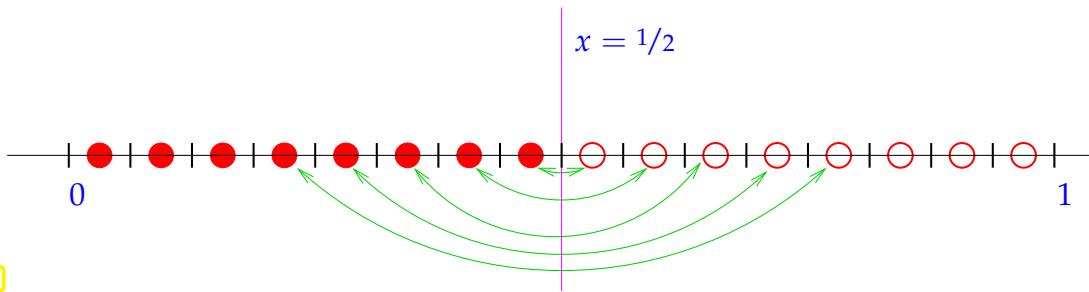


Fig. 207

Chebychev expansion for Chebychev interpolants  
 $\Updownarrow$   
 Trigonometric interpolation at equidistant points → Section 5.6.3

Trigonometric interpolation at equidistant points can be done very efficiently by means of FFT-based algorithms, see Code 5.6.3.3. We can also apply these for the computation of Chebychev expansion coefficients.

Details: (6.1.3.36)  $\iff$  square linear system of equations for unknown  $\beta_j$ .

$$q\left(\frac{k}{2(n+1)} + \frac{1}{4(n+1)}\right) = \sum_{j=-n}^{n+1} \left( \beta_j \exp\left(-\frac{2\pi j}{4(n+1)}\right) \right) \exp\left(-\frac{2\pi i}{2n+1} kj\right) = z_k .$$

$\Updownarrow$

$$\sum_{j=0}^{2n+1} \left( \beta_j \exp\left(-\frac{2\pi i(j-n)}{4(n+1)}\right) \right) \underbrace{\exp\left(-\frac{2\pi i}{2n+1} kj\right)}_{=\omega_{2(n+1)}^{kj} !} = \exp\left(-\pi i \frac{nk}{n+1}\right) z_k , \quad k = 0, \dots, 2n+1 .$$

$\Updownarrow$

$$\mathbf{F}_{2(n+1)} \mathbf{c} = \mathbf{b} \quad \text{with} \quad \mathbf{c} = \left[ \beta_j \exp\left(-\frac{2\pi j}{4(n+1)}\right) \right]_{j=0}^{2n+1} , \quad \mathbf{b} = \left[ \exp\left(-\pi i \frac{nk}{n+1}\right) z_k \right]_{k=0}^{2n+1} . \quad (6.1.3.37)$$

$(2n+2) \times (2n+2)$  Fourier matrix, see (4.2.0.13)

► solve (6.1.3.37) with inverse discrete Fourier transform, see 4.2:

asymptotic complexity  $\mathcal{O}(n \log n)$  ( $\rightarrow$  Section 4.3)

Note: by symmetry of  $\mathbf{z}$   $\Rightarrow \beta_{2n+1} = 0$ , cf. (6.1.3.34)!

#### C++ code 6.1.3.38: Efficient computation of Chebychev expansion coefficient of Chebychev interpolant

```

2 // efficiently compute coefficients  $\alpha_j$  in the Chebychev expansion
3 //  $p = \sum_{j=0}^n \alpha_j T_j$  of  $p \in \mathcal{P}_n$  based on values  $y_k$ ,
4 //  $k = 0, \dots, n$ , in Chebychev nodes  $t_k$ ,  $k = 0, \dots, n$ 
5 // IN: values  $y_k$  passed in  $y$ 
6 // OUT: coefficients  $\alpha_j$ 

```

```

7  VectorXd chebexp(const VectorXd& y) {
8      const int n = y.size() - 1;           // degree of polynomial
9      const std::complex<double> M_I(0, 1); // imaginary unit
10     // create vector z, see (6.1.3.36)
11     VectorXcd b(2*(n + 1));
12     const std::complex<double> om = -M_I*(M_PI*n)/((double)(n+1));
13     for (int j = 0; j <= n; ++j) {
14         b(j) = std::exp(om*double(j))*y(j); // this cast to double is necessary!!
15         b(2*n+1-j) = std::exp(om*double(2*n+1-j))*y(j);
16     }
17
18     // Solve linear system (6.1.3.37) with effort O(n log n)
19     Eigen::FFT<double> fft; // EIGEN's helper class for DFT
20     VectorXcd c = fft.inv(b); // -> c = ifft(z), inverse fourier transform
21     // recover  $\beta_j$ , see (6.1.3.37)
22     VectorXd beta(c.size());
23     const std::complex<double> sc = M_PI_2/(n + 1)*M_I;
24     for (unsigned j = 0; j < c.size(); ++j)
25         beta(j) = (std::exp(sc*double(-n+j))*c[j]).real();
26     // recover  $\alpha_j$ , see (6.1.3.34)
27     VectorXd alpha = 2*beta.segment(n,n); alpha(0) = beta(n);
28     return alpha;
29 }
```

**Remark 6.1.3.39 (Chebychev representation of built-in functions)** Computers use approximation by sums of Chebychev polynomials in the computation of functions like  $\log, \exp, \sin, \cos, \dots$ . The evaluation by means of Clenshaw algorithm according to Code 6.1.3.31 is more efficient and stable than the approximation by Taylor polynomials.

## 6.2 Mean Square Best Approximation

There is a particular family of norms for which the best approximant of a function  $f$  in a finite dimensional function space  $V_N$ , that is, the element of  $V_N$  that is closest to  $f$  with respect to that particular norm can actually be computed. It turns out that this computation boils down to solving a kind of least squares problem, similar to the least squares problems in  $\mathbb{K}^n$  discussed in Chapter 3.

### 6.2.1 Abstract Theory

Concerning mean square best approximation it is useful to learn an abstract framework first into which the concrete examples can be fit later.

#### 6.2.1.1 Mean Square Norms

Mean square norms generalize the Euclidean norm on  $\mathbb{K}^n$ , see [NS02, Sect. 4.4]. In a sense, they endow a vector space with a geometry and give a meaning to concepts like “orthogonality”.

**Definition 6.2.1.1. (Semi-)inner product [NS02, Sect. 4.4]**

Let  $V$  be a vector space over the field  $\mathbb{K}$ . A mapping  $b : V \times V \rightarrow \mathbb{K}$  is called an **inner product** on  $V$ , if it satisfies

(i)  $b$  is **linear** in the first argument:  $b(\alpha v + \beta w, u) = \alpha b(v, u) + \beta b(w, u)$  for all  $\alpha, \beta \in \mathbb{K}$ ,  $v, w \in V$ ,

(ii)  $b$  is (anti-)**symmetric**:  $b(v, w) = \overline{b(w, v)}$  ( $\overline{\phantom{x}} \hat{=} \text{complex conjugation}$ ),

(iii)  $b$  is **positive definite**:  $v \neq 0 \Leftrightarrow b(v, v) > 0$ .

$b$  is a **semi-inner product**, if it still complies with (i) and (ii), but is only positive semi-definite:  $b(v, v) \geq 0$  for all  $v \in V$ .

notation: usually we write  $(\cdot, \cdot)_V$  for an inner product on the vector space  $V$ .

**Definition 6.2.1.2. Orthogonality**

Let  $V$  be a vector space equipped with a (semi-)inner product  $(\cdot, \cdot)_V$ . Any two elements  $v$  and  $w$  of  $V$  are called **orthogonal**, if  $(v, w)_V = 0$ . We write  $v \perp w$ .

notation: If  $W \subset V$  is a (closed) subspace:  $v \perp W \Leftrightarrow (v, w)_V = 0 \forall w \in W$ .

**Theorem 6.2.1.3. Mean square (semi-)norm/Inner product (semi-)norm**

If  $(\cdot, \cdot)_V$  is a (semi-)inner product ( $\rightarrow$  Def. 6.2.1.1) on the vector space  $V$ , then

$$\|v\|_V := \sqrt{(v, v)_V}$$

defines a (semi-)norm ( $\rightarrow$  Def. 1.5.5.4) on  $V$ , the **mean square (semi-)norm**/ inner product (semi-)norm induced by  $(\cdot, \cdot)_V$ .

**§6.2.1.4 (Examples for mean square norms)**

♦ The **Euclidean norm** on  $\mathbb{K}^n$  induced by the dot product (Euclidean inner product).

$$(\mathbf{x}, \mathbf{y})_{\mathbb{K}^n} := \sum_{j=1}^n (\mathbf{x}_j \overline{\mathbf{y}_j}) \quad [\text{"Mathematical indexing"!}] \quad \mathbf{x}, \mathbf{y} \in \mathbb{K}^n .$$

♦ The  **$L^2$ -norm** (5.2.4.5) on  $C^0([a, b])$  induced by the  $L^2([a, b])$  inner product

$$(f, g)_{L^2([a, b])} := \int_a^b f(\tau) \overline{g(\tau)} d\tau , \quad f, g \in C^0([a, b]) . \quad (6.2.1.5)$$

**6.2.1.2 Normal Equations**

Mean square best approximation = best approximation in a mean square norm

From § 3.1.1.6 we know that in Euclidean space  $\mathbb{K}^n$  the best approximation of vector  $\mathbf{x} \in \mathbb{K}^n$  in a subspace  $V \subset \mathbb{K}^n$  is unique and given by the **orthogonal projection** of  $\mathbf{x}$  onto  $V$ . Now we generalize this to vector spaces equipped with inner products.

$X \triangleq$  a vector space over  $\mathbb{K} = \mathbb{R}$ , equipped with a mean square semi-norm  $\|\cdot\|_X$  induced by a semi-inner product  $(\cdot, \cdot)_X$ , see Thm. 6.2.1.3.

It can be an infinite dimensional function space, e.g.,  $X = C^0([a, b])$ .

$V \triangleq$  a finite-dimensional subspace of  $X$ , with basis  $\mathfrak{B}_V := \{b_1, \dots, b_N\} \subset V$ ,  $N := \dim V$ .

### Assumption 6.2.1.6.

The semi-inner product  $(\cdot, \cdot)_X$  is a genuine inner product ( $\rightarrow$  Def. 6.2.1.1) on  $V$ , that is, it is positive definite:  $(v, v)_X > 0 \forall v \in V \setminus \{0\}$ .

Now we give a *formula* for the element  $q$  of  $V$ , which is nearest to a given element  $f$  of  $X$  with respect to the norm  $\|\cdot\|_X$ . This is a genuine generalization of Thm. 3.1.2.1.

### Theorem 6.2.1.7. Mean square norm best approximation through normal equations

Given any  $f \in X$  there is a unique  $q \in V$  such that

$$\|f - q\|_X = \inf_{p \in V} \|f - p\|_X .$$

Its coefficients  $\gamma_j, j = 1, \dots, N$ , with respect to the basis  $\mathfrak{B}_V := \{b_1, \dots, b_N\}$  of  $V$  ( $q = \sum_{j=1}^N \gamma_j b_j$ ) are the unique solution of the **normal equations**

$$\mathbf{M}[\gamma_j]_{j=1}^N = [(f, b_j)]_{j=1}^N, \quad \mathbf{M} := \begin{bmatrix} (b_1, b_1)_X & \dots & (b_1, b_N)_X \\ \vdots & & \vdots \\ (b_N, b_1)_X & \dots & (b_N, b_N)_X \end{bmatrix} \in \mathbb{K}^{N,N}. \quad (6.2.1.8)$$

*Proof.* (*inspired by Rem. 3.1.2.5*) We first show that  $\mathbf{M}$  is s.p.d. ( $\rightarrow$  Def. 1.1.2.6). Symmetry is clear from the definition and the symmetry of  $(\cdot, \cdot)_X$ . That  $\mathbf{M}$  is even positive definite follows from

$$\mathbf{x}^\top \mathbf{M} \mathbf{x} = \sum_{k=1}^N \sum_{j=1}^N \xi_k \bar{\xi}_j (b_k, b_j)_X = \left\| \sum_{j=1}^N \xi_j b_j \right\|_X^2 > 0 , \quad (6.2.1.9)$$

if  $\mathbf{x} := [\xi_j]_{j=1}^N \neq \mathbf{0} \Leftrightarrow \sum_{j=1}^N \xi_j b_j \neq \mathbf{0}$ , since  $\|\cdot\|_X$  is a *norm* on  $V$  by Ass. 6.2.1.6.

Now, writing  $\mathbf{c} := [\gamma_j]_{j=1}^N \in \mathbb{K}^N$ ,  $\mathbf{b} := [(f, b_j)]_{j=1}^N \in \mathbb{K}^N$ , and using the basis representation

$$q = \sum_{j=1}^N \gamma_j b_j ,$$

we find

$$\Phi(\mathbf{c}) := \|f - q\|_X^2 = \|f\|_X^2 - 2\mathbf{b}^\top \mathbf{c} + \mathbf{c}^\top \mathbf{M} \mathbf{c} .$$

Applying the differentiation rules from Ex. 8.4.1.12 to  $\Phi : \mathbb{K}^N \rightarrow \mathbb{R}$ ,  $\mathbf{c} \mapsto \Phi(\mathbf{c})$ , we obtain

$$\text{grad } \Phi(\mathbf{c}) = 2(\mathbf{M}\mathbf{c} - \mathbf{b}), \quad (6.2.1.10)$$

$$\mathbf{H}\Phi(\mathbf{c}) = 2\mathbf{M}. \quad (\text{independent of } \mathbf{c}) \quad (6.2.1.11)$$

Since  $\mathbf{M}$  is s.p.d., the unique solution of  $\text{grad } \Phi(\mathbf{c}) = \mathbf{M}\mathbf{c} - \mathbf{b} = 0$  yields the unique global minimizer of  $\Phi$ ; the Hessian  $2\mathbf{M}$  is s.p.d. everywhere!  $\square$

The unique  $q$  from Thm. 6.2.1.7 is called the **best approximant** of  $f$  in  $V$ .

### Corollary 6.2.1.12. Best approximant by orthogonal projection

If  $q$  is the best approximant of  $f$  in  $V$ , then  $f - q$  is **orthogonal** to every  $p \in V$ :

$$(f - q, p)_X = 0 \quad \forall p \in V \Leftrightarrow f - q \perp V.$$

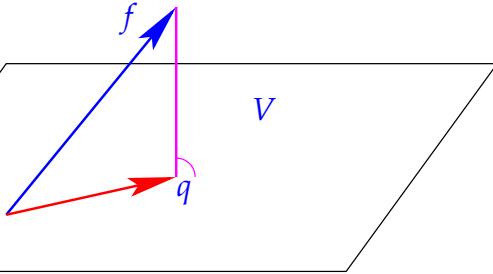


Fig. 208

The message of Cor. 6.2.1.12:

- the best approximation error  $f - q$  for  $f \in X$  in  $V$  is **orthogonal** to the subspace  $V$ .

See § 3.1.1.6 for related discussion in Euclidean space  $\mathbb{K}^n$ .

**Remark 6.2.1.13 (Connection with linear least squares problems Chapter 3)** In Section 3.1.1 we introduced the concept of least squares solutions of overdetermined linear systems of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m > n$ , see Def. 3.1.1.1. Thm. 3.1.2.1 taught that the normal equations  $\mathbf{A}^\top \mathbf{Ax} = \mathbf{A}^\top \mathbf{b}$  give the least squares solution, if  $\text{rank}(\mathbf{A}) = n$ .

In fact, Thm. 3.1.2.1 and the above Thm. 6.2.1.7 agree if  $X = \mathbb{K}^n$  (Euclidean space) and  $V = \text{Span}\{\mathbf{a}_1, \dots, \mathbf{a}_n\}$ , where  $\mathbf{a}_j \in \mathbb{R}^m$  are the columns of  $\mathbf{A}$  and  $N = n$ .  $\square$

### 6.2.1.3 Orthonormal Bases

In the setting of Section 6.2.1.2 we may ask: Which choice of basis  $\mathcal{B} = \{b_1, \dots, b_N\}$  of  $V \subset X$  renders the normal equations (6.2.1.8) particularly simple? Answer: A basis  $\mathcal{B}$ , for which  $(b_k, b_j)_X = \delta_{kj}$  ( $\delta_{kj}$  the Kronecker symbol), because this will imply  $\mathbf{M} = \mathbf{I}$  for the coefficient matrix of the normal equations.

### Definition 6.2.1.14. Orthonormal basis

A subset  $\{b_1, \dots, b_N\}$  of an  $N$ -dimensional vector space  $V$  with inner product ( $\rightarrow$  Def. 6.2.1.1)  $(\cdot, \cdot)_V$  is an **orthonormal basis (ONB)**, if  $(b_k, b_j)_V = \delta_{kj}$ .

A basis of  $\{b_1, \dots, b_N\}$  of  $V$  is called **orthogonal**, if  $(b_k, b_j)_V = 0$  for  $k \neq j$ .

**Corollary 6.2.1.15. ONB representation of best approximant**

If  $\{b_1, \dots, b_N\}$  is an orthonormal basis ( $\rightarrow$  Def. 6.2.1.14) of  $V \subset X$ , then the best approximant  $q := \operatorname{argmin}_{p \in V} \|f - p\|_X$  of  $f \in X$  has the representation

$$q = \sum_{j=1}^N (f, b_j)_X b_j. \quad (6.2.1.16)$$

**§6.2.1.17 (Gram-Schmidt orthonormalization)** From Section 1.5.1 we already know how to compute orthonormal bases: The algorithm from § 1.5.1.1 can be run in the framework of any vector space  $V$  endowed with an inner product  $(\cdot, \cdot)_V$  and induced mean square norm  $\|\cdot\|_V$ .

```

1:  $b_1 := \frac{p_1}{\|p_1\|_V}$  % 1st output vector
2: for  $j = 2, \dots, k$  do
   { % Orthogonal projection
3:    $b_j := p_j$ 
4:   for  $\ell = 1, 2, \dots, j-1$  do (6.2.1.18)
5:     {  $b_j \leftarrow b_j - (p_j, b_\ell)_V b_\ell$  }
6:   if ( $b_j = 0$ ) then STOP
7:   else {  $b_j \leftarrow \frac{b_j}{\|b_j\|_V}$  }
   }
```

**Theorem 6.2.1.19. Gram-Schmidt orthonormalization**

When supplied with  $k \in \mathbb{N}$  linearly independent vectors  $p_1, \dots, p_k \in V$  in a vector space with inner product  $(\cdot, \cdot)_V$ , Algorithm (6.2.1.18) computes vectors  $b_1, \dots, b_k$  with

$$(b_\ell, b_j)_V = \delta_{\ell j}, \quad \ell, j \in \{1, \dots, k\}, \\ \text{Span}\{b_1, \dots, b_k\} = \text{Span}\{p_1, \dots, p_k\}$$

for all  $\ell \in \{1, \dots, k\}$ .

This suggests the following alternative approach to the computation of the mean square best approximant  $q$  in  $V$  of  $f \in X$ :

- ① Orthonormalize a basis  $\{b_1, \dots, b_N\}$  of  $V$ ,  $N := \dim V$ , using Gram-Schmidt algorithm (6.2.1.18).
- ② Compute  $q$  according to (6.2.1.16).

Number of inner products to be evaluated:  $O(N^2)$  for  $N \rightarrow \infty$ .

## 6.2.2 Polynomial Mean Square Best Approximation

Now we apply the results of Section 6.2.1 in the following setting:

$X$ : function space  $C^0([a, b])$ ,  $-\infty < a < b < \infty$ , of  $\mathbb{R}$ -values continuous functions,  
 $V$ : space  $\mathcal{P}_m$  of polynomials of degree  $\leq m$ .

**Remark 6.2.2.1 (Inner products on spaces  $\mathcal{P}_m$  of polynomials)** To match the abstract framework of Section 6.2.1 we need to find (semi-)inner products on  $C^0([a, b])$  that supply positive definite inner products on  $\mathcal{P}_m$ . The following options are commonly considered:

- ◆ On any interval  $[a, b]$  we can use the  $L^2([a, b])$ -inner product  $(\cdot, \cdot)_{L^2([a,b])}$ , defined in (6.2.1.5).
- ◆ Given a positive weight function  $w : [a, b] \rightarrow \mathbb{R}$ ,  $w(t) > 0$  for all  $t \in [a, b]$ , we can consider the

weighted  $L^2$ -inner product on the interval  $[a, b]$

$$(f, g)_{w, [a, b]} := \int_a^b w(\tau) f(\tau) \overline{g}(\tau) d\tau . \quad (6.2.2.2)$$

- ◆ For  $n \geq m$  and  $n + 1$  distinct points collected in the set  $\mathcal{T} := \{t_0, t_1, \dots, t_n\} \subset [a, b]$  we can use the **discrete  $L^2$ -inner product**

$$(f, g)_{\mathcal{T}} := \sum_{j=0}^n f(t_j) \overline{g(t_j)} . \quad (6.2.2.3)$$

Since a polynomial of degree  $\leq m$  must be zero everywhere, if it vanishes in at least  $m + 1$  distinct points,  $(f, g)_{\mathcal{T}}$  is positive definite on  $\mathcal{P}_m$ .

For all these inner products on  $\mathcal{P}_m$  holds

$$(\{t \mapsto tf(t)\}, g)_X = (f, \{t \mapsto tg(t)\})_X , \quad f, g \in C^0([a, b]) , \quad (6.2.2.4)$$

that is, multiplication with the independent variable can be shifted to the other function inside the inner product.

☞ notation: Note that we have to plug a function into the slots of the inner products; this is indicated by the notation  $\{t \mapsto \dots\}$ .

#### Assumption 6.2.2.5. Self-adjointness of multiplication operator

We assume the inner product  $(\cdot, \cdot)_X$  to satisfy (6.2.2.4).

The ideas of Section 6.2.1.3 that center around the use of orthonormal bases can also be applied to polynomials.

#### Definition 6.2.2.6. Orthonormal polynomials → Def. 6.2.1.14

Let  $(\cdot, \cdot)_X$  be an inner product on  $\mathcal{P}_m$ . A sequence  $r_0, r_1, \dots, r_m$  provides **orthonormal polynomials** (ONPs) with respect to  $(\cdot, \cdot)_X$ , if

$$r_\ell \in \mathcal{P}_\ell , \quad (r_k, r_\ell)_X = \delta_{k\ell} , \quad \ell, k \in \{0, \dots, m\} . \quad (6.2.2.7)$$

The polynomials are just **orthogonal**, if  $(r_k, r_\ell)_X = 0$  for  $k \neq \ell$ .

By virtue of Thm. 6.2.1.19 orthonormal polynomials can be generated by applying Gram-Schmidt orthonormalization from § 6.2.1.17 to the ordered basis of monomials  $\{t \mapsto t^j\}_{j=0}^m$ .

**Lemma 6.2.2.8. Uniqueness of orthonormal polynomials**

The sequence of orthonormal polynomials from Def. 6.2.2.6 is unique up to signs, supplies an  $(\cdot, \cdot)_X$ -orthonormal basis ( $\rightarrow$  Def. 6.2.1.14) of  $\mathcal{P}_m$ , and satisfies

$$\text{Span}\{r_0, \dots, r_k\} = \mathcal{P}_k, \quad k \in \{0, \dots, m\}. \quad (6.2.2.9)$$

*Proof.* Comparing Def. 6.2.1.14 and (6.2.2.7) the ONB-property of  $\{r_0, \dots, r_m\}$  is immediate. Then (6.2.2.7) follows from dimensional considerations.

$r_0$  must be a constant, which, up to sign, is fixed by the normalization condition  $\|r_0\|_X = 1$ .

$\mathcal{P}_{k-1} \subset \mathcal{P}_k$  has co-dimension 1 so that there is a unit “vector” in  $\mathcal{P}_k$ , which is orthogonal to  $\mathcal{P}_{k-1}$  and unique up to sign.  $\square$

**§6.2.2.10 (Orthonormal polynomials by orthogonal projection)** Let  $r_0, \dots, r_m \in \mathcal{T}_p$  be a sequence of orthonormal polynomials according to Def. 6.2.2.6. From (6.2.2.7) we conclude that  $r_k \in \mathcal{P}_k$  has leading coefficient  $\neq 0$ .

Hence  $s_k(t) := t \cdot r_k(t)$  is a polynomial of degree  $k+1$  with leading coefficient  $\neq 0$ , that is  $s_k \in \mathcal{P}_{k+1} \setminus \mathcal{P}_k$ . Therefore,  $r_{k+1}$  can be obtained by *orthogonally projecting*  $s_k$  onto  $\mathcal{P}_k$  plus normalization, cf. Lines 4-5 of Algorithm (6.2.1.18):

$$r_{k+1} = \pm \frac{\tilde{r}_{k+1}}{\|\tilde{r}_{k+1}\|_X}, \quad \tilde{r}_{k+1} = s_k - \sum_{j=0}^k (s_k, r_j)_X r_j. \quad (6.2.2.11)$$

Straightforward computations confirm that  $r_{k+1} \perp \mathcal{P}_k$ .

The sum in (6.2.2.11) collapses to two terms! In fact, since  $(r_k, q)_X = 0$  for all  $q \in \mathcal{P}_{k-1}$ , by Ass. 6.2.2.5

$$(s_k, r_j)_X = (\{t \mapsto tr_k(t)\}, r_j)_X \stackrel{(6.2.2.4)}{=} (r_k, \{t \mapsto tr_j\})_X = 0, \quad \text{if } j < k-1,$$

because in this case  $\{t \mapsto tr_j\} \in \mathcal{P}_{k-1}$ . As a consequence (6.2.2.11) reduces to the **3-term recursion**

$$\begin{aligned} r_{k+1} &= \pm \frac{\tilde{r}_{k+1}}{\|\tilde{r}_{k+1}\|_X}, \\ \tilde{r}_{k+1} &= s_k - (\{t \mapsto tr_k\}, r_k)_X r_k - (\{t \mapsto tr_k\}, r_{k-1})_X r_{k-1} \end{aligned}, \quad k = 1, \dots, m-1. \quad (6.2.2.12)$$

The recursion starts with  $r_{-1} := 0$ ,  $r_0 = \{t \mapsto 1/\|1\|_X\}$ .  $\square$

The 3-term recursion (6.2.2.12) can be recast in various ways. Forgoing normalization the next theorem presents one of them.

**Theorem 6.2.2.13. 3-term recursion for orthogonal polynomials**

Given **any** inner product  $(\cdot, \cdot)_X$  on  $\mathcal{P}_m$ ,  $m \in \mathbb{N}$ , define  $p_{-1} := 0$ ,  $p_0 = 1$ , and

$$p_{k+1}(t) := (t - \alpha_{k+1})p_k(t) - \beta_k p_{k-1}(t), \quad k = 0, 1, \dots, m-1, \\ \text{with } \alpha_{k+1} := \frac{(\{t \mapsto tp_k(t)\}, p_k)_X}{\|p_k\|_X^2}, \quad \beta_k := \frac{\|p_k\|_X^2}{\|p_{k-1}\|_X^2}. \quad (6.2.2.14)$$

Then ♦  $p_k \in \mathcal{P}_k$  has leading coefficient = 1, and  
♦  $\{p_0, p_1, \dots, p_m\}$  is an **orthogonal basis** of  $\mathcal{P}_m$ .

*Proof.* (by rather straightforward induction) We first confirm, thanks to the definition of  $\alpha_1$ ,

$$(p_0, p_1)_X = (p_0, \{t \mapsto (t - \alpha_1)p_0(t)\})_X = (p_0, \{t \mapsto tp_0(t)\})_X - \alpha_1(p_0, p_0)_X = 0.$$

For the induction step we assume that the assertion is true for  $p_0, \dots, p_k$  and observe that for  $p_{k+1}$  according to (6.2.2.14) we have

$$(p_k, p_{k+1})_X = (p_k, \{t \mapsto (t - \alpha_{k+1})p_k(t)\} - \beta_k p_{k-1})_X \\ = (p_k, \{t \mapsto tp_k(t)\})_X - \underbrace{\alpha_{k+1}(p_k, p_k)_X}_{=(\{t \mapsto tp_k(t)\}, p_k)_X} - \underbrace{\beta_k(p_k, p_{k-1})_X}_{=0} = 0,$$

$$(p_{k-1}, p_{k+1})_X = (p_{k-1}, \{t \mapsto (t - \alpha_{k+1})p_k(t)\} - \beta_k p_{k-1})_X \\ = (p_{k-1}, \{t \mapsto tp_k(t)\})_X - \alpha_{k+1}(p_{k-1}, p_k)_X - \beta_k(p_{k-1}, p_{k-1})_X = 0,$$

$$(p_\ell, p_{k+1})_X = (p_\ell, \{t \mapsto (t - \alpha_{k+1})p_k(t)\} - \beta_k p_{k-1})_X \\ = \underbrace{(p_\ell, \{t \mapsto tp_k(t)\})_X}_{=0} - \underbrace{\alpha_{k+1}(p_\ell, p_k)_X}_{0} - \underbrace{\beta_k(p_\ell, p_{k-1})_X}_{0} = 0, \quad \ell = 0, \dots, k-2.$$

This amounts to the assertion of orthogonality for  $k+1$ . Above, several inner products vanish because of the induction hypothesis!

□

**Remark 6.2.2.15 ( $L^2([-1, 1])$ -orthogonal polynomials)** An important inner product on  $C^0[-1, 1]$  is the  **$L^2$ -inner product**, see (6.2.1.5)

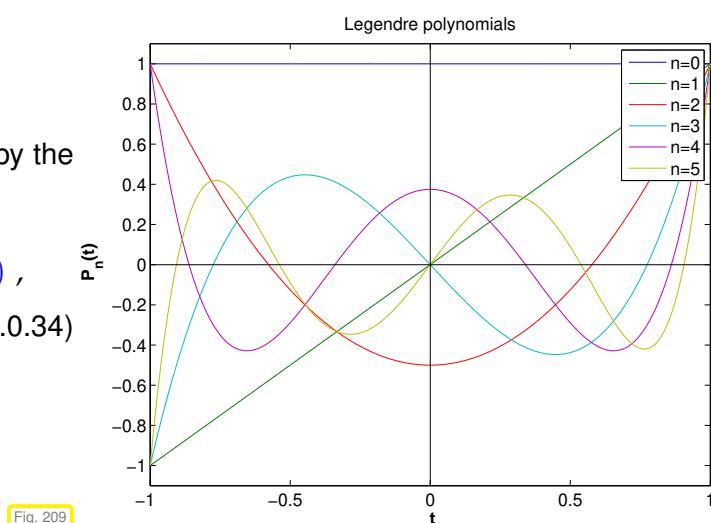
$$(f, g)_{L^2([-1, 1])} := \int_{-1}^1 f(\tau) \bar{g}(\tau) d\tau, \quad f, g \in C^0([-1, 1]).$$

It is a natural question what is the unique sequence of  $L^2([-1, 1])$ -orthonormal polynomials. Their rather simple characterization will be discussed in the sequel.

The **Legendre polynomials**  $P_n$  can be defined by the 3-term recursion

$$P_{n+1}(t) := \frac{2n+1}{n+1} t P_n(t) - \frac{n}{n+1} P_{n-1}(t), \quad (7.3.0.34)$$

$$P_0 := 1, \quad P_1(t) := t.$$



$P_k \in \mathcal{P}_k$  is immediate and so is the parity

$$P_k(t) = P_k(-t) \quad \text{if } k \text{ is even,} \quad P_k(t) = -P_k(-t) \quad \text{if } k \text{ is odd.} \quad (6.2.2.16)$$

Orthogonality,  $(P_k, P_m)_{L^2(\Omega)} = 0$  or  $m \neq k$ , as well as  $\|P_k\|_{L^2([-1,1])}^2 = 2^{2/2n+1}$  can be proved by induction based on (6.2.2.16). This implies that the  $L^2([-1,1])$ -orthonormal polynomials are  $r_k = \sqrt{k+1/2} P_k$  and their 3-term recursion from Thm. 6.2.2.13 reads

$$r_{n+1}(t) = \frac{\sqrt{4(n+1)^2 - 1}}{n+1} t r_n(t) - \frac{n}{n+1} \sqrt{\frac{4(n+1)^2 - 1}{4n^2 - 1}} u_{n-1}(t), \quad (6.2.2.17)$$

$$r_{-1} := 0, \quad r_0 = \frac{1}{2}\sqrt{2}. \quad (6.2.2.18)$$

↓

**§6.2.2.19 (Discrete orthogonal polynomials)** Since they involve integrals, weighted  $L^2$ -inner products (6.2.2.2) are not accessible computationally, unless one resigns to approximation, see Chapter 7 for corresponding theory and techniques.

Therefore, given a point set  $\mathcal{T} := \{t_0, t_1, \dots, t_n\}$ , we focus on the associated discrete  $L^2$ -inner product

$$(f, g)_X := (f, g)_{\mathcal{T}} := \sum_{j=0}^n f(t_j) \overline{g(t_j)}, \quad f, g \in C^0([a, b]), \quad (6.2.2.3)$$

which is positive definite on  $\mathcal{P}_n$  and satisfies Ass. 6.2.2.5.

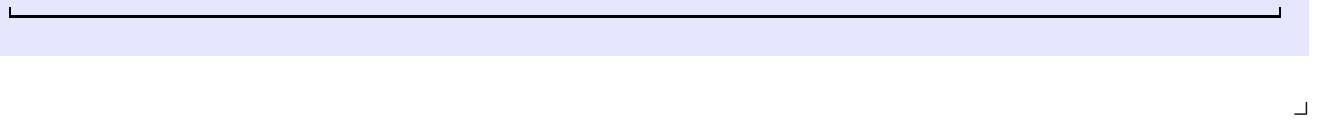
The polynomials  $p_k$  generated by the 3-term recursion (6.2.2.14) from Thm. 6.2.2.13 are then called **discrete orthogonal polynomials**. The following C++ code computes the recursion coefficients  $\alpha_k$  and  $\beta_k$ ,  $k = 1, \dots, n-1$ .

#### C++ code 6.2.2.20: Computation of weights in 3-term recursion for discrete orthogonal polynomials

```

2 // Computation of coefficients alpha, beta from 6.2.2.13
3 // IN : t = points in the definition of the discrete L^2-inner product
4 // n = maximal index desired
5 // alpha, beta are used to save coefficients of recursion
6 void coeffortho(const VectorXd& t, const unsigned n, VectorXd& alpha, VectorXd&
    beta) {
7     const unsigned m = t.size(); // maximal degree of orthogonal polynomial
8     alpha = VectorXd( std::min(n-1, m-2) + 1 );
9     beta = VectorXd( std::min(n-1, m-2) + 1 );
10    alpha(0) = t.sum() / m;
11    // initialization of recursion; we store only the values of
12    // the polynomials at the points in T
13    VectorXd p0,
14        p1 = VectorXd::Ones(m),
15        p2 = t - alpha(0)*VectorXd::Ones(m);
16    for (unsigned k = 0; k < std::min(n-1, m-2); ++k) {
17        p0 = p1; p1 = p2;
18        // 3-term recursion (6.2.2.14)
19        alpha(k+1) = p1.dot(t.cwiseProduct(p1)) / p1.squaredNorm();
20        beta(k) = p1.squaredNorm() / p0.squaredNorm();
21        p2 = (t - alpha(k+1)*VectorXd::Ones(m)).cwiseProduct(p1) - beta(k)*p0;
22    }
23}

```



**§6.2.2.21 (Polynomial fitting)** Given a point set  $\mathcal{T} := \{t_0, t_1, \dots, t_n\} \subset [a, b]$ , and a function  $f : [a, b] \rightarrow \mathbb{K}$ , we may seek to approximate  $f$  by its polynomial best approximant with respect to the discrete  $L^2$ -norm  $\|\cdot\|_{\mathcal{T}}$  induced by the discrete  $L^2$ -inner product (6.2.2.3).

### Definition 6.2.2.22. Fitted polynomial

Given a point set  $\mathcal{T} := \{t_0, t_1, \dots, t_n\} \subset [a, b]$ , and a function  $f : [a, b] \rightarrow \mathbb{K}$  we call

$$q_k := \operatorname{argmin}_{p \in \mathcal{P}_k} \|f - p\|_{\mathcal{T}}, \quad k \in \{0, \dots, n\},$$

the **fitting polynomial** to  $f$  on  $\mathcal{T}$  of degree  $k$ .

The stable and efficient computation of fitting polynomials can rely on combining Thm. 6.2.2.13 with Cor. 6.2.1.15:

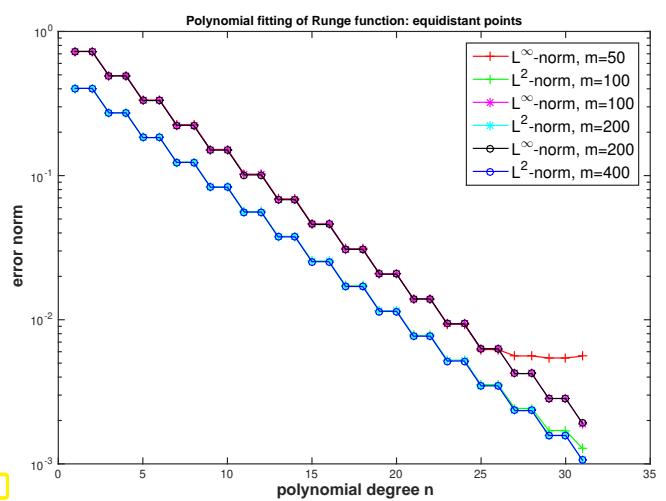
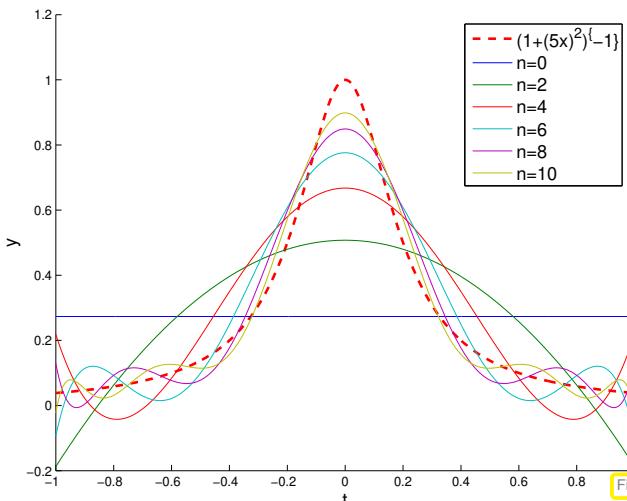
- ① (Pre-)compute the weights  $\alpha_\ell$  and  $\beta_\ell$  for the 3-term recursion (6.2.2.14).
- ② (Pre-)compute the values of the orthogonal polynomials  $p_k$  at desired evaluation points  $x_i \in \mathbb{R}$ ,  $i = 1, \dots, N$ .
- ③ Compute the inner product  $(f, p_\ell)_X$ ,  $\ell = 0, \dots, k$ , and use (6.2.1.16) to linearly combine the vectors  $[p_\ell(x_i)]_{i=1}^N$ ,  $\ell = 0, \dots, k$ .

This yields  $[q(x_i)]_{i=1}^N$ ,  $q \in \mathcal{P}_k$  the fitting polynomial.

**EXAMPLE 6.2.2.23 (Approximation by discrete polynomial fitting)** We use equidistant points  $\mathcal{T} := \{t_k = -1 + k \frac{2}{m}, k = 0, \dots, m\} \subset [-1, 1]$ ,  $m \in \mathbb{N}$  to compute fitting polynomials ( $\rightarrow$  Def. 6.2.2.22) for two different functions.

We monitor the  $L^2$ -norm and  $L^\infty$ -norm of the approximation error, both norms approximated by sampling in  $\xi_j = -1 + \frac{j}{500}$ ,  $j = 0, \dots, 1000$ .

- ①  $f(t) = (1 + (5t)^2)^{-1}$ ,  $I = [-1, 1] \rightarrow$  Ex. 6.1.2.10, analytic in complex neighborhood of  $[-1, 1]$ :



➤ We observe exponential convergence ( $\rightarrow$  Def. 6.1.2.7) in the polynomial degree  $n$ .

②  $f(t) = \max\{0, 1 - 2 * |x + \frac{1}{4}|^2\}$ ,  $f$  only in  $C^0([-1, 1])$ :

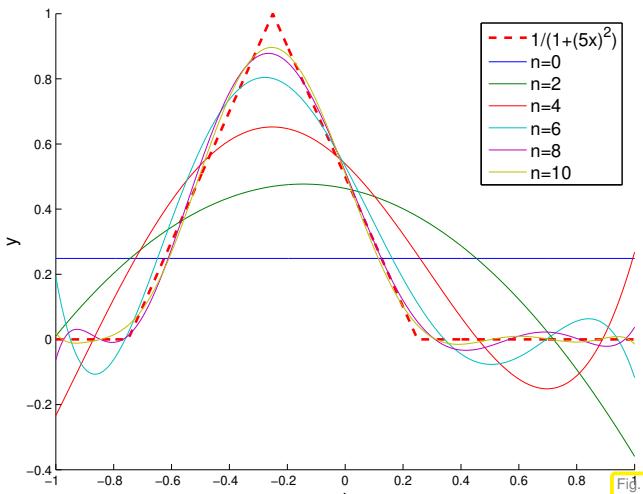
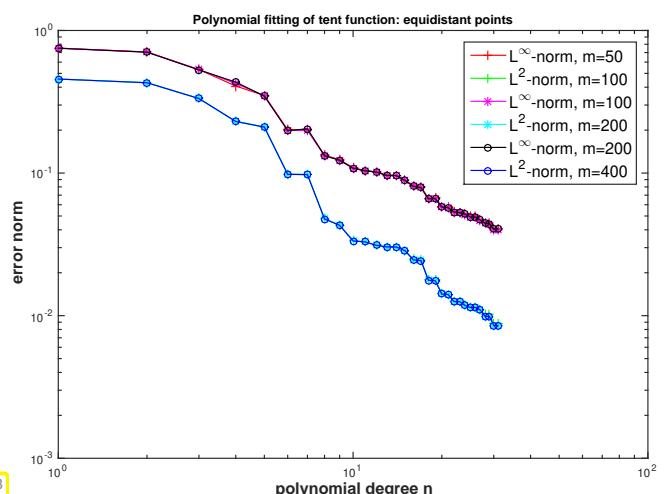


Fig. 212



➤ We observe only algebraic convergence ( $\rightarrow$  Def. 6.1.2.7 in the polynomial degree  $n$  (for  $n \ll m$ !).

③ “bump function”

$$f(t) = \max\{\cos(4\pi|t + \frac{1}{4}|), 0\}.$$

➤ Merely  $f \in C^1([-1, 1])$

Doubly logarithmic plot suggests “asymptotic” algebraic convergence.

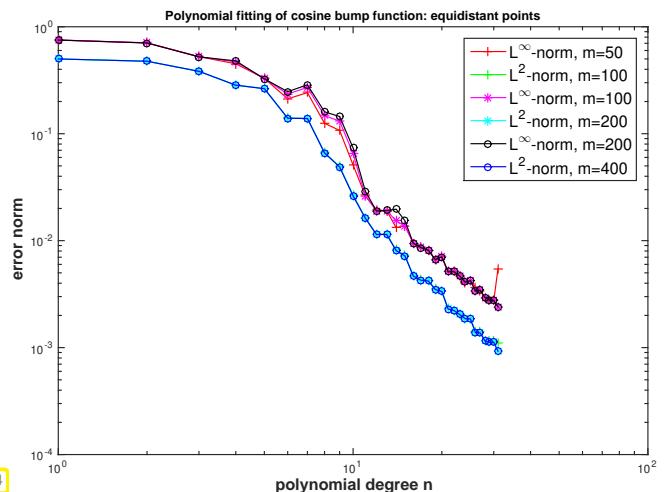


Fig. 214

## 6.3 Uniform Best Approximation

**§6.3.0.1 (The alternation theorem)** Given an interval  $[a, b]$  we seek a best approximant of a function  $f \in C^0([a, b])$  in the space  $\mathcal{P}_n$  of polynomials of degree  $\leq n$  with respect to the supremum norm  $\|\cdot\|_{L^\infty([a,b])}$ :

$$q \in \operatorname{argmin}_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty(I)}.$$

The results of Section 6.2.1 cannot be applied because the supremum norm is not induced by an inner product on  $\mathcal{P}_n$ .

Theory provides us with surprisingly precise necessary and sufficient conditions to be satisfied by the polynomial  $L^\infty([a, b])$ -best approximant  $q$ .

**Theorem 6.3.0.2. Chebychev alternation theorem**

Given  $f \in C^0[a, b]$ ,  $a < b$ , and a polynomial degree  $n \in \mathbb{N}$ , a polynomial  $q \in \mathcal{P}_n$  satisfies

$$q = \operatorname{argmin}_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty(I)}$$

if and only if there exist  $n + 2$  points  $a \leq \xi_0 < \xi_1 < \dots < \xi_{n+1} \leq b$  such that

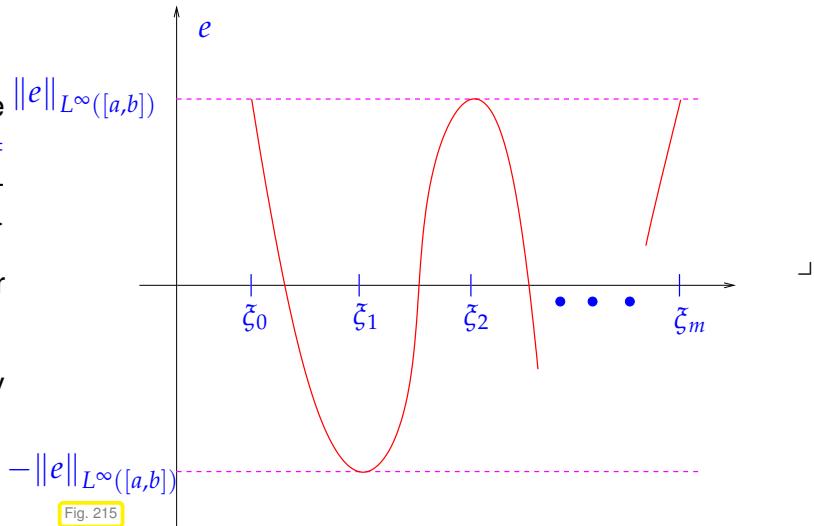
$$\begin{aligned} |e(\xi_j)| &= \|e\|_{L^\infty([a,b])}, \quad j = 0, \dots, n+1, \\ e(\xi_j) &= -e(\xi_{j+1}), \quad j = 0, \dots, n, \end{aligned}$$

where  $e := f - q$  denotes the approximation error.

Visualization of the behavior of the  $\|e\|_{L^\infty([a,b])}$ -best approximation error  $e := f - q$  according to the Chebychev alternation theorem Thm. 6.3.0.2. ▷

The extrema of the approximation error are sometimes called **alternants**.

Compare with the shape of the Chebychev polynomials → Def. 6.1.3.3.



**§6.3.0.3 (Remez algorithm)** The widely used iterative algorithm (**Remez algorithm**) for finding an  $L^\infty$ -best approximant is motivated by the alternation theorem. The idea is to determine successively better approximations of the set of alternants:  $\mathcal{A}^{(0)} \rightarrow \mathcal{A}^{(1)} \rightarrow \dots$ ,  $\#\mathcal{A}^{(l)} = n + 2$ .

Key is the observation that, due to the alternation theorem, the polynomial  $L^\infty([a, b])$ -best approximant  $q$  will satisfy (one of the) **interpolation conditions**

$$q(\xi_k) \pm (-1)^k \delta = f(\xi_k), \quad k = 0, \dots, n+1, \quad \delta := \|f - q\|_{L^\infty([a,b])}. \quad (6.3.0.4)$$

- ① Initial guess  $\mathcal{A}^{(0)} := \{\xi_0^{(0)} < \xi_1^{(0)} < \dots < \xi_n^{(0)} < \xi_{n+1}^{(0)}\} \subset [a, b]$  “arbitrary”, for instance extremal points of the Chebychev polynomial  $T_{n+1}$ , → Def. 6.1.3.3, so so-called **Chebychev alternants**

$$\xi_j^{(0)} = \frac{1}{2}(a + b) + \frac{1}{2}(b - a) \cos\left(\frac{j}{n+1}\pi\right), \quad j = 0, \dots, n+1. \quad (6.3.0.5)$$

- ② Given approximate alternants  $\mathcal{A}^{(l)} := \{\xi_0^{(l)} < \xi_1^{(l)} < \dots < \xi_n^{(l)} < \xi_{n+1}^{(l)}\} \subset [a, b]$  determine  $q \in \mathcal{P}_n$  and a deviation  $\delta \in \mathbb{R}$  satisfying the extended interpolation condition

$$q(\xi_k^{(l)}) + (-1)^k \delta = f(\xi_k^{(l)}), \quad k = 0, \dots, n+1. \quad (6.3.0.6)$$

After choosing a basis for  $\mathcal{P}_n$ , this is  $(n + 2) \times (n + 2)$  linear system of equations, cf. § 5.1.0.15.

- ③ Choose  $\mathcal{A}^{(l+1)}$  as the set of extremal points of  $f - q$ , truncated in case more than  $n + 2$  of these exist.

These extreme can be located approximately by sampling on a fine grid covering  $[a, b]$ ). If the derivative of  $f \in C^1([a, b])$  is available, too, then search for zeros of  $(f - p)'$  using the secant method from § 8.3.2.21.

- ④ If  $\|f - q\|_{L^\infty([a,b])} \leq \text{TOL} \cdot \|d\|_{L^\infty([a,b])}$  STOP, else GOTO ②.  
( $\text{TOL}$  is a prescribed relative tolerance.)

#### C++ code 6.3.0.7: Remez algorithm for uniform polynomial approximation on an interval

```

2 // IN : f = handle to the Function, point evaluation
3 // df = handle to the Derivative of f, point evaluation
4 // a, b = interval boundaries
5 // d = degree of polynomial
6 // c will be saved to save the coefficients of the interpolant in
7 // monomial basis
8 template <class Function, class Derivative>
9 void remez(const Function &f, const Derivative &df, const double a,
10           const double b, const unsigned d, const double tol, VectorXd &c) {
11   const unsigned n = 8 * d;                                // number of sampling points
12   VectorXd xtab = VectorXd::LinSpaced(n, a, b); // points of sampling grid
13   VectorXd ftab = feval(f, xtab); // function values at sampling grid
14   double fsupn =
15     ftab.cwiseAbs().maxCoeff(); // approximate supremum norm of f
16   VectorXd dftab = feval(df, xtab); // derivative values at sampling grid
17
18   // The vector xe stores the current guess for the alternants
19   // initial guess is Chebychev alternant (6.3.0.5)
20   const double h = M_PI / (d + 1);
21   VectorXd xe(d + 2);
22   for (unsigned i = 0; i < d + 2; ++i) {
23     xe(i) = (a + b) / 2. + (a - b) / 2. * std::cos(h * i);
24   }
25
26   VectorXd fxe = feval(f, xe); // f evaluated at alternants
27   const unsigned maxit = 10;
28   // Main iteration loop of Remez algorithm
29   for (unsigned k = 0; k < maxit; ++k) {
30     // Interpolation at d+2 points xe with deviations
31     // ±δ Algorithm uses monomial basis, which is not
32     // optimal
33     MatrixXd V = vander(xe), A(d + 2, d + 2);
34     // build Matrix A, LSE
35     A.block(0, 0, d + 2, d + 1) = V.block(0, 1, d + 2, d + 1);
36     for (unsigned r = 0; r < d + 2; ++r) {
37       A(r, d + 1) = std::pow(-1, r);
38     }
39
40     c = A.lu().solve(fxe); // solve for coefficients of polynomial q
41     VectorXd cd(
42       d); // to compute monomial coefficients of derivative q'
43     for (unsigned i = 0; i < d; ++i) {
44       cd(i) = (d - i) * c(i);
45     }
46
47     // Find initial guesses for the inner extremes by sampling
48     // track sign changes of the derivative of the approximation error

```

```

49 |     VectorXd deltab = polyval(cd, xtab) - dftab ,
50 |     s = deltab.head(n - 1).cwiseProduct(deltab.tail(n - 1));
51 |     VectorXd ind = findNegative(s), // ind = find(s < 0)
52 |     xx0 = select(xtab, ind); // approximate zeros of e'
53 |     const unsigned nx = ind.size(); // number of approximate zeros
54 |
55 |     if (nx < d) { // too few extrema; bail out
56 |         std::cerr << "Too few extrema!\n";
57 |         return;
58 |     }
59 |
60 |     // Secant method to determine zeros of derivative of approximation
61 |     // error
62 |     VectorXd F0 = polyval(cd, xx0) - feval(df, xx0);
63 |     // initial guesses from shifting sampling points
64 |     VectorXd xx1 = xx0 + (b - a) / (2 * n) * VectorXd::Ones(nx.size()),
65 |     F1 = polyval(cd, xx1) - feval(df, xx1);
66 |     // Main loop of the secant method
67 |     while (F1.cwiseAbs().minCoeff() > 1e-12) {
68 |         VectorXd xx2 = xx1 - (F1.cwiseQuotient(F1 - F0)).cwiseProduct(xx1 - xx0);
69 |         xx0 = xx1;
70 |         xx1 = xx2;
71 |         F0 = F1;
72 |         F1 = polyval(cd, xx1) - feval(df, xx1);
73 |     }
74 |
75 |     // Determine new approximation for alternants; store in xe
76 |     // If too many zeros of the derivative  $(f - p)'$ 
77 |     // have been found, select those, where the deviation is maximal
78 |     if (nx == d) {
79 |         xe = VectorXd(xx0.size() + 2);
80 |         xe << a, xx0, b;
81 |     } else if (nx == d + 1) {
82 |         xe = VectorXd(xx0.size() + 1);
83 |         if (xx0.minCoeff() - a > b - xx0.maxCoeff())
84 |             xe << a, xx0;
85 |     } else
86 |         xe << xx0, b;
87 |     } else if (nx == d + 2) {
88 |         xe = xx0;
89 |     } else {
90 |         VectorXd del = (polyval(c.head(d + 1), xx0) - feval(f, xx0)).cwiseAbs(),
91 |         ind = sort_indices(del);
92 |         xe = select(xx0, ind.tail(d + 2));
93 |     }
94 |
95 |     // Deviation in sampling points and approximate alternants
96 |     fxe = feval(f, xe);
97 |     VectorXd del(xe.size() + 2);
98 |     del << polyval(c.head(d + 1), a * VectorXd::Ones(1)) -
99 |         ftab(0) * VectorXd::Ones(1),
100 |         polyval(c.head(d + 1), xe) - fxe,
101 |         polyval(c.head(d + 1), b * VectorXd::Ones(1)) -
102 |         ftab(ftab.size() - 1) * VectorXd::Ones(1);
103 |     // Approximation of supremum norm of approximation error
104 |     const double dev = del.cwiseAbs().maxCoeff();
105 |     // Termination of Remez iteration
106 |     if (dev < tol * fsupn) {
107 |         break;
108 |     }

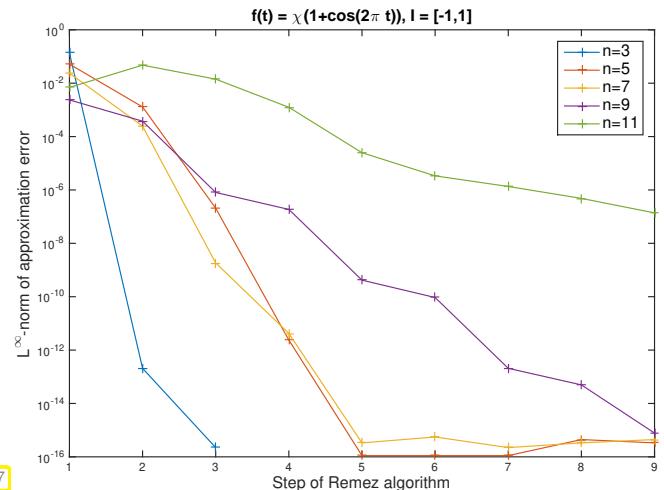
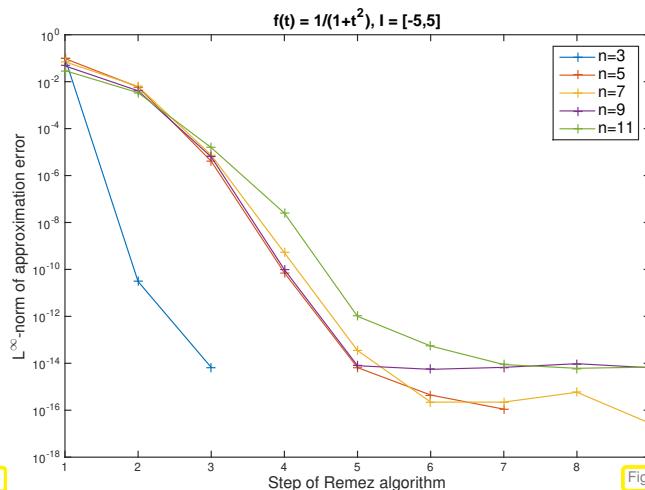
```

```

109 }
110 VectorXd tmp = c.head(d + 1);
111 c = tmp;
112 }
```

**EXPERIMENT 6.3.0.8 (Convergence of Remez algorithm)** We examine the convergence of the Remez algorithm from Code 6.3.0.7 for two different functions:

- $f(t) = (1+t^2)^{-1}$ ,  $I = [-5, 5] \rightarrow$  Bsp. 6.1.2.10
- $f(t) = \begin{cases} \frac{1}{2}(1+\cos(2\pi t)) & , \text{ if } |t| < \frac{1}{2} \\ 0 & \text{else.} \end{cases}$ ,  $I = [-1, 1]$



Convergence in both cases; faster convergence observed for smooth function, for which machine precision is reached after a few steps.

## 6.4 Approximation by Trigonometric Polynomials

Task: approximation of a continuous **1-periodic** function

$$f \in C^0(\mathbb{R}) , \quad f(t+1) = f(t) \quad \forall t \in \mathbb{R} .$$

Dubious: approximation by global polynomials, because they are not 1-periodic  
 ➤ the global polynomial Lagrange interpolant will lack an essential structural property.

The natural space for approximating **generic periodic** functions is a space of **trigonometric polynomials** with the same period.

Remember from Def. 5.6.1.1: space of 1-periodic trigonometric polynomials of degree  $2n$ ,  $n \in \mathbb{N}$ ,

$$\mathcal{P}_{2n}^T = \text{Span} \left\{ t \mapsto 1, t \mapsto \sin(2\pi t), t \mapsto \cos(2\pi t), t \mapsto \sin(4\pi t), t \mapsto \cos(4\pi t), \dots, t \mapsto \sin(2\pi nt), t \mapsto \cos(2\pi nt) \right\} \quad (6.4.0.1a)$$

$$= \text{Span} \{ t \mapsto \exp(2\pi i k t) : k = -n, \dots, n \} . \quad (6.4.0.1b)$$

Both sets of functions provide a basis of  $\mathcal{P}_{2n}^T$  (when considered as a space of  $\mathbb{C}$ -valued functions on  $\mathbb{R}$ ).

### 6.4.1 Approximation by Trigonometric Interpolation



Idea: Adapt the policy of **approximation by interpolation** from § 6.0.0.6

Now: Employ **trigonometric interpolation** from Section 5.6  
into space  $\mathcal{P}_{2n}^T$  of 1-periodic **trigonometric polynomials** → Def. 5.6.1.1.

#### Recall: Trigonometric interpolation → Section 5.6

Given nodes  $t_0 < t_1 < \dots < t_{2n}$ ,  $t_k \in [0, 1[$ , and values  $y_k \in \mathbb{R}$ ,  $k = 0, \dots, 2n$  find

$$q \in \mathcal{P}_{2n}^T := \text{Span}\{t \mapsto \cos(2\pi jt), t \mapsto \sin(2\pi jt)\}_{j=0}^n, \quad (6.4.1.2)$$

$$\text{with } q(t_k) = y_k \text{ for all } k = 0, \dots, 2n. \quad (6.4.1.3)$$

Terminology:  $\mathcal{P}_{2n}^T \hat{=} \text{space of trigonometric polynomials of degree } 2n$ .

From Section 5.6 remember a few more facts about trigonometric polynomials and trigonometric interpolation:

- ◆ Cor. 5.6.1.6: Dimension of the space of trigonometric polynomials:  $\dim \mathcal{P}_{2n}^T = 2n + 1$
- ◆ Trigonometric interpolation can be reduced to polynomial interpolation on the unit circle  $S^1 \subset \mathbb{C}$  in the complex plane, see (5.6.1.5).
  - ▶ existence & uniqueness of trigonometric interpolant  $q$  satisfying (6.4.1.2) and (6.4.1.3)
- ◆ Code 5.6.3.3: Efficient FFT-based algorithms for trigonometric interpolation in **equidistant** nodes  $t_k = \frac{k}{2n+1}$ ,  $k = 0, \dots, 2n$ .

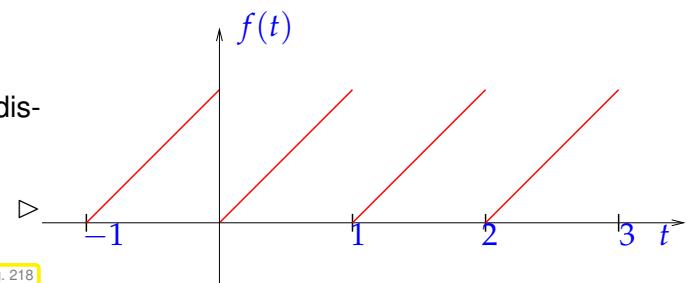
The relationship of trigonometric interpolation and polynomial interpolation on the unit circle suggests a uniform distribution of nodes for general trigonometric interpolation.

#### Nodes for function approximation by trigonometric interpolation

Trigonometric approximation of generic **1-periodic** continuous functions  $\in C^0(\mathbb{R})$  in  $\mathcal{P}_{2n}^T$  usually relies on **equidistant** interpolation nodes  $t_k = \frac{k}{2n+1}$ ,  $k = 0, \dots, 2n$ .

notation: **trigonometric interpolation operator** in  $2n + 1$  equidistant nodes  $t_k = \frac{k}{2n+1}$ ,  $k = 0, \dots, 2n$

$$\mathbf{T}_n : C^0([0, 1]) \rightarrow \mathcal{P}_{2n}^T, \quad \mathbf{T}_n(f)(t_k) = f(t_k) \quad \forall k \in \{0, \dots, 2n\}. \quad (6.4.1.5)$$



Note that a function  $f \in C^0([0, 1])$  can spawn a discontinuous 1-periodic function on  $\mathbb{R}$

A prominent example is the “sawtooth function”

## 6.4.2 Trigonometric Interpolation Error Estimates

Our focus will be on the asymptotic behavior of

$$\|f - T_n f\|_{L^\infty([0,1])} \quad \text{and} \quad \|f - T_n f\|_{L^2([0,1])} \quad \text{as} \quad n \rightarrow \infty,$$

for functions  $f : [0,1] \rightarrow \mathbb{C}$  with different smoothness properties. To begin with we perform an empirical study.

**EXPERIMENT 6.4.2.1 (Interpolation error: trigonometric interpolation)** Now we study the asymptotic behavior of the error of equidistant trigonometric interpolation as  $n \rightarrow \infty$  in a numerical experiment for functions with different smoothness properties.

#1 Step function:  $f(t) = 0$  for  $|t - \frac{1}{2}| > \frac{1}{4}$ ,  $f(t) = 1$  for  $|t - \frac{1}{2}| \leq \frac{1}{4}$

#2  $C^\infty$  periodic function:  $f(t) = \frac{1}{\sqrt{1 + \frac{1}{2} \sin(2\pi t)}}$ .

#3 “wedge function”:  $f(t) = |t - \frac{1}{2}|$

Approximate computation of norms of interpolation errors on equidistant grid with 4096 points.

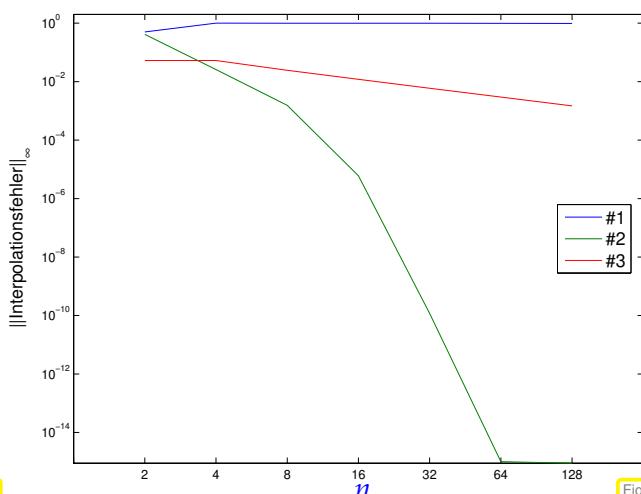


Fig. 220

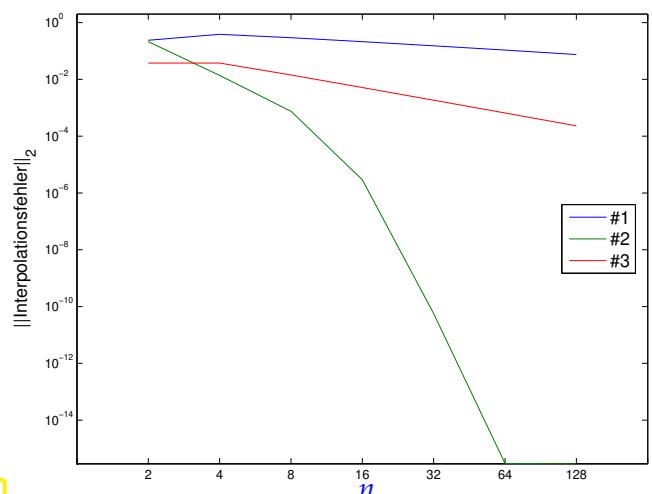


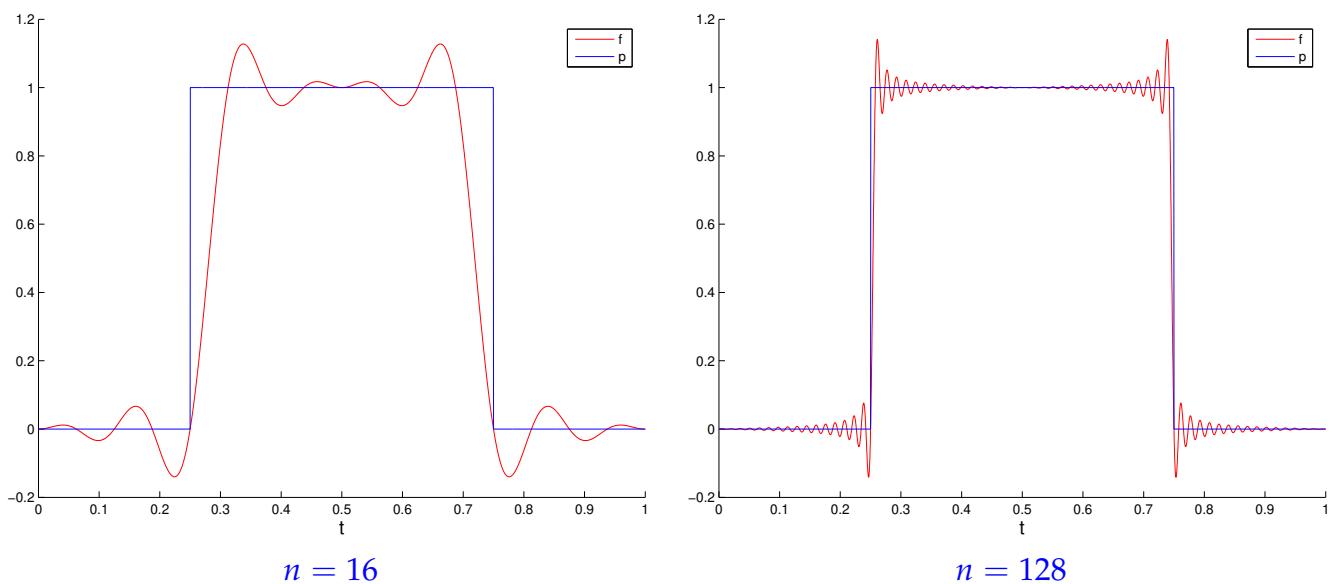
Fig. 220

Observation:

- Function #1: no convergence in  $L^\infty$ -norm, algebraic convergence in  $L^2$ -norm
- Function #3: algebraic convergence in both norms
- Function #2: exponential convergence in both norms

We conclude that in this experiment higher smoothness of  $f$  leads to faster convergence of the trigonometric interpolant.  $\square$

**EXPERIMENT 6.4.2.2 (Gibbs phenomenon)** Of course the smooth trigonometric interpolants of step function fail to converge in  $L^\infty$ -norm in Exp. 6.4.2.1. Moreover, they will not even converge “visually” to the step function, which becomes manifest by a closer inspection of the interpolants.

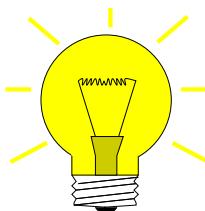


Observation: overshooting in neighborhood of discontinuity: **Gibbs phenomenon**

**§6.4.2.3 (Fourier series** → (4.2.5.6)) From (6.4.0.1b) we learn that the space of trigonometric polynomials  $\mathcal{P}_{2n}^T$  is spanned by the  **$2n + 1$  Fourier modes**  $t \mapsto \exp(2\pi i kt)$  of “lowest frequency”, that is, for  $-n \leq k \leq n$ . In Section 4.2.5 we learned that every function  $f : [0, 1] \rightarrow \mathbb{C}$  with finite  $L^2([0, 1])$ -norm can be expanded in a Fourier series (→ Thm. 4.2.5.28)

$$f(t) = \sum_{k \in \mathbb{Z}} \widehat{f}_k \exp(-2\pi i kt) \quad \text{in } L^2([0, 1]) , \quad \widehat{f}_k := \int_0^1 f(t) \exp(2\pi i kt) dt .$$

(A limit in  $L^2([0, 1])$  means that  $\left\| f - \sum_{k=-M}^M \widehat{f}_k \exp(-2\pi i k \cdot) \right\|_{L^2([0, 1])} \rightarrow 0$  for  $M \rightarrow \infty$ . Also remember the notation  $\widehat{f}_k$  for the **Fourier coefficients**.)



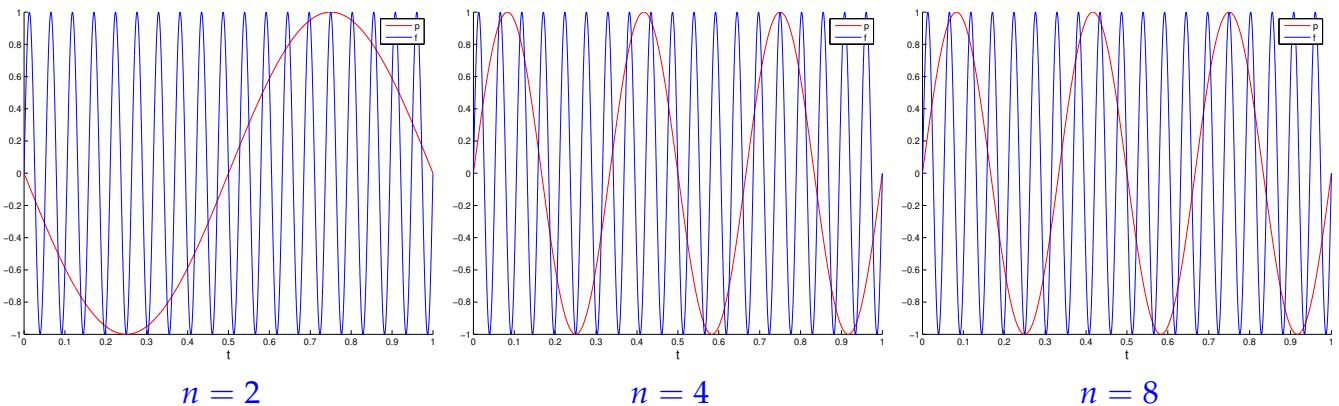
Idea: Study trigonometric interpolation error for interpolants in Fourier series representation.

**§6.4.2.4 (Aliasing)** We study the action of the trigonometric interpolation operator  $T_n$  from (6.4.1.5) on individual Fourier modes  $\mu_k(t) := \exp(-2\pi k t)$ ,  $t \in \mathbb{R}$ ,  $k \in \mathbb{Z}$ . Due to the 1-periodicity of  $t \mapsto \exp(2\pi i t)$  we find for every node  $t_j := \frac{j}{2n+1}$ ,  $j = 0, \dots, 2n$ ,

$$\mu_k(t_j) = \exp(-2\pi i k \frac{j}{2n+1}) = \exp(-2\pi i (k - \ell(2n+1)) \frac{j}{2n+1}) = \mu_{k-\ell(2n+1)}(t_j) \quad \forall \ell \in \mathbb{Z} .$$

When sampled on the node set  $\mathcal{T}_n := \{t_0, \dots, t_{2n}\}$  all the Fourier modes  $\mu_{k-\ell(2n+1)}$ ,  $\ell \in \mathbb{Z}$ , yield the same values. Thus the trigonometric cannot distinguish them! This phenomenon is called **aliasing**.

Aliasing demonstrated for  $f(t) = \sin(2\pi \cdot 19t) = \text{Im}(\exp(2\pi i 19t))$  for different node sets.



The “low-frequency” sine waves plotted in red coincide with  $f$  on the node set  $\mathcal{T}_n$ .

Since  $\mathbf{T}_n \mu_k = \mu_k$  for  $k = -n, \dots, n$ , that is, for  $\mu_k \in \mathcal{P}_{2n}^T$ , the aliasing effect yields

$$\mathbf{T}_n \mu_k = \mu_{\tilde{k}} \quad , \quad \tilde{k} \in \{-n, \dots, n\} , \quad k - \tilde{k} \in (2n+1)\mathbb{Z} \quad [\tilde{k} := k \pmod{(2n+1)}] . \quad (6.4.2.5)$$

$(\tilde{n} = n, \tilde{n+1} = -n, \tilde{-n} = -n, \tilde{-n-1} = n, \tilde{2n} = -1, \text{etc.})$

► Trigonometric interpolation by  $\mathbf{T}_n$  folds all Fourier modes (“frequencies”) to the finite range  $\{-n, \dots, n\}$ .

↓

From (6.4.2.5), by linearity of  $\mathbf{T}_n$ , we obtain for  $f : [0, 1] \rightarrow \mathbb{C}$  in Fourier series representation

$$f(t) = \sum_{j=-\infty}^{\infty} \hat{f}_j \mu_j(t) \quad \blacktriangleright \quad \mathbf{T}_n(f)(t) = \sum_{j=-n}^n \gamma_j \mu_j(t) , \quad \gamma_j = \sum_{\ell=-\infty}^{\ell} \hat{f}_{j+\ell(2n+1)} . \quad (6.4.2.6)$$

We can read the trigonometric polynomial  $\mathbf{T}_n f \in \mathcal{P}_{2n}^T$  as a Fourier series with non-zero coefficients only in the index range  $\{-n, \dots, n\}$ . Thus, for the Fourier coefficients of the trigonometric interpolation error  $E(t) := f(t) - \mathbf{T}_n f(t)$  we find from (6.4.2.6)

$$\widehat{E}_j = \begin{cases} - \sum_{\ell \in \mathbb{Z} \setminus \{0\}} \hat{f}_{j+\ell(2n+1)} & , \text{ if } j \in \{-n, \dots, n\} , \\ \hat{f}_j & , \text{ if } |j| > n , \end{cases} \quad j \in \mathbb{Z} . \quad (6.4.2.7)$$

Since we have (sufficient smoothness of  $f$  assumed)

$$E(t) := f(t) - \mathbf{T}_n f(t) = \sum_{j=-\infty}^{\infty} \widehat{E}_j e^{-2\pi i jt} , \quad (6.4.2.8)$$

we conclude, in particular from the isometry property asserted in Thm. 4.2.5.28,

$$\|f - \mathbf{T}_n f\|_{L^2([0,1])}^2 = \sum_{j=-n}^n \left| \sum_{\ell \in \mathbb{Z} \setminus \{0\}} \hat{f}_{j+\ell(2n+1)} \right|^2 + \sum_{|j| > n} |\hat{f}_j|^2 , \quad (6.4.2.9)$$

$$\|f - \mathbf{T}_n f\|_{L^\infty([0,1])} \leq \sum_{j=-n}^n \sum_{\ell \in \mathbb{Z} \setminus \{0\}} |\hat{f}_{j+\ell(2n+1)}| + \sum_{|j| > n} |\hat{f}_j| . \quad (6.4.2.10)$$

In order to estimate these norms of the trigonometric interpolation error we need quantitative information about the decay of the Fourier coefficients  $\hat{f}_j$  as  $|j| \rightarrow \infty$ .

**§6.4.2.11 (Fourier expansions of derivatives)** For 1-periodic  $c \in C^0(\mathbb{R})$  with integrable derivative  $\dot{c} := \frac{d}{dt}c$  we find by integration by parts (boundary terms cancel due to periodicity)

$$\widehat{c}_j = \int_0^1 c(t) e^{2\pi i jt} dt = -2\pi i j \cdot \int_0^1 \dot{c}(t) e^{2\pi i jt} dt = (-2\pi i j) \widehat{\dot{c}}_j, \quad j \in \mathbb{Z}.$$

We can also arrive at this formula by (formal) term-wise differentiation of the Fourier series:

$$c(t) = \sum_{j=-\infty}^{\infty} \widehat{c}_j e^{-2\pi i jt} \implies \dot{c}(t) = \sum_{j=-\infty}^{\infty} \underbrace{(-2\pi i j) \widehat{c}_j}_{=\widehat{c}_j} e^{-2\pi i jt}. \quad (6.4.2.12)$$

### Lemma 6.4.2.13. Fourier coefficients of derivatives

For the Fourier coefficients of the derivatives a 1-periodic function  $f \in C^{k-1}(\mathbb{R})$ ,  $k \in \mathbb{N}$ , with integrable  $k$ -th derivative  $f^{(k)}$  holds

$$\widehat{f^{(k)}}_j = (-2\pi i)^k \widehat{f}_j, \quad j \in \mathbb{Z}.$$

□

**§6.4.2.14 (Fourier coefficients and smoothness)** From Lemma 6.4.2.13 and the trivial estimates ( $|\exp(2\pi i t)| = 1$ )

$$|\widehat{f}_j| \leq \int_0^1 |f(t)| dt \leq \|f\|_{L^\infty([0,1])} \quad \forall j \in \mathbb{Z}, \quad (6.4.2.15)$$

we conclude that  $((2\pi|j|)^m \widehat{f}_j)_{j \in \mathbb{Z}}$ ,  $m \in \mathbb{N}$ , is bounded, provided that  $f \in C^{m-1}(\mathbb{R})$  with integrable  $m$ -th derivative

### Lemma 6.4.2.16. Decay of Fourier coefficients

If  $f \in C^{k-1}(\mathbb{R})$  with integrable  $k$ -th derivative, then  $\widehat{f}_j = O(j^{-k})$  for  $|j| \rightarrow \infty$

### Decay of Fourier coefficients and smoothness

The smoother a periodic function the faster the decay of its Fourier coefficients

The isometry property of Thm. 4.2.5.28 also yields for  $f \in C^{k-1}(\mathbb{R})$  with  $f^{(k)} \in L^2(\text{unitintv})$  that

$$\|f^{(k)}\|_{L^2([0,1])}^2 = (2\pi)^k \sum_{j=-\infty}^{\infty} j^{2k} |\widehat{f}_j|^2. \quad (6.4.2.18)$$

□

We can now combine the identity (6.4.2.18) with (6.4.2.9) and obtain an interpolation error estimate in  $L^2([0,1])$ -norm.

**Theorem 6.4.2.19.**  $L^2$ -error estimate for trigonometric interpolation

If  $f \in C^{k-1}(\mathbb{R})$ ,  $k \in \mathbb{N}$ , with square-integrable  $k$ -th derivative ( $f^{(k)} \in L^2([0, 1])$ ), then

$$\|f - T_n f\|_{L^2([0, 1])} \leq \sqrt{1 + c_k} n^{-k} \|f^{(k)}\|_{L^2([0, 1])}, \quad (6.4.2.20)$$

with  $c_k = 2 \sum_{\ell=1}^{\infty} (2\ell - 1)^{-2k} < \infty$ .

Thm. 6.4.2.19 confirms *algebraic convergence* of the  $L^2$ -norm of the trigonometric interpolation error for functions with *limited smoothness*. Higher rates can be expected for smoother functions, which we have also found in cases #1 and #3 in Exp. 6.4.2.1.

### 6.4.3 Trigonometric Interpolation of Analytic Periodic Functions

In § 6.1.2.31 we saw that we can expect exponential decay of the maximum norm of polynomial interpolation errors in the case of “very smooth” interpolands. To capture this property of functions we resorted to the notion of *analytic* functions, as defined in Def. 6.1.2.32. Since trigonometric interpolation is closely connected to polynomial interpolation (on the unit circle  $S^1$ , see Section 5.6.2), it is not surprising that analyticity of interpolands will also involve exponential convergence of trigonometric interpolants. This result will be established in this section.

In case #2 of Exp. 6.4.2.1 we already saw an instance of exponential convergence for an analytic interpoland. A more detailed study follows.

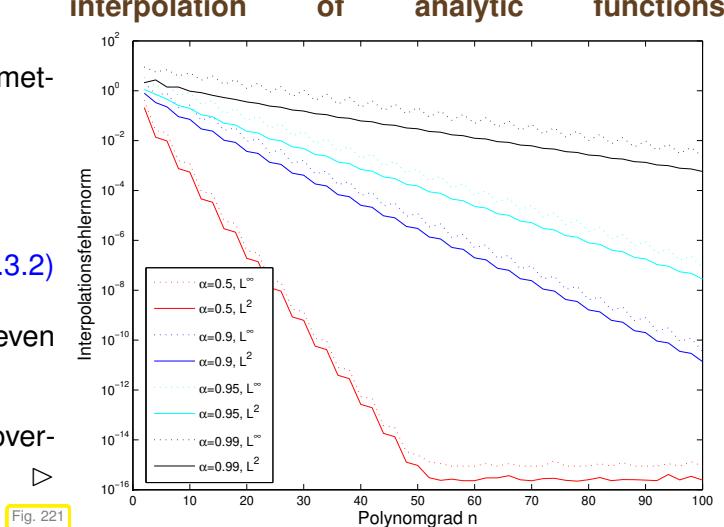
#### EXPERIMENT 6.4.3.1 (Trigonometric interpolation of analytic functions)

We study the convergence of equidistant trigonometric interpolation for the interpoland

$$f(t) = \frac{1}{\sqrt{1 - \alpha \sin(2\pi t)}} \quad \text{on } I = [0, 1]. \quad (6.4.3.2)$$

For  $0 \leq \alpha < 1$  we have  $f \in C^\infty(\mathbb{R})$ , and  $f$  is even analytic ( $\rightarrow$  Def. 6.1.2.32).

Approximative computations of error norms by “oversampling” in 4096 points.  $\triangleright$



$\gg$  Observation: exponential convergence in  $n$ , faster for smaller  $\alpha$   $\triangleright$

**§6.4.3.3 (Decay of Fourier coefficients of analytic functions)** Lemma 6.4.2.16 asserts algebraic decay of the Fourier coefficients of functions with limited smoothness. As analytic 1-periodic functions are “infinitely smooth”, they will always belong to  $C^\infty(\mathbb{R})$ , we expect a stronger result in this case. In fact, we can conclude *exponential decay* of the Fourier coefficients.

**Theorem 6.4.3.4. Exponential decay of Fourier coefficients of analytic functions**

If  $f : \mathbb{R} \rightarrow \mathbb{C}$  is 1-periodic and has an **analytic extension** to the strip

$$\tilde{S} := \{z \in \mathbb{C} : -\eta \leq \operatorname{Im} z \leq \eta\}, \quad \text{for some } \eta > 0,$$

then its Fourier coefficients decay according to

$$|\hat{f}_k| \leq q^{|k|} \cdot \|f\|_{L^\infty(\tilde{S})} \quad \forall k \in \mathbb{Z} \quad \text{with} \quad q := \exp(-2\pi\eta) \in ]0, 1[. \quad (6.4.3.5)$$

*Proof.* ①: A first variant of the proof uses techniques from complex analysis:

Let  $f : \mathbb{R} \mapsto \mathbb{C}$  be 1-periodic with analytic extension to the (closed) strip  $S := \{z \in \mathbb{C} : -\eta \leq \operatorname{Im} z \leq \eta\}, \eta > 0$ .

➤  $g_r(t) := f(t + ir)$  is 1-periodic and smooth for  $-\eta \leq r \leq \eta$

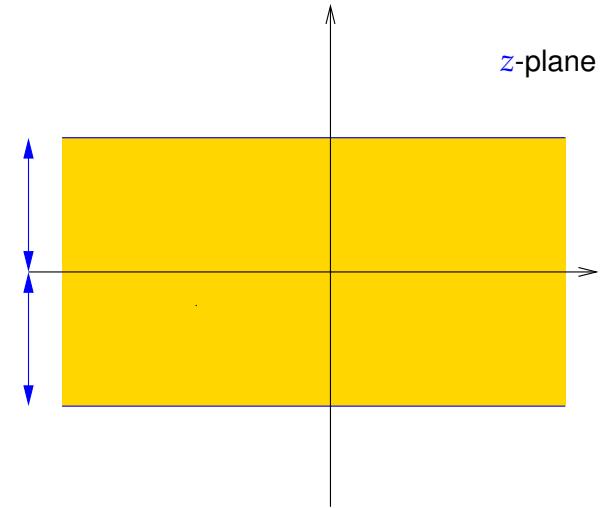


Fig. 222

We compute the Fourier coefficients of  $g_r$

$$\hat{g}_{rk} = \int_0^1 f(t + ir) e^{-2\pi i kt} dt = \int_0^1 f(t) e^{2\pi i k(t - ir)} dt = e^{-2\pi rk} \hat{f}_k,$$

and arrive at the bound

$$|\hat{f}_k| \leq e^{-2\pi rk} \underbrace{\max\{|f(t + iur)|, 0 \leq t \leq 1\}}_{\|g_r\|_{L^\infty([0,1])}} \quad \forall k \in \mathbb{Z}, \quad \forall -\eta^- \leq r \leq \eta^+.$$

The assertion of the theorem follows directly.

②: A second variant of the proof merely relies on classical calculus.

By the assumptions of the theorem the function  $f : \mathbb{R} \rightarrow \mathbb{C}$  can be expanded into a **power series** with radius of convergence  $\eta$  around at every  $y \in \mathbb{R}$

$$\forall y \in \mathbb{R} : \exists c_n(y) \in \mathbb{C} : f(x) = \sum_{n=0}^{\infty} c_n(y) (x - y)^n \quad \forall x : |x - y| \leq \eta. \quad (6.4.3.6)$$

This implies  $|c_n(y)|\eta^n \leq C(y)$  and, since  $[0, 1]$  is compact:

$$\exists C > 0 : |c_n(y)|\eta^n \leq C \quad \forall n \in \mathbb{N}_0, \forall y \in \mathbb{R}. \quad (6.4.3.7)$$

The power series (6.4.3.6) is a Taylor series, which means

$$c_n(y) = \frac{1}{n!} f^{(n)}(y) \stackrel{(6.4.3.7)}{\implies} |f^{(n)}(y)| \leq C n! \eta^{-n} \quad \forall n \in \mathbb{N}_0, \quad y \in \mathbb{R}. \quad (6.4.3.8)$$

We use this estimate to bound the Fourier coefficients of  $\hat{f}$ . Starting from the defining formula (4.2.5.18), we continue with  $n$ -fold integration by parts

$$\widehat{f}_k := \int_0^1 f(t) \exp(2\pi i k t) dt = \frac{(-1)^n}{(2\pi i k)^n} \int_0^1 f^{(n)}(t) \exp(2\pi i k t) dt. \quad (6.4.3.9)$$

We combine this formula with (6.4.3.8) and obtain

$$|\widehat{f}_k| \leq C \frac{n!}{(2\pi |k| \eta)^n} \quad \forall n \in \mathbb{N}, k \in \mathbb{Z}. \quad (6.4.3.10)$$

Next, we use Stirling's formula in the form

$$n! \leq e n^{n+1/2} e^{-n}, \quad n \in \mathbb{N},$$

which gives

$$|\widehat{f}_k| \leq C e \frac{n^{n+1/2}}{(2\pi |k|)^n} e^{-n} \quad \forall n \in \mathbb{N}, k \in \mathbb{Z}.$$

We can also “interpolate” and replace  $n$  with a real number.

$$\Rightarrow |\widehat{f}_k| \leq C e \frac{r^{r+1/2}}{(2\pi |k|)^r} e^{-r} \quad \forall r \geq 1, k \in \mathbb{Z}.$$

Finally, we set  $r := 2\pi |k| \eta$  and arrive at

$$|\widehat{f}_k| \leq C e \sqrt{2\pi k \eta} \underbrace{\exp(-2\pi \eta)}_{=:q}, \quad k \in \mathbb{Z}.$$

This bound is slightly less explicit than that obtain in ①.  $\square$

Knowing exponential decay of the Fourier coefficients, the geometric sum formula can be used to extract estimates from (6.4.2.9) and (6.4.2.10):

### Lemma 6.4.3.11. Interpolation error estimates for exponentially decaying Fourier coefficients

If  $f : \mathbb{R} \rightarrow \mathbb{C}$  is 1-periodic and  $f \in L^2([0, 1])$ , then

$$\exists M > 0, \rho \in ]0, 1[ : |\widehat{f}(k)| \leq M \rho^{|k|} \Rightarrow \begin{aligned} \|f - T_n(f)\|_{L^2([0, 1])} &\leq M \frac{\rho^{n/2}}{1 - \rho^n} \frac{2\sqrt{2}}{\sqrt{1 - \rho^2}}, \\ \|f - T_n(f)\|_{L^\infty([0, 1])} &\leq 4M \frac{\rho^{n/2}}{1 - \rho}. \end{aligned}$$

This estimate can be combined with the result of Thm. 6.4.3.4 and gives the main result of this section:

**Theorem 6.4.3.12. Exponential convergence of trigonometric interpolation for analytic interpolands**

If  $f : \mathbb{R} \rightarrow \mathbb{C}$  is 1-periodic and possesses an *analytic extension* to the strip

$$\bar{S} := \{z \in \mathbb{C} : -\eta \leq \operatorname{Im} z \leq \eta\}, \quad \text{for some } \eta > 0,$$

then there is  $C_\eta > 0$  depending only on  $\eta$  such that

$$\|f - T_n f\|_* \leq C_\eta e^{-\pi \eta n} \|f\|_{L^\infty(\bar{S})}, \quad n \in \mathbb{N}, \quad (* = L^2([0,1]), L^\infty([0,1])) . \quad (6.4.3.13)$$

The speed of exponential convergence clearly depends on the width  $\eta$  of the “strip of analyticity”  $\bar{S}$ .

**§6.4.3.14 (Convergence of trigonometric interpolation for analytic interpolands)**

Explanation of the observation made in Exp. 6.4.3.1, cf. Rem. 6.1.3.22:

Similar to Chebychev interpolants, also trigonometric interpolants converge exponentially fast, if the interpoland  $f$  is 1-periodic *analytic* ( $\rightarrow$  Def. 6.1.2.32) in a strip around the real axis in  $\mathbb{C}$ , see Thm. 6.4.3.12 for details.

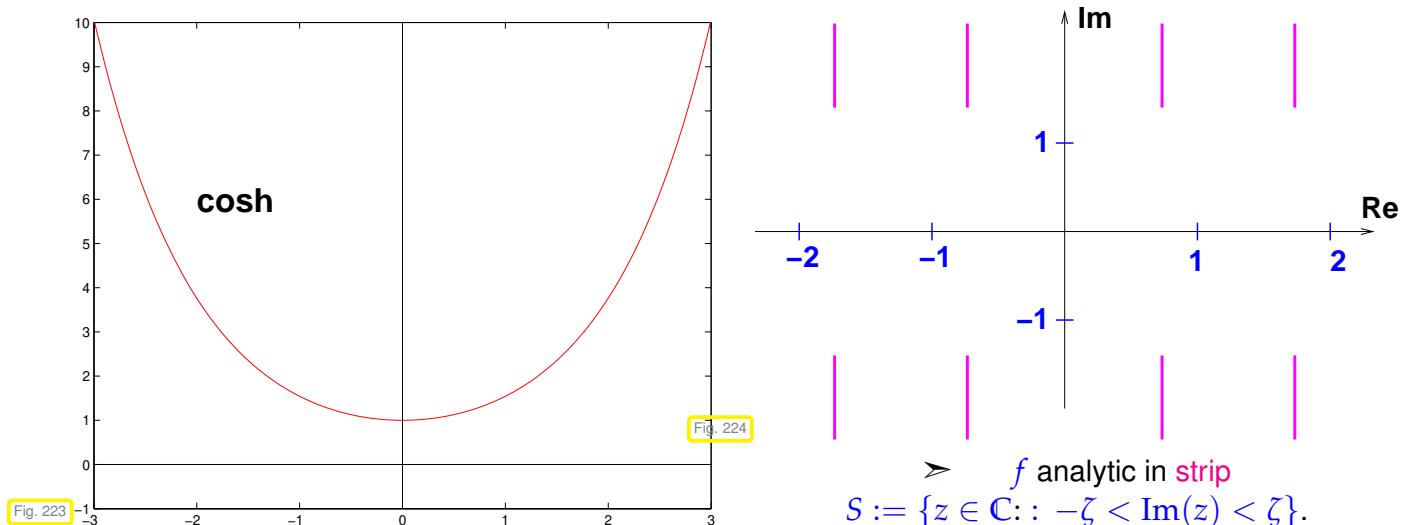
Note that  $f$  from (6.4.3.2) is holomorphic, where

$$1 + \alpha \sin(2\pi z) \notin \mathbb{R}_0^- \Leftrightarrow \sin(2\pi z) = \sin(2\pi x) \cosh(2\pi y) + i \cos(2\pi x) \sinh(2\pi y) \notin [-\infty, -1 - \frac{1}{\alpha}] ,$$

because the square root is holomorphic in  $\mathbb{C} \setminus \mathbb{R}_0^-$ .

Domain of analyticity of  $f$ :

$$\mathbb{C} \setminus \bigcup_{k \in \mathbb{Z}} \left( \frac{k}{2} + \frac{1}{4} + i(\mathbb{R} \setminus [-\zeta, \zeta]) \right), \quad \zeta \in \mathbb{R}^+, \quad \cosh(2\pi\zeta) = 1 + \frac{1}{\alpha} .$$



> As  $\alpha$  decreases the strip of analyticity becomes wider, since  $x \rightarrow \cosh(x)$  is increasing for  $x > 0$ .  $\square$

## 6.5 Approximation by Piecewise Polynomials

Recall some alternatives to interpolation by global polynomials discussed in Chapter 5:

- ◆ piecewise linear/quadratic interpolation → Section 5.3.2, Ex. 5.3.2.4,
  - ◆ cubic Hermite interpolation → Section 5.4,
  - ◆ (cubic) spline interpolation → Section 5.5.1.
- ☞ All these interpolation schemes rely on **piecewise polynomials** (of different global smoothness)

Focus in this section: function approximation by *piecewise polynomial* interpolants

**§6.5.0.1 (Grid/mesh)** The attribute “piecewise” refers to *partitioning of the interval* on which we aim to approximate. In the case of data interpolation the natural choice was to use intervals defined by interpolation nodes. Yet we already saw exceptions in the case of shape-preserving interpolation by means of quadratic splines, see Section 5.5.3.

In the case of function approximation based on an interpolation scheme the additional freedom to choose the interpolation nodes suggests that those be decoupled from the partitioning.



Idea: use **piecewise polynomials** with respect to a **grid/mesh**

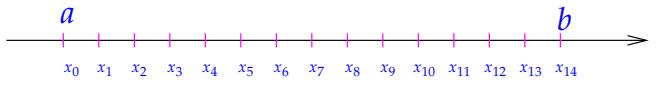
$$\mathcal{M} := \{a = x_0 < x_1 < \dots < x_{m-1} < x_m = b\} \quad (6.5.0.2)$$

to approximate function  $f : [a, b] \mapsto \mathbb{R}, a < b$ .

Borrowing from terminology for splines, *cf.* Def. 5.5.0.1, the underlying mesh for piecewise polynomial approximation is sometimes called the “knot set”.

Terminology:

- ◆  $x_j \hat{=} \text{nodes}$  of the mesh  $\mathcal{M}$ ,
- ◆  $[x_{j-1}, x_j[ \hat{=} \text{intervals/cells}$  of the mesh,
- ◆  $h_{\mathcal{M}} := \max_j |x_j - x_{j-1}| \hat{=} \text{mesh width}$ ,
- ◆ If  $x_j = a + jh \hat{=} \text{equidistant}$  (uniform) mesh with meshwidth  $h > 0$



**Remark 6.5.0.3 (Local approximation by piecewise polynomials)** We will see that most approximation schemes relying on piecewise polynomials are local in the sense that finding the approximant on a cell of the mesh relies only on a fixed number of function evaluations in a neighborhood of the cell.

- $\mathcal{O}(1)$  computational effort to find interpolant on  $[x_{j-1}, x_j]$  (independent of  $m$ )
- $\mathcal{O}(m)$  computational effort to determine piecewise polynomial approximant for  $m \rightarrow \infty$  ('fine meshes')

Contrast this with the computational cost of computing global polynomial interpolants, which will usually be  $\mathcal{O}(n^2)$  for polynomial degree  $n \rightarrow \infty$ . □

### 6.5.1 Piecewise Polynomial Lagrange Interpolation

Given: interval  $[a, b] \subset \mathbb{R}$  endowed with mesh  $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_{m-1} < x_m = b\}$ .

Recall theory of polynomial interpolation  $\rightarrow$  Section 5.2.2:  $n+1$  data points needed to fix interpolating polynomial, see Thm. 5.2.2.7.

► Approach to *local* Lagrange interpolation ( $\rightarrow$  (5.2.2.2)) of  $f \in C(I)$  on mesh  $\mathcal{M}$

#### General local Lagrange interpolation on a mesh (PPLIP)

- ① Choose *local degree*  $n_j \in \mathbb{N}_0$  for each cell of the mesh,  $j = 1, \dots, m$ .
- ② Choose set of *local* interpolation nodes

$$\mathcal{T}^j := \{t_0^j, \dots, t_{n_j}^j\} \subset I_j := [x_{j-1}, x_j], \quad j = 1, \dots, m,$$

for each mesh cell/grid interval  $I_j$ .

- ③ Define *piecewise polynomial* (PP) interpolant  $s : [x_0, x_m] \rightarrow \mathbb{K}$ :

$$s_j := s|_{I_j} \in \mathcal{P}_{n_j} \quad \text{and} \quad s_j(t_i^j) = f(t_i^j) \quad i = 0, \dots, n_j, \quad j = 1, \dots, m. \quad (6.5.1.2)$$

Owing to Thm. 5.2.2.7,  $s_j$  is well defined.

#### Corollary 6.5.1.3. Piecewise polynomials Lagrange interpolation operator

The mapping  $f \mapsto s$  defines a *linear* operator  $\mathbf{L}_{\mathcal{M}} : C^0([a, b]) \mapsto C_{\mathcal{M}, \text{pw}}^0([a, b])$  in the sense of Def. 5.1.0.19.

Obviously,  $\mathbf{L}_{\mathcal{M}}$  depends on  $\mathcal{M}$ , the local degrees  $n_j$ , and the sets  $\mathcal{T}_j$  of local interpolation points (the latter two are suppressed in notation).

#### Corollary 6.5.1.4. Continuous local Lagrange interpolants

If the local degrees  $n_j$  are at least 1 and the local interpolation nodes  $t_k^j, j = 1, \dots, m, k = 0, \dots, n_j$ , for local Lagrange interpolation satisfy

$$t_{n_j}^j = t_0^{j+1} \quad \forall j = 1, \dots, m-1 \quad \Rightarrow \quad s \in C^0([a, b]), \quad (6.5.1.5)$$

then the piecewise polynomial Lagrange interpolant according to (6.5.1.2) is *continuous* on  $[a, b]$ :  $s \in C^0([a, b])$ .

Focus:

*asymptotic* behavior of (some norm of) interpolation error

$$\|f - \mathbf{L}_{\mathcal{M}}f\| \leq CT(N) \quad \text{for } N \rightarrow \infty, \quad (6.5.1.6)$$

where  $N := \sum_{j=1}^m (n_j + 1)$ .

The decay of the bound  $T(N)$  will characterize the type of convergence:

- algebraic convergence or exponential convergence, see Section 6.1.2, Def. 6.1.2.7.

But why do we choose this strange number  $N$  as parameter when investigating the approximation error?

Because, by Thm. 5.2.1.2, it agrees with the dimension of the space of discontinuous, piecewise polynomials functions

$$\{q : [a, b] \rightarrow \mathbb{R} : q|_{I_j} \in \mathcal{P}_{n_j} \forall j = 1, \dots, m\} !$$

This dimension tells us the number of real parameters we need to describe the interpolant  $s$ , that is, the “information cost” of  $s$ .  $N$  is also proportional to the number of interpolation conditions, which agrees with the number of  $f$ -evaluations needed to compute  $s$  (why only proportional in general?).

Special case: uniform polynomial degree  $n_j = n$  for all  $j = 1, \dots, m$ .

► Then we may aim for estimates  $\|f - I_M f\| \leq CT(h_M)$  for  $h_M \rightarrow 0$

in terms of meshwidth  $h_M$ .

Terminology: investigations of this kind are called the study of  **$h$ -convergence**.

### EXAMPLE 6.5.1.7 ( $h$ -convergence of piecewise polynomial interpolation)

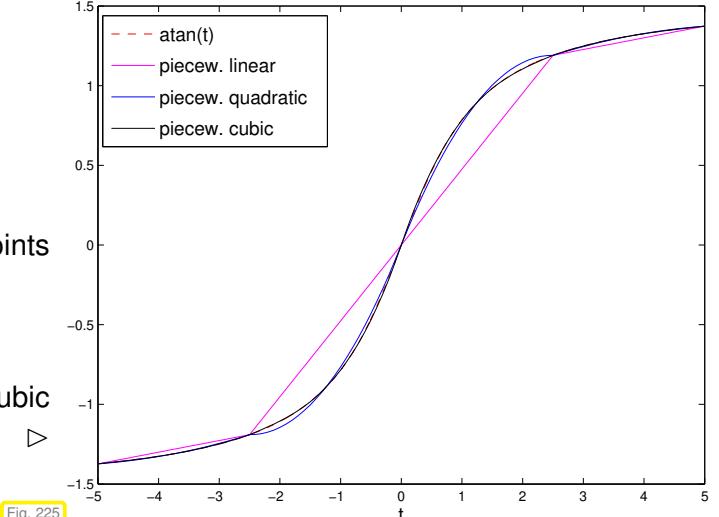
Compare Exp. 5.3.1.6:

$$f(t) = \arctan t, I = [-5, 5]$$

$$\text{Grid } \mathcal{M} := \{-5, -\frac{5}{2}, 0, \frac{5}{2}, 5\}$$

Local interpolation nodes equidistant in  $I_j$ , endpoints included, (6.5.1.5) satisfied.

Plots of the piecewise linear, quadratic and cubic polynomial interpolants

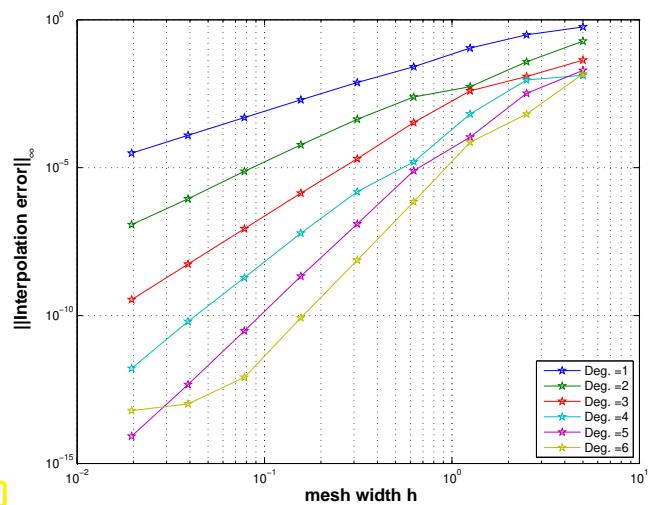
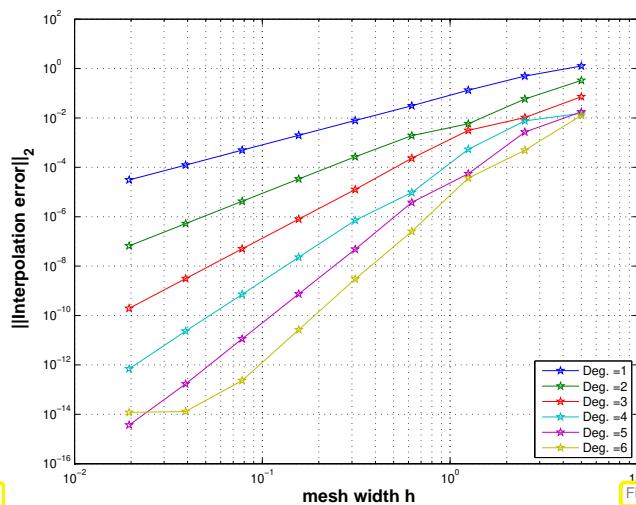


- Sequence of (equidistant) meshes:  $\mathcal{M}_i := \{-5 + j 2^{-i} 10\}_{j=0}^{2^i}, i = 1, \dots, 6$ .
- Equidistant local interpolation nodes (endpoints of grid intervals included).

Monitored: interpolation error in (approximate)  $L^\infty$ - and  $L^2$ -norms, see (6.1.3.21), (6.1.3.20)

$$\|g\|_{L^\infty([-5, 5])} \approx \max_{j=0, \dots, 1000} |g(-5 + j/100)|,$$

$$\|g\|_{L^2([-5, 5])} \approx \sqrt{\frac{1}{1000}} \cdot \left( \frac{1}{2} g(-5)^2 + \sum_{j=1}^{999} |g(-5 + j/100)|^2 + \frac{1}{2} g(5)^2 \right)^{1/2}.$$



Observation: Algebraic convergence ( $\rightarrow$  Def. 6.1.2.7) for meshwidth  $h_M \rightarrow 0$

(nearly linear error norm graphs in doubly logarithmic scale, see Rem. 6.1.2.9)

Observation: rate of algebraic convergence increases with polynomial degree  $n$

Rates  $\alpha$  of algebraic convergence  $O(h_M^\alpha)$  of norms of interpolation error:

$n$	1	2	3	4	5	6
w.r.t. $L^2$ -norm	1.9957	2.9747	4.0256	4.8070	6.0013	5.2012
w.r.t. $L^\infty$ -norm	1.9529	2.8989	3.9712	4.7057	5.9801	4.9228

- Higher polynomial degree provides faster algebraic decrease of interpolation error norms. Empiric evidence for rates  $\alpha = p + 1$

Here: rates estimated by linear regression ( $\rightarrow$  Ex. 3.1.1.3) based on PYTHON's `polyfit` and the interpolation errors for meshwidth  $h \leq 10 \cdot 2^{-5}$ . This was done in order to avoid erratic "preasymptotic", that is, for large meshwidth  $h$ , behavior of the error.

The bad rates for  $n = 6$  are probably due to the impact of roundoff, because the norms of the interpolation error had dropped below machine precision, see Fig. 226, 227. □

**§6.5.1.8 (Approximation error estimates for piecewise polynomial Lagrange interpolation)** The observations made in Ex. 6.5.1.7 are easily explained by applying the polynomial interpolation error estimates of Section 6.1.2, for instance

$$\|f - \mathcal{L}_T f\|_{L^\infty(I)} \leq \frac{\|f^{(n+1)}\|_{L^\infty(I)}}{(n+1)!} \max_{t \in I} |(t - t_0) \cdots (t - t_n)|. \quad (6.1.2.19)$$

locally on the mesh intervals  $[x_{j-1}, x_j]$ ,  $j = 1, \dots, m$ : for constant polynomial degree  $n = n_j$ ,  $j = 1, \dots, m$ , we get for  $f \in C^{n+1}([x_0, x_m])$  (smoothness requirement)

$$(6.1.2.19) \Rightarrow \|f - s\|_{L^\infty([x_0, x_m])} \leq \frac{h_M^{n+1}}{(n+1)!} \|f^{(n+1)}\|_{L^\infty([x_0, x_m])}, \quad (6.5.1.9)$$

with mesh width  $h_M := \max\{|x_j - x_{j-1}| : j = 1, \dots, m\}$ . □

Another special case: fixed mesh  $M$ , uniform polynomial degree  $n$

Study estimates  $\|f - \mathcal{I}_M f\| \leq CT(n)$  for  $n \rightarrow \infty$ .

Terminology:

investigation of  $p$ -convergence

**EXAMPLE 6.5.1.10 ( $p$ -convergence of piecewise polynomial interpolation)** We study  $p$ -convergence in the setting of Ex. 6.5.1.7.

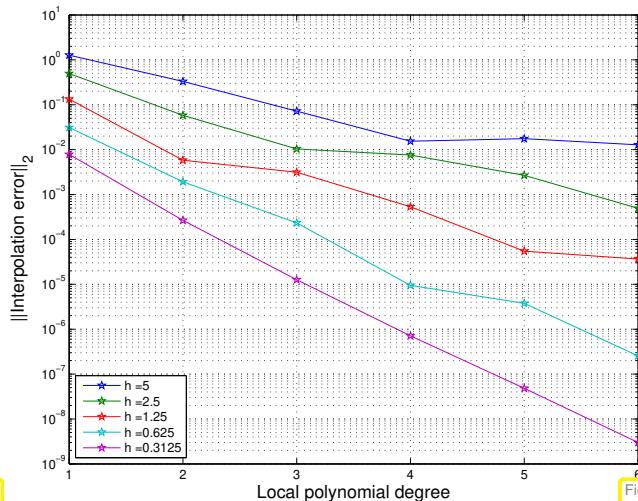


Fig. 228

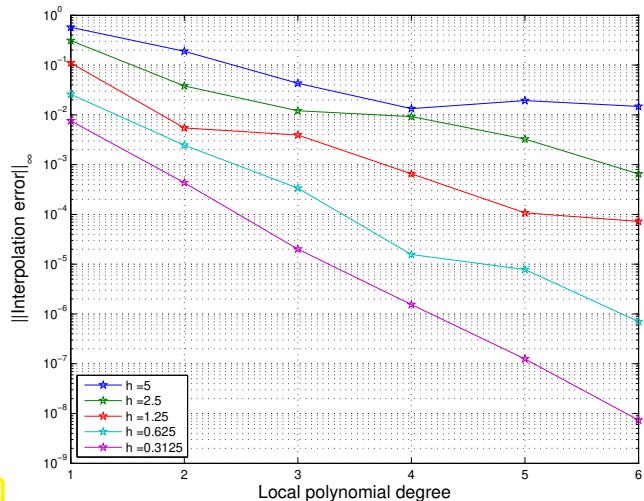


Fig. 229

Observation: (apparent) exponential convergence in polynomial degree

Note: in the case of  $p$ -convergence the situation is the same as for standard polynomial interpolation, see 6.1.2.

In this example we deal with an analytic function, see Rem. 6.1.3.22. Though equidistant local interpolation nodes are used cf. Ex. 6.1.2.10, the mesh intervals seems to be small enough that even in this case exponential convergence prevails. □

## 6.5.2 Cubic Hermite Interpolation: Error Estimates

See Section 5.4 for definition and algorithms for cubic Hermite interpolation of data points, with a focus on shape preservation, however. If the derivative  $f'$  of the interpoland  $f$  is available (in procedural form), then it can be used to fix local cubic polynomials by prescribing point values and derivative values in the endpoints of grid intervals.

**Definition 6.5.2.1. Piecewise cubic Hermite interpolant (with exact slopes) → Def. 5.4.1.1**

Given  $f \in C^1([a, b])$  and a mesh  $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_{m-1} < x_m = b\}$  the **piecewise cubic Hermite interpolant** (with exact slopes)  $s : [a, b] \rightarrow \mathbb{R}$  is defined as

$$s|_{[x_{j-1}, x_j]} \in \mathcal{P}_3, \quad j = 1, \dots, m, \quad s(x_j) = f(x_j), \quad s'(x_j) = f'(x_j), \quad j = 0, \dots, m.$$

Clearly, the piecewise cubic Hermite interpolant is continuously differentiable:  $s \in C^1([a, b])$ , cf. Cor. 5.4.1.2.

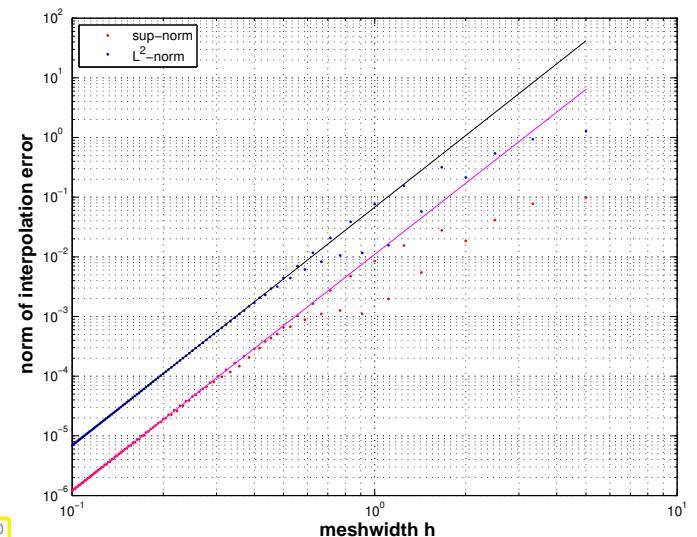
**EXPERIMENT 6.5.2.2 (Convergence of Hermite interpolation with exact slopes)** In this experiment we study the  $h$ -convergence of Cubic Hermite interpolation for a smooth function.

Piecewise cubic Hermite interpolation of

$$f(x) = \arctan(x).$$

- ◆ domain:  $I = (-5, 5)$
- ◆ mesh  $\mathcal{T} = \{-5 + h j\}_{j=0}^n \subset I, h = \frac{10}{n}$ ,
- ◆ exact slopes  $c_i = f'(t_i), i = 0, \dots, n$

► algebraic convergence  $O(h^4)$



Approximate computation of error norms analogous to Ex. 6.5.1.7.

The observation made in Exp. 6.5.2.2 matches the theoretical prediction of the rate of algebraic convergence for cubic Hermite interpolation with exact slopes for a smooth function.

### Theorem 6.5.2.3. Convergence of approximation by cubic Hermite interpolation

Let  $s$  be the cubic Hermite interpolant of  $f \in C^4([a, b])$  on a mesh  $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_{m-1} < x_m = b\}$  according to Def. 6.5.2.1. Then

$$\|f - s\|_{L^\infty([a, b])} \leq \frac{1}{4!} h_{\mathcal{M}}^4 \|f^{(4)}\|_{L^\infty([a, b])},$$

with the meshwidth  $h_{\mathcal{M}} := \max_j |x_j - x_{j-1}|$ .

In Section 5.4.2 we saw variants of cubic Hermite interpolation, for which the slopes  $c_j = s'(x_j)$  were computed from the values  $y_j$  in preprocessing step. Now we study the use of such a scheme for approximation.

### EXPERIMENT 6.5.2.4 (Convergence of Hermite interpolation with averaged slopes)

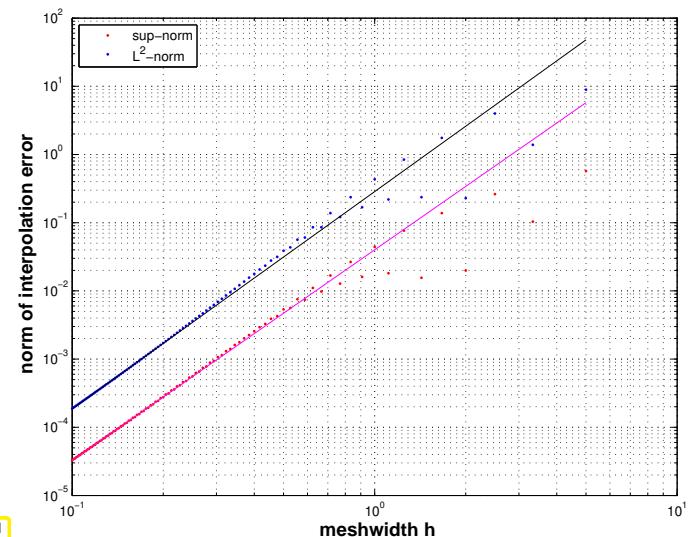
Piecewise cubic Hermite interpolation of

$$f(x) = \arctan(x).$$

- ◆ domain:  $I = (-5, 5)$
- ◆ equidistant mesh  $\mathcal{T}$  in  $I$ , see Exp. 6.5.2.2,
- ◆ averaged local slopes, see (5.4.1.8)

► algebraic convergence  $O(h^3)$  in meshwidth

Code → GITLAB



We observe lower rate of algebraic convergence compared to the use of exact slopes due to averaging (5.4.1.8). From the plot we deduce  $\mathcal{O}(h^3)$  asymptotic decay of the  $L^2$ - and  $L^\infty$ -norms of the approximation error for meshwidth  $h \rightarrow 0$ .

]

### 6.5.3 Cubic Spline Interpolation: Error Estimates [Han02, Ch. 47]

Recall concept and algorithms for cubic spline interpolation from Section 5.5.1. As an interpolation scheme it can also serve as the foundation for an approximation scheme according to § 6.0.0.6: the mesh will double as knot set, see Def. 5.5.0.1. Cubic spline interpolation is not local as we saw in § 5.5.2.6. Nevertheless, cubic spline interpolants can be computed with an effort of  $\mathcal{O}(m)$  as elaborated in § 5.5.1.3.

**EXPERIMENT 6.5.3.1 (Approximation by complete cubic spline interpolants)** We take  $I = [-1, 1]$  and rely on an equidistant mesh (knot set)  $\mathcal{M} := \{-1 + \frac{2}{n}j\}_{j=0}^n, n \in \mathbb{N} \rightarrow$  meshwidth  $h = 2/n$ .

We study  $h$ -convergence of **complete** ( $\rightarrow$  § 5.5.1.9) cubic spline interpolation, where the slopes at the endpoints of the interval are made to agree with the derivatives of the interpoland at these points. As interpolands we consider

$$f_1(t) = \frac{1}{1 + e^{-2t}} \in C^\infty(I) , \quad f_2(t) = \begin{cases} 0 & , \text{ if } t < -\frac{2}{5} , \\ \frac{1}{2}(1 + \cos(\pi(t - \frac{3}{5}))) & , \text{ if } -\frac{2}{5} < t < \frac{3}{5} , \\ 1 & \text{otherwise.} \end{cases} \in C^1(I) .$$

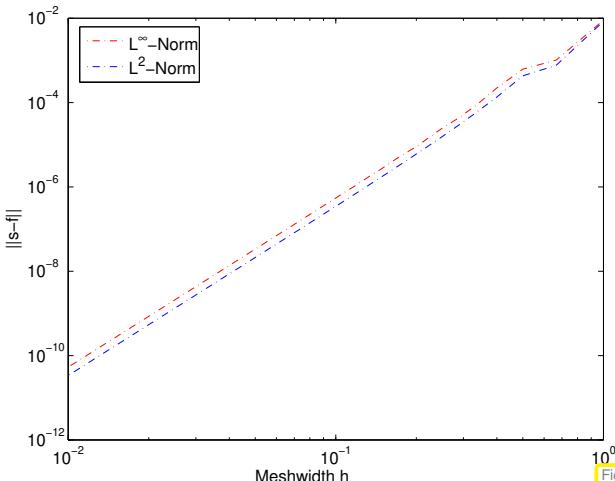


Fig. 232

$$\|f_1 - s\|_{L^\infty([-1,1])} = \mathcal{O}(h^4)$$

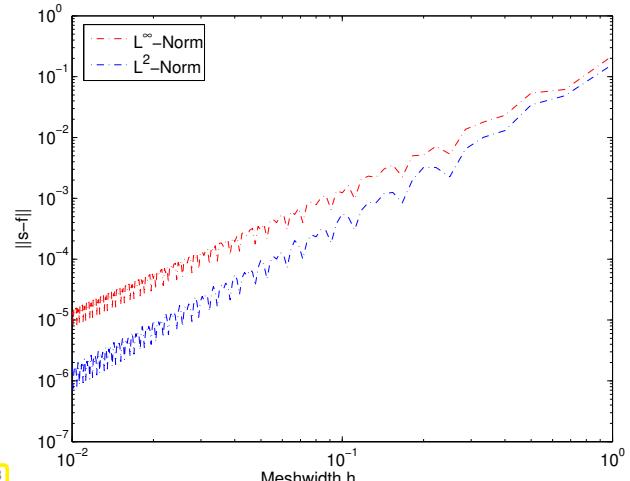


Fig. 233

$$\|f_2 - s\|_{L^\infty([-1,1])} = \mathcal{O}(h^2)$$

The codes used to run this experiment can be accessed through → [GITLAB](#).

We observe algebraic order of convergence in  $h$  with empiric rate approximately given by  $\min\{1 + \text{regularity of } f, 4\}$ . ]

We remark that there is the following theoretical result [HM76], [DR08, Rem. 9.2]:

$$f \in C^4([t_0, t_n]) \quad \|f - s\|_{L^\infty([t_0, t_n])} \leq \frac{5}{384} h^4 \|f^{(4)}\|_{L^\infty([t_0, t_n])} . \quad (6.5.3.2)$$

## Summary and Learning Outcomes

This chapter is meant to impart the following knowledge and skills.

- You should be able to extract the (asymptotic) convergence of approximation errors from empirical data.
- You should know a few relevant norms on spaces of functions.
- You should be able to construct an approximation scheme on an arbitrary interval from an interpolation scheme on a fixed interval.
- You should recall bounds for the pointwise approximation error of polynomial interpolation for functions in  $C^r$  and analytic functions.
- You should be familiar with Chebychev interpolation: rationale, definition, and algorithms.
- You should know about trigonometric interpolation and the behavior of the associated pointwise approximation errors.
- You should be able to predict the convergence of piecewise polynomial interpolation in terms of meshwidth  $h \rightarrow 0$ .

# Bibliography

- [CB95] Q. Chen and I. Babuska. “Approximate optimal points for polynomial interpolation of real functions in an interval and in a triangle”. In: *Comp. Meth. Appl. Mech. Engr.* 128 (1995), pp. 405–417 (cit. on p. 407).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 396, 446).
- [Dav75] P.J. Davis. *Interpolation and Approximation*. New York: Dover, 1975 (cit. on pp. 388, 390).
- [DY10] L. Demanet and L. Ying. *On Chebyshev interpolation of analytic functions*. Online notes. 2010.
- [DH03] P. Deuflhard and A. Hohmann. *Numerical Analysis in Modern Scientific Computing*. Vol. 43. Texts in Applied Mathematics. Springer, 2003 (cit. on p. 405).
- [HM76] C.A. Hall and W.W. Meyer. “Optimal error bounds for cubic spline interpolation”. In: *J. Approx. Theory* 16 (1976), pp. 105–122 (cit. on p. 446).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 396, 404, 405, 412, 446).
- [JWZ19] Peter Jantsch, Clayton G. Webster, and Guannan Zhang. “On the Lebesgue constant of weighted Leja points for Lagrange interpolation on unbounded domains”. In: *IMA J. Numer. Anal.* 39.2 (2019), pp. 1039–1057. DOI: [10.1093/imanum/dry002](https://doi.org/10.1093/imanum/dry002).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on pp. 416, 417).
- [Ran00] R. Rannacher. *Einführung in die Numerische Mathematik*. Vorlesungsskriptum Universität Heidelberg. 2000 (cit. on p. 397).
- [Rem84] R. Remmert. *Funktionentheorie I*. Grundwissen Mathematik 5. Berlin: Springer, 1984 (cit. on p. 400).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on pp. 387, 392, 397).
- [Tad86] Eitan Tadmor. “The Exponential Accuracy of Fourier and Chebyshev Differencing Methods”. In: *SIAM Journal on Numerical Analysis* 23.1 (1986), pp. 1–10. DOI: [10.1137/0723001](https://doi.org/10.1137/0723001).
- [TT10] Rodney Taylor and Vilmos Totik. “Lebesgue constants for Leja points”. In: *IMA J. Numer. Anal.* 30.2 (2010), pp. 462–486. DOI: [10.1093/imanum/drn082](https://doi.org/10.1093/imanum/drn082).
- [Ver86] P. Vértesi. “On the optimal Lebesgue constants for polynomial interpolation”. In: *Acta Math. Hungaria* 47.1-2 (1986), pp. 165–178 (cit. on p. 407).
- [Ver90] P. Vértesi. “Optimal Lebesgue constant for Lagrange interpolation”. In: *SIAM J. Numer. Anal.* 27.5 (1990), pp. 1322–1331 (cit. on p. 407).

# Chapter 7

## Numerical Quadrature



*Supplementary literature.* Numerical quadrature is covered in [Han02, p. VII] and [DR08, Ch. 10].

Numerical quadrature deals with the *approximate* numerical evaluation of integrals  $\int_{\Omega} f(\mathbf{x}) d\mathbf{x}$  for a given (closed) integration domain  $\Omega \subset \mathbb{R}^d$ . Thus, the underlying problem in the sense of § 1.5.5.1 is the mapping

$$I : \left\{ \begin{array}{rcl} C^0(\Omega) & \rightarrow & \mathbb{R} \\ f & \mapsto & \int_{\Omega} f(\mathbf{x}) d\mathbf{x} \end{array} \right., \quad (7.0.0.1)$$

with data space  $X := C^0(\Omega)$  and result space  $Y := \mathbb{R}$ .

If  $f$  is complex-valued or vector-valued, then so is the integral. The methods presented in this chapter can immediately be generalized to this case by componentwise application.

**§7.0.0.2 (Integrands in procedural form)** The integrand  $f$ , a continuos function  $f : \Omega \subset \mathbb{R}^d \mapsto \mathbb{R}$  should not be thought of as given by an analytic expression, but as given in **procedural form**, cf. Rem. 5.1.0.8.

For instance, in C++ the integrand is provided through a “functor” data type with an evaluation operator **double operator** (`Point &`) or a corresponding member function, see Code 5.1.0.9 for an example, or a lambda function (to Section 0.3.3).

- General methods for numerical quadrature should rely only on finitely many *point evaluations* of the integrand. ↓

In this chapter the focus is on the special case  $d = 1$ :  $\Omega = [a, b]$  (an interval) (Multidimensional numerical quadrature is substantially more difficult, unless  $\Omega$  is tensor-product domain, a multi-dimensional box. Multidimensional numerical quadrature will be treated in the course “Numerical methods for partial differential equations”).

### Remark 7.0.0.3 (Importance of numerical quadrature)

- ☞ Numerical quadrature methods are key building blocks for so-called variational methods for the numerical treatment of partial differential equations. A prominent example is the **finite element method**. They are also a pivotal device for the numerical solution of integral equations and in computational statistics. ↓

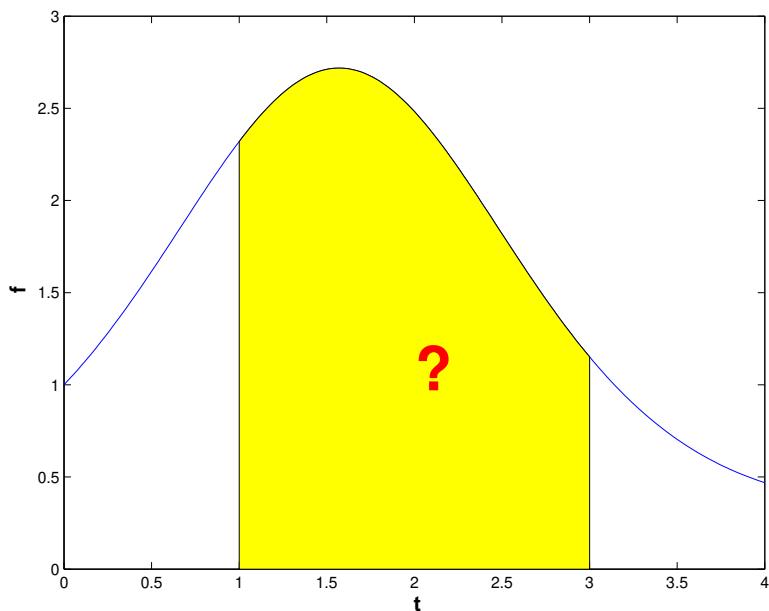


Fig. 234

**EXAMPLE 7.0.0.4 (Heating production in electrical circuits)** In Ex. 2.1.0.3 we learned about the nodal analysis of electrical circuits. Its application to a non-linear circuit will be discussed in Ex. 8.0.0.1, which will reveal that every computation of currents and voltages can be rather time-consuming. In this example we consider a non-linear circuit in quasi-stationary operation (capacities and inductances are ignored). Then the computation of branch currents and nodal voltages entails solving a non-linear system of equations.

Now assume time-harmonic periodic excitation with period  $T > 0$ .

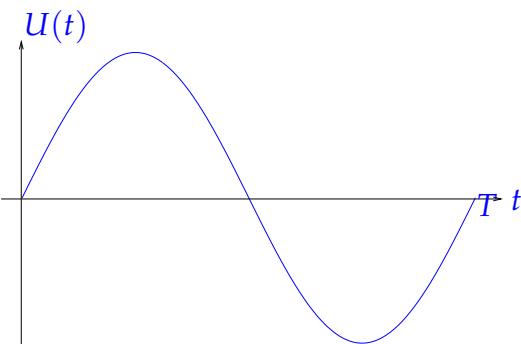


Fig. 235

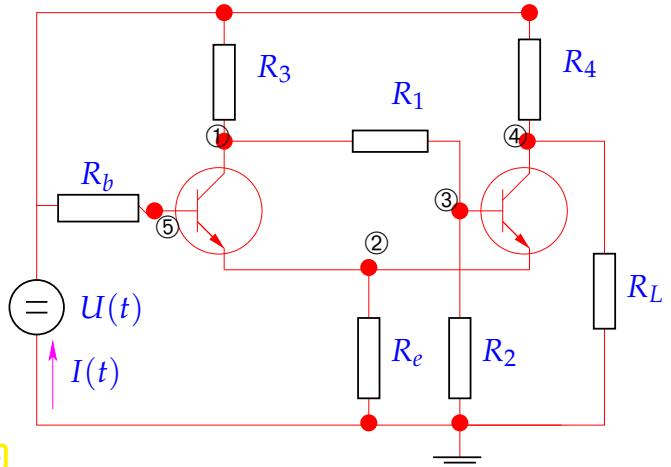


Fig. 236

The goal is to compute the energy dissipated by the circuit, which is equal to the energy injected by the voltage source. This energy can be obtained by integrating the power  $P(t) = U(t)I(t)$  over period  $[0, T]$ :

$$W_{\text{therm}} = \int_0^T U(t)I(t) dt, \quad \text{where } I = I(U).$$

**double**  $I(\mathbf{double} U)$  involves solving non-linear system of equations, see Ex. 8.0.0.1!

This is a typical example where “point evaluation” by solving the non-linear circuit equations is the only way to gather information about the integrand. □

## Contents

7.1 Quadrature Formulas . . . . .	451
7.2 Polynomial Quadrature Formulas . . . . .	454
7.3 Gauss Quadrature . . . . .	456
7.4 Composite Quadrature . . . . .	470
7.5 Adaptive Quadrature . . . . .	478

## 7.1 Quadrature Formulas

Quadrature formulas realize the approximation of an integral through finitely many point evaluations of the integrand.

### Definition 7.1.0.1. Quadrature formula/quadrature rule

An  $n$ -point quadrature formula/quadrature rule (QR) on  $[a, b]$  provides an approximation of the value of an integral through a *weighted sum* of point values of the integrand:

$$\int_a^b f(t) dt \approx Q_n(f) := \sum_{j=1}^n w_j^n f(c_j^n) . \quad (7.1.0.2)$$

Terminology:  $w_j^n$  : quadrature weights  $\in \mathbb{R}$   
 $c_j^n$  : quadrature nodes  $\in [a, b]$  (also called quadrature points)

Obviously (7.1.0.2) is compatible with integrands  $f$  given in *procedural form* as `double f(double t)`, compare § 7.0.0.2.

### C++-code 7.1.0.3: C++ template implementing generic quadrature formula

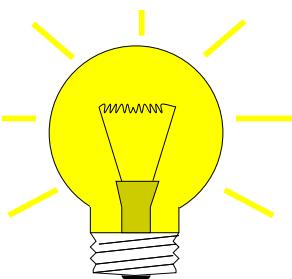
```

2 // Generic numerical quadrature routine implementing (7.1.0.2):
3 // f is a handle to a function, e.g. as lambda function
4 // c, w pass quadrature nodes c_j ∈ [a,b], and weights w_j ∈ ℝ
5 // in a Eigen::VectorXd
6 template <class Function>
7 double quadformula(Function& f, const VectorXd& c, const VectorXd& w) {
8     const std::size_t n = c.size();
9     double I = 0;
10    for (std::size_t i = 0; i < n; ++i) I += w(i)*f(c(i));
11    return I;
12 }
```

A single invocation costs  $n$  point evaluations of the integrand plus  $n$  additions and multiplications.

**Remark 7.1.0.4 (Transformation of quadrature rules)** In the setting of function approximation by polynomials we learned in Rem. 6.1.1.14 that an approximation schemes for any interval could be obtained from an approximation scheme on a single *reference interval* ( $[-1, 1]$  in Rem. 6.1.1.14) by means of *affine pullback*, see (6.1.1.18). A similar affine transformation technique makes it possible to derive quadrature formula for an arbitrary interval from a single quadrature formula on a reference interval.

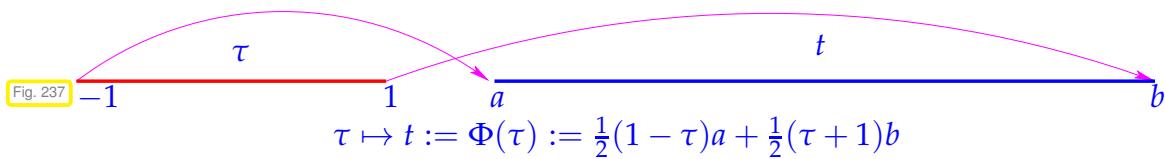
Given: quadrature formula  $(\widehat{c}_j, \widehat{w}_j)_{j=1}^n$  on *reference interval*  $[-1, 1]$



Idea: transformation formula for integrals

$$\int_a^b f(t) dt = \frac{1}{2}(b-a) \int_{-1}^1 \widehat{f}(\tau) d\tau , \quad (7.1.0.5)$$

$$\widehat{f}(\tau) := f\left(\frac{1}{2}(1-\tau)a + \frac{1}{2}(\tau+1)b\right) .$$



Note that  $\hat{f}$  is the affine pullback  $\Phi^* f$  of  $f$  to  $[-1, 1]$  as defined in Eq. (6.1.1.16).

► quadrature formula for general interval  $[a, b], a, b \in \mathbb{R}$ :

$$\int_a^b f(t) dt \approx \frac{1}{2}(b-a) \sum_{j=1}^n \widehat{w}_j \hat{f}(\hat{c}_j) = \sum_{j=1}^n w_j f(c_j) \quad \text{with} \quad c_j = \frac{1}{2}(1-\hat{c}_j)a + \frac{1}{2}(1+\hat{c}_j)b, \\ w_j = \frac{1}{2}(b-a)\widehat{w}_j.$$

In words, the nodes are just mapped through the affine transformation  $c_j = \Phi(\hat{c}_j)$ , the weights are scaled by the ratio of lengths of  $[a, b]$  and  $[-1, 1]$ .

► A 1D quadrature formula on arbitrary intervals can be specified by providing its weights  $\widehat{w}_j$ /nodes  $\hat{c}_j$  for the integration domain  $[-1, 1]$  (reference interval). Then the above transformation is assumed.

Another common choice for the reference interval:  $[0, 1]$ , pay attention! □

**Remark 7.1.0.6 (Tables of quadrature rules)** In many codes families of quadrature rules are used to control the quadrature error. Usually, suitable sequences of weights  $w_j^n$  and nodes  $c_j^n$  are precomputed and stored in tables up to sufficiently large values of  $n$ . A possible interface could be the following:

```
struct QuadTab {
    template <typename VecType>
    static void getrule(int n, VecType &c, VecType &w,
                        double a=-1.0, double b=1.0);
}
```

Calling the method `getrule()` fills the vectors `c` and `w` with the nodes and the weights for a desired  $n$ -point quadrature on  $[a, b]$  with  $[-1, 1]$  being the default reference interval. For `VecType` we may assume the basic functionality of `Eigen::VectorXd`. □

**§7.1.0.7 (Quadrature by approximation schemes)** Every approximation scheme  $A : C^0([a, b]) \rightarrow V$ ,  $V$  a space of “simple functions” on  $[a, b]$ , see § 6.0.0.5, gives rise to a method for numerical quadrature according to

$$\int_a^b f(t) dt \approx Q_A(f) := \int_a^b (Af)(t) dt. \quad (7.1.0.8)$$

As explained in § 6.0.0.6 every interpolation scheme  $I_T$  based on the node set  $T = \{t_0, t_1, \dots, t_n\} \subset [a, b]$  ( $\rightarrow$  § 5.1.0.6) induces an approximation scheme, and, hence, also a quadrature scheme on  $[a, b]$ :

$$\int_a^b f(t) dt \approx \int_a^b I_T[f(t_0), \dots, f(t_n)]^\top(t) dt. \quad (7.1.0.9)$$

### Lemma 7.1.0.10. Quadrature formulas from linear interpolation schemes

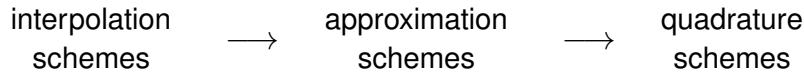
Every linear interpolation operator  $I_T$  according to Def. 5.1.0.19 spawns a quadrature formula ( $\rightarrow$  Def. 7.1.0.1) by (7.1.0.9).

*Proof.* Writing  $\mathbf{e}_j$  for the  $j$ -th unit vector of  $\mathbb{R}^n$ , we have by linearity

$$\int_a^b \mathbf{I}_{\mathcal{T}}[f(t_0), \dots, f(t_n)]^\top dt = \sum_{j=0}^n f(t_j) \underbrace{\int_a^b (\mathbf{I}_{\mathcal{T}}(\mathbf{e}_j))(t) dt}_{\text{weight } w_j}. \quad (7.1.0.11)$$

Hence, we have arrived at an  $n+1$ -point quadrature formula with nodes  $t_j$ , whose weights are the integrals of the **cardinal interpolants** for the interpolation scheme  $\mathcal{T}$ .  $\square$

Summing up, we have found



↓

**§7.1.0.12 (Convergence of numerical quadrature)** In general the quadrature formula (7.1.0.2) will only provide an approximate value for the integral.

➤ For a generic integrand we will encounter a non-vanishing

$$\text{quadrature error} \quad E_n(f) := \left| \int_a^b f(t) dt - Q_n(f) \right|$$

As in the case of function approximation by interpolation Section 6.1.2, our focus will on the **asymptotic** behavior of the quadrature error as a function of the number  $n$  of point evaluations of the integrand.

Therefore consider **families** of quadrature rules  $\{Q_n\}_n$  ( $\rightarrow$  Def. 7.1.0.1) described by

- ◆ quadrature weights  $\{w_j^n, j = 1, \dots, n\}_{n \in \mathbb{N}}$  and
- ◆ quadrature nodes  $\{c_j^n, j = 1, \dots, n\}_{n \in \mathbb{N}}$ .

We study the *asymptotic* behavior of the **quadrature error**  $E(n)$  for  $n \rightarrow \infty$

As in the case of interpolation errors in § 6.1.2.5 we make the usual qualitative distinction, see Def. 6.1.2.7:

- ▷ algebraic convergence  $E(n) = O(n^{-p})$ , rate  $p > 0$
- ▷ exponential convergence  $E(n) = O(q^n)$ ,  $0 \leq q < 1$

Note that the number  $n$  of nodes agrees with the number of  $f$ -evaluations required for the evaluation of the quadrature formula. This is usually used as a *measure for the cost* of computing  $Q_n(f)$ .

Therefore, in the sequel, we consider the quadrature error as a function of  $n$ .  $\downarrow$

**§7.1.0.13 (Quadrature error from approximation error)** Bounds for the maximum norm of the approximation error of an approximation scheme directly translate into estimates of the quadrature error of the induced quadrature scheme (7.1.0.8):

$$\left| \int_a^b f(t) dt - QA(f) \right| \leq \int_a^b |f(t) - A(f)(t)| dt \leq |b-a| \|f - A(f)\|_{L^\infty([a,b])}. \quad (7.1.0.14)$$

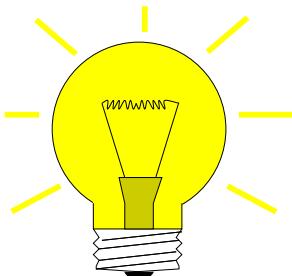
Hence, the various estimates derived in Section 6.1.2 and Section 6.1.3.2 give us quadrature error estimates “for free”. More details will be given in the next section.  $\downarrow$

## 7.2 Polynomial Quadrature Formulas

Now we specialize the general recipe of § 7.1.0.7 for approximation schemes based on global polynomials, the Lagrange approximation scheme as introduced in Section 6.1, Def. 6.1.2.1.



*Supplementary literature.* This topic is discussed in [DR08, Sect. 10.2].



Idea: replace integrand  $f$  with  $p_{n-1} := l_{\mathcal{T}} \in \mathcal{P}_{n-1}$  = polynomial Lagrange interpolant of  $f$  ( $\rightarrow$  Cor. 5.2.2.8) for given node set  $\mathcal{T} := \{t_0, \dots, t_{n-1}\} \subset [a, b]$

$$\Rightarrow \int_a^b f(t) dt \approx Q_n(f) := \int_a^b p_{n-1}(t) dt . \quad (7.2.0.1)$$

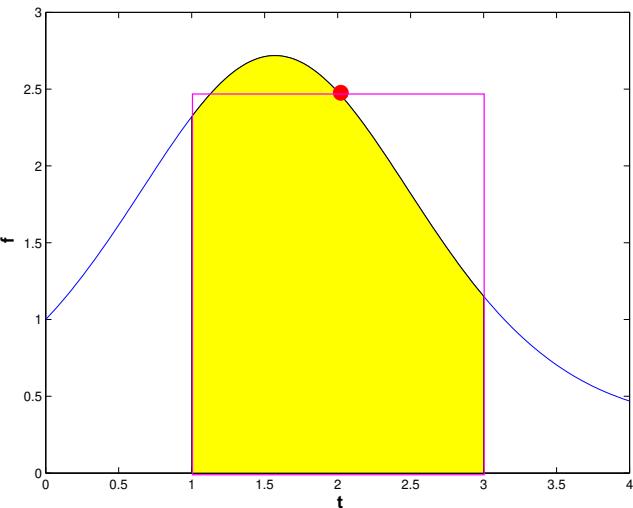
The cardinal interpolants for Lagrange interpolation are the Lagrange polynomials (5.2.2.4)

$$L_i(t) := \prod_{\substack{j=0 \\ j \neq i}}^{n-1} \frac{t - t_j}{t_i - t_j}, \quad i = 0, \dots, n-1 \quad \xrightarrow{(5.2.2.6)} \quad p_{n-1}(t) = \sum_{i=0}^{n-1} f(t_i) L_i(t) .$$

Then (7.1.0.11) amounts to the  $n$ -point quadrature formula

$$\int_a^b p_{n-1}(t) dt = \sum_{i=0}^{n-1} f(t_i) \int_a^b L_i(t) dt \quad \xrightarrow{\text{nodes } c_i = t_{i-1}, \text{ weights } w_i := \int_a^b L_i(t) dt} \quad (7.2.0.2)$$

### EXAMPLE 7.2.0.3 (Midpoint rule)



[midpoint rule,  $n = 1$ ] The midpoint rule is (7.2.0.2) for  $n = 1$  and  $t_0 = \frac{1}{2}(a + b)$ . It leads to the 1-point quadrature formula

$$\int_a^b f(t) dt \approx Q_{mp}(f) = (b - a)f\left(\frac{1}{2}(a + b)\right) .$$

“midpoint”

▷ the area under the graph of  $f$  is approximated by the area of a rectangle.

Fig. 238

**EXAMPLE 7.2.0.4 (Newton-Cotes formulas  $\rightarrow$  [Han02, Ch. 38])** The  $n := m + 1$ -point Newton-Cotes formulas arise from Lagrange interpolation in equidistant nodes (6.1.2.3) in the integration interval  $[a, b]$ :

► Equidistant quadrature nodes

$$t_j := a + hj, \quad h := \frac{b - a}{m}, \quad j = 0, \dots, m:$$

The weights for the interval  $[0, 1]$  can be found, e.g., by symbolic computation using MAPLE: the following MAPLE function expects the polynomial degree as input argument.

```
> newtoncotes := m -> factor(int(interp([seq(i/n, i=0..m)],  
[seq(f(i/n), i=0..m)], z), z=0..1)):
```

Weights on general intervals  $[a, b]$  can then be deduced by the affine transformation rule as explained in Rem. 7.1.0.4.

- $n = 2$ : **Trapezoidal rule** (integrate linear interpolant of integrand in endpoints)

```
> trapez := newtoncotes(1);
```

$$\begin{aligned}\widehat{Q}_{\text{trp}}(f) &:= \frac{1}{2}(f(0) + f(1)) \quad (7.2.0.5) \\ \left( \int_a^b f(t) dt \approx \frac{b-a}{2}(f(a) + f(b)) \right)\end{aligned}$$

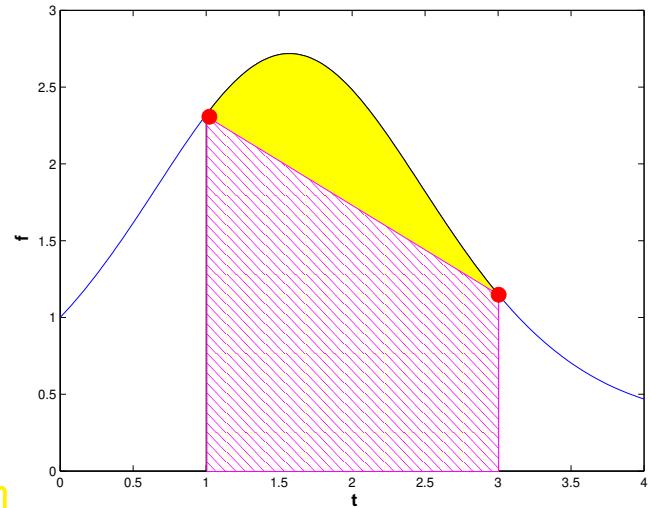


Fig. 239

- $n = 3$ : **Simpson rule**

```
> simpson := newtoncotes(2);
```

$$\frac{h}{6} \left( f(0) + 4f(\frac{1}{2}) + f(1) \right) \quad \left( \int_a^b f(t) dt \approx \frac{b-a}{6} \left( f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right) \right) \quad (7.2.0.6)$$

- $n = 5$ : **Milne rule**

```
> milne := newtoncotes(4);
```

$$\begin{aligned}\frac{1}{90} h \left( 7f(0) + 32f(\frac{1}{4}) + 12f(\frac{1}{2}) + 32f(\frac{3}{4}) + 7f(1) \right) \\ \left( \frac{b-a}{90} (7f(a) + 32f(a + (b-a)/4) + 12f(a + (b-a)/2) + 32f(a + 3(b-a)/4) + 7f(b)) \right)\end{aligned}$$

- $n = 7$ : **Weddle rule**

```
> weddle := newtoncotes(6);
```

$$\begin{aligned}\frac{1}{840} h (41f(0) + 216f(\frac{1}{6}) + 27f(\frac{1}{3}) + 272f(\frac{1}{2}) \\ + 27f(\frac{2}{3}) + 216f(\frac{5}{6}) + 41f(1))\end{aligned}$$

- $n \geq 8$ : quadrature formulas with *negative weights*

```
> newtoncotes(8);
```

$$\frac{1}{28350} h (989 f(0) + 5888 f(\tfrac{1}{8}) - 928 f(\tfrac{1}{4}) + 10496 f(\tfrac{3}{8}) \\ - 4540 f(\tfrac{1}{2}) + 10496 f(\tfrac{5}{8}) - 928 f(\tfrac{3}{4}) + 5888 f(\tfrac{7}{8}) + 989 f(1))$$

From Ex. 6.1.2.10 we know that the approximation error incurred by Lagrange interpolation in equidistant nodes can blow up even for analytic functions. This blow-up can also infect the quadrature error of Newton-Cotes formulas for large  $n$ , which renders them essentially useless. In addition they will be marred by large (in modulus) and negative weights, which compromises numerical stability ( $\rightarrow$  Def. 1.5.5.19)

↓

**Remark 7.2.0.7 (Clenshaw-Curtis quadrature rules [Tre08])** The considerations of Section 6.1.3 confirmed the superiority of the “optimal” Chebychev nodes (6.1.3.10) for globally polynomial Lagrange interpolation. This suggests that we use these nodes also for numerical quadrature with weights given by (7.2.0.2). This yields the so-called **Clenshaw-Curtis rules** with the following rather desirable property:

**Theorem 7.2.0.8. Positivity of Clenshaw-Curtis weights**

The weights  $w_j^n$ ,  $j = 1, \dots, n$ , for every  $n$ -point Clenshaw-Curtis rule are positive.

The weights of any  $n$ -point Clenshaw-Curtis rule can be computed with a computational effort of  $O(n \log n)$  using FFT.

**§7.2.0.9 (Error estimates for polynomial quadrature)** As a concrete application of § 7.1.0.13, (7.1.0.14) we use the  $L^\infty$ -bound (6.1.2.19) for Lagrange interpolation

$$\|f - L_T f\|_{L^\infty(I)} \leq \frac{\|f^{(n+1)}\|_{L^\infty(I)}}{(n+1)!} \max_{t \in I} |(t-t_0) \cdots (t-t_n)|. \quad (6.1.2.19)$$

to conclude for any  $n$ -point quadrature rule based on polynomial interpolation:

$$f \in C^n([a, b]) \Rightarrow \left| \int_a^b f(t) dt - Q_n(f) \right| \leq \frac{1}{n!} (b-a)^{n+1} \|f^{(n)}\|_{L^\infty([a, b])}. \quad (7.2.0.10)$$

Much sharper estimates for Clenshaw-Curtis rules ( $\rightarrow$  Rem. 7.2.0.7) can be inferred from the interpolation error estimate (6.1.3.14) for Chebychev interpolation. For functions with *limited smoothness algebraic convergence* of the quadrature error for Clenshaw-Curtis quadrature follows from (6.1.3.18). For integrands that possess an *analytic* extension to the complex plane in a neighborhood of  $[a, b]$ , we can conclude *exponential* convergence from (6.1.3.24).

## 7.3 Gauss Quadrature



*Supplementary literature.* Gauss quadrature is discussed in detail in [Han02, Ch. 40-41],

[DR08, Sect.10.3]

How to gauge the “quality” of an  $n$ -point quadrature formula  $Q_n$  without testing it for specific integrands? The next definition gives an answer.

#### Definition 7.3.0.1. Order of a quadrature rule

The **order** of quadrature rule  $Q_n : C^0([a, b]) \rightarrow \mathbb{R}$  is defined as

$$\text{order}(Q_n) := \max\{m \in \mathbb{N}_0 : Q_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_m\} + 1, \quad (7.3.0.2)$$

that is, as the maximal degree  $+1$  of polynomials for which the quadrature rule is guaranteed to be exact.

**§7.3.0.3 (Invariance of order under (affine) transformation)** First we note a simple consequence of the invariance of the polynomial space  $\mathcal{P}_n$  under affine pullback, see Lemma 6.1.1.17.

#### Corollary 7.3.0.4. Invariance of order under affine transformation

An affine transformation of a quadrature rule according to Rem. 7.1.0.4 does not change its order.

□

**§7.3.0.5 (Order of polynomial quadrature rules)** Further, by construction all polynomial  $n$ -point quadrature rules possess order at least  $n$ .

#### Theorem 7.3.0.6. Sufficient order conditions for quadrature rules

An  $n$ -point quadrature rule on  $[a, b]$  ( $\rightarrow$  Def. 7.1.0.1)

$$Q_n(f) := \sum_{j=1}^n w_j f(c_j), \quad f \in C^0([a, b]),$$

with nodes  $t_j \in [a, b]$  and weights  $w_j \in \mathbb{R}$ ,  $j = 1, \dots, n$ , has **order  $\geq n$** , if and only if

$$w_j = \int_a^b L_{j-1}(t) dt, \quad j = 1, \dots, n,$$

where  $L_k$ ,  $k = 0, \dots, n-1$ , is the  $k$ -th **Lagrange polynomial** (5.2.2.4) associated with the ordered node set  $\{c_1, c_2, \dots, c_n\}$ .

*Proof.* The conclusion of the theorem is a direct consequence of the facts that

$$\mathcal{P}_{n-1} = \text{Span}\{L_0, \dots, L_{n-1}\} \quad \text{and} \quad L_k(t_j) = \delta_{k+1,j}, \quad k+1, j \in \{1, \dots, n\} :$$

just plug a Lagrange polynomial into the quadrature formula. □

By construction (7.2.0.2) polynomial  $n$ -point quadrature formulas (7.2.0.1) exact for  $f \in \mathcal{P}_{n-1} \Rightarrow$   $n$ -point polynomial quadrature formula has **at least** order  $n$ . □

**Remark 7.3.0.7 (Choice of (local) quadrature weights)** Thm. 7.3.0.6 provides a concrete formula for quadrature weights, which guaranteed order  $n$  for an  $n$ -point quadrature formula. Yet evaluating integrals

of Lagrange polynomials may be cumbersome. Here we give a general recipe for finding the weights  $w_j$  according to Thm. 7.3.0.6 without dealing with Lagrange polynomials.

Given: arbitrary nodes  $c_1, \dots, c_n$  for  $n$ -point (local) quadrature formula on  $[a, b]$

From Def. 7.3.0.1 we immediately conclude the following procedure: If  $p_0, \dots, p_{n-1}$  is any basis of  $\mathcal{P}_n$ , then, thanks to the linearity of the integral and quadrature formulas,

$$Q_n(p_j) = \int_a^b p_j(t) dt \quad \forall j = 0, \dots, n-1 \quad \Leftrightarrow \quad Q_n \text{ has order } \geq n. \quad (7.3.0.8)$$

➤  $n \times n$  linear system of equations, see (7.3.0.15) for an example:

$$\begin{bmatrix} p_0(c_1) & \dots & p_0(c_n) \\ \vdots & & \vdots \\ p_{n-1}(c_1) & \dots & p_{n-1}(c_n) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} \int_a^b p_0(t) dt \\ \vdots \\ \int_a^b p_{n-1}(t) dt \end{bmatrix}. \quad (7.3.0.9)$$

For instance, for the computation of quadrature weights, one may choose the monomial basis  $p_j(t) = t^j$ . □

**EXAMPLE 7.3.0.10 (Orders of simple polynomial quadrature formulas)** From the order rule for polynomial quadrature rule we immediately conclude the orders of simple representatives.

$n$	Order
1	midpoint rule
2	trapezoidal rule (7.2.0.5)
3	Simpson rule (7.2.0.6)
4	$\frac{3}{8}$ -rule
5	Milne rule

The orders for even  $n$  surpass the predictions of Thm. 7.3.0.6 by 1, which can be verified by straightforward computations; following Def. 7.3.0.1 check the exactness of the quadrature rule on  $[0, 1]$  (this is sufficient → Cor. 7.3.0.4) for monomials  $\{t \mapsto t^k\}$ ,  $k = 0, \dots, q-1$ , which form a basis of  $\mathcal{P}_q$ , where  $q$  is the order that is to be confirmed: essentially one has to show

$$Q(\{t \mapsto t^k\}) = \sum_{j=1}^n w_j c_j^k = \frac{1}{k+1}, \quad k = 0, \dots, q-1, \quad (7.3.0.11)$$

where  $Q \doteq$  quadrature rule on  $[0, 1]$  given by (7.1.0.2).

For the Simpson rule (7.2.0.6) we can also confirm order 4 with symbolic calculations in MAPLE:

```
> rule := 1/3*h*(f(2*h)+4*f(h)+f(0))
> err := taylor(rule - int(f(x), x=0..2*h), h=0, 6);
```

$$err := \left( \frac{1}{90} (D^{(4)})(f)(0)h^5 + O(h^6), h, 6 \right)$$

➤ Composite Simpson rule possesses order 4, indeed! □

**§7.3.0.12 (Maximal order of  $n$ -point quadrature rules)** Natural question: Can an  $n$ -point quadrature formula achieve an order  $> n$ ?

A negative result limits the maximal order that can be achieved:

**Theorem 7.3.0.13. Maximal order of  $n$ -point quadrature rule**

*The maximal order of an  $n$ -point quadrature rule is  $2n$ .*

*Proof.* Consider a generic  $n$ -point quadrature rule according to Def. 7.1.0.1

$$Q_n(f) := \sum_{j=1}^n w_j^n f(c_j^n), \quad (7.1.0.2)$$

We build a polynomial of degree  $2n$  that cannot be integrated exactly by  $Q_n$ . We choose polynomial

$$q(t) := (t - c_1)^2 \cdots (t - c_n)^2 \in \mathcal{P}_{2n}.$$

On the one hand,  $q$  is strictly positive almost everywhere, which means

$$\int_a^b q(t) dt > 0.$$

On the other hand, we find a different value

$$Q_n(q) = \sum_{j=1}^n w_j^n \underbrace{q(c_j^n)}_{=0} = 0.$$

□

Can we at least find  $n$ -point rules with order  $2n$ ?

Heuristics: A quadrature formula has order  $m \in \mathbb{N}$  already, if it is exact for  $m$  polynomials  $\in \mathcal{P}_{m-1}$  that form a basis of  $\mathcal{P}_{m-1}$  (recall Thm. 5.2.1.2).

↑  
An  $n$ -point quadrature formula has  $2n$  “degrees of freedom” ( $n$  node positions,  $n$  weights).



It might be possible to achieve order  $2n = \dim \mathcal{P}_{2n-1}$   
(“No. of equations = No. of unknowns”)

**EXAMPLE 7.3.0.14 (2-point quadrature rule of order 4)** Necessary & sufficient conditions for order 4, cf. (7.3.0.9), integrate the functions of the monomial basis of  $\mathcal{P}_3$  exactly:

$$Q_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_3 \Leftrightarrow Q_n(\{t \mapsto t^q\}) = \frac{1}{q+1} (b^{q+1} - a^{q+1}), \quad q = 0, 1, 2, 3.$$



4 equations for weights  $w_j$  and nodes  $c_j, j = 1, 2$  ( $a = -1, b = 1$ ), cf. Rem. 7.3.0.7

$$\begin{aligned} \int_{-1}^1 1 dt &= 2 = 1w_1 + 1w_2, & \int_{-1}^1 t dt &= 0 = c_1 w_1 + c_2 w_2 \\ \int_{-1}^1 t^2 dt &= \frac{2}{3} = c_1^2 w_1 + c_2^2 w_2, & \int_{-1}^1 t^3 dt &= 0 = c_1^3 w_1 + c_2^3 w_2. \end{aligned} \quad (7.3.0.15)$$

Solve using MAPLE:

```
> eqns := {seq(int(x^k, x=-1..1) = w[1]*xi[1]^k+w[2]*xi[2]^k, k=0..3)};  
> sols := solve(eqns, indets(eqns, name)):  
> convert(sols, radical);
```

➤ weights & nodes:  $\{w_2 = 1, w_1 = 1, c_1 = 1/3\sqrt{3}, c_2 = -1/3\sqrt{3}\}$

► quadrature formula (order 4):  $\int_{-1}^1 f(x) dx \approx f\left(\frac{1}{\sqrt{3}}\right) + f\left(-\frac{1}{\sqrt{3}}\right)$  (7.3.0.16)

□

□

**§7.3.0.17 (Construction of  $n$ -point quadrature rules with maximal order  $2n$ )** First we search for *necessary conditions* that have to be met by the nodes, if an  $n$ -point quadrature rule has order  $2n$ .

Optimist's **assumption**:  $\exists$  family of  $n$ -point quadrature formulas on  $[-1, 1]$

$$Q_n(f) := \sum_{j=1}^n w_j^n f(c_j^n) \approx \int_{-1}^1 f(t) dt, \quad w_j \in \mathbb{R}, n \in \mathbb{N},$$

of order  $2n \Leftrightarrow$  exact for polynomials  $\in \mathcal{P}_{2n-1}$ . (7.3.0.18)

Define  $\bar{P}_n(t) := (t - c_1^n) \cdots (t - c_n^n), \quad t \in \mathbb{R} \Rightarrow \bar{P}_n \in \mathcal{P}_n$ .

Note:  $\bar{P}_n$  has leading coefficient = 1.

By assumption on the order of  $Q_n$  we know that for any  $q \in \mathcal{P}_{n-1}$

$$\int_{-1}^1 \underbrace{q(t)\bar{P}_n(t)}_{\in \mathcal{P}_{2n-1}} dt \stackrel{(7.3.0.18)}{=} \sum_{j=1}^n w_j^n q(c_j^n) \underbrace{\bar{P}_n(c_j^n)}_{=0} = 0.$$

We conclude  $L^2([-1, 1])$ -orthogonality:  $\int_{-1}^1 q(t)\bar{P}_n(t) dt = 0 \quad \forall q \in \mathcal{P}_{n-1}$ . (7.3.0.19)

↗  $L^2([-1, 1])$ -inner product of  $q$  and  $\bar{P}_n$ , see (6.2.1.5).

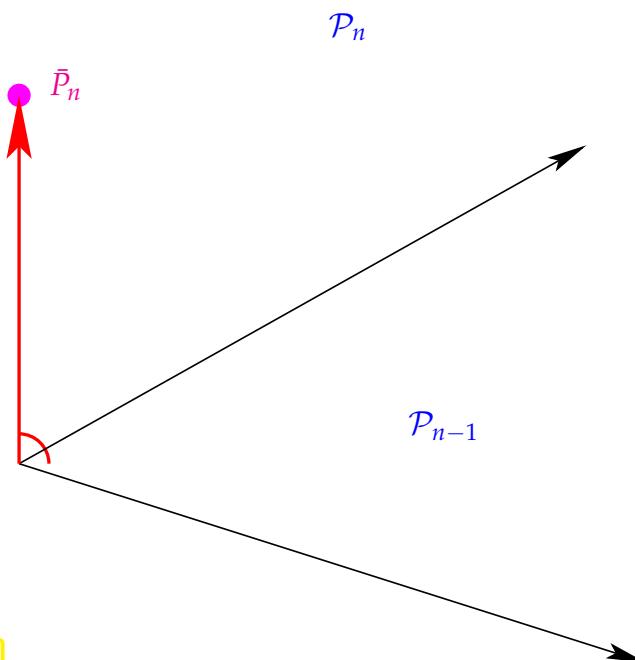


Fig. 240

We can also give algebraic arguments for existence and uniqueness of  $\bar{P}_n$ . Switching to a monomial representation of  $\bar{P}_n$

$$\bar{P}_n(t) = t^n + \alpha_{n-1}t^{n-1} + \cdots + \alpha_1t + \alpha_0 ,$$

by linearity of the integral, (7.3.0.19) is equivalent to

$$\sum_{j=0}^{n-1} \alpha_j \int_{-1}^1 t^j t^\ell dt = - \int_{-1}^1 t^n t^\ell dt , \quad \ell = 0, \dots, n-1 .$$

This is a linear system of equations  $\mathbf{A}[\alpha_j]_{j=0}^{n-1} = \mathbf{b}$  with a symmetric, positive definite ( $\rightarrow$  Def. 1.1.2.6) coefficient matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$ . The  $\mathbf{A}$  is positive definite can be concluded from

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \int_{-1}^1 \left( \sum_{j=0}^{n-1} (\mathbf{x})_j t^j \right)^2 dt > 0 \quad , \text{ if } \mathbf{x} \neq 0 .$$

Hence,  $\mathbf{A}$  is regular and the coefficients  $\alpha_j$  are uniquely determined. Thus there is *only one*  $n$ -point quadrature rule of order  $2n$ .

The nodes of an  $n$ -point quadrature formula of **order  $2n$** , if it exists, must coincide with the unique zeros of the polynomials  $\bar{P}_n \in \mathcal{P}_n \setminus \{\mathbf{0}\}$  satisfying (7.3.0.19).

↓

### Remark 7.3.0.20 (Gram-Schmidt orthogonalization of polynomials)

Recall:  $(f, g) \mapsto \int_a^b f(t)g(t) dt$  is an inner product on  $C^0([a, b])$ , the  **$L^2$ -inner product**, see Rem. 6.2.2.1, [NS02, Sect. 4.4, Ex. 2], [Gut09, Ex. 6.5]

- Treat space of polynomials  $\mathcal{P}_n$  as a vector space equipped with an inner product.
- As we have seen in Section 6.2.2, abstract techniques for vector spaces with inner product can be applied to polynomials, for instance **Gram-Schmidt orthogonalization**, cf. § 6.2.1.17, [NS02, Thm. 4.8], [Gut09, Alg. 6.1].

Now carry out the abstract Gram-Schmidt orthogonalization according to Algorithm (6.2.1.18) and recall Thm. 6.2.1.19: in a vector space  $\mathcal{V}$  with inner product  $(\cdot, \cdot)_V$  orthogonal vectors  $q_0, q_1, \dots$  spanning the same subspaces as the linearly independent vectors  $v_0, v_1, \dots$  are constructed recursively via

$$q_{n+1} := v_{n+1} - \sum_{k=0}^n \frac{(v_{n+1}, q_k)_V}{(q_k, q_k)_V} q_k , \quad q_0 := v_0 . \quad (7.3.0.21)$$

➤ Construction of  $\bar{P}_n$  by **Gram-Schmidt orthogonalization** of monomial basis  $\{1, t, t^2, \dots, t^{n-1}\}$  (the  $v_k$ s!) of  $\mathcal{P}_{n-1}$  w.r.t.  $L^2([-1, 1])$ -inner product through the recursion

$$\bar{P}_0(t) := 1 , \quad \bar{P}_{n+1}(t) = t^n - \sum_{k=0}^n \frac{\int_{-1}^1 t^n \bar{P}_k(t) dt}{\int_{-1}^1 \bar{P}_k^2(t) dt} \cdot \bar{P}_k(t) \quad (7.3.0.22)$$

Note:  $\bar{P}_n$  has leading coefficient = 1  $\Rightarrow$   $\bar{P}_n$  uniquely defined (up to sign) by (7.3.0.22). □

The considerations so far only reveal *necessary conditions* on the nodes of an  $n$ -point quadrature rule of order  $2n$ :

They do by no means confirm the existence of such rules, but offer a clear hint on how to construct them:

**Theorem 7.3.0.23. Existence of  $n$ -point quadrature formulas of order  $2n$**

Let  $\{\bar{P}_n\}_{n \in \mathbb{N}_0}$  be a family of non-zero polynomials that satisfies

- $\bar{P}_n \in \mathcal{P}_n$ ,
- $\int_{-1}^1 q(t) \bar{P}_n(t) dt = 0$  for all  $q \in \mathcal{P}_{n-1}$  ( $L^2([-1, 1])$ -orthogonality),
- The set  $\{c_j^n\}_{j=1}^m$ ,  $m \leq n$ , of real zeros of  $\bar{P}_n$  is contained in  $[-1, 1]$ .

Then the quadrature rule ( $\rightarrow$  Def. 7.1.0.1)  $Q_n(f) := \sum_{j=1}^m w_j^n f(c_j^n)$

with weights chosen according to Thm. 7.3.0.6 provides a quadrature formula of order  $2n$  on  $[-1, 1]$ .

*Proof.* Conclude from the orthogonality of the  $\bar{P}_n$  that  $\{\bar{P}_k\}_{k=0}^n$  is a basis of  $\mathcal{P}_n$  and

$$\int_{-1}^1 h(t) \bar{P}_n(t) dt = 0 \quad \forall h \in \mathcal{P}_{n-1} . \quad (7.3.0.24)$$

Recall division of polynomials with remainder (Euclid's algorithm  $\rightarrow$  Course "Diskrete Mathematik"): for any  $p \in \mathcal{P}_{2n-1}$

$$p(t) = h(t) \bar{P}_n(t) + r(t) , \quad \text{for some } h \in \mathcal{P}_{n-1}, r \in \mathcal{P}_{n-1} . \quad (7.3.0.25)$$

Apply this representation to the integral:

$$\int_{-1}^1 p(t) dt = \underbrace{\int_{-1}^1 h(t) \bar{P}_n(t) dt}_{=0 \text{ by (7.3.0.24)}} + \int_{-1}^1 r(t) dt \stackrel{(*)}{=} \sum_{j=1}^m w_j^n r(c_j^n) , \quad (7.3.0.26)$$

(\*) by choice of weights according to Rem. 7.3.0.7  $Q_n$  is exact for polynomials of degree  $\leq n-1$ !

By choice of nodes as zeros of  $\bar{P}_n$  using (7.3.0.24):

$$\sum_{j=1}^m w_j^n p(c_j^n) \stackrel{(7.3.0.25)}{=} \sum_{j=1}^m w_j^n h(c_j^n) \underbrace{\bar{P}_n(c_j^n)}_{=0} + \sum_{j=1}^m w_j^n r(c_j^n) \stackrel{(7.3.0.26)}{=} \int_{-1}^1 p(t) dt.$$

□

**§7.3.0.27 (Legendre polynomials and Gauss-Legendre quadrature)** The family of polynomials  $\{\bar{P}_n\}_{n \in \mathbb{N}_0}$  are so-called **orthogonal polynomials** w.r.t. the  $L^2([-1, 1])$ -inner product, see Def. 6.2.2.6. We have made use of orthogonal polynomials already in Section 6.2.2.  $L^2([-1, 1])$ -orthogonal polynomials play a key role in analysis.

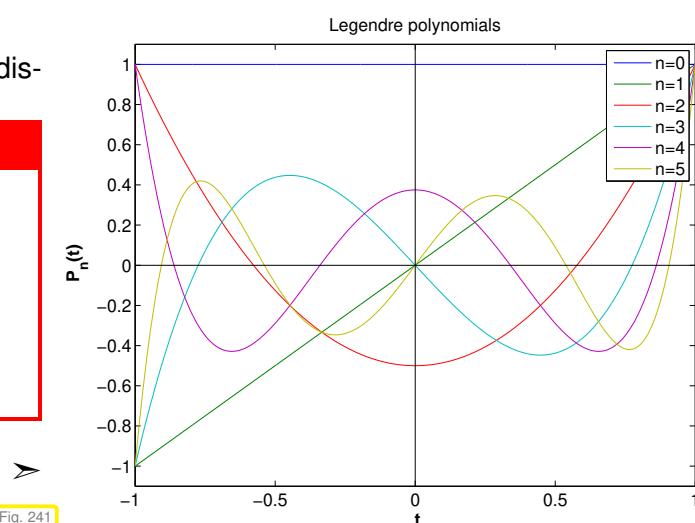
The  $L^2([-1, 1])$ -orthogonal are those already discussed in Rem. 6.2.2.15:

**Definition 7.3.0.28. Legendre polynomials**

The  $n$ -th **Legendre polynomial**  $P_n$  is defined by

- $P_n \in \mathcal{P}_n$ ,
- $\int_{-1}^1 P_n(t) q(t) dt = 0 \quad \forall q \in \mathcal{P}_{n-1}$ ,
- $P_n(1) = 1$ .

Legendre polynomials  $P_0, \dots, P_5$



Notice: the polynomials  $\bar{P}_n$  defined by (7.3.0.22) and the Legendre polynomials  $P_n$  of Def. 7.3.0.28 (merely) differ by a constant factor!



Gauss points  $\xi_j^n =$  zeros of Legendre polynomial  $P_n$

Note: the above considerations, recall (7.3.0.19), show that the nodes of an  $n$ -point quadrature formula of order  $2n$  on  $[-1, 1]$  must agree with the zeros of  $L^2([-1, 1])$ -orthogonal polynomials.



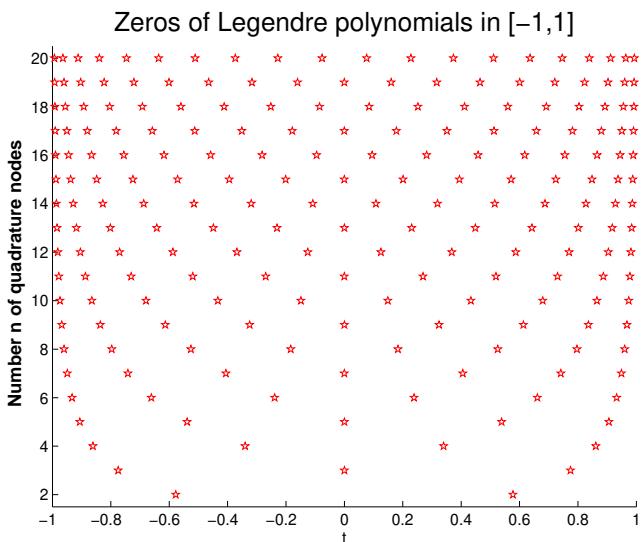
$n$ -point quadrature formulas of order  $2n$  are unique

This is not surprising in light of “ $2n$  equations for  $2n$  degrees of freedom”.



We are not done yet: the zeros of  $\bar{P}_n$  from (7.3.0.22) may lie outside  $[-1, 1]$ . In principle  $\bar{P}_n$  could also have less than  $n$  real zeros.

The next lemma shows that all this cannot happen.



▷ Obviously:

**Lemma 7.3.0.29. Zeros of Legendre polynomials**

$P_n$  has  $n$  distinct zeros in  $[-1,1]$ .

Zeros of Legendre polynomials = Gauss points

*Proof.* (indirect) Assume that  $P_n$  has only  $m < n$  zeros  $\zeta_1, \dots, \zeta_m$  in  $[-1,1]$  at which it changes sign. Define

$$\begin{aligned} q(t) := \prod_{j=1}^m (t - \zeta_j) &\Rightarrow qP_n \geq 0 \quad \text{or} \quad qP_n \leq 0. \\ &\Rightarrow \int_{-1}^1 q(t)P_n(t) dt \neq 0. \end{aligned}$$

As  $q \in \mathcal{P}_{n-1}$ , this contradicts (7.3.0.24). □

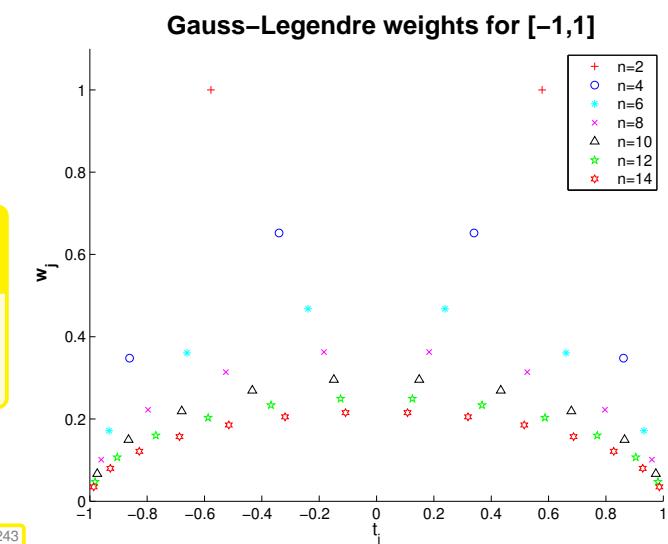
### Definition 7.3.0.30. Gauss-Legendre quadrature formulas

The  $n$ -point Quadrature formulas whose nodes, the **Gauss points**, are given by the zeros of the  $n$ -th Legendre polynomial (→ Def. 7.3.0.28), and whose weights are chosen according to Thm. 7.3.0.6, are called **Gauss-Legendre quadrature formulas**.

Obviously ▷

### Lemma 7.3.0.31. Positivity of Gauss-Legendre quadrature weights

The weights of the Gauss-Legendre quadrature formulas are positive.



*Proof.* Writing  $\xi_j^n$ ,  $j = 1, \dots, n$ , for the nodes (Gauss points) of the  $n$ -point Gauss-Legendre quadrature

formula,  $n \in \mathbb{N}$ , we define

$$q_k(t) = \prod_{\substack{j=1 \\ j \neq k}}^n (t - \xi_j^n)^2 \Rightarrow q_k \in \mathcal{P}_{2n-2}.$$

This polynomial is integrated exactly by the quadrature rule: since  $q_k(\xi_j^n) = 0$  for  $j \neq k$

$$0 < \int_{-1}^1 q(t) dt = w_k^n \underbrace{q(\xi_k^n)}_{>0},$$

where  $w_j^n$  are the quadrature weights.  $\square$

**Remark 7.3.0.32 (3-Term recursion for Legendre polynomials)** From Thm. 6.2.2.13 we learn the orthogonal polynomials satisfy the 3-term recursion (6.2.2.14), see also (7.3.0.34). To keep this chapter self-contained we derive it independently for Legendre polynomials.

Note: the polynomials  $\bar{P}_n$  from (7.3.0.22) are uniquely characterized by the two properties (try a proof!)

- (i)  $\bar{P}_n \in \mathcal{P}_n$  with leading coefficient 1:  $\bar{P}(t) = t^n + \dots$ ,
- (ii)  $\int_{-1}^1 \bar{P}_k(t) \bar{P}_j(t) dt = 0$ , if  $j \neq k$  ( $L^2([-1, 1])$ -orthogonality).

➤ we get the same polynomials  $\bar{P}_n$  by another Gram-Schmidt orthogonalization procedure, cf. (7.3.0.21) and § 6.2.2.10:

$$\bar{P}_{n+1}(t) = t \bar{P}_n(t) - \sum_{k=0}^n \frac{\int_{-1}^1 \tau \bar{P}_n(\tau) \bar{P}_k(\tau) d\tau}{\int_{-1}^1 \bar{P}_k^2(\tau) d\tau} \cdot \bar{P}_k(t)$$

By the orthogonality property (7.3.0.24) the sum collapses, since

$$\int_{-1}^1 \tau \bar{P}_n(\tau) \bar{P}_k(\tau) d\tau = \int_{-1}^1 \bar{P}_n(\tau) \underbrace{(\tau \bar{P}_k(\tau))}_{\in \mathcal{P}_{k+1}} d\tau = 0,$$

if  $k+1 < n$ :

$$\bar{P}_{n+1}(t) = t \bar{P}_n(t) - \frac{\int_{-1}^1 \tau \bar{P}_n(\tau) \bar{P}_n(\tau) d\tau}{\int_{-1}^1 \bar{P}_n^2(\tau) d\tau} \cdot \bar{P}_n(t) - \frac{\int_{-1}^1 \tau \bar{P}_n(\tau) \bar{P}_{n-1}(\tau) d\tau}{\int_{-1}^1 \bar{P}_{n-1}^2(\tau) d\tau} \cdot \bar{P}_{n-1}(t). \quad (7.3.0.33)$$

After rescaling (tedious!): **3-term recursion** for Legendre polynomials

$$P_{n+1}(t) := \frac{2n+1}{n+1} t P_n(t) - \frac{n}{n+1} P_{n-1}(t), \quad P_0 := 1, \quad P_1(t) := t. \quad (7.3.0.34)$$

Reminder (→ Section 6.1.3.1): we have a similar 3-term recursion (6.1.3.5) for Chebychev polynomials. Coincidence? Of course not, nothing in mathematics holds “by accident”. By Thm. 6.2.2.13 3-term recursions are a distinguishing feature of so-called families of **orthogonal polynomials**, to which the Chebychev polynomials belong as well, spawned by Gram-Schmidt orthogonalization with respect to a weighted  $L^2$ -inner product, however, see [Han02, p. VI].

- Efficient and *stable* evaluation of Legendre polynomials by means of 3-term recursion (7.3.0.34), cf. the analogous algorithm for Chebychev polynomials given in Code 6.1.3.6.

### C++-code 7.3.0.35: computing Legendre polynomials

```

2 // returns the values of the first n - 1 legendre polynomials
3 // in point x as columns of the matrix L
4 void legendre(const unsigned n, const VectorXd& x, MatrixXd& L) {
5     L = MatrixXd::Ones(n,n); //  $p_0(x) = 1$ 
6     L.col(1) = x; //  $p_1(x) = x$ 
7     for (unsigned j = 1; j < n - 1; ++j) {
8         //  $p_{j+1}(x) = \frac{2j+1}{j+1}xp_j(x) - \frac{j}{j+1}p_{j-1}(x)$  Eq. (7.3.0.34)
9         V.col(j+1) = (2.*j+1)/(j+1.)*V.col(j-1).cwiseProduct(x)
10        - j/(j+1.)*V.col(j-1);
11    }
12}

```

Comments on Code 7.3.0.35:

- ☛ return value: matrix  $\mathbf{V}$  with  $(\mathbf{V})_{ij} = P_i(x_j)$
- ☛ line 2: takes into account initialization of Legendre 3-term recursion (7.3.0.34)

**Remark 7.3.0.36 (Computing Gauss nodes and weights)** There are several efficient ways to find the Gauss points. Here we discuss an intriguing connection with an eigenvalue problem.

Compute nodes/weights of Gaussian quadrature by solving an eigenvalue problem!

([Golub-Welsch algorithm](#) [[Gan+05](#), Sect. 3.5.4], [[Tre08](#), Sect. 1])

In codes Gauss nodes and weights are usually retrieved from tables, cf. Rem. 7.1.0.6.

### C++-code 7.3.0.37: Golub-Welsch algorithm

```

2 struct QuadRule {
3     Eigen::VectorXd nodes, weights;
4 };
5
6 QuadRule gaussquad_(const unsigned n) {
7     QuadRule qr;
8     // Symmetric matrix whose eigenvalues provide Gauss points
9     Eigen::MatrixXd M = Eigen::MatrixXd::Zero(n, n);
10    for (unsigned i = 1; i < n; ++i) { //
11        const double b = i / std::sqrt(4. * i * i - 1.);
12        M(i, i - 1) = M(i - 1, i) = b;
13    }
14    // using Eigen's built-in solver for symmetric eigenvalue problems
15    Eigen::SelfAdjointEigenSolver<Eigen::MatrixXd> eig(M);
16
17    qr.nodes = eig.eigenvalues(); //
18    qr.weights = 2 * eig.eigenvectors().topRows<1>().array().pow(2); //
19    return qr;
20}

```

Justification: rewrite 3-term recurrence (7.3.0.34) for scaled Legendre polynomials  $\tilde{P}_n = \frac{1}{\sqrt{n+1/2}} P_n$

$$t\tilde{P}_n(t) = \underbrace{\frac{n}{\sqrt{4n^2-1}}}_{=: \beta_n} \tilde{P}_{n-1}(t) + \underbrace{\frac{n+1}{\sqrt{4(n+1)^2-1}}}_{=: \beta_{n+1}} \tilde{P}_{n+1}(t). \quad (7.3.0.38)$$

For fixed  $t \in \mathbb{R}$  (7.3.0.38) can be expressed as

$$t \begin{bmatrix} \tilde{P}_0(t) \\ \tilde{P}_1(t) \\ \vdots \\ \tilde{P}_{n-1}(t) \end{bmatrix} = \underbrace{\begin{bmatrix} 0 & \beta_1 & & & \\ \beta_1 & 0 & \beta_2 & & \\ & \beta_2 & \ddots & \ddots & \\ & & \ddots & \ddots & \ddots \\ & & & 0 & \beta_{n-1} \\ & & & & \beta_{n-1} & 0 \end{bmatrix}}_{=:J_n \in \mathbb{R}^{n,n}} \begin{bmatrix} \tilde{P}_0(t) \\ \tilde{P}_1(t) \\ \vdots \\ \tilde{P}_{n-1}(t) \end{bmatrix} + \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \beta_n \tilde{P}_n(t) \end{bmatrix}$$

►  $\tilde{P}_n(\xi) = 0 \Leftrightarrow \xi \mathbf{p}(\xi) = J_n \mathbf{p}(\xi)$  (an eigenvalue problem!).

► The zeros of  $P_n$  can be obtained as the  $n$  real eigenvalues of the symmetric tridiagonal matrix  $J_n \in \mathbb{R}^{n,n}$ !

This matrix  $J_n$  is initialized in Line 10–Line 13 of Code 7.3.0.37. The computation of the weights in Line 18 of Code 7.3.0.37 is explained in [Gan+05, Sect. 3.5.4].

**§7.3.0.39 (Quadrature error and best approximation error)** The positivity of the weights  $w_j^n$  for all  $n$ -point Gauss-Legendre and Clenshaw-Curtis quadrature rules has important consequences.

**Theorem 7.3.0.40. Quadrature error estimate for quadrature rules with positive weights**

For every  $n$ -point quadrature rule  $Q_n$  as in (7.1.0.2) of order  $q \in \mathbb{N}$  with weights  $w_j \geq 0$ ,  $j = 1, \dots, n$  the quadrature error satisfies

$$E_n(f) := \left| \int_a^b f(t) dt - Q_n(f) \right| \leq 2|b-a| \underbrace{\inf_{p \in \mathcal{P}_{q-1}} \|f - p\|_{L^\infty([a,b])}}_{\text{best approximation error}} \quad \forall f \in C^0([a,b]). \quad (7.3.0.41)$$

*Proof.* The proof runs parallel to the derivation of (6.1.2.30). Writing  $E_n(f)$  for the quadrature error, the left hand side of (7.3.0.41), we find by the definition Def. 7.3.0.1 of the order of a quadrature rule

$$\begin{aligned} E_n(f) = E_n(f - p) &\leq \left| \int_a^b (f - p)(t) dt \right| + \left| \sum_{j=1}^n w_j (f - p)(c_j) \right| \\ &\leq |b-a| \|f - p\|_{L^\infty([a,b])} + \left( \sum_{j=1}^n |w_j| \right) \|f - p\|_{L^\infty([a,b])}. \end{aligned} \quad (7.3.0.42)$$

Since the quadrature rule is exact for constants and  $w_j \geq 0$

$$\sum_{j=1}^n |w_j| = \sum_{j=1}^n w_j = |b-a|,$$

which finishes the proof. □

Drawing on best approximation estimates from Section 6.1.1 and Rem. 6.1.3.22, we immediately get results about the asymptotic decay of the quadrature error for  $n$ -point Gauss-Legendre and Clenshaw-Curtis quadrature as  $n \rightarrow \infty$ :

$$\begin{aligned} f \in C^r([a, b]) &\Rightarrow E_n(f) \rightarrow 0 \text{ algebraically with rate } r, \\ f \in C^\infty([a, b]) \text{ "analytically extensible"} &\Rightarrow E_n(f) \rightarrow 0 \text{ exponentially,} \end{aligned}$$

as  $n \rightarrow \infty$ , see Def. 6.1.2.7 for type of convergence.

Appealing to Thm. 6.1.1.11 and Rem. 6.1.1.19, and (6.1.2.19), the dependence of the constants on the length of the integration interval can be quantified for integrands with limited smoothness.

#### Lemma 7.3.0.43. Quadrature error estimates for $C^r$ -integrands

For every  $n$ -point quadrature rule  $Q_n$  as in (7.1.0.2) of order  $q \in \mathbb{N}$  with weights  $w_j \geq 0$ ,  $j = 1, \dots, n$  we find that the quadrature error  $E_n(f)$  for an integrand  $f \in C^r([a, b])$ ,  $r \in \mathbb{N}_0$ , satisfies

$$\text{in the case } q \geq r: \quad E_n(f) \leq C q^{-r} |b-a|^{r+1} \|f^{(r)}\|_{L^\infty([a,b])}, \quad (7.3.0.44)$$

$$\text{in the case } q < r: \quad E_n(f) \leq \frac{|b-a|^{q+1}}{q!} \|f^{(q)}\|_{L^\infty([a,b])}, \quad (7.3.0.45)$$

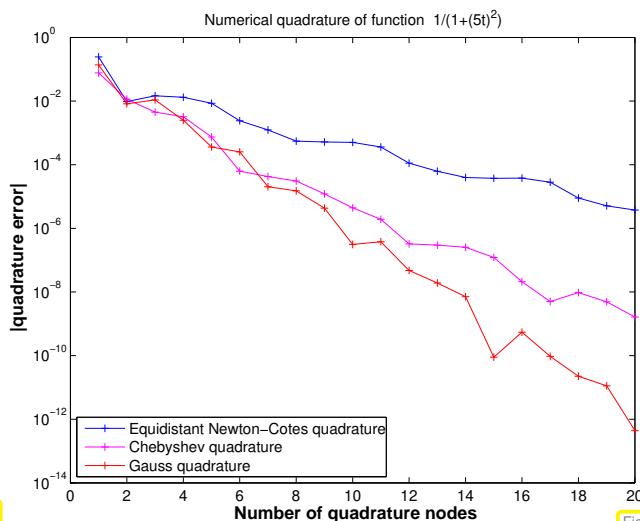
with a constant  $C > 0$  independent of  $n$ ,  $f$ , and  $[a, b]$ .

Please note the different estimates depending on whether the smoothness of  $f$  (as described by  $r$ ) or the order of the quadrature rule is the “limiting factor”.  $\dashv$

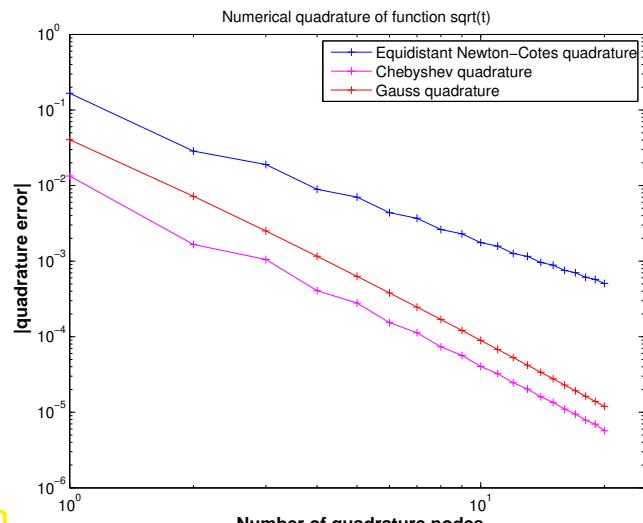
#### EXAMPLE 7.3.0.46 (Convergence of global quadrature rules)

We examine three families of global polynomial ( $\rightarrow$  Thm. 7.3.0.6) quadrature rules: Newton-Cotes formulas, Gauss-Legendre rules, and Clenshaw-Curtis rules. We record the convergence of the quadrature errors for the interval  $[0, 1]$  and two different functions

1.  $f_1(t) = \frac{1}{1+(5t)^2}$ , an analytic function, see Rem. 6.1.2.41,
2.  $f_2(t) = \sqrt{t}$ , merely continuous, derivatives singular in  $t = 0$ .



quadrature error,  $f_1(t) := \frac{1}{1+(5t)^2}$  on  $[0, 1]$



quadrature error,  $f_2(t) := \sqrt{t}$  on  $[0, 1]$

Asymptotic behavior of quadrature error  $\epsilon_n := \left| \int_0^1 f(t) dt - Q_n(f) \right|$  for " $n \rightarrow \infty$ ":

$\Rightarrow$

exponential convergence  $\epsilon_n \approx O(q^n)$ ,  $0 < q < 1$ , for  $C^\infty$ -integrand  $f_1$ : Newton-Cotes quadrature :  $q \approx 0.61$ , Clenshaw-Curtis quadrature :  $q \approx 0.40$ , Gauss-Legendre quadrature :  $q \approx 0.27$

- algebraic convergence  $\epsilon_n \approx O(n^{-\alpha})$ ,  $\alpha > 0$ , for integrand  $f_2$  with singularity at  $t = 0 \rightsquigarrow$  Newton-Cotes quadrature :  $\alpha \approx 1.8$ , Clenshaw-Curtis quadrature :  $\alpha \approx 2.5$ , Gauss-Legendre quadrature :  $\alpha \approx 2.7$

**Remark 7.3.0.47 (Removing a singularity by transformation)** From Ex. 7.3.0.46 teaches us that a lack of smoothness of the integrand can thwart exponential convergence and severely limits the rate of algebraic convergence of a global quadrature rule for  $n \rightarrow \infty$ .

Idea: recover integral with smooth integrand by “analytic preprocessing”

Here is an example:

For a general but smooth  $f \in C^\infty([0, b])$  compute  $\int_0^b \sqrt{t}f(t) dt$  via a quadrature rule, e.g.,  $n$ -point Gauss-Legendre quadrature on  $[0, b]$ . Due to the presence of a square-root singularity at  $t = 0$  the direct application of  $n$ -point Gauss-Legendre quadrature will result in a rather slow algebraic convergence of the quadrature error as  $n \rightarrow \infty$ , see Ex. 7.3.0.46.

Trick: Transformation of integrand by substitution rule:

$$\text{substitution } s = \sqrt{t}: \quad \int_0^b \sqrt{t}f(t) dt = \int_0^{\sqrt{b}} 2s^2 f(s^2) ds. \quad (7.3.0.48)$$

Then: Apply Gauss-Legendre quadrature rule to smooth integrand

**Remark 7.3.0.49 (The message of asymptotic estimates)** There is one blot on most  $n$ -asymptotic estimates obtained from Thm. 7.3.0.40: the bounds usually involve quantities like norms of higher derivatives of the interpolant that are elusive in general, in particular for integrands given only in procedural form, see § 7.0.0.2. Such unknown quantities are often hidden in “generic constants  $C$ ”. Can we extract useful information from estimates marred by the presence of such constants?

For fixed integrand  $f$  let us assume sharp algebraic convergence (in  $n$ ) with rate  $r \in \mathbb{N}$  of the quadrature error  $E_n(f)$  for a family of  $n$ -point quadrature rules:

$$E_n(f) = O(n^{-r}) \stackrel{\text{sharp}}{\implies} E_n(f) \approx Cn^{-r}, \quad (7.3.0.50)$$

with a “generic constant  $C > 0$ ” independent of  $n$ .

Goal: Reduction of the quadrature error by a factor of  $\rho > 1$

Which (minimal) increase in the number  $n$  of quadrature points accomplishes this?

$$\frac{Cn_{\text{old}}^{-r}}{Cn_{\text{new}}^{-r}} = \rho \Leftrightarrow n_{\text{new}} : n_{\text{old}} = \sqrt[r]{\rho}. \quad (7.3.0.51)$$

In the case of *algebraic convergence* with rate  $r \in \mathbb{R}$  a reduction of the quadrature error by a factor of  $\rho$  is bought by an increase of the number of quadrature points by a factor of  $\rho^{1/r}$ .

(7.3.0.44)  $\geq$  gains in accuracy are “cheaper” for smoother integrands!

Now assume *sharp exponential convergence* (in  $n$ ) of the quadrature error  $E_n(f)$  for a family of  $n$ -point quadrature rules,  $0 \leq q < 1$ :

$$E_n(f) = O(q^n) \xrightarrow{\text{sharp}} E_n(f) \approx Cq^n , \quad (7.3.0.52)$$

with a “generic constant  $C > 0$ ” *independent* of  $n$ .

Error reduction by a factor  $\rho > 1$  results from

$$\frac{Cq^{n_{\text{old}}}}{Cq^{n_{\text{new}}}} = \rho \Leftrightarrow n_{\text{new}} - n_{\text{old}} = -\frac{\log \rho}{\log q} .$$

In the case of *exponential convergence* (7.3.0.52) a *fixed* increase of the number of quadrature points by  $-\log \rho : \log q$  results in a reduction of the quadrature error by a factor of  $\rho > 1$ .

## 7.4 Composite Quadrature

In 6, Section 6.5.1 we studied approximation by piecewise polynomial interpolants. A similar idea underlies the so-called composite quadrature rules on an interval  $[a, b]$ . Analogously to piecewise polynomial techniques they start from a grid/mesh

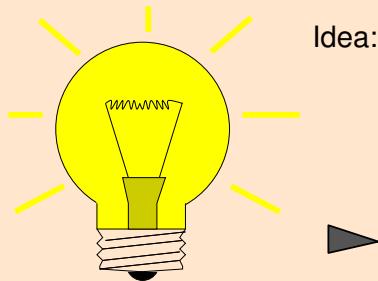
$$\mathcal{M} := \{a = x_0 < x_1 < \dots < x_{m-1} < x_m = b\} \quad (6.5.0.2)$$

and appeal to the trivial identity

$$\int_a^b f(t) dt = \sum_{j=1}^m \int_{x_{j-1}}^{x_j} f(t) dt . \quad (7.4.0.1)$$

On each mesh interval  $[x_{j-1}, x_j]$  we then use a *local quadrature rule*, which may be one of the polynomial quadrature formulas from 7.2.

## General construction of composite quadrature rules



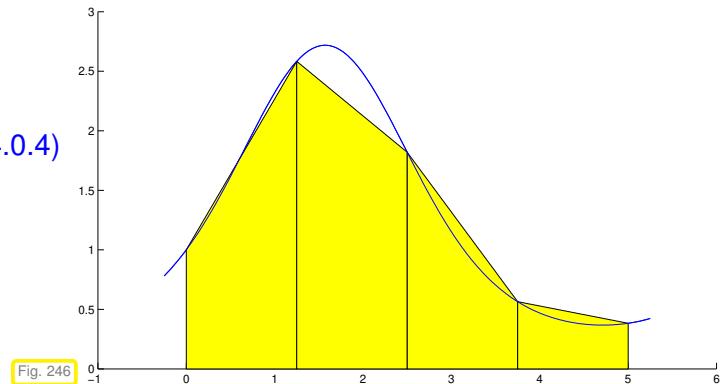
- Idea:
- Partition integration domain  $[a, b]$  by a mesh/grid ( $\rightarrow$  Section 6.5)  
 $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$
  - Apply quadrature formulas from Section 7.2, Section 7.3 **locally** on mesh intervals  $I_j := [x_{j-1}, x_j], j = 1, \dots, m$ , and sum up.

composite quadrature rule

### EXAMPLE 7.4.0.3 (Simple composite polynomial quadrature rules)

Composite trapezoidal rule, cf. (7.2.0.5)

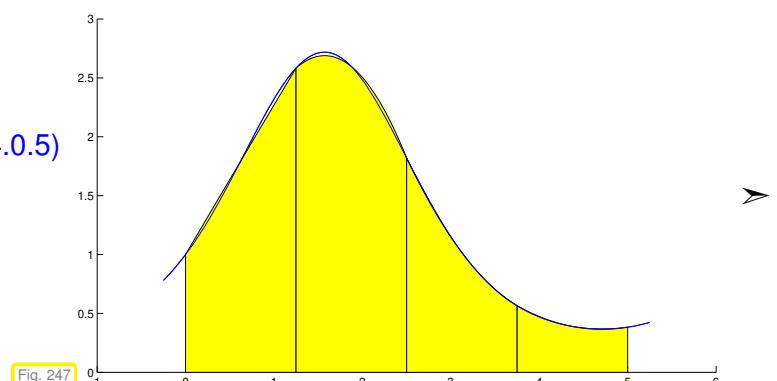
$$\int_a^b f(t) dt = \frac{1}{2}(x_1 - x_0)f(a) + \sum_{j=1}^{m-1} \frac{1}{2}(x_{j+1} - x_{j-1})f(x_j) + \frac{1}{2}(x_m - x_{m-1})f(b). \quad (7.4.0.4)$$



➤ arising from piecewise linear interpolation of  $f$ .

Composite Simpson rule, cf. (7.2.0.6)

$$\int_a^b f(t) dt = \frac{1}{6}(x_1 - x_0)f(a) + \sum_{j=1}^{m-1} \frac{1}{6}(x_{j+1} - x_{j-1})f(x_j) + \sum_{j=1}^m \frac{2}{3}(x_j - x_{j-1})f(\frac{1}{2}(x_j + x_{j-1})) + \frac{1}{6}(x_m - x_{m-1})f(b). \quad (7.4.0.5)$$



related to piecewise quadratic Lagrangian interpolation.

Formulas (7.4.0.4), (7.4.0.5) directly suggest efficient implementation with minimal number of  $f$ -evaluations.

#### C++-code 7.4.0.6: Equidistant composite trapezoidal rule (7.4.0.4)

```

1 #include <cassert>
2
3 // N-interval equidistant trapezoidal rule
4 template <class Function>
5 double trapezoidal(Function &f, const double a, const double b,
6                     const unsigned N) {
7     double l = 0;

```

```

8  const double h = (b - a) / N; // interval length
9
10 for (unsigned i = 0; i < N; ++i) {
11     // rule:  $T = (b - a)/2 * (f(a) + f(b))$ ,
12     // apply on N intervals:  $[a + i \cdot h, a + (i+1) \cdot h]$ ,  $i=0..(N-1)$ 
13     l += h / 2 * (f(a + i * h) + f(a + (i + 1) * h));
14 }
15 return l;
16 }

17 // Alternative implementation of n-point equidistant trapezoidal rule
18 template <typename Functor>
19 double equidTrapezoidalRule(Functor &&f, double a, double b, unsigned int n) {
20     assert(n>=2);
21     const double h = (b - a)/(n-1);
22     double t = a + h, s = 0.0;
23     for (unsigned int i = 1; i < n-1; t += h, ++i)
24         s += f(t);
25     return (0.5 * h * f(a) + 0.5 * h * f(b) + h * s);
26 }
27 }
```

### C++-code 7.4.0.7: Equidistant composite Simpson rule (7.4.0.5)

```

1 template <class Function>
2 double simpson(Function& f, const double a, const double b, const unsigned N) {
3     double l = 0;
4     const double h = (b - a)/N; // intervall length
5
6     for (unsigned i = 0; i < N; ++i) {
7         // rule:  $S = (b - a)/6 * (f(a) + 4*f(0.5*(a + b)) + f(b))$ 
8         // apply on  $[a + i \cdot h, a + (i+1) \cdot h]$ 
9         l += h/6 * (f(a + i * h) + 4*f(a + (i+0.5)*h) + f(a + (i+1)*h));
10    }
11
12    return l;
13 }
```

In both cases the function object passed in  $f$  must provide an evaluation operator **double operator (double) const**. □

**Remark 7.4.0.8 (Composite quadrature and piecewise polynomial interpolation)** Composite quadrature scheme based on local polynomial quadrature can usually be understood as “quadrature by approximation schemes” as explained in § 7.1.0.7. The underlying approximation schemes belong to the class of general local Lagrangian interpolation schemes introduced in Section 6.5.1.

In other words, many composite quadrature schemes arise from replacing the integrand by a piecewise interpolating polynomial, see Fig. 246 and Fig. 247 and compare with Fig. 225. □

To see the main rationale behind the use of composite quadrature rules recall Lemma 7.3.0.43: for a polynomial quadrature rule (7.2.0.1) of order  $q$  with positive weights and  $f \in C^r([a,b])$  the quadrature error shrinks with the  $\min\{r,q\} + 1$ -st power of the length  $|b - a|$  of the integration domain! Hence, applying polynomial quadrature rules to small mesh intervals should lead to a small overall quadrature error.

**§7.4.0.9 (Quadrature error estimate for composite polynomial quadrature rules)** Assume a composite quadrature rule  $\mathcal{Q}$  on  $[x_0, x_m] = [a, b]$ ,  $b > a$ , based on  $n_j$ -point local quadrature rules  $\mathcal{Q}_{n_j}^j$  with **positive weights** (e.g. local Gauss-Legendre quadrature rules or local Clenshaw-Curtis quadrature rules) and of *fixed* orders  $q_j \in \mathbb{N}$  on each mesh interval  $[x_{j-1}, x_j]$ . From Lemma 7.3.0.43 recall the estimate for  $f \in C^r([x_{j-1}, x_j])$

$$\left| \int_{x_{j-1}}^{x_j} f(t) dt - Q_{n_j}^j(f) \right| \leq C |x_j - x_{j-1}|^{\min\{r, q_j\}+1} \|f^{(\min\{r, q_j\})}\|_{L^\infty([x_{j-1}, x_j])}. \quad (7.2.0.10)$$

with  $C > 0$  independent of  $f$  and  $j$ . For  $f \in C^r([a, b])$ , summing up these bounds we get for the global quadrature error

$$\left| \int_{x_0}^{x_m} f(t) dt - Q(f) \right| \leq C \sum_{j=1}^m h_j^{\min\{r, q_j\}+1} \|f^{(\min\{r, q_j\})}\|_{L^\infty([x_{j-1}, x_j])},$$

with local meshwidths  $h_j = x_j - x_{j-1}$ . If  $q_j = q$ ,  $q \in \mathbb{N}$ , for all  $j = 1, \dots, m$ , then, as  $\sum_j h_j = b - a$ ,

$$\left| \int_{x_0}^{x_m} f(t) dt - Q(f) \right| \leq C h_M^{\min\{q, r\}} |b - a| \|f^{(\min\{q, r\})}\|_{L^\infty([a, b])}, \quad (7.4.0.10)$$

with (global) meshwidth  $h_M := \max_j h_j$ .

(7.4.0.10)  $\longleftrightarrow$  **Algebraic convergence** in no. of  $f$ -evaluations for  $n \rightarrow \infty$

**§7.4.0.11 (Constructing families of composite quadrature rules)** As with polynomial quadrature rules, we study the asymptotic behavior of the quadrature error for families of composite quadrature rules as a function on the total number  $n$  of function evaluations.

As in the case of  $M$ -piecewise polynomial approximation of function ( $\rightarrow$  Section 6.5.1) families of composite quadrature rules can be generated in two different ways:

(I) use a sequence of successively refined meshes  $(M_k = \{x_j^k\}_j)_{k \in \mathbb{N}}$  with  $\#M_k = m(k) + 1$ ,  $m(k) \rightarrow \infty$  for  $k \rightarrow \infty$ , combined with the *same* (transformed,  $\rightarrow$  Rem. 7.1.0.4) local quadrature rule on all mesh intervals  $[x_{j-1}^k, x_j^k]$ . Examples are the composite trapezoidal rule and composite Simpson rule from Ex. 7.4.0.3 on sequences of equidistant meshes.

➤ **h-convergence**

(II) On a *fixed* mesh  $M = \{x_j\}_{j=0}^m$ , on each cell use the *same* (transformed) local quadrature rule taken from a *sequence* of polynomial quadrature rules of increasing order.

➤ **p-convergence**

**EXPERIMENT 7.4.0.12 (Quadrature errors for composite quadrature rules)** Composite quadrature rules based on

- trapezoidal rule (7.2.0.5) ➤ local order 2 (exact for linear functions, see Ex. 7.3.0.10),
- Simpson rule (7.2.0.6) ➤ local order 4 (exact for cubic polynomials, see Ex. 7.3.0.10)

on equidistant mesh  $\mathcal{M} := \{jh\}_{j=0}^n$ ,  $h = 1/n$ ,  $n \in \mathbb{N}$ .

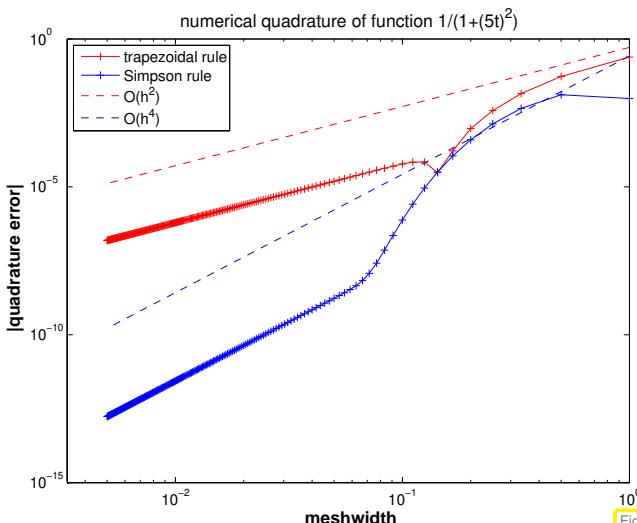


Fig. 248

$$\text{quadrature error, } f_1(t) := \frac{1}{1+(5t)^2} \text{ on } [0, 1]$$

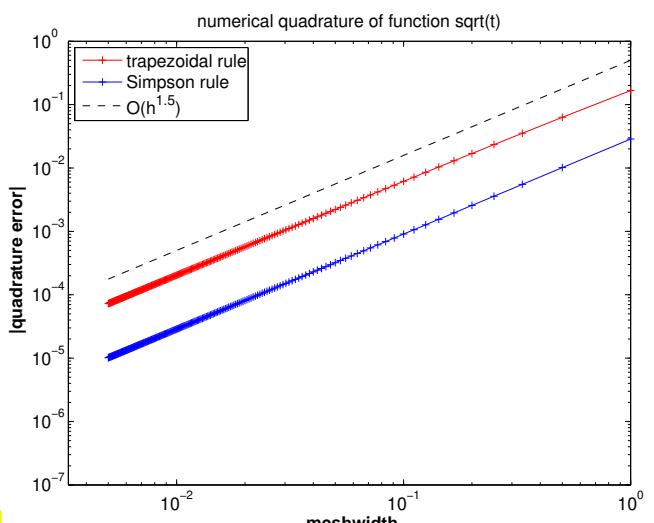


Fig. 249

$$\text{quadrature error, } f_2(t) := \sqrt{t} \text{ on } [0, 1]$$

Asymptotic behavior of quadrature error  $E(n) := \left| \int_0^1 f(t) dt - Q_n(f) \right|$  for meshwidth "h → 0":

- Throughout we observe algebraic convergence  $E(n) = O(h^\alpha)$  of with rate  $\alpha > 0$  for  $h = n^{-1} \rightarrow 0$
- Sufficiently smooth integrand  $f_1$ : trapezoidal rule →  $\alpha = 2$ , Simpson rule →  $\alpha = 4$  !?
- singular integrand  $f_2$ :  $\alpha = 3/2$  for trapezoidal rule & Simpson rule !

(lack of) smoothness of integrand limits convergence!

↓

**Remark 7.4.0.13 (Composite quadrature rules vs. global quadrature rules)** For a fixed integrand  $f \in C^r([a, b])$  of limited smoothness on an interval  $[a, b]$  we compare

- a family of composite quadrature rules based on single local  $\ell$ -point rule (with positive weights) of order  $q$  on a sequence of equidistant meshes  $(\mathcal{M}_k = \{x_j^k\}_{j=0}^{\#\mathcal{M}_k})_{k \in \mathbb{N}}$ ,
- the family of Gauss-Legendre quadrature rules from Def. 7.3.0.30.

We study the asymptotic dependence of the quadrature error on the number  $n$  of function evaluations.

For the composite quadrature rules we have  $n \approx \ell \#\mathcal{M}_k \approx \ell h_M^{-1}$ . Combined with (7.4.0.10), we find for quadrature error  $E_n^{\text{comp}}(f)$  of the composite quadrature rules

$$E_n^{\text{comp}}(f) \leq C_1 n^{-\min\{q, r\}}, \quad (7.4.0.14)$$

with  $C_1 > 0$  independent of  $\mathcal{M} = \mathcal{M}_k$ .

The quadrature errors  $E_n^{\text{GL}}(f)$  of the  $n$ -point Gauss-Legendre quadrature rules are given in Lemma 7.3.0.43, (7.3.0.44):

$$E_n^{\text{GL}}(f) \leq C_2 n^{-r}, \quad (7.4.0.15)$$

with  $C_2 > 0$  independent of  $n$ .

Gauss-Legendre quadrature converges *at least as fast* fixed order composite quadrature on equidistant meshes.

Moreover, Gauss-Legendre quadrature “automatically detects” the smoothness of the integrand, and enjoys fast exponential convergence for analytic integrands.

Use Gauss-Legendre quadrature instead of fixed order composite quadrature on equidistant meshes.

**EXPERIMENT 7.4.0.16 (Empiric convergence of equidistant trapezoidal rule)** Sometimes there are surprises: Now we will witness a convergence behavior of a composite quadrature rule that is much better than predicted by the order of the local quadrature formula.

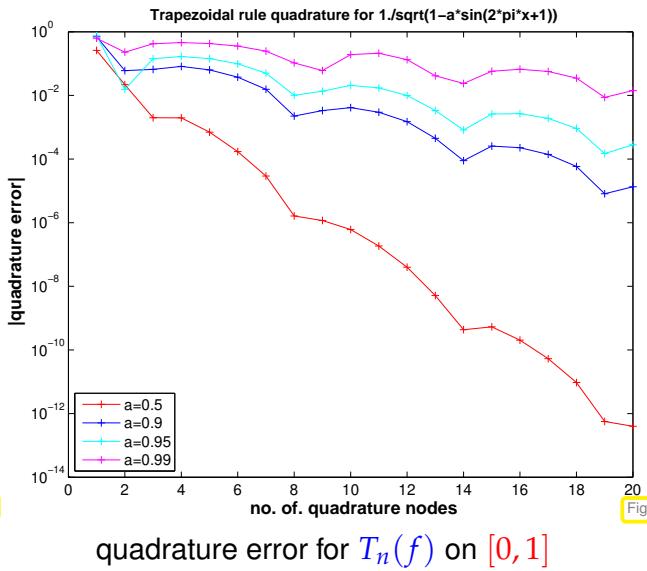
We consider the equidistant trapezoidal rule (order 2), see (7.4.0.4), Code 7.4.0.6

$$\int_a^b f(t) dt \approx T_m(f) := h \left( \frac{1}{2}f(a) + \sum_{k=1}^{m-1} f(kh) + \frac{1}{2}f(b) \right), \quad h := \frac{b-a}{m}. \quad (7.4.0.17)$$

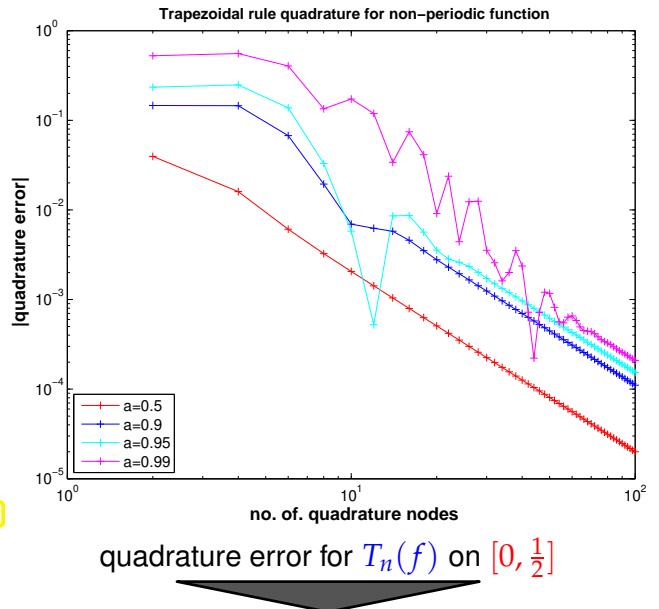
and the 1-periodic smooth (analytic) integrand

$$f(t) = \frac{1}{\sqrt{1-a \sin(2\pi t-1)}}, \quad 0 < a < 1.$$

(As “exact value of integral” we use  $T_{500}$  in the computation of quadrature errors.)



exponential convergence !!



merely algebraic convergence

### §7.4.0.18 (Equidistant trapezoidal rule (for periodic integrands))

In this § we use  $I := [0, 1]$  as a reference interval, cf. Exp. 7.4.0.16. We rely on similar techniques as in Section 5.6, Section 5.6.2. Again, a key tool will be the bijective mapping, see Fig. 178,

$$\Phi_{S^1} : I \rightarrow S^1 := \{z \in \mathbb{C} : |z| = 1\} \quad , \quad t \mapsto z := \exp(2\pi i t), \quad (5.6.2.1)$$

which induces the general pullback, *c.f.* (6.1.1.16),

$$(\Phi_{S^1}^{-1})^* : C^0([0,1]) \rightarrow C^0(S^1) \quad , \quad ((\Phi_{S^1}^{-1})^* f)(z) := f(\Phi_{S^1}^{-1}(z)) \quad , \quad z \in S^1 .$$

If  $f \in C^r(\mathbb{R})$  and **1-periodic**, then  $(\Phi_{S^1}^{-1})^* f \in C^r(S^1)$ . Further,  $\Phi_{S^1}$  maps equidistant nodes on  $I := [0,1]$  to equispaced nodes on  $S^1$ , which are the **roots of unity**:

$$\Phi_{S^1}\left(\frac{j}{n}\right) = \exp(2\pi i \frac{j}{n}) \quad [ \exp(2\pi i \frac{j}{n})^n = 1; ] . \quad (7.4.0.19)$$

Now consider an  $n$ -point polynomial quadrature rule on  $S^1$  based on the set of equidistant nodes  $\mathcal{Z} := \{z_j := \exp(2\pi i \frac{j-1}{n}), j = 1, \dots, n\}$  and defined as

$$Q_n^{S^1}(g) := \int_{S^1} (\mathcal{L}_{\mathcal{Z}} g)(\tau) dS(\tau) = \sum_{j=1}^n w_j^{S^1} g(z_j) , \quad (7.4.0.20)$$

where  $\mathcal{L}_{\mathcal{Z}}$  is the Lagrange interpolation operator ( $\rightarrow$  Def. 6.1.2.1). This means that the weights obey Thm. 7.3.0.6, where the definition (5.2.2.4) of Lagrange polynomials remains the same for complex nodes. By sheer *symmetry*, all the weights have to be the same, which, since the rule will be at least of order 1, means

$$w_j^{S^1} = \frac{2\pi}{n} , \quad j = 1, \dots, n .$$

Moreover, the quadrature rule  $Q_n^{S^1}$  will be of order  $n$ , see Def. 7.3.0.1, that is, it will *integrate polynomials of degree  $\leq n - 1$ , exactly*.

By transformation ( $\rightarrow$  Rem. 7.1.0.4 and pullback (7.4.0.18)),  $Q_n^{S^1}$  induces a quadrature rule on  $I := [0,1]$  by

$$Q_n^I(f) := \frac{1}{2\pi} Q_n^{S^1}((\Phi_{S^1}^{-1})^* f) = \frac{1}{2\pi} \sum_{j=1}^n w_j^{S^1} f(\Phi_{S^1}^{-1}(z_j)) = \sum_{j=1}^n \frac{1}{n} f\left(\frac{j-1}{n}\right) . \quad (7.4.0.21)$$

This is exactly the **equidistant trapezoidal rule** (7.4.0.17), if  $f$  is 1-periodic,  $f(0) = f(1)$ :  $Q_n^I = T_n$ . Hence we arrive at the following estimate for the quadrature error

$$E_n(f) := \left| \int_0^1 f(t) dt - T_n(f) \right| \leq 2\pi \max_{z \in S^1} \left| ((\Phi_{S^1}^{-1})^* f)(z) - (\mathcal{L}_{\mathcal{Z}}(\Phi_{S^1}^{-1})^* f)(z) \right| .$$

Equivalently, one can show that  $T_n$  integrates trigonometric polynomials up to degree  $2n - 1$  exactly:

$$f(t) = e^{2\pi i kt} \quad \Rightarrow \quad \begin{cases} \int_0^1 f(t) dt = \begin{cases} 0 & , \text{if } k \neq 0 , \\ 1 & , \text{if } k = 0 . \end{cases} \\ T_n(f) = \frac{1}{n} \sum_{l=0}^{m-1} e^{\frac{2\pi i l k}{n}} \stackrel{(4.2.0.8)}{=} \begin{cases} 0 & , \text{if } k \notin n\mathbb{Z} , \\ 1 & , \text{if } k \in n\mathbb{Z} . \end{cases} \end{cases}$$

### Lemma 7.4.0.22. Exact quadrature by equidistant trapezoidal rule

The  $n$ -point equidistant trapezoidal quadrature rule

$$\int_a^b f(t) dt \approx T_n(f) := h \left( \frac{1}{2}f(a) + \sum_{k=1}^{n-1} f(kh) + \frac{1}{2}f(b) \right) , \quad h := \frac{b-a}{n} . \quad (7.4.0.17)$$

is exact for trigonometric polynomials ( $\rightarrow$  Section 6.4.1) of degree  $< 2n$ .

- Since the weights of the equidistant trapezoidal rule are clearly positive, by Thm. 7.3.0.40 the asymptotic behavior of the quadrature error can directly be inferred from estimates for equidistant trigonometric interpolation. Such estimates are given, e.g., in Thm. 6.4.3.12 and they confirm exponential convergence for periodic integrands that allow analytic extension, in agreement with the observations made in Exp. 7.4.0.16.

### Numerical Quadrature of *periodic integrands*

Use the equidistant trapezoidal rule for the numerical quadrature of a periodic integrand over its interval of periodicity.

#### Remark 7.4.0.24 (Approximate computation of Fourier coefficients)

Recall from Section 4.2.5: recovery of signal  $(y_k)_{k \in \mathbb{Z}}$  from its Fourier transform  $c(t)$

$$y_j = \int_0^1 c(t) \exp(2\pi i jt) dt . \quad (4.2.5.18)$$

Task: approximate computation of  $y_j$

Recall:  $c(t)$  obtained from  $(y_k)_{k \in \mathbb{Z}}$  through Fourier series

$$c(t) = \sum_{k \in \mathbb{Z}} y_k \exp(-2\pi i kt) . \quad (4.2.5.6)$$

►  $c(t)$  smooth & 1-periodic for finite/rapidly decaying  $(y_k)_{k \in \mathbb{Z}}$ .

Exp. 7.4.0.16 ►

use **equidistant trapezoidal rule** (7.4.0.17)  
for approximate evaluation of integral in (4.2.5.18).

☞ Boils down to inverse DFT (4.2.0.20); hardly surprising in light of the derivation of (4.2.5.18) in Section 4.2.5.

#### C++-code 7.4.0.25: DFT-based approximate computation of Fourier coefficients

```

1 # include <iostream>
2 # include <vector>
3 # include <complex>
4 # include <unsupported/Eigen/FFT>
5
6 template <class Function>
7 void fourcoeffcomp(std::vector<std::complex<double>>& y, Function& c, const unsigned
8   m, const unsigned ovsmpl = 2) {
9   // Compute the Fourier coefficients  $y_{-m}, \dots, y_m$  of the function
10  //  $c : [0, 1] \rightarrow \mathbb{C}$  using an oversampling factor ovsmpl.
11  //  $c$  must be a handle to a function  $\theta(t)$ , e.g. a lambda function
12  const unsigned N = (2*m + 1)*ovsmpl; // number of quadrature points
13  const double h = 1./N;
14
15  // evaluate function in N points
16  std::vector<std::complex<double>> c_eval(N);
17  for (unsigned i = 0; i < N; ++i) {
18    c_eval[i] = c(i*h);
19  }
20
21  // inverse discrete fourier transformation

```

```

21 Eigen::FFT<double> fft;
22 std::vector<std::complex<double>> z;
23 fft.inv(z, c_eval);
24
25 // Undo oversampling and wrapping of Fourier coefficient array
26 // -> y contains same values as z but in a different order:
27 // y = [z(N-m+1:N), z(1:m+1)]
28 y = std::vector<std::complex<double>>();
29 y.reserve(N);
30 for (unsigned i = N - m; i < N; ++i) {
31     y.push_back(z[i]);
32 }
33 for (unsigned j = 0; j < m + 1; ++j) {
34     y.push_back(z[j]);
35 }
36

```

## 7.5 Adaptive Quadrature

Hitherto, we have just “blindly” applied quadrature rules for the approximate evaluation of  $\int_a^b f(t) dt$ , oblivious of any properties of the integrand  $f$ . This led us to the conclusion of Rem. 7.4.0.13 that Gauss-Legendre quadrature (→ Def. 7.3.0.30) should be preferred to composite quadrature rules (→ Section 7.4) in general. Now the composite quadrature rule will partly be rehabilitated, because they offer the flexibility to *adjust the quadrature rule to the integrand*, a policy known as adaptive quadrature.

### Adaptive numerical quadrature

The policy of **adaptive quadrature** approximates  $\int_a^b f(t) dt$  by a quadrature formula (7.1.0.2), whose nodes  $c_j^n$  are chosen depending on the integrand  $f$ .

We distinguish

- (I) **a priori** adaptive quadrature: the nodes are *fixed* before the evaluation of the quadrature formula, taking into account external information about  $f$ , and
- (II) **a posteriori** adaptive quadrature: the node positions are chosen or improved based on information gleaned *during the computation* inside a loop. It terminates when sufficient accuracy has been reached.

In this section we will chiefly discuss a posteriori adaptive quadrature for composite quadrature rules (→ Section 7.4) based on a single local quadrature rule (and its transformation).



*Supplementary literature.* [DH03, Sect. 9.7]

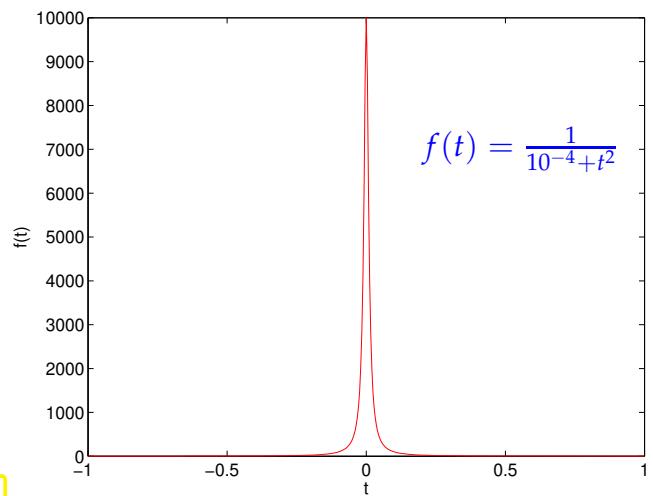
**EXAMPLE 7.5.0.2 (Rationale for adaptive quadrature)** This example presents an extreme case. We consider the composite trapezoidal rule (7.4.0.4) on a mesh  $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$  and for the integrand  $f(t) = \frac{1}{10^{-4} + t^2}$  on  $[-1, 1]$ .

$f$  is a spike-like function

▷

Intuition: quadrature nodes should cluster around 0, whereas hardly any are needed close to the endpoints of the integration interval, where the function has very small (in modulus) values.

➤ Use locally refined mesh !



A quantitative justification can appeal to (7.2.0.10) and the resulting bound for the local quadrature error (for  $f \in C^2([a, b])$ ):

$$\int_{x_{k-1}}^{x_k} f(t) dt - \frac{1}{2}(f(x_{k-1}) + f(x_k)) \leq \frac{1}{8} h_k^3 \|f''\|_{L^\infty([x_{k-1}, x_k])}, \quad h_k := x_k - x_{k-1}. \quad (7.5.0.3)$$

➤ Suggests the use of small mesh intervals, where  $|f''|$  is large !

**§7.5.0.4 (Goal: equidistribution of errors)** The ultimate but elusive goal is to find a mesh with a minimal number of cells that just delivers a quadrature error below a prescribed threshold. A more practical goal is to adjust the local meshwidths  $h_k := x_k - x_{k-1}$  in order to achieve a minimal sum of local error bounds. This leads to the constrained minimization problem:

$$\sum_{k=1}^m h_k^3 \|f''\|_{L^\infty([x_{k-1}, x_k])} \rightarrow \min \quad \text{s.t.} \quad \sum_{k=1}^m h_k = b - a. \quad (7.5.0.5)$$

#### Lemma 7.5.0.6.

Let  $f : \mathbb{R}_0^+ \rightarrow \mathbb{R}_0^+$  be a convex function with  $f(0) = 0$  and  $x > 0$ . Then the constrained minimization problem: seek  $\zeta_1, \dots, \zeta_m \in \mathbb{R}_0^+$  such that

$$\sum_{k=1}^m f(\zeta_k) \rightarrow \min \quad \text{and} \quad \sum_{k=1}^m \zeta_k = x, \quad (7.5.0.7)$$

has the solution  $\zeta_1 = \zeta_2 = \dots = \zeta_m = \frac{x}{m}$ .

This means that we should strive for equal bounds  $h_k^3 \|f''\|_{L^\infty([x_{k-1}, x_k])}$  for all mesh cells.

#### Error equidistribution principle

The mesh for a posteriori adaptive composite numerical quadrature should be chosen to achieve equal contributions of all mesh intervals to the quadrature error

As indicated above, guided by the equidistribution principle, the improvement of the mesh will be done gradually in an iteration. The change of the mesh in each step is called **mesh adaptation** and there are two fundamentally different ways to do it:

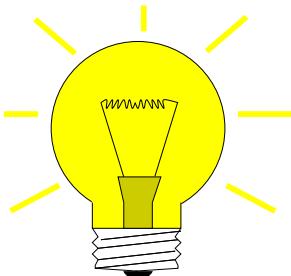
- (I) by **moving** nodes, keeping their total number, but making them cluster where mesh intervals should be small, or
- (II) by **adding** nodes, where mesh intervals should be small (**mesh refinement**).

Algorithms for a posteriori adaptive quadrature based on mesh refinement usually have the following structure:

#### Adaptation loop for numerical quadrature

- (1) **ESTIMATE**: based on available information compute an approximation for the quadrature error on every mesh interval.
- (2) **CHECK TERMINATION**: if total error sufficient small  $\rightarrow$  **STOP**
- (3) **MARK**: single out mesh intervals with the largest or above average error contributions.
- (4) **REFINE**: add node(s) inside the marked mesh intervals. GOTO (1)

**§7.5.0.10 (Adaptive multilevel quadrature)** We now see a concrete algorithm based on the two composite quadrature rules introduced in Ex. 7.4.0.3.



Idea: local error estimation by comparing local results of two quadrature formulas  $Q_1, Q_2$  of *different* order  $\rightarrow$  local error estimates

heuristics:  $\text{error}(Q_2) \ll \text{error}(Q_1) \Rightarrow \text{error}(Q_1) \approx Q_2(f) - Q_1(f)$ .

Here:  $Q_1 = \text{trapezoidal rule (order 2)} \leftrightarrow Q_2 = \text{Simpson rule (order 4)}$

Given: initial mesh  $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$

① (Error estimation)

For  $I_k = [x_{k-1}, x_k], k = 1, \dots, m$  (midpoints  $p_k := \frac{1}{2}(x_{k-1} + x_k)$ )

$$\text{EST}_k := \left| \underbrace{\frac{h_k}{6}(f(x_{k-1}) + 4f(p_k) + f(x_k))}_{\text{Simpson rule}} - \underbrace{\frac{h_k}{4}(f(x_{k-1}) + 2f(p_k) + f(x_k))}_{\text{trapezoidal rule on split mesh interval}} \right|. \quad (7.5.0.11)$$

② (Check termination)

Simpson rule on  $\mathcal{M} \Rightarrow$  intermediate approximation  $I \approx \int_a^b f(t) dt$

$$\text{If } \sum_{k=1}^m \text{EST}_k \leq \text{RTOL} \cdot I \quad (\text{RTOL} := \text{prescribed relative tolerance}) \Rightarrow \text{STOP} \quad (7.5.0.12)$$

**③** (Marking)

Marked intervals:  $\mathcal{S} := \{k \in \{1, \dots, m\} : \text{EST}_k \geq \eta \cdot \frac{1}{m} \sum_{j=1}^m \text{EST}_j\}, \quad \eta \approx 0.9.$  (7.5.0.13)

**④** (Local mesh refinement)

new mesh:  $\mathcal{M}^* := \mathcal{M} \cup \{p_k := \frac{1}{2}(x_{k-1} + x_k) : k \in \mathcal{S}\}.$  (7.5.0.14)

Then continue with step ① and mesh  $\mathcal{M} \leftarrow \mathcal{M}^*.$

The following C++ code give a (non-optimal) recursive implementation

### C++-code 7.5.0.15: $h$ -adaptive numerical quadrature

```

2 // Adaptive multilevel quadrature of a function passed in f.
3 // The vector M passes the positions of current quadrature nodes
4 template <class Function>
5 double adaptquad(Function& f, VectorXd& M, double rtol, double atol) {
6     const std::size_t n = M.size(); // number of nodes
7     VectorXd h = M.tail(n-1)-M.head(n-1), // distance of quadrature nodes
8         mp = 0.5*(M.head(n-1)+M.tail(n-1)); // midpoints
9     // Values of integrand at nodes and midpoints
10    VectorXd fx(n), fm(n-1);
11    for (unsigned i = 0; i < n; ++i) fx(i) = f(M(i));
12    for (unsigned j = 0; j < n - 1; ++j) fm(j) = f(mp(j));
13    // trapezoidal rule (7.4.0.4)
14    const VectorXd trp_loc = 1./4*h.cwiseProduct(fx.head(n-1)+2*fm+fx.tail(n-1));
15    // Simpson rule (7.4.0.5)
16    const VectorXd simp_loc = 1./6*h.cwiseProduct(fx.head(n-1)+4*fm+fx.tail(n-1));
17
18    // Simpson approximation for the integral value
19    double I = simp_loc.sum();
20    // local error estimate (7.5.0.11)
21    const VectorXd est_loc = (simp_loc-trp_loc).cwiseAbs();
22    // estimate for quadrature error
23    const double err_tot = est_loc.sum();
24
25    // STOP: Termination based on (7.5.0.12)
26    if (err_tot > rtol*std::abs(I) && err_tot > atol) { //
27        // find cells where error is large
28        std::vector<double> new_cells;
29        for (unsigned i = 0; i < est_loc.size(); ++i) {
30            // MARK by criterion (7.5.0.13) & REFINER by (7.5.0.14)
31            if (est_loc(i) > 0.9/(n-1)*err_tot) {
32                // new quadrature point = midpoint of interval with large error
33                new_cells.push_back(mp(i));
34            }
35
36            // create new set of quadrature nodes
37            // (necessary to convert std::vector to Eigen vector)
38            Eigen::Map<VectorXd> tmp(new_cells.data(), new_cells.size());
39            VectorXd new_M(M.size() + tmp.size());
40            new_M << M, tmp; // concatenate old cells and new cells
41            // nodes of a mesh are supposed to be sorted
42            std::sort(new_M.data(), new_M.data() + new_M.size());
43            I = adaptquad(f, new_M, rtol, atol); // recursion
44        }
45        return I;
46    }
}

```

Comments on Code 7.5.0.15:

- Arguments:  $f \hat{=} \text{handle}$  to function  $f$ ,  $M \hat{=} \text{initial mesh}$ ,  $\text{rtol} \hat{=} \text{relative tolerance for termination}$ ,  $\text{atol} \hat{=} \text{absolute tolerance for termination}$ , necessary in case the exact integral value  $= 0$ , which renders a relative tolerance meaningless.
- Line 7: compute lengths of mesh-intervals  $[x_{j-1}, x_j]$ ,
- Line 8: store positions of midpoints  $p_j$ ,
- Line 9: evaluate function (vector arguments!),
- Line 13: local composite trapezoidal rule (7.4.0.4),

- Line 15: local simpson rule (7.2.0.6),
- Line 18: value obtained from composite simpson rule is used as intermediate approximation for integral value,
- Line 20: difference of values obtained from local composite trapezoidal rule ( $\sim Q_1$ ) and local simpson rule ( $\sim Q_2$ ) is used as an estimate for the local quadrature error.
- Line 22: estimate for global error by summing up **moduli** of local error contributions,
- Line 26: terminate, once the estimated total error is below the relative or absolute error threshold,
- Line 43 otherwise, add midpoints of mesh intervals with large error contributions according to (7.5.0.14) to the mesh and continue.

**C++-code 7.5.0.16: Call of adaptquad():**

```

1 # include "./adaptquad.hpp"
2 # include <iostream>
3 # include <cmath>
4
5 int main () {
6     auto f =[](double x) { return std::exp(-x*x); };
7     VectorXd M(4);
8     M << -100, 0.1, 0.5, 100;
9     std::cout << "Sqrt(Pi) - Int_{-100}^{100} exp(-x*x) dx = ";
10    std::cout << adaptquad(f, M, 1e-10, 1e-12) - std::sqrt(M_PI) << "\n";
11    return 0;
12 }
```

**Remark 7.5.0.17 (Estimation of “wrong quadrature error”?)** In Code 7.5.0.15 we use the higher order quadrature rule, the Simpson rule of order 4, to compute an approximate value for the integral. This is reasonable, because it would be foolish not to use this information after we have collected it for the sake of error estimation.

Yet, according to our heuristics, what `est_loc` and `est_tot` give us are estimates for the error of the second-order trapezoidal rule, which we do not use for the actual computations.

However, experience teaches that

`est_loc` gives useful (for the sake of mesh refinement) information about the *distribution* of the error of the Simpson rule, though it fails to capture its *size*.

Therefore, the termination criterion of Line 26 may not be appropriate!

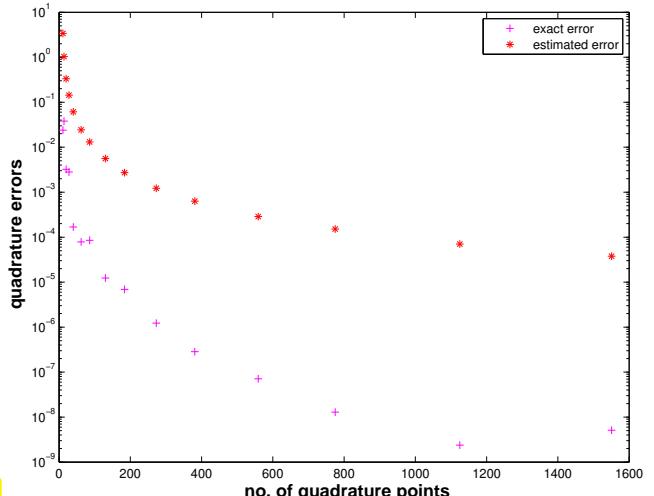
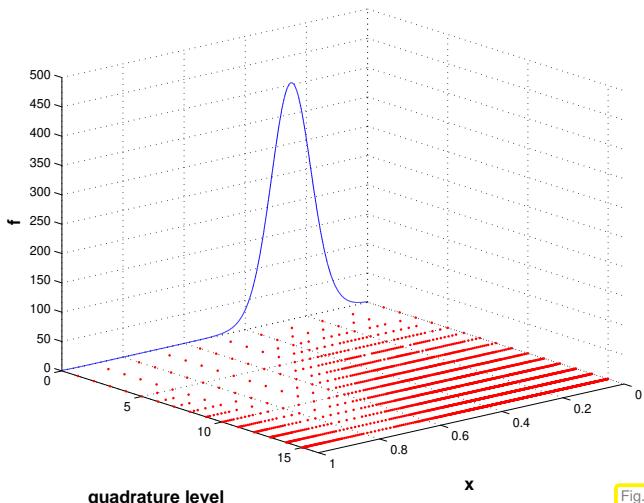
**EXPERIMENT 7.5.0.18 ( $h$ -adaptive numerical quadrature)** In this numerical test we investigate whether the adaptive technique from § 7.5.0.10 produces an appropriate distribution of integration nodes. We do this for different functions.

◆ approximate  $\int_0^1 \exp(6 \sin(2\pi t)) dt$ , initial mesh  $\mathcal{M}_0 = \{j/10\}_{j=0}^{10}$

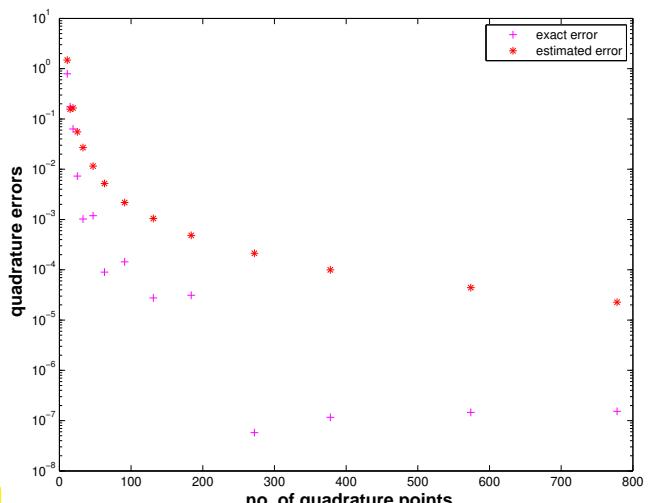
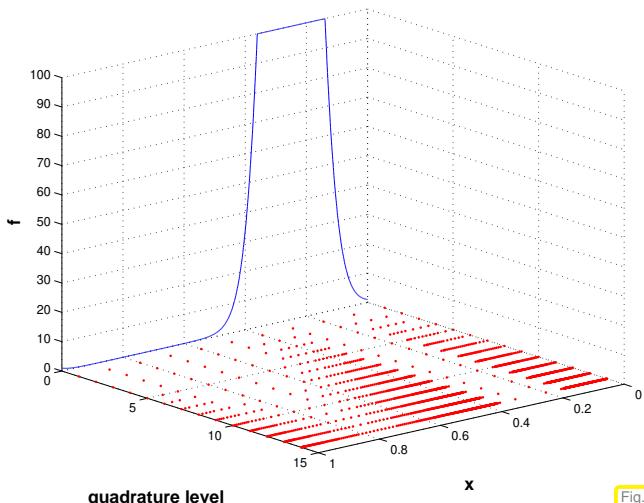
Algorithm: adaptive quadrature, Code 7.5.0.15 with tolerances  $\text{rtol} = 10^{-6}$ ,  $\text{abstol} = 10^{-10}$

We monitor the distribution of quadrature points during the adaptive quadrature and the true and estimated quadrature errors. The “exact” value for the integral is computed by composite Simpson rule on an

equidistant mesh with  $10^7$  intervals.



- ◆ approximate  $\int_0^1 \min\{\exp(6 \sin(2\pi t)), 100\} dt$ , initial mesh as above



Observation:

- Adaptive quadrature locally decreases meshwidth where integrand features variations or kinks.
- Trend for estimated error mirrors behavior of true error.
- Overestimation may be due to taking the modulus in (7.5.0.11)

However, the important piece of information we want to extract from  $\text{EST}_k$  is about the *distribution* of the quadrature error.

### Remark 7.5.0.19 (Adaptive quadrature in PYTHON)

`q = scipy.integrate.quad(fun,a,b,tol):`

adaptive multigrid quadrature  
(local low order quadrature formulas)

`q = scipy.integrate.quadrature(fun,a,b,tol):`

adaptive Gauss-Lobatto quadrature

## Learning Outcomes

- ◆ You should know what is a quadrature formula and terminology connected with it,
- ◆ You should be able to transform quadrature formulas to arbitrary intervals.
- ◆ You should understand how interpolation and approximation schemes spawn quadrature formulas and how quadrature errors are connected to interpolation/approximation errors.
- ◆ You should be able to compute the weights of polynomial quadrature formulas.
- ◆ You should know the concept of order of a quadrature rule and why it is invariant under (affine) transformation
- ◆ You should remember the maximal and minimal order of polynomial quadrature rules.
- ◆ You should know the order of the  $n$ -point Gauss-Legendre quadrature rule.
- ◆ You should understand why Gauss-Legendre quadrature converges exponentially for integrands that can be extended analytically and algebraically for integrands with limited smoothness.
- ◆ You should be able to apply regularizing transformations to integrals with non-smooth integrands.
- ◆ You should know about asymptotic convergence of the  $h$ -version of composite quadrature.
- ◆ You should know the principles of adaptive composite quadrature.

# Bibliography

- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 449, 454, 456).
- [DH03] P. Deuflhard and A. Hohmann. *Numerical Analysis in Modern Scientific Computing*. Vol. 43. Texts in Applied Mathematics. Springer, 2003 (cit. on p. 478).
- [Gan+05] M. Gander, W. Gander, G. Golub, and D. Gruntz. *Scientific Computing: An introduction using MATLAB*. Springer, 2005 (cit. on pp. 466, 467).
- [GLR07] Andreas Glaser, Xiangtao Liu, and Vladimir Rokhlin. “A fast algorithm for the calculation of the roots of special functions”. In: *SIAM J. Sci. Comput.* 29.4 (2007), pp. 1420–1438. DOI: [10.1137/06067016X](https://doi.org/10.1137/06067016X).
- [Gut09] M.H. Gutknecht. *Lineare Algebra*. Lecture Notes. SAM, ETH Zürich, 2009 (cit. on p. 461).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 449, 454, 456, 465).
- [Joh08] S.G. Johnson. *Notes on the convergence of trapezoidal-rule quadrature*. MIT online course notes, <http://math.mit.edu/~stevenj/trapezoidal.pdf>. 2008.
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on p. 461).
- [Tre08] Lloyd N. Trefethen. “Is Gauss quadrature better than Clenshaw-Curtis?” In: *SIAM Rev.* 50.1 (2008), pp. 67–87. DOI: [10.1137/060659831](https://doi.org/10.1137/060659831) (cit. on pp. 456, 466).
- [TWXX] Lloyd N. Trefethen and J. A. C. Weideman. *THE EXPONENTIALLY CONVERGENT TRAPEZOIDAL RULE*. XX.

# Chapter 8

## Iterative Methods for Non-Linear Systems of Equations

### EXAMPLE 8.0.0.1 (Non-linear electric circuit)

Non-linear systems naturally arise in mathematical models of electrical circuits, once non-linear circuit elements are introduced. This generalizes Ex. 2.1.0.3, where the current-voltage relationship for all circuit elements was the simple proportionality (2.1.0.5) (of the complex amplitudes  $\mathbf{U}$  and  $\mathbf{I}$ ).

As an example we consider the

Schmitt trigger circuit



Its key non-linear circuit element is the NPN bipolar junction transistor:

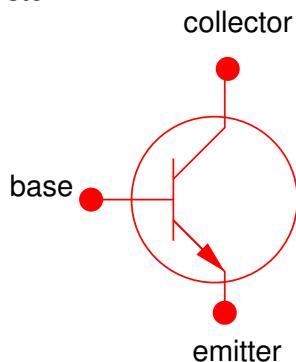
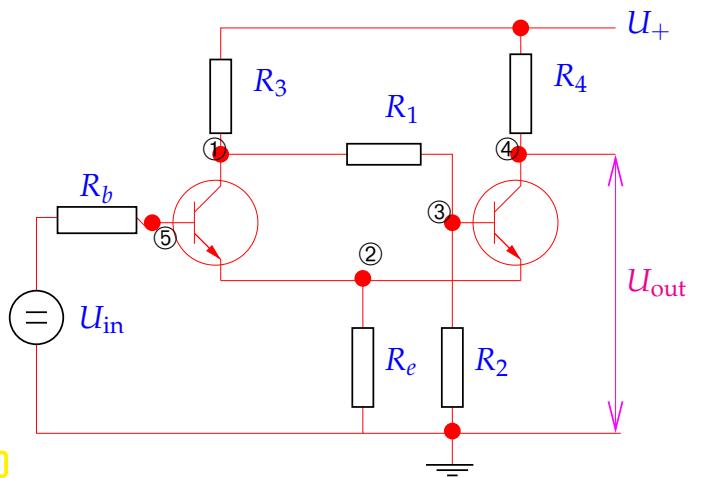


Fig. 257



A transistor has three ports: emitter, collector, and base. Transistor models give the port currents as functions of the applied voltages, for instance the Ebers-Moll model (large signal approximation):

$$\begin{aligned} I_C &= I_S \left( e^{\frac{U_{BE}}{U_T}} - e^{\frac{U_{BC}}{U_T}} \right) - \frac{I_S}{\beta_R} \left( e^{\frac{U_{BC}}{U_T}} - 1 \right) = I_C(U_{BE}, U_{BC}), \\ I_B &= \frac{I_S}{\beta_F} \left( e^{\frac{U_{BE}}{U_T}} - 1 \right) + \frac{I_S}{\beta_R} \left( e^{\frac{U_{BC}}{U_T}} - 1 \right) = I_B(U_{BE}, U_{BC}), \\ I_E &= I_S \left( e^{\frac{U_{BE}}{U_T}} - e^{\frac{U_{BC}}{U_T}} \right) + \frac{I_S}{\beta_F} \left( e^{\frac{U_{BE}}{U_T}} - 1 \right) = I_E(U_{BE}, U_{BC}). \end{aligned} \quad (8.0.0.2)$$

$I_C$ ,  $I_B$ ,  $I_E$ : current in collector/base/emitter,

$U_{BE}$ ,  $U_{BC}$ : potential drop between base-emitter, base-collector.

The parameters have the following meanings:  $\beta_F$  is the forward common emitter current gain (20 to 500),  $\beta_R$  is the reverse common emitter current gain (0 to 20),  $I_S$  is the reverse saturation current (on the order of  $10^{-15}$  to  $10^{-12}$  amperes),  $U_T$  is the thermal voltage (approximately 26 mV at 300 K).

The circuit of Fig. 257 has 5 nodes ①–⑤ with unknown nodal potentials. Kirchhoffs law (2.1.0.4) plus the constitutive relations gives an equation for each of them.

► **Non-linear system of equations** from nodal analysis, static case ( $\rightarrow$  Ex. 2.1.0.3):

$$\begin{aligned} \textcircled{1} : \quad & R_3(U_1 - U_+) + R_1(U_1 - U_3) + I_B(U_5 - U_1, U_5 - U_2) = 0, \\ \textcircled{2} : \quad & R_e U_2 + I_E(U_5 - U_1, U_5 - U_2) + I_E(U_3 - U_4, U_3 - U_2) = 0, \\ \textcircled{3} : \quad & R_1(U_3 - U_1) + I_B(U_3 - U_4, U_3 - U_2) = 0, \\ \textcircled{4} : \quad & R_4(U_4 - U_+) + I_C(U_3 - U_4, U_3 - U_2) = 0, \\ \textcircled{5} : \quad & R_b(U_5 - U_{\text{in}}) + I_B(U_5 - U_1, U_5 - U_2) = 0. \end{aligned} \quad (8.0.0.3)$$

5 equations  $\leftrightarrow$  5 unknowns  $U_1, U_2, U_3, U_4, U_5$

Formally:  $(8.0.0.3) \longleftrightarrow F(\mathbf{u}) = 0$  with a function  $F : \mathbb{R}^5 \rightarrow \mathbb{R}^5$

**Remark 8.0.0.4 (General non-linear systems of equations)** A **non-linear system of equations** is a concept almost *too abstract to be useful*, because it covers an extremely wide variety of problems. Nevertheless in this chapter we will mainly look at “generic” methods for such systems. This means that every method discussed may take a good deal of fine-tuning before it will really perform satisfactorily for a given non-linear system of equations.

**§8.0.0.5 (Generic/general non-linear system of equations)** Let us try to describe the “problem” of having to solve a non-linear system of equations, where the concept of a “problem” was first introduced in § 1.5.5.1.

Given: function  $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ ,  $n \in \mathbb{N}$



Possible meanings:  $\square$   $F$  is known as an analytic expression.

$\square$   $F$  is merely available in **procedural form** allowing point evaluations.

Here,  $D$  is the domain of definition of the function  $F$ , which cannot be evaluated for  $\mathbf{x} \notin D$ .

Sought: solution(s)  $\mathbf{x} \in D$  of **non-linear equation**  $F(\mathbf{x}) = 0$

Note:  $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n \leftrightarrow$  “same number of equations and unknowns”

In contrast to the situation for linear systems of equations ( $\rightarrow$  Thm. 2.2.1.4), the class of non-linear systems is far too big to allow a general theory:

There are no general results existence & uniqueness of solutions of  $F(\mathbf{x}) = 0$

## Contents

<b>8.1 Iterative Methods . . . . .</b>	<b>489</b>
8.1.1 Speed of Convergence . . . . .	491
8.1.2 Termination Criteria/Stopping Rules . . . . .	495
<b>8.2 Fixed-Point Iterations . . . . .</b>	<b>498</b>
8.2.1 Consistent Fixed-Point Iterations . . . . .	499
8.2.2 Convergence of Fixed-Point Iterations . . . . .	501
<b>8.3 Finding Zeros of Scalar Functions . . . . .</b>	<b>507</b>

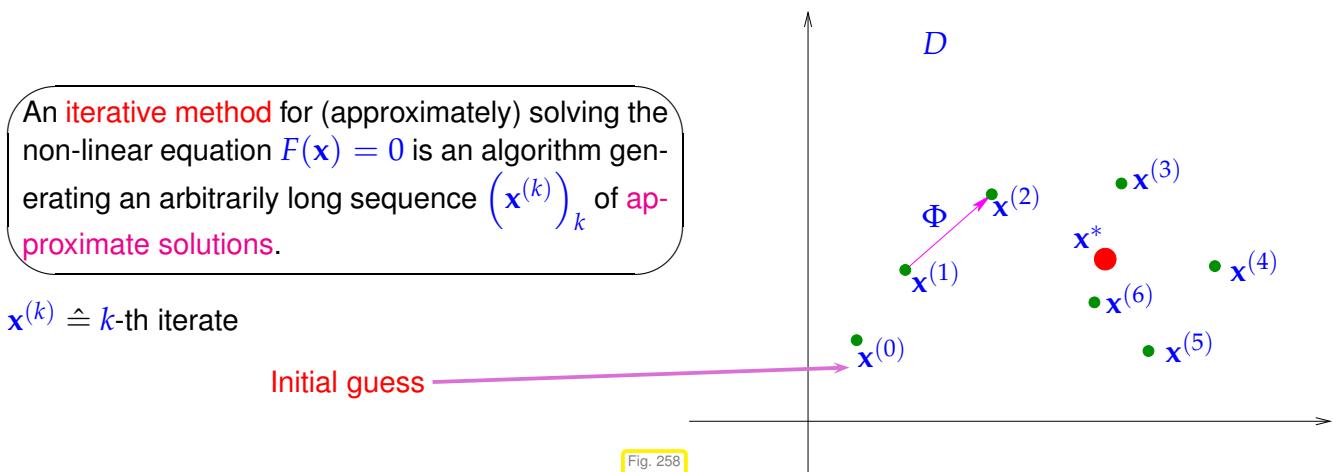
8.3.1	Bisection . . . . .	507
8.3.2	Model Function Methods . . . . .	509
8.3.3	Asymptotic Efficiency of Iterative Methods for Zero Finding . . . . .	520
<b>8.4</b>	<b>Newton's Method in <math>\mathbb{R}^n</math></b> . . . . .	<b>522</b>
8.4.1	The Newton Iteration . . . . .	522
8.4.2	Convergence of Newton's Method . . . . .	532
8.4.3	Termination of Newton Iteration . . . . .	534
8.4.4	Damped Newton Method . . . . .	536
8.4.5	Quasi-Newton Method . . . . .	539
<b>8.5</b>	<b>Non-linear Least Squares [DR08, Ch. 6]</b> . . . . .	<b>544</b>
8.5.1	(Damped) Newton Method . . . . .	546
8.5.2	Gauss-Newton Method . . . . .	547
8.5.3	Trust Region Method (Levenberg-Marquardt Method) . . . . .	550

## 8.1 Iterative Methods

**Remark 8.1.0.1 (Necessity of iterative approximation)** Gaussian elimination ( $\rightarrow$  Section 2.3) provides an algorithm that, if carried out in exact arithmetic (no roundoff errors), computes the solution of a linear system of equations with a *finite* number of elementary operations. However, linear systems of equations represent an exceptional case, because it is hardly ever possible to solve general systems of non-linear equations using only finitely many elementary operations.  $\square$

### §8.1.0.2 (Generic iterations)

All methods for general non-lineare systems of equations are *iterative* in the sense that they will usually yield only *approximate solutions* whenever they terminate after finite time.



**§8.1.0.3 ((Stationary)  $m$ -point iterative method)** All the iterative methods discussed below fall in the class of (stationary)  $m$ -point,  $m \in \mathbb{N}$ , iterative methods, for which the iterate  $\mathbf{x}^{(k)}$  depends on  $F$  and the  $m$  most recent iterates  $\mathbf{x}^{(k-1)}, \dots, \mathbf{x}^{(k-m)}$ , e.g.,

$$\mathbf{x}^{(k)} = \underbrace{\Phi_F(\mathbf{x}^{(k-1)}, \dots, \mathbf{x}^{(k-m)})}_{\text{iteration function for } m\text{-point method}} \quad (8.1.0.4)$$

Terminology:  $\Phi_F$  is called the *iteration function*.

Note: The *initial guess(es)*  $\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(m-1)} \in \mathbb{R}^n$  have to be provided.  $\square$

**§8.1.0.5 (Key issues with iterative methods)** When applying an iterative method to solve a non-linear system of equations  $F(\mathbf{x}) = \mathbf{0}$ , the following issues arise:

- ◆ **Convergence**: Does the sequence  $(\mathbf{x}^{(k)})_k$  converge to a limit:  $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$ ?
- ◆ **Consistency**: Does the limit, if it exists, provide a solution of the non-linear system of equations:  $F(\mathbf{x}^*) = \mathbf{0}$ ?
- ◆ **Speed** of convergence: How “fast” does  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|$  ( $\|\cdot\|$  a suitable norm on  $\mathbb{R}^N$ ) decrease for increasing  $k$ ?

More formal definitions can be given:

#### Definition 8.1.0.6. Convergence of iterative methods

An iterative method **converges** (for fixed initial guess(es)) : $\Leftrightarrow$   $\mathbf{x}^{(k)} \xrightarrow{k \rightarrow \infty} \mathbf{x}^*$  and  $F(\mathbf{x}^*) = \mathbf{0}$ .

#### Definition 8.1.0.7. Consistency of iterative methods

A stationary  $m$ -point iterative method is **consistent** with the non-linear system of equations  $F(\mathbf{x}) = \mathbf{0}$ , if and only if

$$\Phi_F(\mathbf{x}^*, \dots, \mathbf{x}^*) = \mathbf{x}^* \iff F(\mathbf{x}^*) = \mathbf{0}.$$

For a consistent stationary iterative method we can study the **error** of the iterates  $\mathbf{x}^{(k)}$  defined as:

$$\mathbf{e}^{(k)} := \mathbf{x}^{(k)} - \mathbf{x}^*.$$

Unfortunately, convergence may critically depend on the choice of initial guesses. The property defined next weakens this dependence:

#### Definition 8.1.0.8. Local and global convergence → [Han02, Def. 17.1]

As stationary  $m$ -point iterative method **converges locally** to  $\mathbf{x}^* \in \mathbb{R}^n$ , if there is a neighborhood  $U \subset D$  of  $\mathbf{x}^*$ , such that

$$\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(m-1)} \in U \Rightarrow \mathbf{x}^{(k)} \text{ well defined} \wedge \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$$

where  $(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$  is the (infinite) sequence of iterates.

If  $U = D$ , the iterative method is **globally convergent**.

Illustration of local convergence

(Only initial guesses “sufficiently close” to  $\mathbf{x}^*$  guarantee convergence.)

Unfortunately, the neighborhood  $U$  is rarely known a priori. It may also be very small.

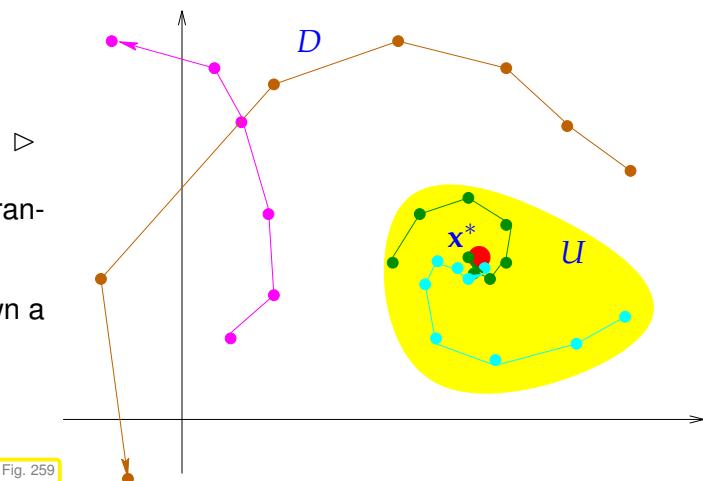


Fig. 259

Our goal: Given a non-linear system of equations, find iterative methods that converge (locally) to a solution of  $F(\mathbf{x}) = 0$ .

Two general questions: How to measure, describe, and predict the speed of convergence?  
When to terminate the iteration?

### 8.1.1 Speed of Convergence

Here and in the sequel,  $\|\cdot\|$  designates a generic vector norm on  $\mathbb{R}^n$ , see Def. 1.5.5.4. Any occurring matrix norm is induced by this vector norm, see Def. 1.5.5.10.

It is important to be aware which statements depend on the choice of norm and which do not!

“Speed of convergence”  $\leftrightarrow$  decrease of norm (see Def. 1.5.5.4) of iteration error

#### Definition 8.1.1.1. Linear convergence

A sequence  $\mathbf{x}^{(k)}$ ,  $k = 0, 1, 2, \dots$ , in  $\mathbb{R}^n$  converges linearly to  $\mathbf{x}^* \in \mathbb{R}^n$ ,

$$\exists 0 < L < 1: \quad \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq L \|\mathbf{x}^{(k)} - \mathbf{x}^*\| \quad \forall k \in \mathbb{N}_0 .$$

Terminology: The least upper bound for  $L$  gives the rate of convergence:

$$\text{rate} = \sup_{k \in \mathbb{N}_0} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|}, \quad \mathbf{x}^* := \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} . \quad (8.1.1.2)$$

#### Remark 8.1.1.3 (Impact of choice of norm)

Fact of convergence of iteration is	independent	of choice of norm
Fact of linear convergence	depends	on choice of norm
Rate of linear convergence	depends	on choice of norm

The first statement is a consequence of the equivalence of all norms on the finite dimensional vector space  $\mathbb{K}^n$ :

#### Definition 8.1.1.4. Equivalence of norms

Two norms  $\|\cdot\|_a$  and  $\|\cdot\|_b$  on a vector space  $V$  are equivalent if

$$\exists \underline{C}, \bar{C} > 0: \underline{C}\|v\|_a \leq \|v\|_b \leq \bar{C}\|v\|_a \quad \forall v \in V.$$

#### Theorem 8.1.1.5. Equivalence of all norms on finite dimensional vector spaces

If  $\dim V < \infty$  all norms ( $\rightarrow$  Def. 1.5.5.4) on  $V$  are equivalent ( $\rightarrow$  Def. 8.1.1.4).

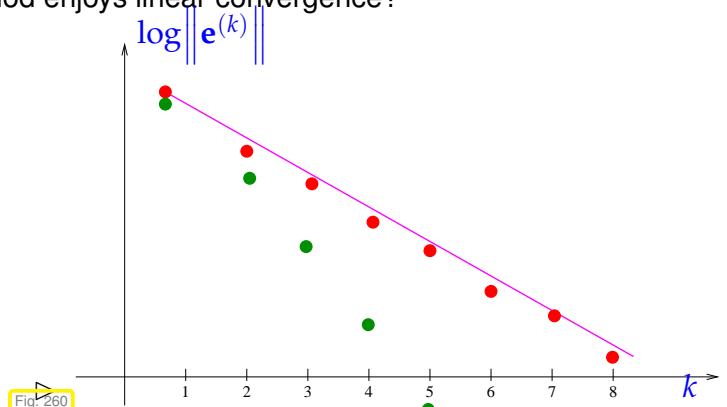
**Remark 8.1.1.6 (Detecting linear convergence)** Often we will study the behavior of a consistent iterative method for a model problem in a numerical experiments and measure the norms of the iteration errors  $\mathbf{e}^{(k)} := \mathbf{x}^{(k)} - \mathbf{x}^*$ . How can we tell that the method enjoys linear convergence?

norms of iteration errors  
↑  
~ on straight line in lin-log plot

$$\|\mathbf{e}^{(k)}\| \leq L^k \|\mathbf{e}^{(0)}\|,$$

$$\log(\|\mathbf{e}^{(k)}\|) \leq k \log(L) + \log(\|\mathbf{e}^{(0)}\|).$$

- $\hat{=}$  linear cvg., ●  $\hat{=}$  faster than linear cvg.



Let us abbreviate the error norm in step  $k$  by  $\epsilon_k := \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$ . In the case of linear convergence (see Def. 8.1.1.1) assume (with  $0 < L < 1$ )

$$\epsilon_{k+1} \approx L\epsilon_k \Rightarrow \log \epsilon_{k+1} \approx \log L + \log \epsilon_k \Rightarrow \log \epsilon_k \approx k \log L + \log \epsilon_0. \quad (8.1.1.7)$$

We conclude that  $\log L < 0$  determines the slope of the graph in lin-log error chart.

Related: guessing time complexity  $\mathcal{O}(n^\alpha)$  of an algorithm from measurements, see § 1.4.1.6.

Note the green dots ● in Fig. 260: Any “faster” convergence also qualifies as linear convergence in the strict sense of the definition. However, whenever this term is used, we tacitly imply, that no “faster convergence” prevails and that the estimates in (8.1.1.7) are sharp.  $\square$

**EXAMPLE 8.1.1.8 (Linearly convergent iteration)** We consider the iteration ( $n = 1$ ):

$$x^{(k+1)} = x^{(k)} + \frac{\cos x^{(k)} + 1}{\sin x^{(k)}}.$$

In the C++ code Code 8.1.1.9  $x$  has to be initialized with the different values for  $x_0$ .

Note: The final iterate  $x^{(15)}$  replaces the exact solution  $x^*$  in the computation of the rate of convergence.

#### C++ code 8.1.1.9: Simple fixed point iteration in 1D

```
2 void fpit(double x0, VectorXd &rates, VectorXd &err)
3 {
4     const unsigned int N = 15;
```

```

5   double x = x0; // initial guess
6   VectorXd y(N);
7
8   for (int i=0; i<N; ++i) {
9     x = x + (cos(x)+1)/sin(x);
10    y(i) = x;
11  }
12  err.resize(N); rates.resize(N);
13  err = y-VectorXd::Constant(N,x);
14  rates = err.bottomRows(N-1).cwiseQuotient(err.topRows(N-1));
15 }
```

$k$	$x^{(0)} = 0.4$		$x^{(0)} = 0.6$		$x^{(0)} = 1$	
	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$
2	3.3887	0.1128	3.4727	0.4791	2.9873	0.4959
3	3.2645	0.4974	3.3056	0.4953	3.0646	0.4989
4	3.2030	0.4992	3.2234	0.4988	3.1031	0.4996
5	3.1723	0.4996	3.1825	0.4995	3.1224	0.4997
6	3.1569	0.4995	3.1620	0.4994	3.1320	0.4995
7	3.1493	0.4990	3.1518	0.4990	3.1368	0.4990
8	3.1454	0.4980	3.1467	0.4980	3.1392	0.4980

Rate of convergence  $\approx 0.5$

Plot of modulus of iteration errors for different initial guesses (see above table)

Observation:

Linear convergence as in Def. 8.1.1.1



w error graphs = straight lines in lin-log scale  
 $\rightarrow$  Rem. 8.1.1.6

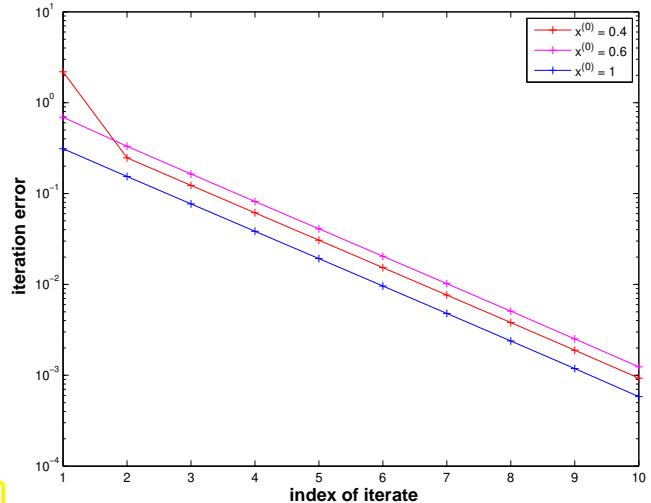


Fig. 261

There are notions of convergence that guarantee a much faster (asymptotic) decay of norm of the iteration error than linear convergence from Def. 8.1.1.1.

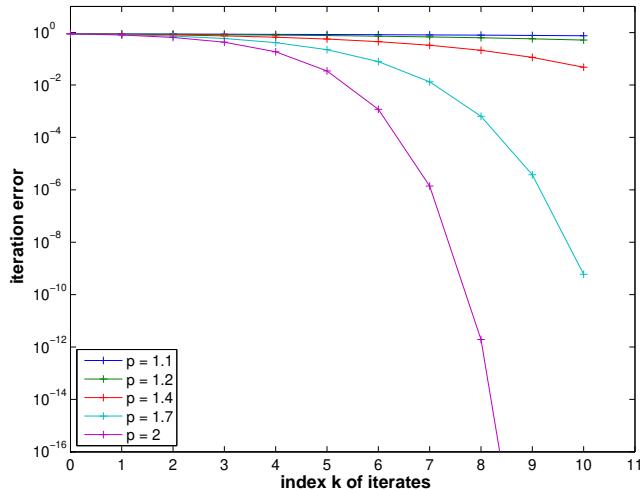
**Definition 8.1.1.10. Order of convergence** → [Han02, Sect. 17.2], [DR08, Def. 5.14], [QSS00, Def. 6.1]

A convergent sequence  $\mathbf{x}^{(k)}$ ,  $k = 0, 1, 2, \dots$ , in  $\mathbb{R}^n$  with limit  $\mathbf{x}^* \in \mathbb{R}^n$  converges with order  $p$ , if

$$\exists C > 0: \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq C \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p \quad \forall k \in \mathbb{N}_0, \quad (8.1.1.11)$$

and, in addition,  $C < 1$  in the case  $p = 1$  (linear convergence → Def. 8.1.1.1).

Of course, the order  $p$  of convergence of an iterative method refers to the largest possible  $p$  in the definition, that is, the error estimate will in general not hold, if  $p$  is replaced with  $p + \epsilon$  for any  $\epsilon > 0$ , cf. Rem. 1.4.1.3.



▷ Qualitative error graphs for convergence of order  $p$   
(lin-log scale)

In the case of convergence of order  $p$  ( $p > 1$ ) (see Def. 8.1.1.10):

$$\begin{aligned} \epsilon_{k+1} \approx C\epsilon_k^p &\Rightarrow \log \epsilon_{k+1} = \log C + p \log \epsilon_k \Rightarrow \log \epsilon_{k+1} = \log C \sum_{l=0}^k p^l + p^{k+1} \log \epsilon_0 \\ &\Rightarrow \log \epsilon_{k+1} = -\frac{\log C}{p-1} + \left( \frac{\log C}{p-1} + \log \epsilon_0 \right) p^{k+1}. \end{aligned}$$

In this case, the error graph is a concave power curve (for sufficiently small  $\epsilon_0$  !)

**Remark 8.1.1.12 (Detecting order of convergence)** How can we guess the order of convergence (→ Def. 8.1.1.10) from tabulated error norms measured in a numerical experiment?

Abbreviate by  $\epsilon_k := \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$  the norm of the iteration error.

$$\text{assume } \epsilon_{k+1} \approx C\epsilon_k^p \Rightarrow \log \epsilon_{k+1} \approx \log C + p \log \epsilon_k \Rightarrow \frac{\log \epsilon_{k+1} - \log \epsilon_k}{\log \epsilon_k - \log \epsilon_{k-1}} \approx p.$$

➤ monitor the quotients  $(\log \epsilon_{k+1} - \log \epsilon_k) / (\log \epsilon_k - \log \epsilon_{k-1})$  over several steps of the iteration. ↴

**EXAMPLE 8.1.1.13 (quadratic convergence = convergence of order 2)** From your analysis course [Str09, Bsp. 3.3.2(iii)] recall the famous iteration for computing  $\sqrt{a}$ ,  $a > 0$ :

$$x^{(k+1)} = \frac{1}{2}(x^{(k)} + \frac{a}{x^{(k)}}) \Rightarrow |x^{(k+1)} - \sqrt{a}| = \frac{1}{2x^{(k)}}|x^{(k)} - \sqrt{a}|^2. \quad (8.1.1.14)$$

By the arithmetic-geometric mean inequality (AGM)  $\sqrt{ab} \leq \frac{1}{2}(a+b)$  we conclude:  $x^{(k)} > \sqrt{a}$  for  $k \geq 1$ . Therefore estimate from (8.1.1.14) means that the sequence from (8.1.1.14) converges with order 2 to  $\sqrt{a}$ .

Note:  $x^{(k+1)} < x^{(k)}$  for all  $k \geq 2$  ➤  $(x^{(k)})_{k \in \mathbb{N}_0}$  converges as a decreasing sequence that is bounded from below (→ analysis course)

Numerical experiment: iterates for  $a = 2$ :

$k$	$x^{(k)}$	$e^{(k)} := x^{(k)} - \sqrt{2}$	$\log \frac{ e^{(k)} }{ e^{(k-1)} } : \log \frac{ e^{(k-1)} }{ e^{(k-2)} }$
0	2.0000000000000000000	0.58578643762690485	
1	1.5000000000000000000	0.08578643762690485	
2	1.41666666666666652	0.00245310429357137	1.850
3	1.41421568627450966	0.00000212390141452	1.984
4	1.41421356237468987	0.000000000000159472	2.000
5	1.41421356237309492	0.00000000000000022	0.630

Note the **doubling** of the number of correct digits in each step !

[impact of roundoff !]

The doubling of the number of significant digits for the iterates holds true for any quadratically convergent iteration:

Recall from Rem. 1.5.3.4 that the relative error ( $\rightarrow$  Def. 1.5.3.3) tells the number of significant digits. Indeed, denoting the relative error in step  $k$  by  $\delta_k$ , we have in the case of quadratic convergence.

$$\begin{aligned} x^{(k)} = x^*(1 + \delta_k) &\Rightarrow x^{(k)} - x^* = \delta_k x^*. \\ \Rightarrow |x^* \delta_{k+1}| &= |x^{(k+1)} - x^*| \leq C|x^{(k)} - x^*|^2 = C|x^* \delta_k|^2 \\ &\Rightarrow |\delta_{k+1}| \leq C|x^*| \delta_k^2. \end{aligned} \quad (8.1.1.15)$$

Note:  $\delta_k \approx 10^{-\ell}$  means that  $x^{(k)}$  has  $\ell$  significant digits.

Also note that if  $C \approx 1$ , then  $\delta_k = 10^{-\ell}$  and (8.1.1.13) implies  $\delta_{k+1} \approx 10^{-2\ell}$ .  $\square$

## 8.1.2 Termination Criteria/Stopping Rules



*Supplementary literature.* Also discussed in [AG11, Sect. 3.1, p. 42].

As remarked above, usually (even without roundoff errors) an iteration will never arrive at an/the exact solution  $x^*$  after finitely many steps. Thus, we can only hope to compute an *approximate* solution by accepting  $x^{(K)}$  as result for some  $K \in \mathbb{N}_0$ . Termination criteria (stopping rules) are used to determine a suitable value for  $K$ .

For the sake of efficiency  $\triangleright$  stop iteration when iteration error is just “small enough” (“small enough” depends on the concrete problem and user demands.)

### §8.1.2.1 (Classification of termination criteria (stopping rules) for iterative solvers for non-linear systems of equations)

A **termination criterion** (stopping rule) is an algorithm deciding in each step of an iterative method whether to **STOP** or to **CONTINUE**.

<b>A priori</b> termination	<b>A posteriori</b> termination
Decision to stop based on information about $F$ and $x^{(0)}$ , made before starting iteration.	Beside $x^{(0)}$ and $F$ , also current and past iterates are used to decide about termination.

A termination criterion for a convergent iteration is deemed **reliable**, if it lets the iteration **CONTINUE**, until the iteration error  $e^{(k)} := x^{(k)} - x^*$ ,  $x^*$  the limit value, satisfies certain conditions (usually imposed before the start of the iteration).  $\square$

**§8.1.2.2 (Ideal termination)** Termination criteria are usually meant to ensure accuracy of the final iterate  $\mathbf{x}^{(K)}$  in the following sense:

$$\begin{aligned}\|\mathbf{x}^{(K)} - \mathbf{x}^*\| &\leq \tau_{\text{abs}}, \quad \tau_{\text{abs}} \hat{=} \text{prescribed (absolute) tolerance.} \\ &\quad \text{or} \\ \|\mathbf{x}^{(K)} - \mathbf{x}^*\| &\leq \tau_{\text{rel}} \|\mathbf{x}^*\|, \quad \tau_{\text{rel}} \hat{=} \text{prescribed (relative) tolerance.}\end{aligned}$$

it seems that the second criterion, asking that the **relative** ( $\rightarrow$  Def. 1.5.3.3) iteration error be below a prescribed threshold, alone would suffice, but the absolute tolerance should be checked, if, by “accident”,  $\|\mathbf{x}^*\| = 0$  is possible. Otherwise, the iteration might fail to terminate at all.

Both criteria enter the “ideal (a posteriori) termination rule”:

$$\text{STOP at step } K = \operatorname{argmin}_{k \in \mathbb{N}_0} \{\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \begin{cases} \tau_{\text{abs}} \\ \tau_{\text{rel}} \|\mathbf{x}^*\| \end{cases} \text{ or } . \quad (8.1.2.3)$$

Obviously, (8.1.2.3) achieves the optimum in terms of efficiency and reliability. As obviously, this termination criterion is not practical, because  $\mathbf{x}^*$  is not known.

□

**§8.1.2.4 (Practical termination criteria for iterations)** The following termination criteria are commonly used in numerical codes:

- ① **A priori termination:** stop iteration after fixed number of steps (possibly depending on  $\mathbf{x}^{(0)}$ ).



Drawback: hard to ensure prescribed accuracy!

(A priori  $\hat{=}$  without actually taking into account the computed iterates, see § 8.1.2.1)

Invoking additional properties of either the non-linear system of equations  $\mathbf{F}(\mathbf{x}) = 0$  or the iteration it is sometimes possible to tell that for sure  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \tau$  for all  $k \geq K$ , though this  $K$  may be (significantly) larger than the optimal termination index from (8.1.2.3), see Rem. 8.1.2.6.

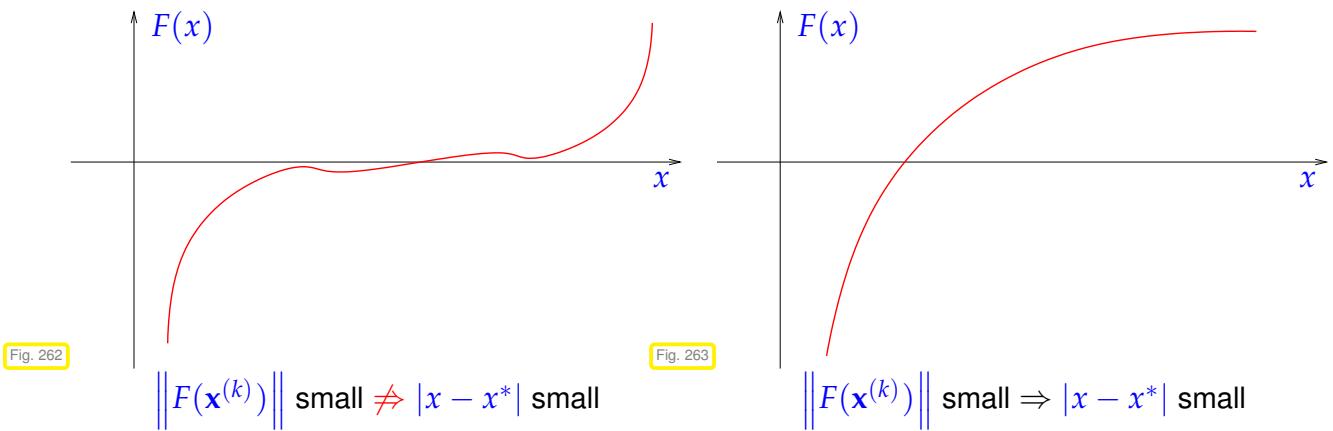
- ② **Residual based** termination: STOP convergent iteration  $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}_0}$ , when

$$\|\mathbf{F}(\mathbf{x}^{(k)})\| \leq \tau, \quad \tau \hat{=} \text{prescribed tolerance} > 0.$$



no guaranteed accuracy

Consider the case  $n = 1$ . If  $F : D \subset \mathbb{R} \rightarrow \mathbb{R}$  is “flat” in the neighborhood of a zero  $x^*$ , then a small value of  $|F(x)|$  does not mean that  $x$  is close to  $x^*$ .



③ **Correction based** termination: STOP convergent iteration  $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}_0}$ , when

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq \begin{cases} \frac{\tau_{\text{abs}}}{\tau_{\text{rel}}} \|\mathbf{x}^{(k+1)}\|, & \tau_{\text{abs}} \text{ prescribed} \\ \tau_{\text{rel}} \|\mathbf{x}^{(k+1)}\|, & \tau_{\text{rel}} \text{ relative tolerances} > 0. \end{cases}$$

Also for this criterion, we have no guarantee that (8.1.2.3) will be satisfied only remotely.

A special variant of correction based termination exploits that  $\mathbb{M}$  is finite! ( $\rightarrow$  Section 1.5.3)

Wait until (convergent) iteration becomes stationary in the discrete set  $\mathbb{M}$  of machine numbers!



possibly grossly inefficient!  
(always computes "up to machine precision")

#### C++ code 8.1.2.5: Square root iteration $\rightarrow$ Ex. 8.1.1.13

```

2 double sqrtit(double a)
3 {
4     double x_old = -1;
5     double x = a;
6     while (x_old != x) {
7         x_old = x;
8         x = 0.5*(x+a/x);
9     }
10    return x;
11 }
```

**Remark 8.1.2.6 (A posteriori termination criterion for linearly convergent iterations**  $\rightarrow$  [DR08, Lemma 5.17, 5.19]) Let us assume that we know that an iteration linearly convergent ( $\rightarrow$  Def. 8.1.1.1) with rate of convergence  $0 < L < 1$ :

The following simple manipulations give an a posteriori termination criterion (for linearly convergent iterations with rate of convergence  $0 < L < 1$ ):

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \stackrel{\triangle\text{-inequ.}}{\leq} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| + \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| + L\|\mathbf{x}^{(k)} - \mathbf{x}^*\|.$$

► Iterates satisfy:

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \frac{L}{1-L} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|. \quad (8.1.2.7)$$

This suggests that we take the right hand side of (8.1.2.7) as a posteriori error bound and use it instead of the inaccessible  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|$  for checking absolute and relative accuracy in (8.1.2.3). The resulting

termination criterium will be **reliable** ( $\rightarrow$  § 8.1.2.1), since we will certainly have achieved the desired accuracy when we stop the iteration.

 Estimating the rate of convergence  $L$  might be difficult.

 Pessimistic estimate for  $L$  will not compromise reliability.

(Using  $\tilde{L} > L$  in (8.1.2.7) still yields a valid *upper bound* for  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|$ .)

↓

**EXAMPLE 8.1.2.8 (A posteriori error bound for linearly convergent iteration)** We revisit the iteration of Ex. 8.1.1.8:

$$x^{(k+1)} = x^{(k)} + \frac{\cos x^{(k)} + 1}{\sin x^{(k)}} \Rightarrow x^{(k)} \rightarrow \pi \text{ for } x^{(0)} \text{ close to } \pi.$$

Observed rate of convergence:  $L = 1/2$

Error and error bound for  $x^{(0)} = 0.4$ :

$k$	$ x^{(k)} - \pi $	$\frac{L}{1-L} x^{(k)} - x^{(k-1)} $	slack of bound
1	2.191562221997101	4.933154875586894	2.741592653589793
2	0.247139097781070	1.944423124216031	1.697284026434961
3	0.122936737876834	0.124202359904236	0.001265622027401
4	0.061390835206217	0.061545902670618	0.000155067464401
5	0.030685773472263	0.030705061733954	0.000019288261691
6	0.015341682696235	0.015344090776028	0.000002408079792
7	0.007670690889185	0.007670991807050	0.000000300917864
8	0.003835326638666	0.003835364250520	0.000000037611854
9	0.001917660968637	0.001917665670029	0.000000004701392
10	0.000958830190489	0.000958830778147	0.000000000587658
11	0.000479415058549	0.000479415131941	0.000000000073392
12	0.000239707524646	0.000239707533903	0.0000000000009257
13	0.000119853761949	0.000119853762696	0.000000000000747
14	0.000059926881308	0.000059926880641	0.000000000000667
15	0.000029963440745	0.000029963440563	0.000000000000181

Hence: the a posteriori error bound is highly accurate in this case!

↓

## 8.2 Fixed-Point Iterations



*Supplementary literature.* The contents of this section are also treated in [DR08, Sect. 5.3],

[QSS00, Sect. 6.3], [AG11, Sect. 3.3]

As before we consider a non-linear system of equations  $F(\mathbf{x}) = 0$ ,  $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ .

1-point stationary iterative methods, see (8.1.0.4), for  $F(\mathbf{x}) = \mathbf{0}$  are also called **fixed point iterations**.

► A fixed point iteration is defined by **iteration function**  $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ :

$$\begin{array}{l} \text{iteration function} \quad \Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n \\ \text{initial guess} \quad \mathbf{x}^{(0)} \in U \end{array} \Rightarrow \boxed{\text{iterates } (\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}: \quad \mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)})} .$$

$\rightarrow$  1-point method, cf. (8.1.0.4)

Here,  $U$  designates the domain of definition of the iteration function  $\Phi$ .

Note that the sequence of iterates need not be well defined:  $\mathbf{x}^{(k)} \notin U$  possible !

### 8.2.1 Consistent Fixed-Point Iterations

Next, we specialize Def. 8.1.0.7 for fixed point iterations:

#### Definition 8.2.1.1. Consistency of fixed point iterations, c.f. Def. 8.1.0.7

A fixed point iteration  $\mathbf{x}^{(k+1)} = \Phi(\mathbf{x}^{(k)})$  is **consistent** with  $F(\mathbf{x}) = \mathbf{0}$ , if, for  $\mathbf{x} \in U \cap D$ ,

$$F(\mathbf{x}) = \mathbf{0} \Leftrightarrow \Phi(\mathbf{x}) = \mathbf{x} .$$

Note: iteration function  $\Phi$  continuous and fixed point iteration (locally) convergent to  $\mathbf{x}^* \in U \Rightarrow \mathbf{x}^*$  is a fixed point of  $\Phi$ .

This is an immediate consequence that for a continuous function limits and function evaluations commute [Str09, Sect. 4.1].

General construction of fixed point iterations that is consistent with  $F(\mathbf{x}) = \mathbf{0}$ :

- ① Rewrite equivalently  $F(\mathbf{x}) = \mathbf{0} \Leftrightarrow \Phi(\mathbf{x}) = \mathbf{x}$  and then
- ② use the fixed point iteration

$$\mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)}) . \quad (8.2.1.2)$$

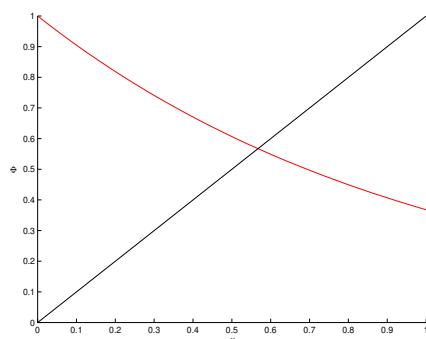
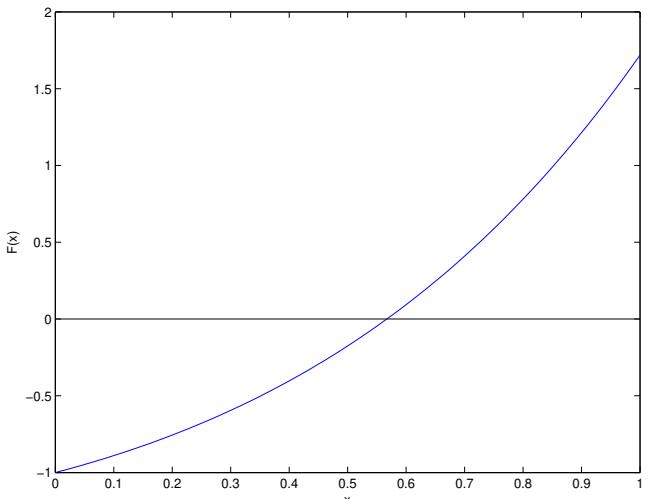
Note: there are *many ways* to transform  $F(\mathbf{x}) = \mathbf{0}$  into a fixed point form !

**EXPERIMENT 8.2.1.3 (Many choices for consistent fixed point iterations)** In this example we construct three different consistent fixed point iteration for a single scalar ( $n = 1$ ) non-linear equation  $F(x) = 0$ . In numerical experiments we will see that they behave very differently.

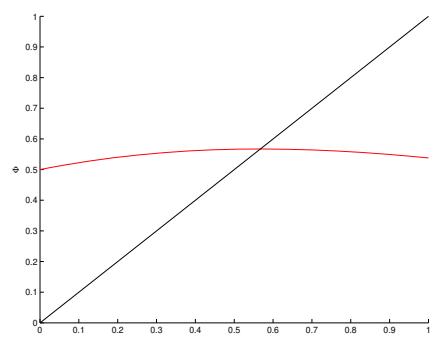
$$F(x) = xe^x - 1, \quad x \in [0, 1].$$

Different fixed point forms:

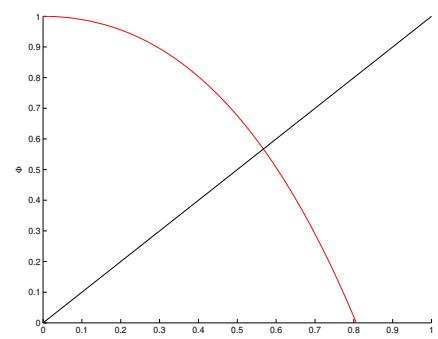
$$\begin{aligned}\Phi_1(x) &= e^{-x}, \\ \Phi_2(x) &= \frac{1+x}{1+e^x}, \\ \Phi_3(x) &= x + 1 - xe^x.\end{aligned}$$



function  $\Phi_1$



function  $\Phi_2$



function  $\Phi_3$

With the same initial guess  $x^{(0)} = 0.5$  for all three fixed point iterations we obtain the following iterates:

$k$	$x^{(k+1)} := \Phi_1(x^{(k)})$	$x^{(k+1)} := \Phi_2(x^{(k)})$	$x^{(k+1)} := \Phi_3(x^{(k)})$
0	0.5000000000000000	0.5000000000000000	0.5000000000000000
1	0.606530659712633	0.566311003197218	0.675639364649936
2	0.545239211892605	0.567143165034862	0.347812678511202
3	0.579703094878068	0.567143290409781	0.855321409174107
4	0.560064627938902	0.567143290409784	-0.156505955383169
5	0.571172148977215	0.567143290409784	0.977326422747719
6	0.564862946980323	0.567143290409784	-0.619764251895580
7	0.568438047570066	0.567143290409784	0.713713087416146
8	0.566409452746921	0.567143290409784	0.256626649129847
9	0.567559634262242	0.567143290409784	0.924920676910549
10	0.566907212935471	0.567143290409784	-0.407422405542253

We can also tabulate the modulus of the iteration error and mark correct digits with red:

$k$	$ x_1^{(k+1)} - x^* $	$ x_2^{(k+1)} - x^* $	$ x_3^{(k+1)} - x^* $
0	0.067143290409784	0.067143290409784	0.067143290409784
1	0.039387369302849	0.000832287212566	0.108496074240152
2	0.021904078517179	0.000000125374922	0.219330611898582
3	0.012559804468284	0.000000000000003	0.288178118764323
4	0.007078662470882	0.000000000000000	0.723649245792953
5	0.004028858567431	0.000000000000000	0.410183132337935
6	0.002280343429460	0.000000000000000	1.186907542305364
7	0.001294757160282	0.000000000000000	0.146569797006362
8	0.000733837662863	0.000000000000000	0.310516641279937
9	0.000416343852458	0.000000000000000	0.357777386500765
10	0.000236077474313	0.000000000000000	0.974565695952037

Observed: linear convergence of  $x_1^{(k)}$ , quadratic convergence of  $x_2^{(k)}$ ,  
no convergence (erratic behavior of  $x_3^{(k)}$ ) ( $x_i^{(0)} = 0.5$  in all cases).  $\square$

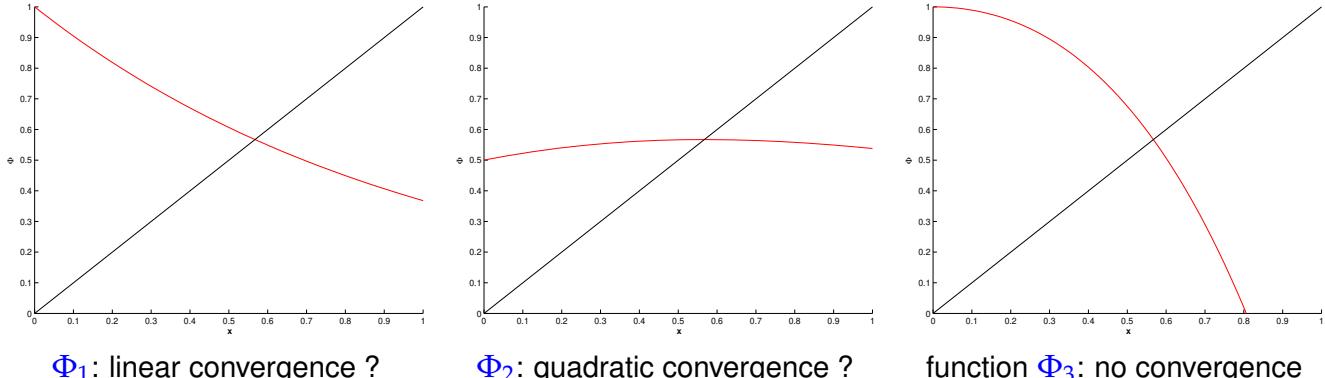
Question: Can we explain/forecast the behaviour of a fixed point iteration?

## 8.2.2 Convergence of Fixed-Point Iterations

In this section we will try to find easily verifiable conditions that ensure convergence (of a certain order) of fixed point iterations. It will turn out that these conditions are surprisingly simple and general.

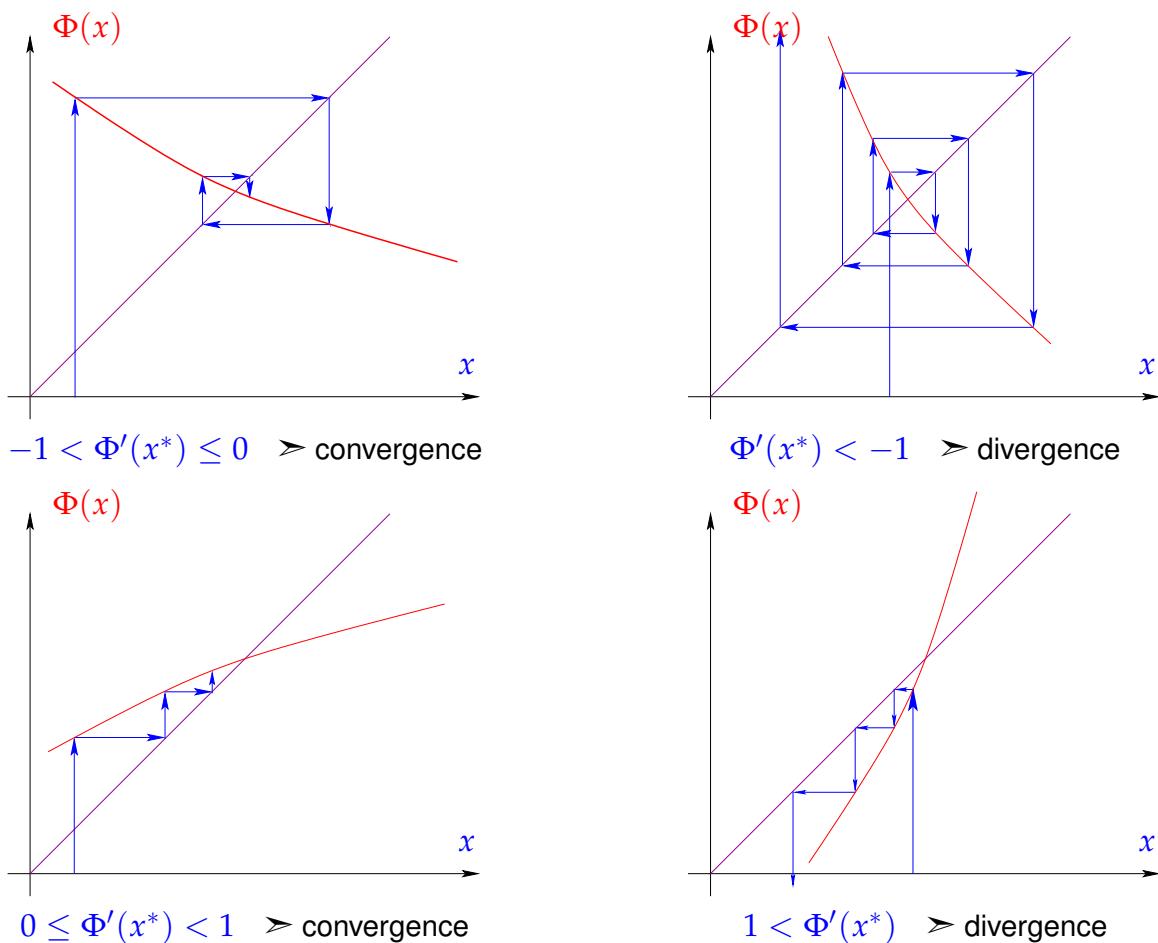
### EXPERIMENT 8.2.2.1 (Exp. 8.2.1.3 revisited)

In Exp. 8.2.1.3 we observed vastly different behavior of different fixed point iterations for  $n = 1$ . Is it possible to predict this from the shape of the graph of the iteration functions?



**Remark 8.2.2.2 (Visualization of fixed point iterations in 1D)**  
 1D setting ( $n = 1$ ):  $\Phi : \mathbb{R} \mapsto \mathbb{R}$  continuously differentiable,  $\Phi(x^*) = x^*$   
 fixed point iteration:  $x^{(k+1)} = \Phi(x^{(k)})$

In 1D it is possible to visualize the different convergence behavior of fixed point iterations: In order to construct  $x^{(k+1)}$  from  $x^{(k)}$  one moves vertically to  $(x^{(k)}, x^{(k+1)} = \Phi(x^{(k)}))$ , then horizontally to the angular bisector of the first/third quadrant, that is, to the point  $(x^{(k+1)}, x^{(k+1)})$ . Returning vertically to the abscissa gives  $x^{(k+1)}$ .



Numerical examples for iteration functions  $\Rightarrow$  Exp. 8.2.1.3, iteration functions  $\Phi_1$  and  $\Phi_3$

It seems that the slope of the iteration function  $\Phi$  in the fixed point, that is, in the point where it intersects the bisector of the first/third quadrant, is crucial.  $\square$

Now we investigate rigorously, when a fixed point iteration will lead to a convergent iteration with a particular qualitative kind of convergence according to Def. 8.1.1.10.

### Definition 8.2.2.3. Contractive mapping

$\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$  is **contractive** (w.r.t. norm  $\|\cdot\|$  on  $\mathbb{R}^n$ ), if

$$\exists L < 1: \quad \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in U. \quad (8.2.2.4)$$

A simple consideration: if  $\Phi(\mathbf{x}^*) = \mathbf{x}^*$  (fixed point), then a fixed point iteration induced by a contractive mapping  $\Phi$  satisfies

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = \|\Phi(\mathbf{x}^{(k)}) - \Phi(\mathbf{x}^*)\| \stackrel{(8.2.2.4)}{\leq} L \|\mathbf{x}^{(k)} - \mathbf{x}^*\|,$$

that is, the iteration **converges** (at least) **linearly** ( $\rightarrow$  Def. 8.1.1.1).

Note that

$\Phi$  contractive  $\Rightarrow$   $\Phi$  has **at most one** fixed point.

**Remark 8.2.2.5 (Banach's fixed point theorem)**  $\rightarrow$  [Str09, Satz 6.5.2], [DR08, Satz 5.8]) A key theorem in calculus (also functional analysis):

**Theorem 8.2.2.6. Banach's fixed point theorem**

If  $D \subset \mathbb{K}^n$  ( $\mathbb{K} = \mathbb{R}, \mathbb{C}$ ) closed and bounded and  $\Phi : D \mapsto D$  satisfies

$$\exists L < 1: \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in D,$$

then there is a unique fixed point  $\mathbf{x}^* \in D$ ,  $\Phi(\mathbf{x}^*) = \mathbf{x}^*$ , which is the limit of the sequence of iterates  $\mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)})$  for any  $\mathbf{x}^{(0)} \in D$ .

*Proof.* Proof based on 1-point iteration  $\mathbf{x}^{(k)} = \Phi(\mathbf{x}^{(k-1)})$ ,  $\mathbf{x}^{(0)} \in D$ :

$$\begin{aligned} \|\mathbf{x}^{(k+N)} - \mathbf{x}^{(k)}\| &\leq \sum_{j=k}^{k+N-1} \|\mathbf{x}^{(j+1)} - \mathbf{x}^{(j)}\| \leq \sum_{j=k}^{k+N-1} L^j \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \\ &\leq \frac{L^k}{1-L} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \xrightarrow{k \rightarrow \infty} 0. \end{aligned}$$

$(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$  Cauchy sequence  $\rightarrow$  convergent  $\mathbf{x}^{(k)} \xrightarrow{k \rightarrow \infty} \mathbf{x}^*$ .  
Continuity of  $\Phi \rightarrow \Phi(\mathbf{x}^*) = \mathbf{x}^*$ . Uniqueness of fixed point is evident.

□

]

A simple criterion for a differentiable  $\Phi$  to be contractive:

**Lemma 8.2.2.7. Sufficient condition for local linear convergence of fixed point iteration** → [Han02, Thm. 17.2], [DR08, Cor. 5.12]

If  $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ ,  $\Phi(\mathbf{x}^*) = \mathbf{x}^*$ ,  $\Phi$  differentiable in  $\mathbf{x}^*$ , and  $\|\mathbf{D}\Phi(\mathbf{x}^*)\| < 1$ , then the fixed point iteration

$$\mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)}), \tag{8.2.1.2}$$

converges locally and at least linearly.

matrix norm, Def. 1.5.5.10 !

☞ notation:  $\mathbf{D}\Phi(\mathbf{x}) \triangleq$  **Jacobian** (ger.: Jacobi-Matrix) of  $\Phi$  at  $\mathbf{x} \in D$  → [Str09, Sect. 7.6]

$$\mathbf{D}\Phi(\mathbf{x}) = \left[ \frac{\partial \Phi_i}{\partial x_j}(\mathbf{x}) \right]_{i,j=1}^n = \begin{bmatrix} \frac{\partial \Phi_1}{\partial x_1}(\mathbf{x}) & \frac{\partial \Phi_1}{\partial x_2}(\mathbf{x}) & \cdots & \cdots & \frac{\partial \Phi_1}{\partial x_n}(\mathbf{x}) \\ \frac{\partial \Phi_2}{\partial x_1}(\mathbf{x}) & & & & \frac{\partial \Phi_2}{\partial x_n}(\mathbf{x}) \\ \vdots & & & & \vdots \\ \frac{\partial \Phi_n}{\partial x_1}(\mathbf{x}) & \frac{\partial \Phi_n}{\partial x_2}(\mathbf{x}) & \cdots & \cdots & \frac{\partial \Phi_n}{\partial x_n}(\mathbf{x}) \end{bmatrix}. \tag{8.2.2.8}$$

“Visualization” of the statement of Lemma 8.2.2.7 in Rem. 8.2.2.2: The iteration converges *locally*, if  $\Phi$  is flat in a neighborhood of  $\mathbf{x}^*$ , it will diverge, if  $\Phi$  is steep there.

*Proof.* (of Lemma 8.2.2.7) By definition of derivative

$$\|\Phi(\mathbf{y}) - \Phi(\mathbf{x}^*) - \mathbf{D}\Phi(\mathbf{x}^*)(\mathbf{y} - \mathbf{x}^*)\| \leq \psi(\|\mathbf{y} - \mathbf{x}^*\|)\|\mathbf{y} - \mathbf{x}^*\|,$$

with  $\psi : \mathbb{R}_0^+ \mapsto \mathbb{R}_0^+$  satisfying  $\lim_{t \rightarrow 0} \psi(t) = 0$ .

Choose  $\delta > 0$  such that

$$L := \psi(t) + \|D\Phi(\mathbf{x}^*)\| \leq \frac{1}{2}(1 + \|D\Phi(\mathbf{x}^*)\|) < 1 \quad \forall 0 \leq t < \delta.$$

By inverse triangle inequality we obtain for fixed point iteration

$$\begin{aligned} \|\Phi(\mathbf{x}) - \mathbf{x}^*\| - \|D\Phi(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*)\| &\leq \psi(\|\mathbf{x} - \mathbf{x}^*\|)\|\mathbf{x} - \mathbf{x}^*\| \\ \Rightarrow \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| &\leq (\psi(t) + \|D\Phi(\mathbf{x}^*)\|)\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq L\|\mathbf{x}^{(k)} - \mathbf{x}^*\|, \end{aligned}$$

if  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\| < \delta$ . □

### Lemma 8.2.2.9. Sufficient condition for linear convergence of fixed point iteration

Let  $U$  be convex and  $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$  be continuously differentiable with

$$L := \sup_{\mathbf{x} \in U} \|D\Phi(\mathbf{x})\| < 1.$$

If  $\Phi(\mathbf{x}^*) = \mathbf{x}^*$  for some interior point  $\mathbf{x}^* \in U$ , then the fixed point iteration  $\mathbf{x}^{(k+1)} = \Phi(\mathbf{x}^{(k)})$  converges to  $\mathbf{x}^*$  at least linearly with rate  $L$ .

Recall:  $U \subset \mathbb{R}^n$  convex  $\Leftrightarrow (t\mathbf{x} + (1-t)\mathbf{y}) \in U$  for all  $\mathbf{x}, \mathbf{y} \in U, 0 \leq t \leq 1$

*Proof.* (of Lemma 8.2.2.9) By the mean value theorem

$$\begin{aligned} \Phi(\mathbf{x}) - \Phi(\mathbf{y}) &= \int_0^1 D\Phi(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x}) d\tau \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(\Phi). \\ \Rightarrow \quad \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\| &\leq L\|\mathbf{y} - \mathbf{x}\|, \\ \Rightarrow \quad \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| &\leq L\|\mathbf{x}^{(k)} - \mathbf{x}^*\|. \end{aligned}$$

We find that  $\Phi$  is contractive on  $U$  with unique fixed point  $\mathbf{x}^*$ , to which  $\mathbf{x}^{(k)}$  converges linearly for  $k \rightarrow \infty$ . □

**Remark 8.2.2.10 (Bound for asymptotic rate of linear convergence)** By asymptotic rate of a linearly converging iteration we mean the contraction factor for the norm of the iteration error that we can expect, when we are already very close to the limit  $\mathbf{x}^*$ .

If  $0 < \|D\Phi(\mathbf{x}^*)\| < 1$ ,  $\mathbf{x}^{(k)} \approx \mathbf{x}^*$  then the (worst) asymptotic rate of linear convergence is  $L = \|D\Phi(\mathbf{x}^*)\|$  □

**EXAMPLE 8.2.2.11 (Multidimensional fixed point iteration)** In this example we encounter the first genuine system of non-linear equations and apply Lemma 8.2.2.9 to it.

$$\begin{aligned} \text{System of equations} &\quad \text{in} && \text{fixed point form:} \\ \left\{ \begin{array}{lcl} x_1 - c(\cos x_1 - \sin x_2) & = & 0 \\ (x_1 - x_2) - c \sin x_2 & = & 0 \end{array} \right. &\Rightarrow& \left\{ \begin{array}{lcl} c(\cos x_1 - \sin x_2) & = & x_1 \\ c(\cos x_1 - 2 \sin x_2) & = & x_2 \end{array} \right. \\ \text{Define: } \Phi \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} &= c \begin{bmatrix} \cos x_1 - \sin x_2 \\ \cos x_1 - 2 \sin x_2 \end{bmatrix} &\Rightarrow& D\Phi \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = -c \begin{bmatrix} \sin x_1 & \cos x_2 \\ \sin x_1 & 2 \cos x_2 \end{bmatrix}. \end{aligned}$$

Choose appropriate norm:  $\|\cdot\| = \infty$ -norm  $\|\cdot\|_\infty$  ( $\rightarrow$  Example 1.5.5.12);

$$\text{if } c < \frac{1}{3} \Rightarrow \|D\Phi(\mathbf{x})\|_\infty < 1 \quad \forall \mathbf{x} \in \mathbb{R}^2,$$

➤ (at least) linear convergence of the fixed point iteration. The existence of a fixed point is also guaranteed, because  $\Phi$  maps into the closed set  $[-3, 3]^2$ . Thus, the Banach fixed point theorem, Thm. 8.2.2.6, can be applied.  $\square$

What about higher order convergence ( $\rightarrow$  Def. 8.1.1.10, cf.  $\Phi_2$  in Ex. 8.2.1.3)? Also in this case we should study the derivatives of the iteration functions in the fixed point (limit point).

We give a refined convergence result only for  $n = 1$  (scalar case,  $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$ ):

**Theorem 8.2.2.12. Taylor's formula**  $\rightarrow$  [Str09, Sect. 5.5]

If  $\Phi : U \subset \mathbb{R} \mapsto \mathbb{R}$ ,  $U$  interval, is  $m + 1$  times continuously differentiable,  $x \in U$

$$\Phi(y) - \Phi(x) = \sum_{k=1}^m \frac{1}{k!} \Phi^{(k)}(x)(y-x)^k + O(|y-x|^{m+1}) \quad \forall y \in U. \quad (8.2.2.13)$$

Now apply Taylor expansion (8.2.2.13) to iteration function  $\Phi$ :

If  $\Phi(x^*) = x^*$  and  $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$  is “sufficiently smooth”, it tells us that

$$x^{(k+1)} - x^* = \Phi(x^{(k)}) - \Phi(x^*) = \sum_{l=1}^m \frac{1}{l!} \Phi^{(l)}(x^*)(x^{(k)} - x^*)^l + O(|x^{(k)} - x^*|^{m+1}). \quad (8.2.2.14)$$

Here we used the Landau symbol  $O(\cdot)$  to describe the local behavior of a remainder term in the vicinity of  $x^*$ .

**Lemma 8.2.2.15. Higher order local convergence of fixed point iterations**

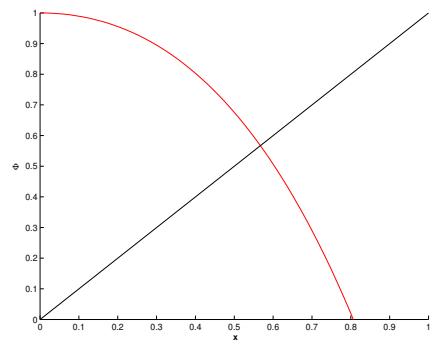
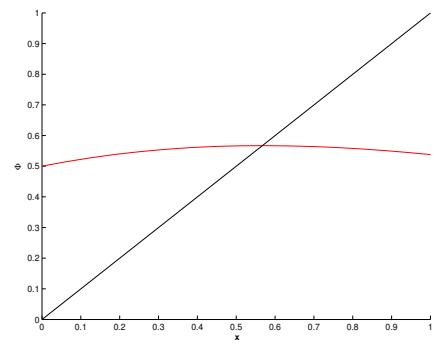
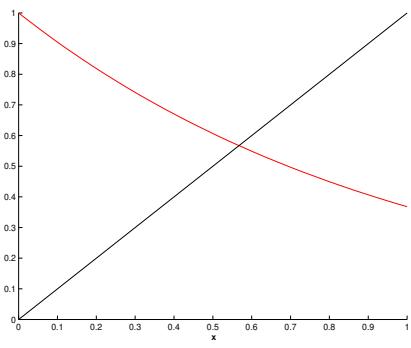
If  $\Phi : U \subset \mathbb{R} \mapsto \mathbb{R}$  is  $m + 1$  times continuously differentiable,  $\Phi(x^*) = x^*$  for some  $x^*$  in the interior of  $U$ , and  $\Phi^{(l)}(x^*) = 0$  for  $l = 1, \dots, m$ ,  $m \geq 1$ , then the fixed point iteration (8.2.1.2) converges **locally** to  $x^*$  with **order  $\geq m + 1$**  ( $\rightarrow$  Def. 8.1.1.10).

*Proof.* For neighborhood  $U$  of  $x^*$

$$(8.2.2.14) \Rightarrow \exists C > 0: |\Phi(y) - \Phi(x^*)| \leq C |y - x^*|^{m+1} \quad \forall y \in U. \\ \delta^m C < 1/2: |x^{(0)} - x^*| < \delta \Rightarrow |x^{(k)} - x^*| < 2^{-k} \delta \Rightarrow \text{local convergence}.$$

Then appeal to (8.2.2.14)  $\square$

**EXPERIMENT 8.2.2.16 (Exp. 8.2.2.1 continued)** Now, Lemma 8.2.2.9 and Lemma 8.2.2.15 permit us a precise prediction of the (asymptotic) convergence we can expect from the different fixed point iterations studied in Exp. 8.2.1.3.



$$\Phi'_2(x) = \frac{1 - xe^x}{(1 + e^x)^2} = 0 \quad , \text{ if } xe^x - 1 = 0 \quad \text{hence quadratic convergence!} .$$

Since  $x^*e^{x^*} - 1 = 0$ , simple computations yield

$$\begin{aligned} \Phi'_1(x) = -e^{-x} &\Rightarrow \Phi'_1(x^*) = -x^* \approx -0.56 \quad \text{hence local linear convergence.} \\ \Phi'_3(x) = 1 - xe^x - e^x &\Rightarrow \Phi'_3(x^*) = -\frac{1}{x^*} \approx -1.79 \quad \text{hence no convergence.} \end{aligned}$$

□

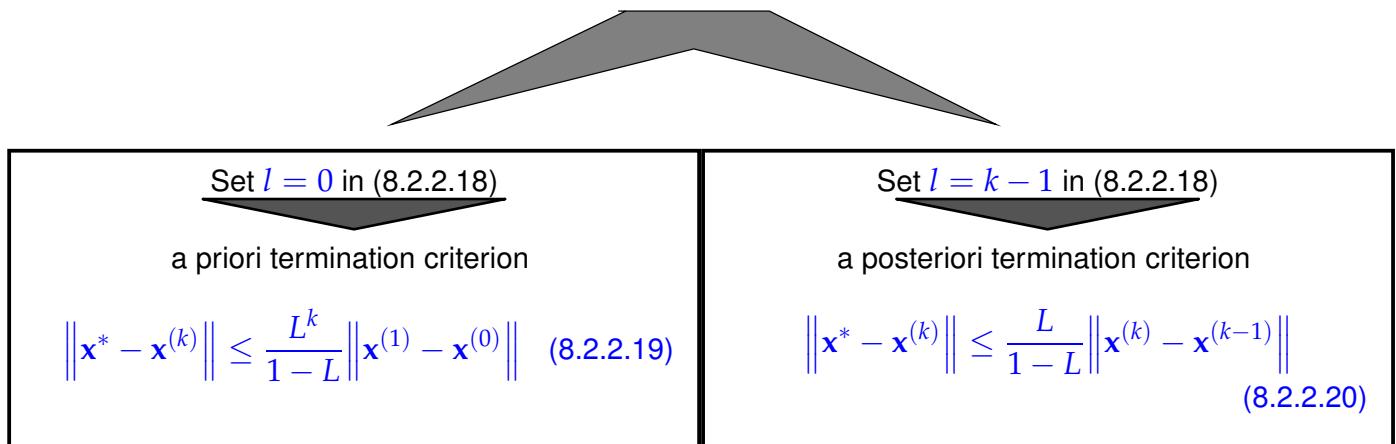
### Remark 8.2.2.17 (Termination criterion for contractive fixed point iteration)

We recall the considerations of Rem. 8.1.2.6 about a termination criterion for contractive fixed point iteration (= linearly convergent fixed point iteration → Def. 8.1.1.1), c.f. (8.2.2.4), with contraction factor (= rate of convergence)  $0 \leq L < 1$ :

$$\begin{aligned} \|\mathbf{x}^{(k+m)} - \mathbf{x}^{(k)}\| &\stackrel{\triangle\text{-ineq.}}{\leq} \sum_{j=k}^{k+m-1} \|\mathbf{x}^{(j+1)} - \mathbf{x}^{(j)}\| \leq \sum_{j=k}^{k+m-1} L^{j-k} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \\ &= \frac{1 - L^m}{1 - L} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq \frac{1 - L^m}{1 - L} L^{k-l} \|\mathbf{x}^{(l+1)} - \mathbf{x}^{(l)}\|. \end{aligned}$$

Hence, for  $m \rightarrow \infty$ , with  $\mathbf{x}^* := \lim_{k \rightarrow \infty} \mathbf{x}^{(k)}$  we find the estimate

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{L^{k-l}}{1 - L} \|\mathbf{x}^{(l+1)} - \mathbf{x}^{(l)}\|. \quad (8.2.2.18)$$



With the same arguments as in Rem. 8.1.2.6 we see that **overestimating  $L$** , that is, using a value for  $L$  that is larger than the true value, still gives **reliable** termination criteria.

However, whereas overestimating  $L$  in (8.2.2.20) will not lead to a severe deterioration of the bound, unless  $L \approx 1$ , using a pessimistic value for  $L$  in (8.2.2.19) will result in a bound way bigger than the true bound, if  $k \gg 1$ . Then the a priori termination criterion (8.2.2.19) will recommend termination many iterations after the accuracy requirements have already been met. This will **thwart** the **efficiency** of the method. □

## 8.3 Finding Zeros of Scalar Functions



*Supplementary literature.* [AG11, Ch. 3] is also devoted to this topic. The algorithm of “bisection” discussed in the next subsection, is treated in [DR08, Sect. 5.5.1] and [AG11, Sect. 3.2].

Now, we focus on scalar case  $n = 1$ :

$F : I \subset \mathbb{R} \mapsto \mathbb{R}$  continuous,  $I$  interval

Sought:

$$x^* \in I: \quad F(x^*) = 0$$

### 8.3.1 Bisection

Idea: use ordering of real numbers & intermediate value theorem [Str09, Sect. 4.6]

Input:  $a, b \in I$  such that

$$F(a)F(b) < 0. \quad (\text{different signs!})$$

►  $\exists x^* \in ]\min\{a, b\}, \max\{a, b\}[:$   
 $F(x^*) = 0,$

as we conclude from the intermediate value theorem.

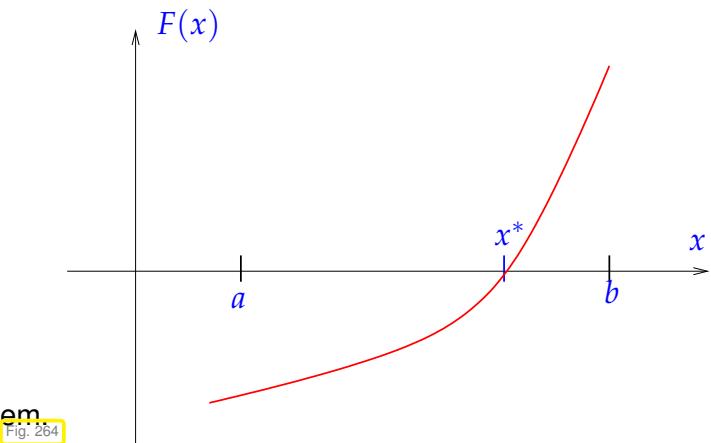


Fig. 264

► Find a sequence of intervals with geometrically decreasing lengths, in each of which  $F$  will change sign.

Such a sequence can easily be found by testing the sign of  $F$  at the midpoint of the current interval, see Code 8.3.1.2.

**§8.3.1.1 (Bisection method)** The following C++ code implements the bisection method for finding the zeros of a function passed through the **function handle**  $F$  in the interval  $[a, b]$  with **absolute tolerance**  $tol$ .

**C++ code 8.3.1.2: Bisection method for solving  $F(x) = 0$  on  $[a, b]$** 

```

2 // Searching zero of F in [a,b] by bisection
3 template <typename Func, typename Scalar>
4 Scalar bisect(Func&& F, Scalar a, Scalar b, Scalar tol)
5 {
6     if (a > b) std::swap(a,b); // sort interval bounds
7     if (F(a)*F(b) > 0) throw "f(a) and f(b) have same sign";
8     static_assert(std::is_floating_point<Scalar>::value,
9                 "Scalar must be a floating point type");
10    int v=F(a) < 0 ? 1 : -1;
11    Scalar x = (a+b)/2; // determine midpoint
12    // termination, relies on machine arithmetic if tol = 0
13    while (b-a > tol) { //
14        assert(a<=x && x<=b); // assert invariant
15        // sgn(f(x)) = sgn(f(b)), then use x as next right boundary
16        if (v*F(x) > 0) b=x;
17        // sgn(f(x)) = sgn(f(a)), then use x as next left boundary
18        else a=x;
19        x = (a+b)/2; // determine next midpoint
20    }
21    return x;
22 }
```

Line 13: the test  $((a < x) \&\& (x < b))$  offers a safeguard against an infinite loop in case  $\text{tol} < \text{resolution of M}$  at zero  $x^*$  (cf. “M-based termination criterion”).

This is also an example for an algorithm that (in the case of  $\text{tol}=0$ ) uses the properties of machine arithmetic to define an a posteriori termination criterion, see Section 8.1.2. The iteration will terminate, when, e.g.,  $a \tilde{+} \frac{1}{2}(b - a) = a$  ( $\tilde{+}$  is the floating point realization of addition), which, by the Ass. 1.5.3.11 can only happen, when

$$\left| \frac{1}{2}(b - a) \right| \leq \text{EPS} \cdot |a| .$$

Since the exact zero is located between  $a$  and  $b$ , this condition implies a relative error  $\leq \text{EPS}$  of the computed zero.  $\square$



Advantages:

- “foolproof”, **robust**: will always terminate with a zero of requested accuracy,
- requires only point evaluations of  $F$ ,
- works with *any* continuous function  $F$ , no derivatives needed.



Drawbacks:

Merely “linear-type” (\*) convergence:  $|x^{(k)} - x^*| \leq 2^{-k}|b - a|$

$$\log_2\left(\frac{|b - a|}{\text{tol}}\right) \text{ steps necessary}$$

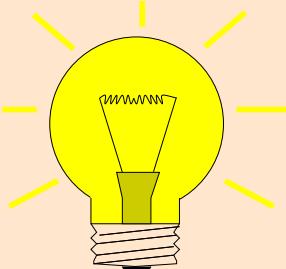
(\*): the convergence of a bisection algorithm is not linear in the sense of Def. 8.1.1.1, because the condition  $\|x^{(k+1)} - x^*\| \leq L \|x^{(k)} - x^*\|$  might be violated at any step of the iteration.

**Remark 8.3.1.3 (Generalized bisection methods)** It is straightforward to combine the bisection idea with more elaborate “model function methods” as they will be discussed in the next section: Instead of stubbornly choosing the midpoint of the probing interval  $[a, b]$  ( $\rightarrow$  Code 8.3.1.2) as next iterate, one may use a refined guess for the location of a zero of  $F$  in  $[a, b]$ .

A method of this type is used by MATLAB’s `fzero` function for root finding in 1D [QSS00, Sect. 6.2.3].  $\square$

### 8.3.2 Model Function Methods

≈ class of iterative methods for finding zeroes of  $F$ : iterate in step  $k+1$  is computed according to the following idea:



Idea: Given recent iterates (approximate zeroes)  
 $x^{(k)}, x^{(k-1)}, \dots, x^{(k-m+1)}, m \in \mathbb{N}$

① replace  $F$  with a  $k$ -dependent model function  $\tilde{F}_k$   
 (based on function values  $F(x^{(k)}), F(x^{(k-1)}), \dots, F(x^{(k-m+1)})$  and,  
 possibly, derivative values  $F'(x^{(k)}), F'(x^{(k-1)}), \dots, F'(x^{(k-m+1)})$ )

②  $x^{(k+1)} := \text{zero of } \tilde{F}_k: \tilde{F}_k(x^{(k+1)}) = 0$   
 (has to be readily available  $\leftrightarrow$  analytic formula)

Distinguish (see § 8.1.0.2 and (8.1.0.4)):

one-point methods :  $x^{(k+1)} = \Phi_F(x^{(k)}), k \in \mathbb{N}$  (e.g., fixed point iteration → Section 8.2)

multi-point methods :  $x^{(k+1)} = \Phi_F(x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}), k \in \mathbb{N}, m = 2, 3, \dots$

#### 8.3.2.1 Newton Method in the Scalar Case



*Supplementary literature.* Newton's method in 1D is discussed in [Han02, Sect. 18.1], [DR08, Sect. 5.5.2], [AG11, Sect. 3.4].

☞ Again we consider the problem of finding zeros of the function  $F: I \subset \mathbb{R} \rightarrow \mathbb{R}$ .

☞ Now we assume that  $F: I \subset \mathbb{R} \mapsto \mathbb{R}$  is continuously differentiable.

Now model function := tangent at  $F$  in  $x^{(k)}$ :

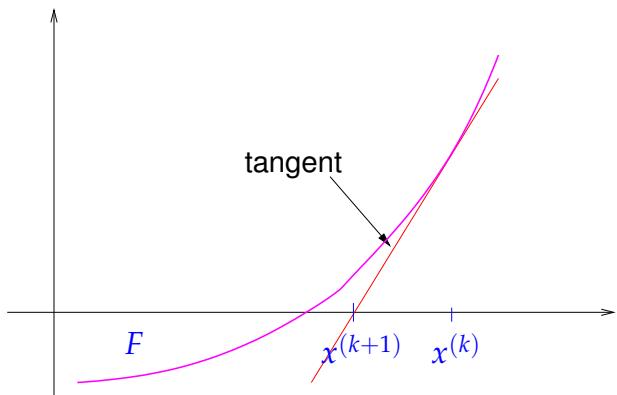
$$\tilde{F}_k(x) := F(x^{(k)}) + F'(x^{(k)})(x - x^{(k)})$$

take  $x^{(k+1)} := \text{zero of tangent}$

We obtain the **Newton iteration**

$$x^{(k+1)} := x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})},$$

(8.3.2.1)



that requires  $F'(x^{(k)}) \neq 0$ .

The following C++ code snippet implements a generic Newton method for zero finding.

- The types **FuncType** and **DervType** must be functor types and provide an evaluation operator **Scalar operator (Scalar) const**.
- The arguments **F** and **DF** must provide functors for **F** and **F'**.

**C++11-code 8.3.2.2: Newton method in the scalar case  $n = 1$** 

```

1 template <typename FuncType, typename DervType, typename Scalar>
2 Scalar newton1D(FuncType &&F, DervType &&DF,
3                   const Scalar &x0, double rtol, double atol)
4 {
5     Scalar s, x = x0;
6     do {
7         s = F(x)/DF(x); // compute Newton correction
8         x -= s;           // compute next iterate
9     }
10    // correction based termination (relative and absolute)
11    while ((std::abs(s) > rtol*std::abs(x)) && (std::abs(s) > atol));
12    return (x);
13 }
```

This code implements a correction-based termination criterion as introduced in § 8.1.2.4, see also § 8.1.2.2 for a discussion of absolute and relative tolerances.

**EXAMPLE 8.3.2.3 (Square root iteration as Newton's method)** In Ex. 8.1.1.13 we learned about the *quadratically convergent* fixed point iteration (8.1.1.14) for the approximate computation of the square root of a positive number. It can be derived as a Newton iteration (8.3.2.1)!

For  $F(x) = x^2 - a$ ,  $a > 0$ , we find  $F'(x) = 2x$ , and, thus, the Newton iteration for finding zeros of  $F$  reads:

$$x^{(k+1)} = x^{(k)} - \frac{(x^{(k)})^2 - a}{2x^{(k)}} = \frac{1}{2}(x^{(k)} + \frac{a}{x^{(k)}}),$$

which is exactly (8.1.1.14). Thus, for this  $F$  Newton's method converges globally with order  $p = 2$ . ✉

**EXAMPLE 8.3.2.4 (Newton method in 1D (→ Exp. 8.2.1.3))** Newton iterations for two different scalar non-linear equations  $F(x) = 0$  with the same solution sets:

$$F(x) = xe^x - 1 \Rightarrow F'(x) = e^x(1+x) \Rightarrow x^{(k+1)} = x^{(k)} - \frac{x^{(k)}e^{x^{(k)}} - 1}{e^{x^{(k)}}(1+x^{(k)})} = \frac{(x^{(k)})^2 + e^{-x^{(k)}}}{1+x^{(k)}}$$

$$F(x) = x - e^{-x} \Rightarrow F'(x) = 1 + e^{-x} \Rightarrow x^{(k+1)} = x^{(k)} - \frac{x^{(k)} - e^{-x^{(k)}}}{1 + e^{-x^{(k)}}} = \frac{1 + x^{(k)}}{1 + e^{x^{(k)}}}$$

Exp. 8.2.1.3 confirms **quadratic convergence** in both cases! (→ Def. 8.1.1.10)

Note that for the computation of its zeros, the function  $F$  in this example can be recast in different forms!

✉

In fact, based on Lemma 8.2.2.15 it is straightforward to show local quadratic convergence of Newton's method to a zero  $x^*$  of  $F$ , provided that  $F'(x^*) \neq 0$ :

Newton iteration (8.3.2.1)  $\hat{=}$  fixed point iteration (→ Section 8.2) with iteration function

$$\Phi(x) = x - \frac{F(x)}{F'(x)} \Rightarrow \Phi'(x) = \frac{F(x)F''(x)}{(F'(x))^2} \Rightarrow \Phi'(x^*) = 0, \text{ if } F(x^*) = 0, F'(x^*) \neq 0.$$

Thus from Lemma 8.2.2.15 we conclude the following result:

### Convergence of Newton's method in 1D

Newton's method locally quadratically converges ( $\rightarrow$  Def. 8.1.1.10) to a zero  $x^*$  of  $F$ , if  $F'(x^*) \neq 0$

#### EXAMPLE 8.3.2.6 (Implicit differentiation of $F$ )

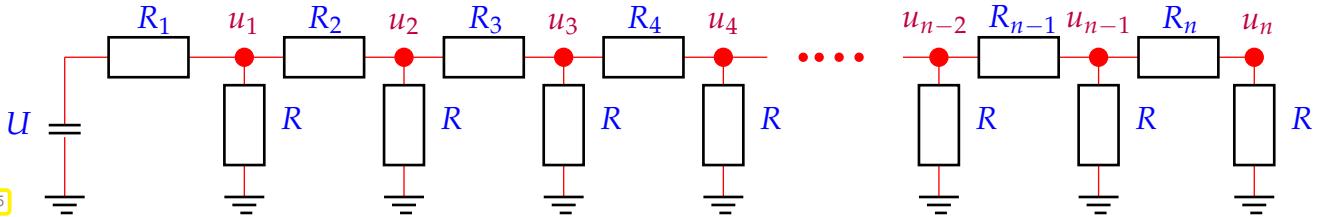


Fig. 265

How do we have to choose the leak resistance  $R$  in the linear circuit displayed in Fig. 265 in order to achieve a prescribed potential at one of the nodes?

The circuit displayed in Fig. 265 is composed of linear resistors only. Thus we can use the nodal analysis of the circuit introduced in Ex. 2.1.0.3 in order to derive a linear system of equations for the nodal potentials  $u_j$  in the nodes represented by  $\bullet$  in Fig. 265. Kirchhoff's current law (2.1.0.4) plus the constitutive relationship  $I = U/R$  for a resistor with resistance  $R$  give

$$\begin{aligned} \text{Node 1: } & \frac{1}{R_1}(u_1 - U) + \frac{1}{R}u_1 + \frac{1}{R_2}(u_1 - u_2) = 0, \\ \text{Node } j: & \frac{1}{R_{j-1}}(u_j - u_{j-1}) + \frac{1}{R}u_j + \frac{1}{R_j}(u_j - u_{j+1}) = 0, \quad j = 2, \dots, n-1, \\ \text{Node } n: & \frac{1}{R_{n-1}}(u_n - u_{n-1}) + \frac{1}{R}u_n = 0. \end{aligned} \quad (8.3.2.7)$$

These  $n$  equations are equivalent to a linear system of equations for the vector  $\mathbf{u} = [u_j]_{j=1}^n \in \mathbb{R}^n$ , which reads in compact notation

$$\left( \mathbf{A} + \frac{1}{R} \cdot \mathbf{I} \right) \mathbf{u} = \mathbf{b}, \quad (8.3.2.8)$$

$$\mathbf{A} = \begin{bmatrix} \frac{1}{R_1} + \frac{1}{R_2} & -\frac{1}{R_2} & & & & \\ -\frac{1}{R_2} & \frac{1}{R_2} + \frac{1}{R_3} & -\frac{1}{R_3} & & & \\ & -\frac{1}{R_3} & \frac{1}{R_3} + \frac{1}{R_4} & -\frac{1}{R_4} & & \\ & & \ddots & \ddots & \ddots & \\ & & & -\frac{1}{R_{n-2}} & \frac{1}{R_{n-2}} + \frac{1}{R_{n-1}} & -\frac{1}{R_{n-1}} \\ & & & & -\frac{1}{R_{n-1}} & \frac{1}{R_n} + \frac{1}{R_{n-1}} -\frac{1}{R_n} \\ & & & & & \frac{1}{R_n} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \frac{U}{R_1} \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Thus the current problem can be formulated as: find  $x \in \mathbb{R}$ ,  $x := R^{-1}$ , such that

$$F(x) = 0 \quad \text{with} \quad F : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto \mathbf{w}^\top (\mathbf{A} + x\mathbf{I})^{-1} \mathbf{b} - 1 \end{cases}, \quad (8.3.2.9)$$

where  $\mathbf{A} \in \mathbb{R}^{n,n}$  is a symmetric, tridiagonal, diagonally dominant matrix,  $\mathbf{w} \in \mathbb{R}^n$  is a unit vector singling out the node of interest, and  $\mathbf{b}$  takes into account the exciting voltage  $U$ .

In order to apply Newton's method to (8.3.2.9), we have to determine the derivative  $F'(x)$  and so by implicit differentiation [Str09, Sect. 7.8], first rewriting  $(\mathbf{u}(x)) \doteq$  vector of nodal potentials as a function of

$$x = R^{-1})$$

$$F(x) = \mathbf{w}^\top \mathbf{u}(x) - 1 \quad , \quad (\mathbf{A} + x\mathbf{I})\mathbf{u}(x) = \mathbf{b} .$$

Then we differentiate the linear system of equations defining  $\mathbf{u}(x)$  on both sides with respect to  $x$  using the product rule (8.4.1.10)

$$\begin{aligned} & (\mathbf{A} + x\mathbf{I})\mathbf{u}(x) = \mathbf{b} \xrightarrow{\frac{d}{dx}} (\mathbf{A} + x\mathbf{I})\mathbf{u}'(x) + \mathbf{u}(x) = \mathbf{0} . \\ \Rightarrow & \mathbf{u}'(x) = -(\mathbf{A} + x\mathbf{I})^{-1}\mathbf{u}(x) . \end{aligned} \quad (8.3.2.10)$$

$$\Rightarrow F'(x) = \mathbf{w}^\top \mathbf{u}'(x) = -\mathbf{w}^\top (\mathbf{A} + x\mathbf{I})^{-1}\mathbf{u}(x) . \quad (8.3.2.11)$$

Thus, the Newton iteration for (8.3.2.9) reads:

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + F'(x^{(k)})^{-1}F(x^{(k)}) \\ &= \frac{\mathbf{w}^\top \mathbf{u}(x^{(k)}) - 1}{\mathbf{w}^\top \mathbf{z}} , \quad (\mathbf{A} + x^{(k)}\mathbf{I})\mathbf{u}(x^{(k)}) = \mathbf{b} , \\ &\quad (\mathbf{A} + x^{(k)}\mathbf{I})\mathbf{z} = \mathbf{u}(x^{(k)}) . \end{aligned} \quad (8.3.2.12)$$

In each step of the iteration we have to solve two linear systems of equations, which can be done with asymptotic effort  $\mathcal{O}(n)$  in this case, because  $\mathbf{A} + x^{(k)}\mathbf{I}$  is tridiagonal.

Note that in a practical application one must demand  $x > 0$ , in addition, because the solution must provide a meaningful conductance (= inverse resistance.)

Also note that bisection ( $\rightarrow$  8.3.1) is a viable alternative to using Newton's method in this case. □

### 8.3.2.2 Special One-Point Methods

Idea underlying other one-point methods: non-linear local approximation

Useful, if *a priori knowledge* about the structure of  $F$  (e.g. about  $F$  being a rational function, see below) is available. This is often the case, because many problems of 1D zero finding are posed for functions given in analytic form with a few parameters.

Prerequisite: Smoothness of  $F$ :  $F \in C^m(I)$  for some  $m > 1$

**EXAMPLE 8.3.2.13 (Halley's iteration  $\rightarrow$  [Han02, Sect. 18.3])** This example demonstrates that non-polynomial model functions can offer excellent approximation of  $F$ . In this example the model function is chosen as a quotient of two linear function, that is, from the simplest class of true rational functions.

Of course, that this function provides a good model function is merely “a matter of luck”, unless you have some more information about  $F$ . Such information might be available from the application context.

Given  $x^{(k)} \in I$ , next iterate := zero of model function:  $h(x^{(k+1)}) = 0$ , where

$$h(x) := \frac{a}{x+b} + c \quad (\text{rational function}) \text{ such that } F^{(j)}(x^{(k)}) = h^{(j)}(x^{(k)}) , \quad j = 0, 1, 2 .$$



$$\frac{a}{x^{(k)} + b} + c = F(x^{(k)}) , \quad -\frac{a}{(x^{(k)} + b)^2} = F'(x^{(k)}) , \quad \frac{2a}{(x^{(k)} + b)^3} = F''(x^{(k)}) .$$



$$x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})} \cdot \frac{1}{1 - \frac{1}{2} \frac{F(x^{(k)}) F''(x^{(k)})}{F'(x^{(k)})^2}}.$$

Halley's iteration for  $F(x) = \frac{1}{(x+1)^2} + \frac{1}{(x+0.1)^2} - 1$ ,  $x > 0$ : and  $x^{(0)} = 0$

$k$	$x^{(k)}$	$F(x^{(k)})$	$x^{(k)} - x^{(k-1)}$	$x^{(k)} - x^*$
1	0.19865959351191	10.90706835180178	-0.19865959351191	-0.84754290138257
2	0.69096314049024	0.94813655914799	-0.49230354697833	-0.35523935440424
3	1.02335017694603	0.03670912956750	-0.33238703645579	-0.02285231794846
4	1.04604398836483	0.00024757037430	-0.02269381141880	-0.00015850652965
5	1.04620248685303	0.00000001255745	-0.00015849848821	-0.00000000804145

Compare with Newton method (8.3.2.1) for the same problem:

$k$	$x^{(k)}$	$F(x^{(k)})$	$x^{(k)} - x^{(k-1)}$	$x^{(k)} - x^*$
1	0.04995004995005	44.38117504792020	-0.04995004995005	-0.99625244494443
2	0.12455117953073	19.62288236082625	-0.07460112958068	-0.92165131536375
3	0.23476467495811	8.57909346342925	-0.11021349542738	-0.81143781993637
4	0.39254785728080	3.63763326452917	-0.15778318232269	-0.65365463761368
5	0.60067545233191	1.42717892023773	-0.20812759505112	-0.44552704256257
6	0.82714994286833	0.46286007749125	-0.22647449053641	-0.21905255202615
7	0.99028203077844	0.09369191826377	-0.16313208791011	-0.05592046411604
8	1.04242438221432	0.00592723560279	-0.05214235143588	-0.00377811268016
9	1.04618505691071	0.00002723158211	-0.00376067469639	-0.00001743798377
10	1.04620249452271	0.00000000058056	-0.00001743761199	-0.00000000037178

Note that Halley's iteration is superior in this case, since  $F$  is a rational function.

! Newton method converges more slowly, but also needs less effort per step (→ Section 8.3.3)

In the previous example Newton's method performed rather poorly. Often its convergence can be boosted by converting the non-linear equation to an equivalent one (that is, one with the same solutions) for another function  $g$ , which is “closer to a linear function”:

Assume  $F \approx \hat{F}$ , where  $\hat{F}$  is invertible with an inverse  $\hat{F}^{-1}$  that can be evaluated with little effort.

$$\Rightarrow g(x) := \hat{F}^{-1}(F(x)) \approx x.$$

Then apply Newton's method to  $g(x)$ , using the formula for the derivative of the inverse of a function

$$\frac{d}{dy}(\hat{F}^{-1})(y) = \frac{1}{\hat{F}'(\hat{F}^{-1}(y))} \Rightarrow g'(x) = \frac{1}{\hat{F}'(g(x))} \cdot F'(x).$$

**EXAMPLE 8.3.2.14 (Adapted Newton method)** As in Ex. 8.3.2.13 we consider

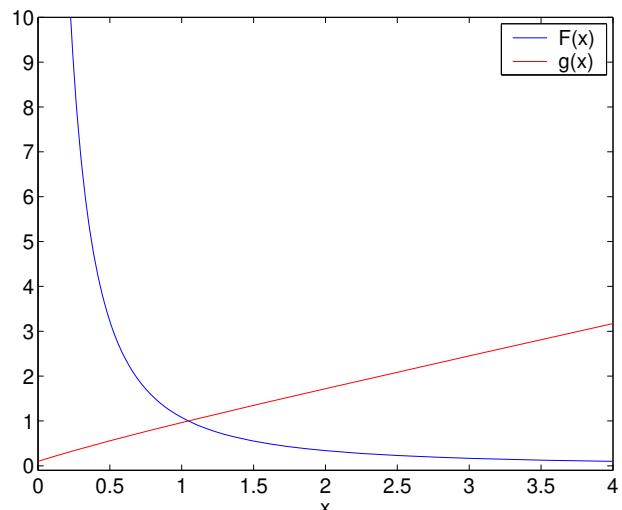
$$F(x) = \frac{1}{(x+1)^2} + \frac{1}{(x+0.1)^2} - 1, \quad x > 0,$$

and try to find its zeros.

Observation:

$$F(x) + 1 \approx 2x^{-2} \text{ for } x \gg 1$$

and so  $g(x) := \frac{1}{\sqrt{F(x)+1}}$  is “almost” linear for  $x \gg 1$ .



Idea: instead of  $F(x) \stackrel{!}{=} 0$  tackle  $g(x) \stackrel{!}{=} 1$  with Newton's method (8.3.2.1).

$$\begin{aligned} x^{(k+1)} &= x^{(k)} - \frac{g(x^{(k)}) - 1}{g'(x^{(k)})} = x^{(k)} + \left( \frac{1}{\sqrt{F(x^{(k)}) + 1}} - 1 \right) \frac{2(F(x^{(k)}) + 1)^{3/2}}{F'(x^{(k)})} \\ &= x^{(k)} + \frac{2(F(x^{(k)}) + 1)(1 - \sqrt{F(x^{(k)}) + 1})}{F'(x^{(k)})}. \end{aligned}$$

Convergence recorded for  $x^{(0)} = 0$ :

$k$	$x^{(k)}$	$F(x^{(k)})$	$x^{(k)} - x^{(k-1)}$	$x^{(k)} - x^*$
1	0.91312431341979	0.24747993091128	0.91312431341979	-0.13307818147469
2	1.04517022155323	0.00161402574513	0.13204590813344	-0.00103227334125
3	1.04620244004116	0.00000008565847	0.00103221848793	-0.00000005485332
4	1.04620249489448	0.0000000000000000	0.00000005485332	-0.0000000000000000

□

For zero finding there is wealth of iterative methods that offer higher order of convergence. One class is discussed next.

**§8.3.2.15 (Modified Newton methods)** Taking the cue from the iteration function of Newton's method (8.3.2.1), we extend it by introducing an extra function  $H$ :

► new fixed point iteration :  $\Phi(x) = x - \frac{F(x)}{F'(x)}H(x)$  with “proper”  $H : I \mapsto \mathbb{R}$ .

Still, every zero of  $F$  is a fixed point of this  $\Phi$ , that is, the fixed point iteration is still consistent ( $\rightarrow$  Def. 8.2.1.1).

Aim: find  $H$  such that the method is of  $p$ -th order. The main tool is Lemma 8.2.2.15, which tells us that we have to ensure  $\Phi^{(\ell)}(x^*) = 0, 1 \leq \ell \leq p-1$ , guarantees local convergence of order  $p$ .

Assume:  $F$  smooth “enough” and  $\exists x^* \in I : F(x^*) = 0, F'(x^*) \neq 0$ . Then we can compute the derivatives of  $\Phi$  appealing to the product rule and quotient rule for derivatives.

$$\Phi = x - uH, \quad \Phi' = 1 - u'H - uH', \quad \Phi'' = -u''H - 2u'H - uH'',$$

$$\text{with } u = \frac{F}{F'} \Rightarrow u' = 1 - \frac{FF''}{(F')^2}, \quad u'' = -\frac{F''}{F'} + 2\frac{F(F'')^2}{(F')^3} - \frac{FF'''}{(F')^2}.$$

$$F(x^*) = 0 \rightarrow u(x^*) = 0, u'(x^*) = 1, u''(x^*) = -\frac{F''(x^*)}{F'(x^*)}.$$

$$\blacktriangleright \Phi'(x^*) = 1 - H(x^*) , \quad \Phi''(x^*) = \frac{F''(x^*)}{F'(x^*)} H(x^*) - 2H'(x^*) . \quad (8.3.2.16)$$

Lemma 8.2.2.15  $\rightarrow$  **Necessary** conditions for local convergence of order  $p$ :

$$p = 2 \text{ (quadratrical convergence): } H(x^*) = 1,$$

$$p = 3 \text{ (cubic convergence): } H(x^*) = 1 \wedge H'(x^*) = \frac{1}{2} \frac{F''(x^*)}{F'(x^*)}.$$

Trial expression:  $H(x) = G(1 - u'(x))$  with “appropriate”  $G$

$$\blacktriangleright \text{fixed point iteration } x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})} G\left(\frac{F(x^{(k)}) F''(x^{(k)})}{(F'(x^{(k)}))^2}\right). \quad (8.3.2.17)$$

### Lemma 8.3.2.18. Cubic convergence of modified Newton methods

If  $F \in C^2(I)$ ,  $F(x^*) = 0$ ,  $F'(x^*) \neq 0$ ,  $G \in C^2(U)$  in a neighbourhood  $U$  of  $0$ ,  $G(0) = 1$ ,  $G'(0) = \frac{1}{2}$ , then the fixed point iteration (8.3.2.17) converge locally cubically to  $x^*$ .

*Proof.* We apply Lemma 8.2.2.15, which tells us that both derivatives from (8.3.2.16) have to vanish. Using the definition of  $H$  we find.

$$H(x^*) = G(0), \quad H'(x^*) = -G'(0)u''(x^*) = G'(0)\frac{F''(x^*)}{F'(x^*)}.$$

Plugging these expressions into (8.3.2.16) finishes the proof.  $\square$

### EXPERIMENT 8.3.2.19 (Application of modified Newton methods)

- $G(t) = \frac{1}{1 - \frac{1}{2}t} \rightarrow$  Halley's iteration ( $\rightarrow$  Ex. 8.3.2.13)
- $G(t) = \frac{2}{1 + \sqrt{1 - 2t}} \rightarrow$  Euler's iteration
- $G(t) = 1 + \frac{1}{2}t \rightarrow$  quadratic inverse interpolation

$k$	$e^{(k)} := x^{(k)} - x^*$		
	Halley	Euler	Quad. Inv.
1	2.81548211105635	3.57571385244736	2.03843730027891
2	1.37597082614957	2.76924150041340	1.02137913293045
3	0.34002908011728	1.95675490333756	0.28835890388161
4	0.00951600547085	1.25252187565405	0.01497518178983
5	0.00000024995484	0.51609312477451	0.00000315361454
6		0.14709716035310	
7		0.00109463314926	
8		0.00000000107549	

### 8.3.2.3 Multi-Point Methods



*Supplementary literature.* The secant method is presented in [Han02, Sect. 18.2], [DR08, Sect. 5.5.3], [AG11, Sect. 3.4].

#### Construction of multi-point iterations in 1D



Idea:

Replace  $F$  with **interpolating polynomial**  
producing interpolatory model function methods

#### §8.3.2.21 (The secant method)

The **secant method** is the simplest representative of model function multi-point methods:

The figure illustrates the geometric idea underlying this 2-point method for zero finding:

$x^{(k+1)}$  is the zero of the secant (red line  $\triangleright$ ) connecting the points  $(x^{(k-1)}, F(x^{(k-1)}))$  and  $(x^k, F(x^k))$  on the graph of  $F$ .

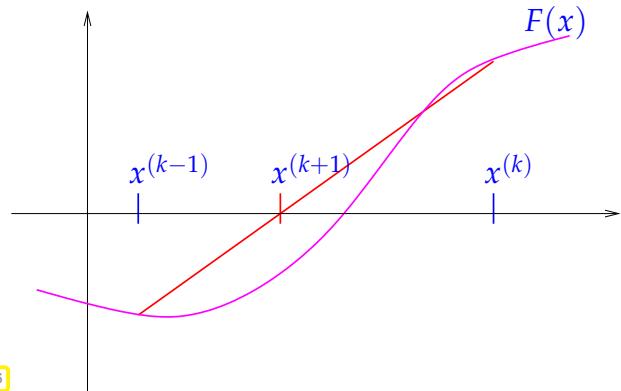


Fig. 266

The secant line is the graph of the function

$$s(x) = F(x^{(k)}) + \frac{F(x^{(k)}) - F(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}(x - x^{(k)}), \quad (8.3.2.22)$$

$$\Rightarrow x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})(x^{(k)} - x^{(k-1)})}{F(x^{(k)}) - F(x^{(k-1)})}. \quad (8.3.2.23)$$

The following C++ code snippet demonstrates the implementation of the abstract secant method for finding the zeros of a function passed through the functor  $F$ .

#### C++ code 8.3.2.24: Secant method for 1D non-linear equation

```

2 // Secant method for solving  $F(x) = 0$  for  $F : D \subset \mathbb{R} \rightarrow \mathbb{R}$ ,
3 // initial guesses  $x_0, x_1$ ,
4 // tolerances  $atol$  (absolute),  $rtol$  (relative)
5 template <typename Func>
6 double secant(double x0, double x1, Func &&F, double rtol, double atol,
7                 unsigned int maxIt) {
8     double fo = F(x0);
9     for (unsigned int i = 0; i < maxIt; ++i) {
10        double fn = F(x1);
11        double s = fn * (x1 - x0) / (fn - fo); // secant correction
12        x0 = x1;
13        x1 = x1 - s;
14        // correction based termination (relative and absolute)
15        if (abs(s) < max(atol, rtol * min(abs(x0), abs(x1))))
16            return x1;
17        fo = fn;
18    }
19    return x1;
20 }
```

This code demonstrates several important features of the secant method:

- Only *one* function evaluation per step
- **no derivatives required!**
- 2-point method: two initial guesses needed

Remember:  $F(x)$  may only be available as output of a (complicated) procedure. In this case it is difficult to find a procedure that evaluates  $F'(x)$ . Thus the significance of methods that do not involve evaluations of derivatives.

**EXPERIMENT 8.3.2.25 (Convergence of secant method)** We empirically examine the convergence of the secant method from Code 8.3.2.24 for a model problem.

$F(x) = xe^x - 1$ , using secant method with initial guesses  $x^{(0)} = 0, x^{(1)} = 5$ .

$k$	$x^{(k)}$	$F(x^{(k)})$	$e^{(k)} := x^{(k)} - x^*$	$\frac{\log  e^{(k+1)}  - \log  e^{(k)} }{\log  e^{(k)}  - \log  e^{(k-1)} }$
2	0.00673794699909	-0.99321649977589	-0.56040534341070	
3	0.01342122983571	-0.98639742654892	-0.55372206057408	24.43308649757745
4	0.98017620833821	1.61209684919288	0.41303291792843	2.70802321457994
5	0.38040476787948	-0.44351476841567	-0.18673852253030	1.48753625853887
6	0.50981028847430	-0.15117846201565	-0.05733300193548	1.51452723840131
7	0.57673091089295	0.02670169957932	0.00958762048317	1.70075240166256
8	0.56668541543431	-0.00126473620459	-0.00045787497547	1.59458505614449
9	0.56713970649585	-0.00000990312376	-0.00000358391394	1.62641838319117
10	0.56714329175406	0.00000000371452	0.00000000134427	
11	0.56714329040978	-0.0000000000000001	-0.0000000000000000	

Rem. 8.1.1.12: the rightmost column of the table provides an estimate for the order of convergence → Def. 8.1.1.10. For further explanations see Rem. 8.1.1.12.

A startling observation: the method seems to have a *fractional* (!) order of convergence, see Def. 8.1.1.10.

**Remark 8.3.2.26 (Fractional order of convergence of secant method)** Indeed, a fractional order of convergence can be proved for the secant method, see [Han02, Sect. 18.2]. Here we give an *asymptotic* argument that holds, if the iterates are already very close to the zero  $x^*$  of  $F$ .

We can write the secant method in the form (8.1.0.4)

$$x^{(k+1)} = \Phi(x^{(k)}, x^{(k-1)}) \quad \text{with} \quad \Phi(x, y) = \Phi(x, y) := x - \frac{F(x)(x - y)}{F(x) - F(y)}. \quad (8.3.2.27)$$

Using  $\Phi$  we find a recursion for the iteration error  $e^{(k)} := x^{(k)} - x^*$ :

$$e^{(k+1)} = \Phi(x^* + e^{(k)}, x^* + e^{(k-1)}) - x^*. \quad (8.3.2.28)$$

Thanks to the asymptotic perspective we may assume that  $|e^{(k)}|, |e^{(k-1)}| \ll 1$  so that we can rely on two-dimensional Taylor expansion around  $(x^*, x^{(*)})$ , cf. [Str09, Satz 7.5.2]:

$$\begin{aligned}\Phi(x^* + h, x^* + k) &= \Phi(x^*, x^*) + \frac{\partial\Phi}{\partial x}(x^*, x^*)h + \frac{\partial\Phi}{\partial y}(x^*, x^*)k + \\ &\quad \frac{1}{2} \frac{\partial^2\Phi}{\partial x^2}(x^*, x^*)h^2 \frac{\partial^2\Phi}{\partial x\partial y}(x^*, x^*)hk \frac{1}{2} \frac{\partial^2\Phi}{\partial y^2}(x^*, x^*)k^2 + R(x^*, h, k), \\ \text{with } |R| &\leq C(h^3 + h^2k + hk^2 + k^3).\end{aligned}\tag{8.3.2.29}$$

Computations invoking the quotient rule and product rule and using  $F(x^*) = 0$  show

$$\Phi(x^*, x^*) = x^*, \quad \frac{\partial\Phi}{\partial x}(x^*, x^*) = \frac{\partial\Phi}{\partial y}(x^*, x^*) = \frac{1}{2} \frac{\partial^2\Phi}{\partial x^2}(x^*, x^*) = \frac{1}{2} \frac{\partial^2\Phi}{\partial y^2}(x^*, x^*) = 0.$$

We may also use MAPLE to find the Taylor expansion (assuming  $F$  sufficiently smooth):

```
> Phi := (x, y) -> x-F(x)*(x-y)/(F(x)-F(y));
> F(s) := 0;
> e2 = normal(mtaylor(Phi(s+e1, s+e0)-s, [e0, e1], 4));
```

➤ truncated error propagation formula (products of three or more error terms ignored)

$$e^{(k+1)} \doteq \frac{1}{2} \frac{F''(x^*)}{F'(x^*)} e^{(k)} e^{(k-1)} = C e^{(k)} e^{(k-1)}. \tag{8.3.2.30}$$

How can we deduce the order of converge from this recursion formula? We try  $e^{(k)} = K(e^{(k-1)})^p$  inspired by the estimate in Def. 8.1.1.10:

$$\begin{aligned}\Rightarrow e^{(k+1)} &= K^{p+1} (e^{(k-1)})^{p^2} \\ \Rightarrow (e^{(k-1)})^{p^2-p-1} &= K^{-p} C \Rightarrow p^2 - p - 1 = 0 \Rightarrow p = \frac{1}{2}(1 \pm \sqrt{5}).\end{aligned}$$

As  $e^{(k)} \rightarrow 0$  for  $k \rightarrow \infty$  we get the order of convergence  $p = \frac{1}{2}(1 + \sqrt{5}) \approx 1.62$  (see Exp. 8.3.2.25 !)

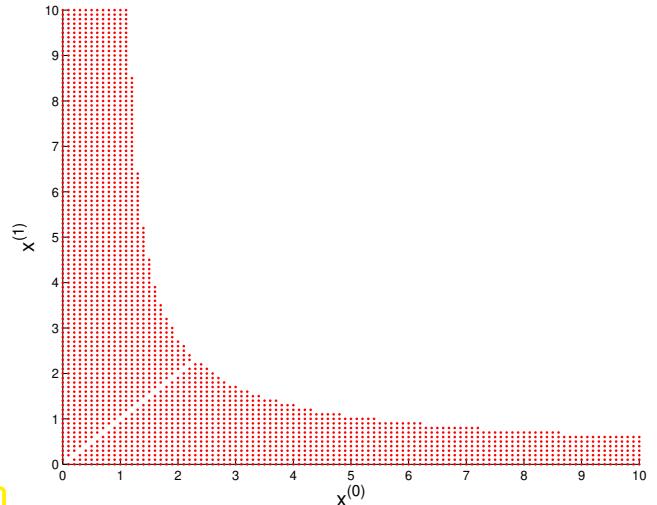
### EXAMPLE 8.3.2.31 (local convergence of the secant method)

Model problem: find zero of

$$F(x) = \arctan(x)$$

•  $\hat{=}$  secant method converges for a pair  $(x^{(0)}, x^{(1)}) \in \mathbb{R}_+^2$  of initial guesses.  $\triangleright$

We observe that the secant method will converge only for initial guesses sufficiently close to  $0$  = local convergence  $\rightarrow$  Def. 8.1.0.8



**§8.3.2.32 (Inverse interpolation)** Another class of multi-point methods: *inverse interpolation*

Assume:

$F : I \subset \mathbb{R} \mapsto \mathbb{R}$  one-to-one (monotone)

$$F(x^*) = 0 \Rightarrow F^{-1}(0) = x^*.$$

- ▶ Interpolate  $F^{-1}$  by polynomial  $p$  of degree  $m - 1$  determined by
- ▶ New approximate zero  $x^{(k+1)} := p(0)$

The graph of  $F^{-1}$  can be obtained by reflecting the graph of  $F$  at the angular bisector.

▷

$$F(x^*) = 0 \Leftrightarrow F^{-1}(0) = x^*$$

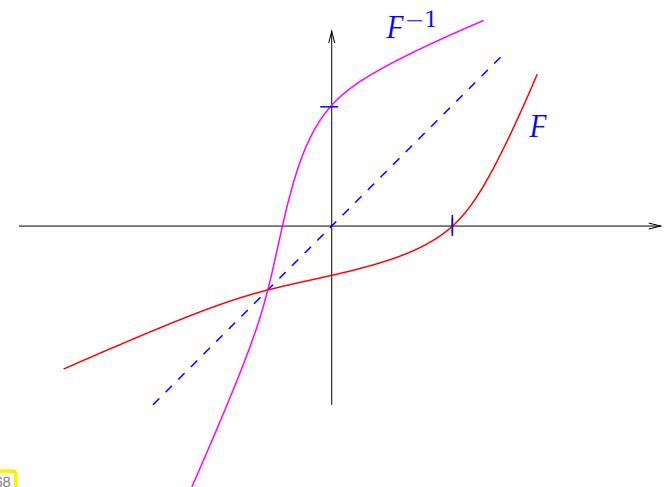


Fig. 268

Case  $m = 2$  (2-point method)

➤ secant method

The interpolation polynomial is a line. In this case we do not get a new method, because the inverse function of a linear function (polynomial of degree 1) is again a polynomial of degree 1.

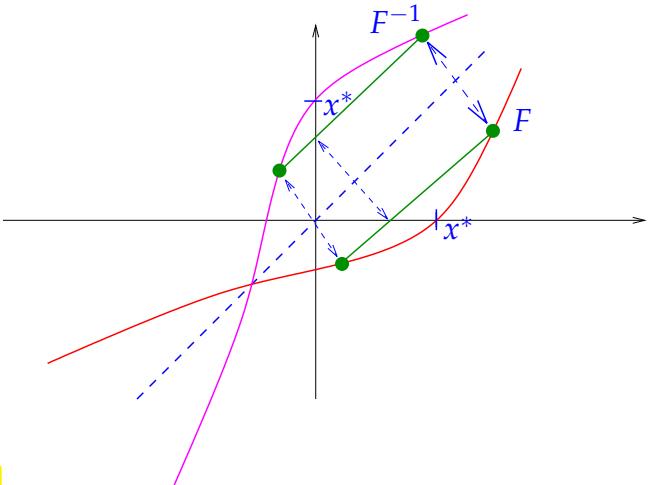


Fig. 269

Case  $m = 3$ : quadratic inverse interpolation, a 3-point method, see [Mol04, Sect. 4.5]

We interpolate the points  $(F(x^{(k)}), x^{(k)}), (F(x^{(k-1)}), x^{(k-1)}), (F(x^{(k-2)}), x^{(k-2)})$  with a parabola (polynomial of degree 2). Note the importance of monotonicity of  $F$ , which ensures that  $F(x^{(k)}), F(x^{(k-1)}), F(x^{(k-2)})$  are mutually different.

MAPLE code:  $p := x \rightarrow a*x^2+b*x+c;$

```
solve({p(f0)=x0, p(f1)=x1, p(f2)=x2}, {a, b, c});  
assign(%); p(0);
```

$$\Rightarrow x^{(k+1)} = \frac{F_0^2(F_1 x_2 - F_2 x_1) + F_1^2(F_2 x_0 - F_0 x_2) + F_2^2(F_0 x_1 - F_1 x_0)}{F_0^2(F_1 - F_2) + F_1^2(F_2 - F_0) + F_2^2(F_0 - F_1)}.$$

$$(F_0 := F(x^{(k-2)}), F_1 := F(x^{(k-1)}), F_2 := F(x^{(k)}), x_0 := x^{(k-2)}, x_1 := x^{(k-1)}, x_2 := x^{(k)})$$

**EXPERIMENT 8.3.2.33 (Convergence of quadratic inverse interpolation)** We test the method for the model problem/initial guesses  $F(x) = xe^x - 1$ ,  $x^{(0)} = 0$ ,  $x^{(1)} = 2.5$ ,  $x^{(2)} = 5$ .

$k$	$x^{(k)}$	$F(x^{(k)})$	$e^{(k)} := x^{(k)} - x^*$	$\frac{\log  e^{(k+1)}  - \log  e^{(k)} }{\log  e^{(k)}  - \log  e^{(k-1)} }$
3	0.08520390058175	-0.90721814294134	-0.48193938982803	
4	0.16009252622586	-0.81211229637354	-0.40705076418392	3.33791154378839
5	0.79879381816390	0.77560534067946	0.23165052775411	2.28740488912208
6	0.63094636752843	0.18579323999999	0.06380307711864	1.82494667289715
7	0.56107750991028	-0.01667806436181	-0.00606578049951	1.87323264214217
8	0.56706941033107	-0.00020413476766	-0.00007388007872	1.79832936980454
9	0.56714331707092	0.00000007367067	0.00000002666114	1.84841261527097
10	0.56714329040980	0.0000000000000003	0.0000000000000001	

Also in this case the numerical experiment hints at a fractional rate of convergence  $p \approx 1.8$ , as in the case of the secant method, see Rem. 8.3.2.26.  $\square$

### 8.3.3 Asymptotic Efficiency of Iterative Methods for Zero Finding

Efficiency is measured by forming the ratio of gain and the effort required to achieve it. For iterative methods for solving  $F(\mathbf{x}) = 0$ ,  $F : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ , this means the following:

Efficiency of an iterative method (for solving $F(\mathbf{x}) = 0$ )	$\leftrightarrow$	computational effort to reach prescribed number of significant digits in result.
---	-------------------	---

Ingredient ①:  $W \triangleq$  computational effort per step

$$(e.g., \quad W \approx \frac{\#\{\text{evaluations of } F\}}{\text{step}} + n \cdot \frac{\#\{\text{evaluations of } F'\}}{\text{step}} + \dots)$$

Ingredient ②: number of steps  $k = k(\rho)$  to achieve relative reduction of error (= gain)

$$\|e^{(k)}\| \leq \rho \|e^{(0)}\|, \quad \rho > 0 \text{ prescribed.} \quad (8.3.3.1)$$

Let us consider an iterative method of order  $p \geq 1$  ( $\rightarrow$  Def. 8.1.1.10). Its error recursion can be converted into expressions (8.3.3.2) and (8.3.3.3) that relate the error norm  $\|e^{(k)}\|$  to  $\|e^{(0)}\|$  and lead to quantitative bounds for the number of steps to achieve (8.3.3.1):

$$\exists C > 0: \|e^{(k)}\| \leq C \|e^{(k-1)}\|^p \quad \forall k \geq 1 \quad (C < 1 \text{ for } p = 1).$$

Assuming  $C \|e^{(0)}\|^{p-1} < 1$  (guarantees convergence!), we find the following minimum number of steps to achieve (8.3.3.1) for sure:

$$p = 1: \quad \|e^{(k)}\| \stackrel{!}{\leq} C^k \|e^{(0)}\| \quad \text{requires} \quad k \geq \frac{\log \rho}{\log C}, \quad (8.3.3.2)$$

$$p > 1: \quad \|e^{(k)}\| \stackrel{!}{\leq} C^{\frac{p^k - 1}{p-1}} \|e^{(0)}\|^{p^k} \quad \text{requires} \quad p^k \geq 1 + \frac{\log \rho}{\log C/p-1 + \log(\|e^{(0)}\|)} \\ \Rightarrow k \geq \log(1 + \frac{\log \rho}{\log L_0}) / \log p, \quad (8.3.3.3)$$

$$L_0 := C^{1/p-1} \|e^{(0)}\| < 1.$$

Now we adopt an asymptotic perspective and ask for a large reduction of the error, that is  $\rho \ll 1$ .

If  $\rho \ll 1$ , then  $\log(1 + \frac{\log \rho}{\log L_0}) \approx \log |\log \rho| - \log |\log L_0| \approx \log |\log \rho|$ . This simplification will be made in the context of asymptotic considerations  $\rho \rightarrow 0$  below.

Notice:

$$|\log \rho| \leftrightarrow \text{Gain in no. of significant digits of } x^{(k)}$$

Measure for efficiency:

$$\text{Efficiency} := \frac{\text{no. of digits gained}}{\text{total work required}} = \frac{|\log \rho|}{k(\rho) \cdot W}$$

(8.3.3.4)

► **asymptotic efficiency for  $\rho \ll 1$**   $\rightarrow$   $|\log \rho| \rightarrow \infty$ ):

$$\text{Efficiency}_{|\rho \ll 1} = \begin{cases} -\frac{\log C}{W} & , \text{if } p = 1 , \\ \frac{\log p}{W} \cdot \frac{|\log \rho|}{\log(|\log \rho|)} & , \text{if } p > 1 . \end{cases} \quad (8.3.3.5)$$

We conclude that

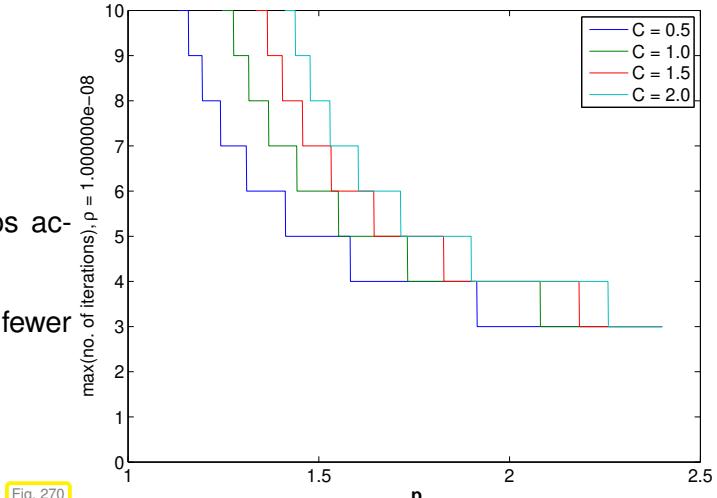
- when requiring high accuracy, linearly convergent iterations should not be used, because their efficiency does not increase for  $\rho \rightarrow 0$ ,
- for method of order  $p > 1$ , the factor  $\frac{\log p}{W}$  offers a gauge for efficiency.

### EXAMPLE 8.3.3.6 (Efficiency of iterative methods)

We choose  $\|e^{(0)}\| = 0.1, \rho = 10^{-8}$ .

The plot displays the number of iteration steps according to (8.3.3.3).

Higher-order method require substantially fewer steps compared to low-order methods.



We compare

Newton's method  $\leftrightarrow$  secant method

in terms of number of steps required for a prescribed guaranteed error reduction, assuming  $C = 1$  in both cases and for  $\|e^{(0)}\| = 0.1$ .

Newton's method requires only marginally fewer steps than the secant method.

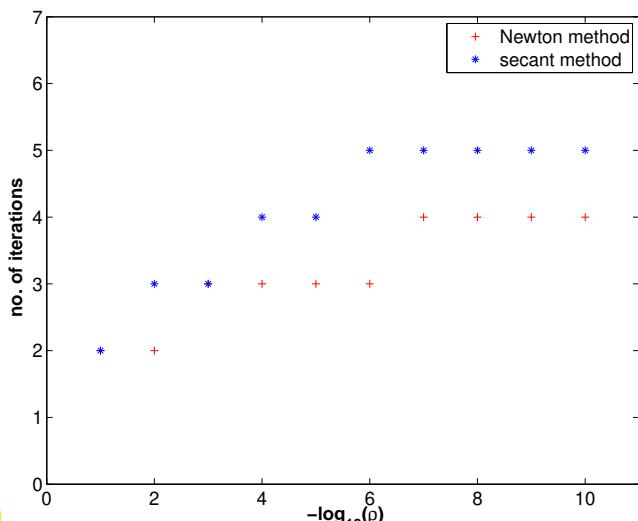


Fig. 271

We draw conclusions from the discussion above and (8.3.3.5):

$$\begin{aligned} W_{\text{Newton}} &= 2W_{\text{secant}}, \\ p_{\text{Newton}} &= 2, p_{\text{secant}} = 1.62 \end{aligned} \Rightarrow \frac{\log p_{\text{Newton}}}{W_{\text{Newton}}} : \frac{\log p_{\text{secant}}}{W_{\text{secant}}} = 0.71 .$$

We set the effort for a step of Newton's method to twice that for a step of the secant method from Code 8.3.2.24, because we need an addition evaluation of  $F'$  in Newton's method.

➤ secant method is more efficient than Newton's method! □

## 8.4 Newton's Method in $\mathbb{R}^n$

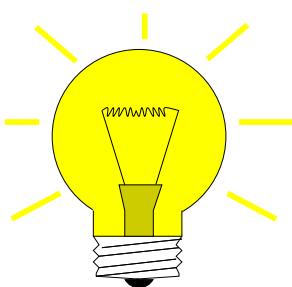
We consider a non-linear system of  $n$  equations with  $n$  unknowns,  $n \in \mathbb{N}$ :

$$\text{for } F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n \text{ find } \mathbf{x}^* \in D: F(\mathbf{x}^*) = 0.$$

We assume:

$F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$  is *continuously differentiable*.

### 8.4.1 The Newton Iteration



Idea ( $\rightarrow$  Section 8.3.2.1):

Given  $\mathbf{x}^{(k)} \in D \succ \mathbf{x}^{(k+1)}$  as zero of affine linear model function

$$F(\mathbf{x}) \approx \tilde{F}_k(\mathbf{x}) := F(\mathbf{x}^{(k)}) + D F(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}) ,$$

$$D F(\mathbf{x}) \in \mathbb{R}^{n,n} = \text{Jacobian}, D F(\mathbf{x}) := \left[ \frac{\partial F_j}{\partial x_k}(\mathbf{x}) \right]_{j,k=1}^n .$$

➤ **Newton iteration:** (generalizes (8.3.2.1) to  $n > 1$ )

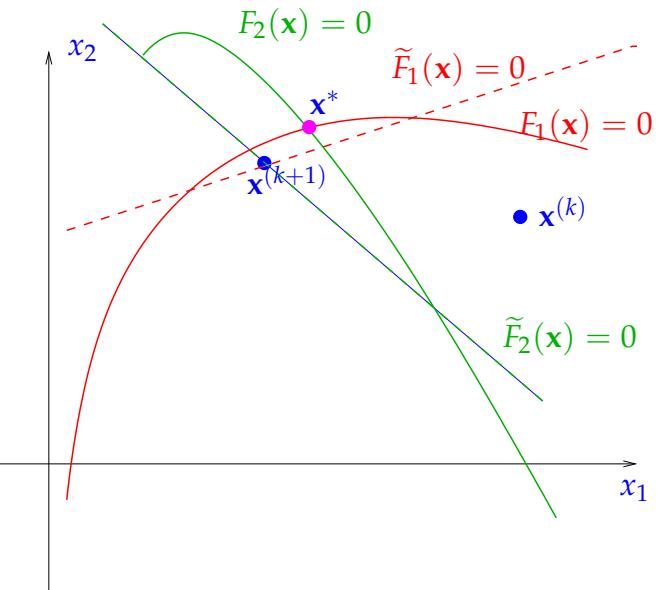
$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - D F(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)}) , \quad [\text{if } D F(\mathbf{x}^{(k)}) \text{ regular}] \quad (8.4.1.1)$$

Terminology:  $-D F(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)})$  is called the **Newton correction**

Illustration of idea of Newton's method for  $n = 2$ : ▷

Sought: intersection point  $\mathbf{x}^*$  of the curves  $F_1(\mathbf{x}) = 0$  and  $F_2(\mathbf{x}) = 0$ .

Idea:  $\mathbf{x}^{(k+1)}$  = the intersection of two straight lines (= zero sets of the components of the model function, cf. Ex. 2.2.2.10) that are approximations of the original curves



The following C++ implementation of the skeleton for Newton's method uses a *correction based* a posteriori termination criterion for the Newton iteration. It stops the iteration if the *relative* size of the Newton correction drops below the prescribed *relative tolerance* `rtol`. If  $\mathbf{x}^* \approx 0$  also the *absolute* size of the Newton correction has to be tested against an *absolute tolerance* `atol` in order to avoid non-termination despite convergence of the iteration,

#### C++11-code 8.4.1.2: Newton's method in C++

```

2  template <typename FuncType, typename JacType, typename VecType>
3  VecType newton(FuncType &&F, JacType &&DFinv, VecType x, const double rtol,
4                  const double atol) {
5      // Note that the vector x passes both the initial guess and also
6      // contains the iterates
7      VecType s(x.size()); // Vector for Newton corrections
8      // Main loop
9      do {
10          s = DFinv(x, F(x)); // compute Newton correction
11          x -= s;             // compute next iterate
12      }
13      // correction based termination (relative and absolute)
14      while ((s.norm() > rtol * x.norm()) && (s.norm() > atol));
15      return x;
16  }
```

☞ Objects of type **FuncType** must feature

```
Eigen::VectorXd operator(const Eigen::VectorXd &x);
```

that evaluates  $F(\mathbf{x})$  ( $Vx \leftrightarrow \mathbf{x}$ ).

☞ Objects of type **JacType** must provide a method

```
Eigen::VectorXd operator (const Eigen::VectorXd &x, const
                           Eigen::VectorXd &f);
```

that computes the Newton correction, that is it returns the solution of a linear system with system matrix  $D F(\mathbf{x})$  ( $\mathbf{x} \leftrightarrow \mathbf{x}$ ) and right hand side  $\mathbf{f} \leftrightarrow \mathbf{f}$ .

☞ The function returns the computed approximate solution of the non-linear system.

The next code demonstrates the invocation of `newton` for a  $2 \times 2$  non-linear system from a code relying on EIGEN. It also demonstrates the use of fixed size eigen matrices and vectors.

#### C++11-code 8.4.1.3: Calling `newton` with EIGEN data types

```

2 void newton2Ddriver(void) {
3     // Function F defined through lambda function
4     auto F = [] (const Eigen::Vector2d &x) {
5         Eigen::Vector2d z;
6         const double x1 = x(0), x2 = x(1);
7         z << x1 * x1 - 2 * x1 - x2 + 1, x1 * x1 + x2 * x2 - 1;
8         return (z);
9     };
10    // Lambda function for the computation of the Newton correction
11    auto DFinv = [] (const Eigen::Vector2d &x, const Eigen::Vector2d &f) {
12        Eigen::Matrix2d J;
13        const double x1 = x(0), x2 = x(1);
14        // Jacobian of F
15        J << 2 * x1 - 2, -1, 2 * x1, 2 * x2;
16        // Solve 2x2 linear system of equations
17        Eigen::Vector2d s = J.lu().solve(f);
18        return (s);
19    };
20    // initial guess
21    Vector2d x0(2., 3.);
22
23    // Invoke Newton's method
24    Vector2d x = newton(F, DFinv, x0, 1E-6, 1E-8);
25    std::cout << "||F(x)|| = " << F(x).norm() << std::endl;
26 }
```



New aspect for  $n \gg 1$  (compared to  $n = 1$ -dimensional case, Section 8.3.2.1):  
 Computation of the Newton correction may be expensive!  
 (because it involves the solution of a LSE, cf. Thm. 2.5.0.2)

**Remark 8.4.1.4 (Affine invariance of Newton method)** An important property of the Newton iteration (8.4.1.1): **affine invariance** → [Deu11, Sect. 1.2.2]

set  $G_{\mathbf{A}}(\mathbf{x}) := \mathbf{A}F(\mathbf{x})$  with regular  $\mathbf{A} \in \mathbb{R}^{n,n}$  so that  $F(\mathbf{x}^*) = 0 \Leftrightarrow G_{\mathbf{A}}(\mathbf{x}^*) = 0$ .



**Affine invariance:** Newton iterations for  $G_{\mathbf{A}}(\mathbf{x}) = 0$  are the same for all regular  $\mathbf{A}$ !

This is a simple computation:

$$\mathbf{D}G(\mathbf{x}) = \mathbf{A}\mathbf{D}F(\mathbf{x}) \Rightarrow \mathbf{D}G(\mathbf{x})^{-1}G(\mathbf{x}) = \mathbf{D}F(\mathbf{x})^{-1}\mathbf{A}^{-1}\mathbf{A}F(\mathbf{x}) = \mathbf{D}F(\mathbf{x})^{-1}F(\mathbf{x}).$$

Use affine invariance as guideline for

- convergence theory for Newton's method: assumptions and results should be affine invariant, too.
- modifying and extending Newton's method: resulting schemes should preserve affine invariance.

In particular, termination criteria for Newton's method should also be affine invariant in the sense that, when applied for  $G_A$  they STOP the iteration at exactly the same step for any choice of  $A$ .  $\square$

The function  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  defining the non-linear system of equations may be given in various formats, as explicit expression or rather implicitly. In most cases,  $D F$  has to be computed *symbolically* in order to obtain concrete formulas for the Newton iteration. We now learn how these symbolic computations can be carried out harnessing advanced techniques of multi-variate calculus.

**§8.4.1.5 (Derivatives and their role in Newton's method)** The reader will probably agree that the derivative of a function  $F : I \subset \mathbb{R} \rightarrow \mathbb{R}$  in  $x \in I$  is a number  $F'(x) \in \mathbb{R}$ , the derivative of a function  $F : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ , in  $x \in D$  a matrix  $D F(x) \in \mathbb{R}^{m,n}$ . However, the nature of a **derivative** in a point is that of a **linear mapping** that approximates  $F$  locally up to second order:

#### Definition 8.4.1.6. Derivative of functions between vector spaces

Let  $V, W$  be finite dimensional vector spaces and  $F : D \subset V \mapsto W$  a sufficiently smooth mapping. The **derivative** (differential)  $D F(x)$  of  $F$  in  $x \in V$  is the *unique* linear mapping  $D F(x) : V \mapsto W$  such that there is a  $\delta > 0$  and a function  $\epsilon : [0, \delta] \rightarrow \mathbb{R}^+$  satisfying  $\lim_{\xi \rightarrow 0} \epsilon(\xi) = 0$  such that

$$\|F(x + h) - F(x) - D F(x)h\| = \epsilon(\|h\|) \quad \forall h \in V, \|h\| < \delta. \quad (8.4.1.7)$$

The vector  $h$  in (8.4.1.7) may be called “direction of differentiation”.

- ☞ Note that  $D F(x)h \in W$  is the vector returned by the linear mapping  $D F(x)$  when applied to  $h \in V$ .
- ☞ In Def. 8.4.1.6  $\|\cdot\|$  can be any norm on  $V$  ( $\rightarrow$  Def. 1.5.5.4).
- ☞ A common shorthand notation for (8.4.1.7) relies on the “little-o” Landau symbol:

$$\|F(x + h) - F(x) - D F(x)h\| = o(h) \quad \text{for } h \rightarrow 0,$$

which designates a remainder term tending to 0 as its arguments tends to 0.

- ☞ Choosing bases of  $V$  and  $W$ ,  $D F(x)$  can be described by the **Jacobian** (matrix) (8.2.2.8), because every linear mapping between finite dimensional vector spaces has a matrix representation after bases have been fixed. Thus, the derivative is usually written as a matrix-valued function on  $D$ .

In the context of the Newton iteration (8.4.1.1) the computation of the Newton correction  $s$  in the  $k+1$ -th step amounts to solving a linear system of equations:

$$s = -D F(x^{(k)})^{-1} F(x^{(k)}) \Leftrightarrow D F(x^{(k)})s = -F(x^{(k)}).$$

Matching this with Def. 8.4.1.6 we see that we need only determine expressions for  $D F(x^{(k)})h$ ,  $h \in V$ , in order to state the LSE yielding the Newton correction. This will become important when applying the “compact” differentiation rules discussed next.  $\square$

**§8.4.1.8 (Differentiation rules → Repetition of basic analysis skills)** Statement of the Newton iteration (8.4.1.1) for  $F : \mathbb{R}^n \mapsto \mathbb{R}^n$  given as analytic expression entails computing the Jacobian  $D F$ . The safe, but tedious way is to use the definition (8.2.2.8) directly and compute the partial derivatives.

To avoid cumbersome component-oriented considerations, it is sometimes useful to know the *rules of multidimensional differentiation*:

Immediate from Def. 8.4.1.6 are the following differentiation rules:

- For  $F : V \mapsto W$  linear, we have  $D F(\mathbf{x}) = F$  for all  $\mathbf{x} \in V$

(For instance, if  $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $F(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{m,n}$ , then  $D F(\mathbf{x}) = \mathbf{A}$  for all  $\mathbf{x} \in \mathbb{R}^n$ .)

- **Chain rule:** For  $F : V \mapsto W$ ,  $G : W \mapsto U$  sufficiently smooth

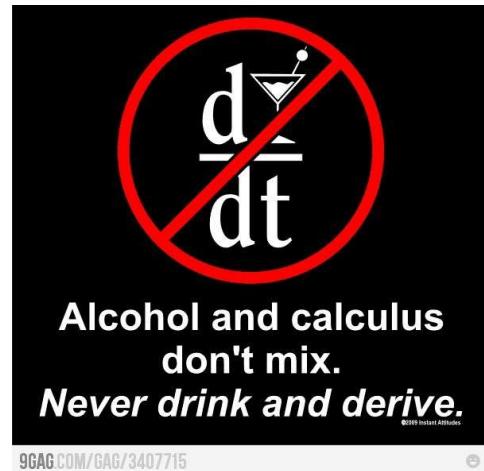
$$D(G \circ F)(\mathbf{x})\mathbf{h} = D G(F(\mathbf{x}))(D F(\mathbf{x}))\mathbf{h}, \quad \mathbf{h} \in V, \mathbf{x} \in D. \quad (8.4.1.9)$$

- **Product rule:**  $F : D \subset V \mapsto W$ ,  $G : D \subset V \mapsto U$  sufficiently smooth,  $b : W \times U \mapsto Z$  bilinear, i.e., linear in each argument:

$$T(\mathbf{x}) = b(F(\mathbf{x}), G(\mathbf{x})) \Rightarrow D T(\mathbf{x})\mathbf{h} = b(D F(\mathbf{x})\mathbf{h}, G(\mathbf{x})) + b(F(\mathbf{x}), D G(\mathbf{x})\mathbf{h}), \quad (8.4.1.10)$$

$$\mathbf{h} \in V, \mathbf{x} \in D.$$

Advice: If you do not feel comfortable with these rules of multidimensional differential calculus, please resort to detailed componentwise/entrywise calculations according to (8.2.2.8) (“pedestrian differentiation”), though they may be tedious.



The first and second derivatives of real-valued functions occur frequently and have special names, see [Str09, Def. 7.3.2] and [Str09, Satz 7.5.3].

#### Definition 8.4.1.11. Gradient and Hessian

For sufficiently smooth  $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}$  the gradient  $\text{grad } F : D \mapsto \mathbb{R}^n$ , and the Hessian (matrix)  $H F(\mathbf{x}) : D \mapsto \mathbb{R}^{n,n}$  are defined as

$$(\text{grad } F(\mathbf{x})^T)\mathbf{h} := D F(\mathbf{x})\mathbf{h}, \quad \mathbf{h}_1^T H F(\mathbf{x})\mathbf{h}_2 := D(D F(\mathbf{x})(\mathbf{h}_1))(\mathbf{h}_2), \quad \mathbf{h}, \mathbf{h}_1, \mathbf{h}_2 \in V.$$

**EXAMPLE 8.4.1.12 (Derivative of a bilinear form)** A simple example:  $\Psi : \mathbb{R}^n \mapsto \mathbb{R}$ ,  $\Psi(\mathbf{x}) := \mathbf{x}^T \mathbf{A} \mathbf{x}$ , with  $\mathbf{A} \in \mathbb{R}^{n,n}$

This is the general matrix representation of a bilinear form on  $\mathbb{R}^n$ .

“High level differentiation”: We apply the product rule (8.4.1.10) with  $F, G = Id$ , which means  $D F(\mathbf{x}) = D G(\mathbf{x}) = I$ , and the bilinear form  $b(\mathbf{x}, \mathbf{y}) := \mathbf{x}^T \mathbf{A} \mathbf{y}$ :

$$\Rightarrow D \Psi(\mathbf{x})\mathbf{h} = \mathbf{h}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{A} \mathbf{h} = \underbrace{(\mathbf{x}^T \mathbf{A}^T + \mathbf{x}^T \mathbf{A})}_{=(\text{grad } \Psi(\mathbf{x}))^T} \mathbf{h},$$

Hence,  $\text{grad } \Psi(\mathbf{x}) = (\mathbf{A} + \mathbf{A}^T)\mathbf{x}$  according to the definition of a gradient ( $\rightarrow$  Def. 8.4.1.11)

“Low level differentiation”: Using the rules of matrix  $\times$  vector multiplication,  $\Psi$  can be written in terms of the

vector components  $x_i, i = 1, \dots, n$ :

$$\Psi(\mathbf{x}) = \sum_{k=1}^n \sum_{j=1}^n (\mathbf{A})_{k,j} x_k x_j = (\mathbf{A})_{i,i} x_i^2 + \sum_{\substack{j=1 \\ j \neq i}}^n (\mathbf{A})_{i,j} x_i x_j + \sum_{\substack{k=1 \\ k \neq i}}^n (\mathbf{A})_{k,i} x_k x_i + \sum_{\substack{k=1 \\ k \neq i}}^n \sum_{\substack{j=1 \\ j \neq i}}^n (\mathbf{A})_{k,j} x_k x_j ,$$

which leads to the partial derivatives

$$\begin{aligned} \frac{\partial \Psi(\mathbf{x})}{\partial x_i}(\mathbf{x}) &= 2(\mathbf{A})_{i,i} x_i + \sum_{\substack{j=1 \\ j \neq i}}^n (\mathbf{A})_{i,j} x_j + \sum_{\substack{k=1 \\ k \neq i}}^n (\mathbf{A})_{k,i} x_k = \sum_{j=1}^n (\mathbf{A})_{i,j} x_j + \sum_{k=1}^n (\mathbf{A})_{k,i} x_k \\ &= (\mathbf{Ax} + \mathbf{A}^\top \mathbf{x})_i . \end{aligned}$$

Of course, both results must agree! □

**EXAMPLE 8.4.1.13 (Derivative of Euclidean norm)** We seek the derivative of the Euclidean norm, that is, of the function  $F(\mathbf{x}) := \|\mathbf{x}\|_2$ ,  $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$  ( $F$  is defined but not differentiable in  $\mathbf{x} = \mathbf{0}$ , just look at the case  $n = 1$ !).

“High level differentiation”: We can write  $F$  as the composition of two functions  $F = G \circ H$  with

$$\begin{aligned} G : \mathbb{R}_+ &\rightarrow \mathbb{R}_+, \quad G(\xi) := \sqrt{\xi}, \\ H : \mathbb{R}^n &\rightarrow \mathbb{R}, \quad H(\mathbf{x}) := \mathbf{x}^\top \mathbf{x} . \end{aligned}$$

Using the rule for the differentiation of bilinear forms from Ex. 8.4.1.12 for the case  $\mathbf{A} = \mathbf{I}$  and basic calculus, we find

$$\begin{aligned} D H(\mathbf{x}) \mathbf{h} &= 2\mathbf{x}^\top \mathbf{h}, \quad \mathbf{x}, \mathbf{h} \in \mathbb{R}^n, \\ D G(\xi) \zeta &= \frac{\zeta}{2\sqrt{\xi}}, \quad \xi > 0, \zeta \in \mathbb{R} . \end{aligned}$$

Finally, the **chain rule** (8.4.1.9) gives

$$\begin{aligned} D F(\mathbf{x}) \mathbf{h} &= D G(H(\mathbf{x}))(D H(\mathbf{x}) \mathbf{h}) = \frac{2\mathbf{x}^\top \mathbf{h}}{2\sqrt{\mathbf{x}^\top \mathbf{x}}} = \frac{\mathbf{x}^\top}{\|\mathbf{x}\|_2} \cdot \mathbf{h} . \tag{8.4.1.14} \\ \text{Def. 8.4.1.11} \Rightarrow \quad \mathbf{grad} F(\mathbf{x}) &= \frac{\mathbf{x}}{\|\mathbf{x}\|_2} . \end{aligned}$$

“Pedestrian differentiation”: For a single component of  $F(\mathbf{x})$  we have

$$\begin{aligned} (F(\mathbf{x}))_i &= \sqrt{x_1^2 + x_2^2 + \cdots + x_i^2 + \cdots + x_n^2}, \\ \Rightarrow \quad \frac{\partial (F(\mathbf{x}))_i}{\partial x_i} &= \frac{1}{2\sqrt{x_1^2 + x_2^2 + \cdots + x_i^2 + \cdots + x_n^2}} \cdot 2x_i \quad [\text{chain rule}] . \\ \Rightarrow \quad \mathbf{grad} F(\mathbf{x}) &= \frac{1}{\|\mathbf{x}\|_2} (x_1, x_2, \dots, x_n)^\top . \end{aligned}$$
□

**§8.4.1.15 (Newton iteration via product rule)** This paragraph explains the use of the general product rule (8.4.1.10) to derive the linear system solved by the Newton correction. It implements the insights from § 8.4.1.5.

We seek solutions of  $F(\mathbf{x}) = \mathbf{0}$  with  $F(\mathbf{x}) := b(G(\mathbf{x}), H(\mathbf{x}))$ , where

- ◆  $V, W$  are some vector spaces (finite- or even infinite-dimensional),
- ◆  $G : D \rightarrow V, H : D \rightarrow W, D \subset \mathbb{R}^n$ , are continuously differentiable in the sense of Def. 8.4.1.6,
- ◆  $b : V \times W \mapsto \mathbb{R}^n$  is bilinear (linear in each argument).

According to the general product rule (8.4.1.10) we have

$$D F(\mathbf{x})\mathbf{h} = b(D G(\mathbf{x})\mathbf{h}, H(\mathbf{x})) + b(G(\mathbf{x}), D H(\mathbf{x})\mathbf{h}), \quad \mathbf{h} \in \mathbb{R}^n. \quad (8.4.1.16)$$

This already defines the linear system of equations to be solved to compute the Newton correction  $\mathbf{s}$

$$b(D G(\mathbf{x}^{(k)})\mathbf{s}, H(\mathbf{x}^{(k)})) + b(G(\mathbf{x}^{(k)}), D H(\mathbf{x}^{(k)})\mathbf{s}) = -b(G(\mathbf{x}^{(k)}), H(\mathbf{x}^{(k)})). \quad (8.4.1.17)$$

Since the left-hand side is linear in  $\mathbf{s}$ , this really represents a square linear system of  $n$  equations. The next example will present a concrete case.  $\square$

**§8.4.1.18 (A “quasi-linear” system of equations)** We call a quasilinear system of equations a non-linear equation of the form

$$\mathbf{A}(\mathbf{x})\mathbf{x} = \mathbf{b} \quad \text{with } \mathbf{b} \in \mathbb{R}^n, \quad (8.4.1.19)$$

where  $\mathbf{A} : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^{n,n}$  is a matrix-valued function. On other words, a quasi-linear system is a “linear system of equations with solution-dependent system matrix”. It incarnates a zero-finding problem for a function  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  with a bilinear structure as introduced in § 8.4.1.15.

For many quasi-linear systems, for which there exist solutions, the fixed point iteration ( $\rightarrow$  Section 8.2)

$$\mathbf{x}^{(k+1)} = \mathbf{A}(\mathbf{x}^{(k)})^{-1}\mathbf{b} \iff \mathbf{A}(\mathbf{x}^{(k)})\mathbf{x}^{(k+1)} = \mathbf{b}, \quad (8.4.1.20)$$

provides a convergent iteration, provided that a good initial guess is available.

We can also reformulate

$$(8.4.1.19) \iff F(\mathbf{x}) = \mathbf{0} \quad \text{with } F(\mathbf{x}) = \mathbf{A}(\mathbf{x})\mathbf{x} - \mathbf{b}.$$

If  $\mathbf{x} \mapsto \mathbf{A}(\mathbf{x})$  is differentiable, the product rule (8.4.1.10) yields

$$D F(\mathbf{x})\mathbf{h} = (D \mathbf{A}(\mathbf{x})\mathbf{h})\mathbf{x} + \mathbf{A}(\mathbf{x})\mathbf{h}, \quad \mathbf{h} \in \mathbb{R}^n. \quad (8.4.1.21)$$

Note that  $D \mathbf{A}(\mathbf{x}^{(k)})$  is a mapping from  $\mathbb{R}^n$  into  $\mathbb{R}^{n,n}$ , which gets  $\mathbf{h}$  as an argument. Then the Newton iteration reads

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{s}, \quad D F(\mathbf{x})\mathbf{s} = (D \mathbf{A}(\mathbf{x}^{(k)})\mathbf{s})\mathbf{x}^{(k)} + \mathbf{A}(\mathbf{x}^{(k)})\mathbf{s} = \mathbf{A}(\mathbf{x}^{(k)})\mathbf{x}^{(k)} - \mathbf{b}. \quad (8.4.1.22)$$

The next example will examine a concrete quasi-linear system of equations.  $\square$

**EXAMPLE 8.4.1.23 (A special quasi-linear system of equations)** We consider the quasi-linear system of equations

$$\mathbf{A}(\mathbf{x})\mathbf{x} = \mathbf{b}, \quad \mathbf{A}(\mathbf{x}) = \begin{bmatrix} \gamma(\mathbf{x}) & 1 & & \\ 1 & \gamma(\mathbf{x}) & 1 & \\ & \ddots & \ddots & \ddots \\ & & \ddots & \ddots \\ & & & 1 & \gamma(\mathbf{x}) & 1 \\ & & & 1 & \gamma(\mathbf{x}) & \\ & & & & 1 & \gamma(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad (8.4.1.24)$$

where  $\gamma(\mathbf{x}) := 3 + \|\mathbf{x}\|_2$  (Euclidean vector norm), the right hand side vector  $\mathbf{b} \in \mathbb{R}^n$  is given and  $\mathbf{x} \in \mathbb{R}^n$  is unknown.

In order to compute the derivative of  $F(\mathbf{x}) := \mathbf{A}(\mathbf{x})\mathbf{x} - \mathbf{b}$  it is advisable to rewrite

$$\mathbf{A}(\mathbf{x})\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{x}\|\mathbf{x}\|_2, \quad \mathbf{T} := \begin{bmatrix} 3 & 1 & & \\ 1 & 3 & 1 & \\ & \ddots & 3 & \ddots & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & 3 & 1 \\ & & & & 1 & 3 \end{bmatrix}.$$

The derivative of the first term is straightforward, because it is linear in  $\mathbf{x}$ , see the discussion following Def. 8.4.1.6.

The “pedestrian” approach to the second term starts with writing it explicitly in components as

$$(\mathbf{x}\|\mathbf{x}\|)_i = x_i \sqrt{x_1^2 + \cdots + x_n^2}, \quad i = 1, \dots, n.$$

Then we can compute the Jacobian according to (8.2.2.8) by taking partial derivatives:

$$\begin{aligned} \frac{\partial}{\partial x_i} (\mathbf{x}\|\mathbf{x}\|)_i &= \sqrt{x_1^2 + \cdots + x_n^2} + x_i \frac{x_i}{\sqrt{x_1^2 + \cdots + x_n^2}}, \\ \frac{\partial}{\partial x_j} (\mathbf{x}\|\mathbf{x}\|)_i &= x_i \frac{x_j}{\sqrt{x_1^2 + \cdots + x_n^2}}, \quad j \neq i. \end{aligned}$$

For the “high level” treatment of the second term  $\mathbf{x} \mapsto \mathbf{x}\|\mathbf{x}\|_2$  we apply the product rule (8.4.1.10), together with (8.4.1.14):

$$\mathbf{D}F(\mathbf{x})\mathbf{h} = \mathbf{T}\mathbf{h} + \|\mathbf{x}\|_2\mathbf{h} + \mathbf{x}\frac{\mathbf{x}^\top \mathbf{h}}{\|\mathbf{x}\|_2} = (\mathbf{A}(\mathbf{x}) + \frac{\mathbf{x}\mathbf{x}^\top}{\|\mathbf{x}\|_2})\mathbf{h}.$$

Thus, in concrete terms the Newton iteration (8.4.1.22) becomes

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left( \mathbf{A}(\mathbf{x}^{(k)}) + \frac{\mathbf{x}^{(k)}(\mathbf{x}^{(k)})^\top}{\|\mathbf{x}^{(k)}\|_2} \right)^{-1} (\mathbf{A}(\mathbf{x}^{(k)})\mathbf{x}^{(k)} - \mathbf{b}).$$

Note that that matrix of the linear system to be solved in each step is a rank-1-modification (2.6.0.16) of the symmetric positive definite tridiagonal matrix  $\mathbf{A}(\mathbf{x}^{(k)})$ , cf. Lemma 2.8.0.12. Thus the Sherman-Morrison-Woodbury formula from Lemma 2.6.0.21 can be used to solve it efficiently.  $\square$

### §8.4.1.25 (Implicit differentiation of bilinear expressions)

- Given are
- ◆ a finite-dimensional vector space  $\mathbf{W}$ ,
  - ◆ a continuously differentiable function  $G : D \subset \mathbb{R}^n \rightarrow \mathbf{V}$  into some vector space  $\mathbf{V}$ ,
  - ◆ a **bilinear** mapping  $b : \mathbf{V} \times \mathbf{W} \rightarrow \mathbf{W}$ .

Let  $F : D \rightarrow \mathbf{W}$  be *implicitly defined* by

$$F(\mathbf{x}) : b(G(\mathbf{x}), F(\mathbf{x})) = \mathbf{b} \in \mathbf{W}. \quad (8.4.1.26)$$

This relationship will provide a valid definition of  $F$  in a neighborhood of  $\mathbf{x}_0 \in W$ , if we assume that there is  $\mathbf{x}_0, \mathbf{z}_0 \in W$  such that  $b(G(\mathbf{x}_0), \mathbf{z}_0) = \mathbf{b}$ , and that the linear mapping  $\mathbf{z} \mapsto b(G(\mathbf{x}_0), \mathbf{z})$  is *invertible*. Then, for  $\mathbf{x}$  close to  $\mathbf{x}_0$ ,  $F(\mathbf{x})$  can be computed by solving a square linear system of equations in  $W$ . In Ex. 8.3.2.6 we already saw an example of an implicitly defined  $F$  for  $W = \mathbb{R}$ .

We want to solve  $F(\mathbf{x}) = \mathbf{0}$  for this implicitly defined  $F$  by means of Newton's method. In order to determine the derivative of  $F$  we resort to **implicit differentiation** [Str09, Sect. 7.8] of the defining equation (8.4.1.26) by means of the general product rule (8.4.1.10). We formally differentiate both sides of (8.4.1.26):

$$b(DG(\mathbf{x})\mathbf{h}, F(\mathbf{x})) + b(G(\mathbf{x}), DF(\mathbf{x})\mathbf{h}) = \mathbf{0} \quad \forall \mathbf{h} \in W, \quad (8.4.1.27)$$

and find that the Newton correction  $\mathbf{s}$  in the  $k+1$ -th Newton step can be computed as follows:

$$\begin{aligned} D F(\mathbf{x}^{(k)})\mathbf{s} &= -F(\mathbf{x}^{(k)}) \Rightarrow b(G(\mathbf{x}^{(k)}), D F(\mathbf{x}^{(k)})\mathbf{s}) = -b(G(\mathbf{x}^{(k)}), F(\mathbf{x}^{(k)})) \\ \stackrel{(8.4.1.27)}{\Rightarrow} b(DG(\mathbf{x}^{(k)})\mathbf{s}, F(\mathbf{x}^{(k)})) &= b(G(\mathbf{x}^{(k)}), F(\mathbf{x}^{(k)})), \end{aligned}$$

which constitutes an  $\dim W \times \dim W$  linear system of equations. The next example discusses a concrete application of implicit differentiation with  $W = \mathbb{R}^{n,n}$ .  $\downarrow$

**EXAMPLE 8.4.1.28 (Derivative of matrix inversion)** We consider matrix inversion as a mapping and (formally) compute its derivative, that is, the derivative of function

$$\text{inv} : \left\{ \begin{array}{ccc} \mathbb{R}_{*}^{n,n} & \rightarrow & \mathbb{R}^{n,n} \\ \mathbf{X} & \mapsto & \mathbf{X}^{-1} \end{array} \right.,$$

where  $\mathbb{R}_{*}^{n,n}$  denotes the (open) set of invertible  $n \times n$ -matrices,  $n \in \mathbb{N}$ .

We apply the technique of **implicit differentiation** from § 8.4.1.25 to the equation

$$\text{inv}(\mathbf{X}) \cdot \mathbf{X} = \mathbf{I}, \quad \mathbf{X} \in \mathbb{R}_{*}^{n,n}. \quad (8.4.1.29)$$

Differentiation on both sides of (8.4.1.29) by means of the product rule (8.4.1.10) yields

$$\begin{aligned} D \text{inv}(\mathbf{X})\mathbf{H} \cdot \mathbf{X} + \text{inv}(\mathbf{X}) \cdot \mathbf{H} &= \mathbf{0}, \quad \mathbf{H} \in \mathbb{R}^{n,n}, \\ \Rightarrow D \text{inv}(\mathbf{X})\mathbf{H} &= -\mathbf{X}^{-1}\mathbf{H}\mathbf{X}^{-1}, \quad \mathbf{H} \in \mathbb{R}^{n,n}. \end{aligned} \quad (8.4.1.30)$$

For  $n = 1$  we get  $D \text{inv}(x)\mathbf{h} = -\frac{\mathbf{h}}{x^2}$ , which recovers the well-known derivative of the function  $x \rightarrow x^{-1}$ .  $\downarrow$

**EXAMPLE 8.4.1.31 (Matrix inversion by means of Newton's method → [Ale12; PS91])** Surprisingly, it is possible to obtain the inverse of a matrix as a the solution of a non-linear system of equations. Thus it can be computed using Newton's method.

Given a regular matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$ , its inverse can be defined as the unique zero of a function:

$$\mathbf{X} = \mathbf{X}^{-1} \iff F(\mathbf{X}) = \mathbf{0} \quad \text{for } F : \left\{ \begin{array}{ccc} \mathbb{R}_{*}^{n,n} & \rightarrow & \mathbb{R}^{n,n} \\ \mathbf{X} & \mapsto & \mathbf{A} - \mathbf{X}^{-1} \end{array} \right..$$

Using (8.4.1.30) we find for the derivative of  $F$  in  $\mathbf{X} \in \mathbb{R}_{*}^{n,n}$

$$D F(\mathbf{X})\mathbf{H} = \mathbf{X}^{-1}\mathbf{H}\mathbf{X}^{-1}, \quad \mathbf{H} \in \mathbb{R}^{n,n}. \quad (8.4.1.32)$$

► The abstract Newton iteration (8.4.1.1) for  $F$  reads

$$\mathbf{X}^{(k+1)} = \mathbf{X}^{(k)} - \mathbf{S}, \quad \mathbf{S} := D F(\mathbf{X}^{(k)})^{-1} F(\mathbf{X}^{(k)}). \quad (8.4.1.33)$$

The **Newton correction**  $\mathbf{S}$  in the  $k$ -th step solves the *linear system of equations*

$$\begin{aligned} \mathbf{D}F(\mathbf{x}^{(k)})\mathbf{S} &\stackrel{(8.4.1.32)}{=} (\mathbf{x}^{(k)})^{-1}\mathbf{S}(\mathbf{x}^{(k)})^{-1} = \mathbf{F}(\mathbf{x}^{(k)}) = \mathbf{A} - (\mathbf{x}^{(k)})^{-1}. \\ \Rightarrow \mathbf{S} &= \mathbf{x}^{(k)}(\mathbf{A} - (\mathbf{x}^{(k)})^{-1})\mathbf{x}^{(k)} = \mathbf{x}^{(k)}\mathbf{A}\mathbf{x}^{(k)} - \mathbf{x}^{(k)}. \end{aligned} \quad (8.4.1.34)$$

in (8.4.1.33)

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{x}^{(k)}\mathbf{A}\mathbf{x}^{(k)} - \mathbf{x}^{(k)}) = \mathbf{x}^{(k)}(2\mathbf{I} - \mathbf{A}\mathbf{x}^{(k)}). \quad (8.4.1.35)$$

To study the convergence of this iteration we derive a recursion for the iteration errors  $\mathbf{E}^{(k)} := \mathbf{x}^{(k)} - \mathbf{A}^{-1}$ :

$$\begin{aligned} \mathbf{E}^{(k+1)} &= \mathbf{x}^{(k+1)} - \mathbf{A}^{-1} \\ &\stackrel{(8.4.1.35)}{=} \mathbf{x}^{(k)}(2\mathbf{I} - \mathbf{A}\mathbf{x}^{(k)}) - \mathbf{A}^{-1} \\ &= (\mathbf{E}^{(k)} + \mathbf{A}^{-1})(2\mathbf{I} - \mathbf{A}(\mathbf{E}^{(k)} + \mathbf{A}^{-1})) - \mathbf{A}^{-1} \\ &= (\mathbf{E}^{(k)} + \mathbf{A}^{-1})(\mathbf{I} - \mathbf{A}\mathbf{E}^{(k)}) - \mathbf{A}^{-1} = -\mathbf{E}^{(k)}\mathbf{A}\mathbf{E}^{(k)}. \end{aligned}$$

For the norm of the iteration error (a matrix norm  $\rightarrow$  Def. 1.5.5.10) we conclude from submultiplicativity (1.5.5.11) a recursive estimate

$$\Rightarrow \|\mathbf{E}^{(k+1)}\| \leq \|\mathbf{E}^{(k)}\|^2 \|\mathbf{A}\|. \quad (8.4.1.36)$$

This holds for any matrix norm according to Def. 1.5.5.10, which is induced by a vector norm. For the relative iteration error we obtain

$$\underbrace{\frac{\|\mathbf{E}^{(k+1)}\|}{\|\mathbf{A}\|}}_{\text{relative error}} \leq \left( \underbrace{\frac{\|\mathbf{E}^{(k)}\|}{\|\mathbf{A}\|}}_{\text{relative error}} \right)^2 \underbrace{\|\mathbf{A}\| \|\mathbf{A}^{-1}\|}_{=\text{cond}(\mathbf{A})}, \quad (8.4.1.37)$$

where the **condition number** is defined in Def. 2.2.2.7.

From (8.4.1.36) we conclude that the iteration will converge ( $\lim_{k \rightarrow \infty} \|\mathbf{E}^{(k)}\| = 0$ ), if

$$\|\mathbf{E}^{(0)}\mathbf{A}\| = \|\mathbf{x}^{(0)}\mathbf{A} - \mathbf{I}\| < 1, \quad (8.4.1.38)$$

which gives a condition on the initial guess  $\mathbf{S}^{(0)}$ . Now let us consider the Euclidean matrix norm  $\|\cdot\|_2$ , which can be expressed in terms of eigenvalues, see Cor. 1.5.5.16. Motivated by this relationship, we use the initial guess  $\mathbf{x}^{(0)} = \alpha \mathbf{A}^\top$  with  $\alpha > 0$  still to be determined.

$$\|\mathbf{x}^{(0)}\mathbf{A} - \mathbf{I}\|_2 = \|\alpha \mathbf{A}^\top \mathbf{A} - \mathbf{I}\|_2 = \alpha \|\mathbf{A}\|_2^2 - 1 \stackrel{!}{<} 1 \Leftrightarrow \alpha < \frac{2}{\|\mathbf{A}\|_2^2},$$

which is a sufficient condition for the initial guess  $\mathbf{x}^{(0)} = \alpha \mathbf{A}^\top$ , in order to make (8.4.1.35) converge. In this case we infer **quadratic convergence** from both (8.4.1.36) and (8.4.1.37).  $\square$

**Remark 8.4.1.39 (Simplified Newton method [DR08, Sect. 5.6.2])** Computing the Newton correction can be expensive owing to the  $\mathcal{O}(n^3)$  asymptotic cost ( $\rightarrow$  § 2.5.0.4) of solving a different large  $n \times n$  linear system of equations in every step of the Newton iteration (8.4.1.1).

We know that the cost of a linear solve can be reduced to  $\mathcal{O}(n^2)$  if the coefficient matrix is available in LU- or QR factorized form, see, e.g., § 2.3.2.15. This motivates the attempt to “samemp\*freeze” the Jacobian in the Newton iteration and use  $\mathbf{D}F(\mathbf{x}^{(0)})$  throughout, which leads to the simplified Newton iteration:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{D}F(\mathbf{x}^{(0)})^{-1}\mathbf{F}(\mathbf{x}^{(k)}), \quad k = 0, 1, \dots.$$

The following C++ function implements a template for this simplified Newton Method and uses the same Jacobian  $\mathbf{D}F(\mathbf{x}^{(0)})$  for all steps, which makes it possible to reuse an LU-decomposition, cf. Rem. 2.5.0.10.

**C++ code 8.4.1.40: Efficient implementation of simplified Newton method**

```

2 // C++ template for simplified Newton method
3 template<typename Func, typename Jac, typename Vec>
4 void simpnewton(Vec& x, Func F, Jac DF, double rtol, double atol)
5 {
6     auto lu = DF(x).lu(); // do LU decomposition once!
7     Vec s;                // Newton correction
8     double ns,nx;         // auxiliary variables for termination control
9     do {
10        s = lu.solve(F(x));
11        x = x-s;           // new iterate
12        ns = s.norm(); nx = x.norm();
13    }
14    // termination based on relative and absolute tolerance
15    while((ns > rtol*nx) && (ns > atol));
16 }

```



Drawback: Switching to the simplified Newton method usually sacrifices the asymptotic quadratic convergence of the Newton method: merely linear convergence can be expected.

**Remark 8.4.1.41 (Numerical Differentiation for computation of Jacobian)** If  $D F(\mathbf{x})$  is not available (e.g. when  $F(\mathbf{x})$  is given only as a procedure) we may resort to approximation by difference quotients:

$$\text{Numerical Differentiation: } \frac{\partial F_i}{\partial x_j}(\mathbf{x}) \approx \frac{F_i(\mathbf{x} + h\vec{e}_j) - F_i(\mathbf{x})}{h} .$$

Caution: Roundoff errors wreak havoc for small  $h$  → Ex. 1.5.4.7 ! Therefore use  $h \approx \sqrt{\text{EPS}}$ .

## 8.4.2 Convergence of Newton's Method

Newton iteration (8.4.1.1)  $\hat{=}$  fixed point iteration ( $\rightarrow$  Section 8.2) with iteration function

$$\Phi(\mathbf{x}) = \mathbf{x} - D F(\mathbf{x})^{-1} F(\mathbf{x}) .$$

[“product rule” (8.4.1.10) :  $D \Phi(\mathbf{x}) = \mathbf{I} - D(\{\mathbf{x} \mapsto D F(\mathbf{x})^{-1}\}) F(\mathbf{x}) - D F(\mathbf{x})^{-1} D F(\mathbf{x})$  ]

$$F(\mathbf{x}^*) = 0 \Rightarrow D \Phi(\mathbf{x}^*) = 0 ,$$

that is, the derivative (Jacobian) of the iteration function of the Newton fixed point iteration vanishes in the limit point. Thus from Lemma 8.2.2.15 we draw the same conclusion as in the scalar case  $n = 1$ , cf. Section 8.3.2.1.

Local quadratic convergence of Newton's method, if  $D F(\mathbf{x}^*)$  regular

**EXPERIMENT 8.4.2.1 (Convergence of Newton's method in 2D)** We study the convergence of Newton's method empirically for  $n = 2$  for

$$F(\mathbf{x}) = \begin{bmatrix} x_1^2 - x_2^4 \\ x_1 - x_2^3 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2 \quad \text{with solution} \quad F\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) = 0 . \quad (8.4.2.2)$$

Jacobian (analytic computation):  $D F(\mathbf{x}) = \begin{bmatrix} \partial_{x_1} F_1(x) & \partial_{x_2} F_1(x) \\ \partial_{x_1} F_2(x) & \partial_{x_2} F_2(x) \end{bmatrix} = \begin{bmatrix} 2x_1 & -4x_2^3 \\ 1 & -3x_2^2 \end{bmatrix}$

Realization of Newton iteration (8.4.1.1):

1. Solve LSE

$$\begin{bmatrix} 2x_1 & -4x_2^3 \\ 1 & -3x_2^2 \end{bmatrix} \Delta \mathbf{x}^{(k)} = -\begin{bmatrix} x_1^2 - x_2^4 \\ x_1 - x_2^3 \end{bmatrix},$$

where  $\mathbf{x}^{(k)} = [x_1, x_2]^T$ .

2. Set  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)}$ .

Monitoring the iteration we obtain the following iterates/error norms:

$k$	$\mathbf{x}^{(k)}$	$\epsilon_k := \ \mathbf{x}^* - \mathbf{x}^{(k)}\ _2$	$\frac{\log \epsilon_{k+1} - \log \epsilon_k}{\log \epsilon_k - \log \epsilon_{k-1}}$
0	$[0.7, 0.7]^T$	4.24e-01	
1	$[0.878500000000000, 1.064285714285714]^T$	1.37e-01	1.69
2	$[1.01815943274188, 1.00914882463936]^T$	2.03e-02	2.23
3	$[1.00023355916300, 1.00015913936075]^T$	2.83e-04	2.15
4	$[1.00000000583852, 1.00000002726552]^T$	2.79e-08	1.77
5	$[0.999999999999998, 1.000000000000000]^T$	2.11e-15	
6	$[1, 1]^T$		

☞ (Some) evidence of quadratic convergence, see Rem. 8.1.1.12. □

There is a sophisticated theory about the convergence of Newton's method. For example one can find the following theorem in [DH03, Thm. 4.10], [Deu11, Sect. 2.1]):

### Theorem 8.4.2.3. Local quadratic convergence of Newton's method

If:

- (A)  $D \subset \mathbb{R}^n$  open and convex,
- (B)  $F : D \mapsto \mathbb{R}^n$  continuously differentiable,
- (C)  $D F(\mathbf{x})$  regular  $\forall \mathbf{x} \in D$ ,
- (D)  $\exists L \geq 0: \|D F(\mathbf{x})^{-1}(D F(\mathbf{x} + \mathbf{v}) - D F(\mathbf{x}))\|_2 \leq L \|\mathbf{v}\|_2 \quad \forall \mathbf{v} \in \mathbb{R}^n, \mathbf{v} + \mathbf{x} \in D, \quad \forall \mathbf{x} \in D$ ,
- (E)  $\exists \mathbf{x}^*: F(\mathbf{x}^*) = 0$  (existence of solution in  $D$ )
- (F) initial guess  $\mathbf{x}^{(0)} \in D$  satisfies  $\rho := \|\mathbf{x}^* - \mathbf{x}^{(0)}\|_2 < \frac{2}{L} \quad \wedge \quad B_\rho(\mathbf{x}^*) \subset D$ .

then the Newton iteration (8.4.1.1) satisfies:

- (i)  $\mathbf{x}^{(k)} \in B_\rho(\mathbf{x}^*) := \{\mathbf{y} \in \mathbb{R}^n, \|\mathbf{y} - \mathbf{x}^*\| < \rho\}$  for all  $k \in \mathbb{N}$ ,
- (ii)  $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$ ,
- (iii)  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2 \leq \frac{L}{2} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^2 \quad (\text{local quadratic convergence})$ .

☞ notation: ball  $B_\rho(\mathbf{z}) := \{\mathbf{x} \in \mathbb{R}^n: \|\mathbf{x} - \mathbf{z}\|_2 \leq \rho\}$

Terminology: (D)  $\hat{=}$  affine invariant Lipschitz condition



Usually, it is hardly possible to verify the assumptions of the theorem for a concrete non-linear system of equations, because neither  $L$  nor  $\mathbf{x}^*$  are known.

- In general: a priori estimates as in Thm. 8.4.2.3 are of little practical relevance.

### 8.4.3 Termination of Newton Iteration

An abstract discussion of ways to stop iterations for solving  $F(\mathbf{x}) = 0$  was presented in Section 8.1.2, with “ideal termination” ( $\rightarrow$  § 8.1.2.2) as ultimate, but unfeasible, goal.

Yet, in 8.4.2 we saw that Newton’s method enjoys (asymptotic) quadratic convergence, which means rapid decrease of the relative error of the iterates, once we are close to the solution, which is exactly the point, when we want to **STOP**. As a consequence, **asymptotically**, the Newton correction (difference of two consecutive iterates) yields rather precise information about the size of the error:

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \ll \|\mathbf{x}^{(k)} - \mathbf{x}^*\| \Rightarrow \|\mathbf{x}^{(k)} - \mathbf{x}^*\| \approx \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|. \quad (8.4.3.1)$$

This suggests the following **correction based** termination criterion:

$$\begin{aligned} \text{STOP, as soon as } & \|\Delta\mathbf{x}^{(k)}\| \leq \tau_{\text{rel}} \|\mathbf{x}^{(k)}\| \quad \text{or} \quad \|\Delta\mathbf{x}^{(k)}\| \leq \tau_{\text{abs}}, \\ \text{with Newton correction } & \Delta\mathbf{x}^{(k)} := D F(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)}). \end{aligned} \quad (8.4.3.2)$$

Here,  $\|\cdot\|$  can be any suitable vector norm,  $\tau_{\text{rel}} \hat{=} \text{relative tolerance}$ ,  $\tau_{\text{abs}} \hat{=} \text{absolute tolerance}$ , see § 8.1.2.2.

➤ quit iterating as soon as  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| = \|D F(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)})\| < \tau \|\mathbf{x}^{(k)}\|$ ,  
with  $\tau$  = tolerance

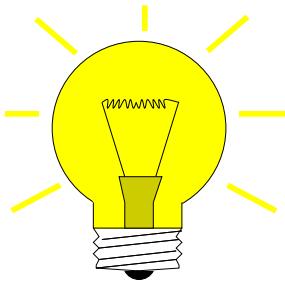
$\rightarrow$  uneconomical: one needless update, because  $\mathbf{x}^{(k)}$  would already be accurate enough.

**Remark 8.4.3.3 (Newton’s iteration; computational effort and termination)** Some facts about the Newton method for solving **large** ( $n \gg 1$ ) non-linear systems of equations:

- ☛ Solving the linear system to compute the Newton correction may be expensive (asymptotic computational effort  $O(n^3)$  for direct elimination  $\rightarrow$  § 2.3.1.5) and accounts for the bulk of numerical cost of a single step of the iteration.
- ☛ In applications only **very few steps** of the iteration will be needed to achieve the desired accuracy due to fast quadratic convergence.
- ▷ The termination criterion (8.4.3.2) computes the last Newton correction  $\Delta\mathbf{x}^{(k)}$  needlessly, because  $\mathbf{x}^{(k)}$  already accurate enough!

Therefore we would like to use an a-posteriori termination criterion that dispenses with computing (and “inverting”) another Jacobian  $D F(\mathbf{x}^{(k)})$  just to tell us that  $\mathbf{x}^{(k)}$  is already accurate enough.  $\downarrow$

**§8.4.3.4 (Termination of Newton iteration based on simplified Newton correction)** Due to fast asymptotic quadratic convergence, we can expect  $D F(\mathbf{x}^{(k-1)}) \approx D F(\mathbf{x}^{(k)})$  during the final steps of the iteration.



Idea: Replace  $D F(\mathbf{x}^{(k)})$  with  $D F(\mathbf{x}^{(k-1)})$  in any correction based termination criterion.

Rationale: LU-decomposition of  $D F(\mathbf{x}^{(k-1)})$  is already available ➤ less effort.

Terminology:

$$\Delta \bar{\mathbf{x}}^{(k)} := D F(\mathbf{x}^{(k-1)})^{-1} F(\mathbf{x}^{(k)}) \triangleq \text{simplified Newton correction}$$

► Economical correction based termination criterion for Newton's method:

$$\begin{aligned} \text{STOP, as soon as } & \|\Delta \bar{\mathbf{x}}^{(k)}\| \leq \tau_{\text{rel}} \|\mathbf{x}^{(k)}\| \quad \text{or} \quad \|\Delta \bar{\mathbf{x}}^{(k)}\| \leq \tau_{\text{abs}}, \\ \text{with simplified Newton correction } & \Delta \bar{\mathbf{x}}^{(k)} := D F(\mathbf{x}^{(k-1)})^{-1} F(\mathbf{x}^{(k)}). \end{aligned} \quad (8.4.3.5)$$

Note that (8.4.3.5) is **affine invariant** → Rem. 8.4.1.4.

Effort: Reuse of LU-factorization (→ Rem. 2.5.0.10) of  $D F(\mathbf{x}^{(k-1)})$  ➤  $\Delta \bar{\mathbf{x}}^{(k)}$  available with  $O(n^2)$  operations

#### C++11 code 8.4.3.6: Generic Newton iteration with termination criterion (8.4.3.5)

```

2  template <typename FuncType, typename JacType, typename VecType>
3  void newton_stc(const FuncType &F, const JacType &DF, VecType &x, double rtol,
4                  double atol) {
5      using scalar_t = typename VecType::Scalar;
6      scalar_t sn;
7      do {
8          auto jacfac = DF(x).lu(); // LU-factorize Jacobian ]
9          x -= jacfac.solve(F(x)); // Compute next iterate
10         // Compute norm of simplified Newton correction
11         sn = jacfac.solve(F(x)).norm();
12     }
13     // Termination based on simplified Newton correction
14     while ((sn > rtol * x.norm()) && (sn > atol));
15 }
```

**Remark 8.4.3.7 (Residual based termination of Newton's method)** If we used the residual based termination criterion

$$\|F(\mathbf{x}^{(k)})\| \leq \tau,$$

then the resulting algorithm would **not** be **affine invariant**, because for  $F(\mathbf{x}) = 0$  and  $\mathbf{A}F(\mathbf{x}) = 0$ ,  $\mathbf{A} \in \mathbb{R}^{n,n}$  regular, the Newton iteration might terminate with different iterates. □

#### Summary: Newton's method



converges *asymptotically* very fast: doubling of number of significant digits in each step



often a very small region of convergence, which requires an initial guess rather close to the solution.

#### 8.4.4 Damped Newton Method

Potentially big problem: Newton method converges quadratically, but only *locally*, which may render it useless, if convergence is guaranteed only for initial guesses very close to exact solution, see also Ex. 8.3.2.31.

In this section we study a method to enlarge the region of convergence, at the expense of quadratic convergence, of course.

**EXAMPLE 8.4.4.1 (Local convergence of Newton's method)** The dark side of local convergence (→ Def. 8.1.0.8): for many initial guesses  $x^{(0)}$  Newton's method will not converge!

In 1D two main causes can be identified:

- ① “Wrong direction” of Newton correction:

$$F(x) = xe^x - 1 \Rightarrow F'(-1) = 0$$

$$\begin{aligned} x^{(0)} < -1 &\Rightarrow x^{(k)} \rightarrow -\infty, \\ x^{(0)} > -1 &\Rightarrow x^{(k)} \rightarrow x^*, \end{aligned}$$

because all Newton corrections for  $x^{(k)} < -1$  make the iterates decrease even further.

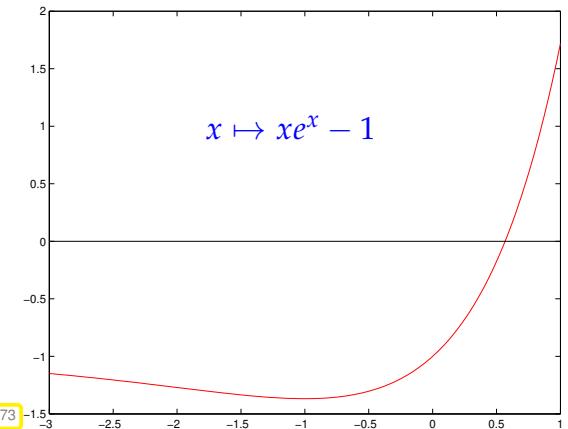


Fig. 273

- ② Newton correction is too large:

$$F(x) = \arctan(ax), \quad a > 0, x \in \mathbb{R}$$

$$\text{with zero } x^* = 0.$$

If  $x^{(k)}$  is located where the function is “flat”, the intersection of the tangents with the  $x$ -axis is “far out”, see Fig. 275.

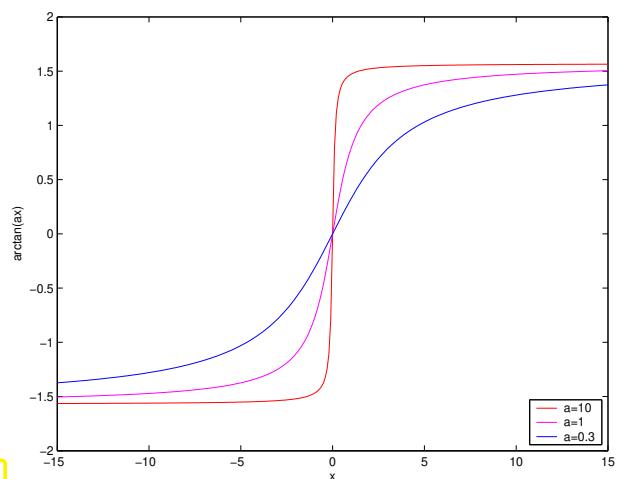
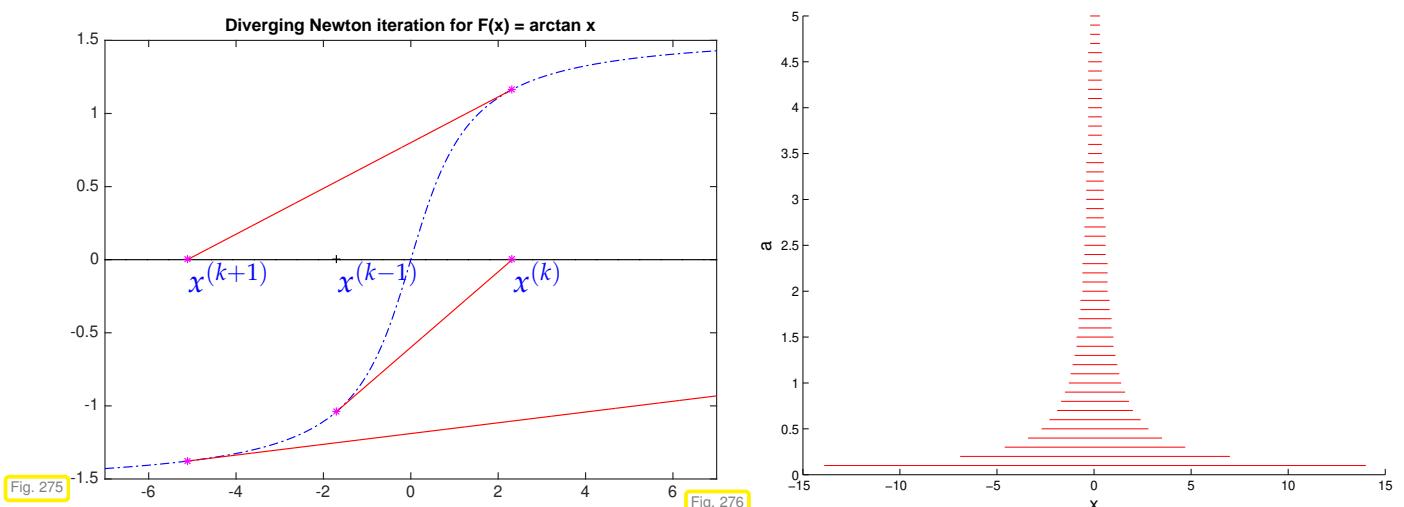
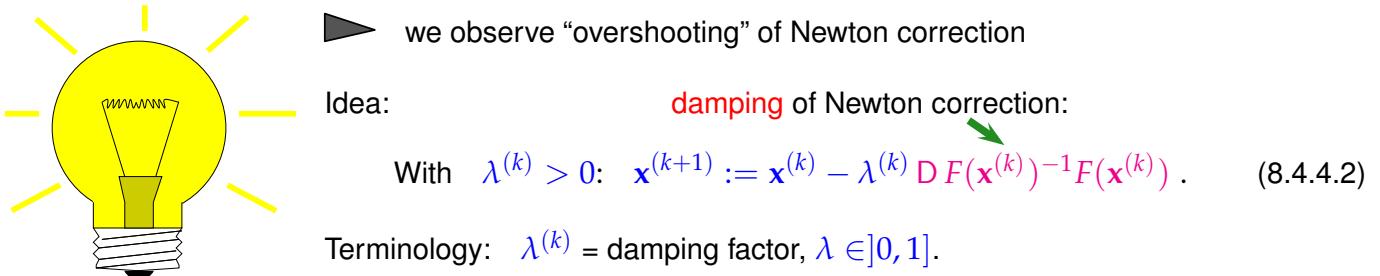


Fig. 274



In Fig. 276 the red zone =  $\{x^{(0)} \in \mathbb{R}, x^{(k)} \rightarrow 0\}$ , domain of initial guesses for which Newton's method converges.

If the Newton correction points in the wrong direction (Item ①), no general remedy is available. If the Newton correction is too large (Item ②), there is an effective cure:



### Affine invariant damping strategy

Choice of damping factor: affine invariant natural monotonicity test (NMT) [Deu11, Ch. 3]:

$$\text{choose “maximal” } 0 < \lambda^{(k)} \leq 1: \quad \left\| \Delta \bar{x}(\lambda^{(k)}) \right\|_2 \leq \left(1 - \frac{\lambda^{(k)}}{2}\right) \left\| \Delta x^{(k)} \right\|_2 \quad (8.4.4.4)$$

where  $\Delta x^{(k)} := D F(x^{(k)})^{-1} F(x^{(k)})$  → current Newton correction ,

$\Delta \bar{x}(\lambda^{(k)}) := D F(x^{(k)})^{-1} F(x^{(k)} + \lambda^{(k)} \Delta x^{(k)})$  → tentative simplified Newton correction .

Heuristics behind control of damping:

- ◆ When the method converges  $\Leftrightarrow$  size of Newton correction decreases  $\Leftrightarrow$  (8.4.4.4) satisfied.
- ◆ In the case of strong damping ( $\lambda^{(k)} \ll 1$ ) the size of the Newton correction cannot be expected to shrink significantly, since iterates do not change much  $\geq$  factor  $(1 - \frac{1}{2}\lambda^{(k)})$  in (8.4.4.4).

Note: As before, reuse of LU-factorization in the computation of  $\Delta x^{(k)}$  and  $\Delta \bar{x}(\lambda^{(k)})$ .

Policy: Reduce damping factor by a factor  $q \in ]0, 1[$  (usually  $q = \frac{1}{2}$ ) until the affine invariant natural monotonicity test (8.4.4.4) passed, see Line 20 in the following C++ code.

**C++ code 8.4.4.5: Generic damped Newton method based on natural monotonicity test**

```

1 template <typename FuncType, typename JacType, typename VecType>
2 void dampnewton(const FuncType &F, const JacType &DF,
3                  VecType &x, double rtol, double atol)
4 {
5     using index_t = typename VecType::Index;
6     using scalar_t = typename VecType::Scalar;
7     const index_t n = x.size(); // No. of unknowns
8     const scalar_t lmin = 1E-3; // Minimal damping factor
9     scalar_t lambda = 1.0; // Initial and actual damping factor
10    VecType s(n), st(n); // Newton corrections
11    VecType xn(n); // Tentative new iterate
12    scalar_t sn, stn; // Norms of Newton corrections
13
14    do {
15        auto jfac = DF(x).lu(); // LU-factorize Jacobian
16        s = jfac.solve(F(x)); // Newton correction
17        sn = s.norm(); // Norm of Newton correction
18        lambda *= 2.0;
19        do {
20            lambda /= 2; // Reduce damping factor
21            if (lambda < lmin) throw "No convergence: lambda -> 0";
22            xn = x - lambda*s; // Tentative next iterate
23            st = jfac.solve(F(xn)); // Simplified Newton correction
24            stn = st.norm();
25        }
26        while (stn > (1-lambda/2)*sn); // Natural monotonicity test
27        x = xn; // Now: xn accepted as new iterate
28        lambda = std::min(2.0*lambda, 1.0); // Try to mitigate damping
29    }
30    // Termination based on simplified Newton correction
31    while ((stn > rtol*x.norm()) && (stn > atol));
32 }
```

The arguments for Code 8.4.4.5 are the same as for Code 8.4.3.6. As termination criterion is uses (8.4.3.5). Note that all calls to `solve` boil down to forward/backward elimination for triangular matrices and incur cost of  $\mathcal{O}(n^2)$  only.

Note: The LU-factorization of the Jacobi matrix  $D F(\mathbf{x}^{(k)})$  is done once per *successful* iteration step and reused for the computation of the simplified Newton correction in Line 23 of the above C++ code.

**EXPERIMENT 8.4.4.6 (Damped Newton method)** We test the damped Newton method for Item ② of Ex. 8.4.4.1, where excessive Newton corrections made Newton's method fail.

$F(\mathbf{x}) = \arctan(\mathbf{x})$ ,  
•  $x^{(0)} = 20$   
•  $q = \frac{1}{2}$   
• LMIN = 0.001  
We observe that damping is effective and asymptotic quadratic convergence is recovered.

$k$	$\lambda^{(k)}$	$\mathbf{x}^{(k)}$	$F(\mathbf{x}^{(k)})$
1	0.03125	0.94199967624205	0.75554074974604
2	0.06250	0.85287592931991	0.70616132170387
3	0.12500	0.70039827977515	0.61099321623952
4	0.25000	0.47271811131169	0.44158487422833
5	0.50000	0.20258686348037	0.19988168667351
6	1.00000	-0.00549825489514	-0.00549819949059
7	1.00000	0.00000011081045	0.00000011081045
8	1.00000	-0.00000000000001	-0.00000000000001

**EXPERIMENT 8.4.4.7 (Failure of damped Newton method)** We examine the effect of damping in the

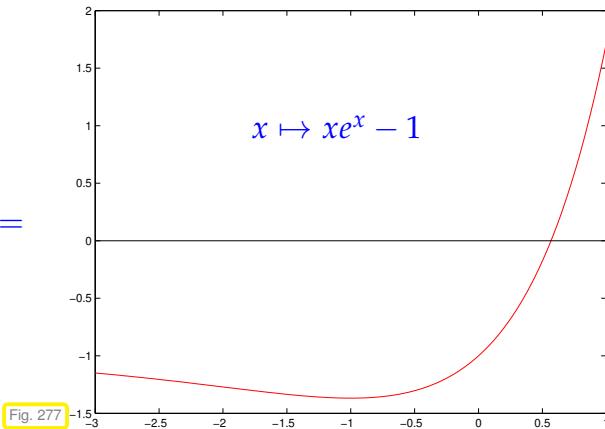
case of Item ❶ of Ex. 8.4.4.1.

- ◆ As in Ex. 8.4.4.1:

$$F(x) = xe^x - 1,$$

- ◆ Initial guess for damped Newton method:  $x^{(0)} = -1.5$

This time the initial guess is to the left of the global minimum of the function.



Observation:

Newton correction pointing in “wrong direction”

► no convergence despite damping

$k$	$\lambda^{(k)}$	$x^{(k)}$	$F(x^{(k)})$
1	0.25000	-4.4908445351690	-1.0503476286303
2	0.06250	-6.1682249558799	-1.0129221310944
3	0.01562	-7.6300006580712	-1.0037055902301
4	0.00390	-8.8476436930246	-1.0012715832278
5	0.00195	-10.5815494437311	-1.0002685596314
Bailed out because of lambda < LMIN !			

## 8.4.5 Quasi-Newton Method



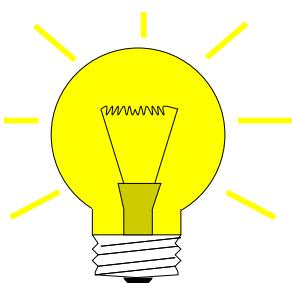
*Supplementary literature.* For related expositions refer to [QSS00, Sect. 7.1.4], [Wer92, p. 2.3.2].

We start with the following question: How can we solve the non-linear system of equations  $F(\mathbf{x}) = 0$ ,  $D : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ , iteratively, in case  $D F(\mathbf{x})$  is not available and numerical differentiation (see Rem. 8.4.1.41) is too expensive?

In 1D ( $n = 1$ ) we can choose among many derivative-free methods that rely on  $F$ -evaluations alone, for instance the secant method (8.3.2.23) from Section 8.3.2.3:

$$x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)}) - F(x^{(k-1)})}{F(x^{(k)}) - F(x^{(k-1)})}. \quad (8.3.2.23)$$

Recall from Rem. 8.3.2.26 that the secant method converges locally with order  $p \approx 1.6$  and beats Newton's method in terms of efficiency (→ Section 8.3.3).

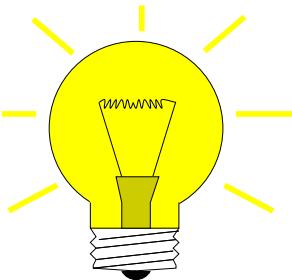


Comparing with (8.3.2.1) we realize that this iteration amounts to a “Newton-type iteration” with the approximation

$$F'(x^{(k)}) \approx \frac{F(x^{(k)}) - F(x^{(k-1)})}{x^{(k)} - x^{(k-1)}} \quad \begin{matrix} \text{“difference quotient”} \\ \text{already computed!} \rightarrow \text{cheap} \end{matrix} \quad (8.4.5.1)$$



Unfortunately, it is not immediate how to generalize the secant method to  $n > 1$ .



Idea: rewrite (8.4.5.1) as a **secant condition** for an approximation  
 $\mathbf{J}_k \approx \mathbf{D} F(\mathbf{x}^{(k)})$ ,  $\mathbf{x}^{(k)} \triangleq$  the current iterate:

$$\mathbf{J}_k(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) = F(\mathbf{x}^{(k)}) - F(\mathbf{x}^{(k-1)}). \quad (8.4.5.2)$$

► Iteration:  $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \mathbf{J}_k^{-1}F(\mathbf{x}^{(k)})$ . (8.4.5.3)



Embarassment of choice: Many different matrices  $\mathbf{J}_k$  fulfill (8.4.5.2)!

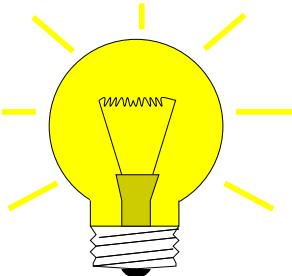


We need extra conditions to fix  $\mathbf{J}_k \in \mathbb{R}^{n,n}$ .

Reasoning: If we assume that  $\mathbf{J}_k$  is a good approximation of  $\mathbf{D} F(\mathbf{x}^{(k)})$ , then it would be foolish not to use the information contained in  $\mathbf{J}_k$  for the construction of  $\mathbf{J}_{k+1}$ .

► Guideline: obtain  $\mathbf{J}_k$  through a “small” **modification** of  $\mathbf{J}_{k-1}$  compliant with (8.4.5.2)

What can “small modification” mean: Demand that  $\mathbf{J}_k$  acts like  $\mathbf{J}_{k-1}$  on the orthogonal complement of the one-dimensional subspace of  $\mathbb{R}^n$  generated by the vector  $\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}$ !



► Broyden’s conditions:  $\mathbf{J}_k \mathbf{z} = \mathbf{J}_{k-1} \mathbf{z} \quad \forall \mathbf{z}: \mathbf{z} \perp (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$ . (8.4.5.4)

► 
$$\mathbf{J}_k := \mathbf{J}_{k-1} + \frac{F(\mathbf{x}^{(k)})(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})^\top}{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_2^2} \quad (8.4.5.5)$$

- ◆ The conditions (8.4.5.2) and (8.4.5.4) uniquely define  $\mathbf{J}_k$ .
- ◆ The update formula (8.4.5.5) means that  $\mathbf{J}_k$  is spawned by a **rank-1-modification** of  $\mathbf{J}_{k-1}$ .



Final form of **Broyden’s quasi-Newton method** for solving  $F(\mathbf{x}) = 0$ :

$$\begin{aligned} \mathbf{x}^{(k+1)} &:= \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)}, \quad \Delta \mathbf{x}^{(k)} := -\mathbf{J}_k^{-1}F(\mathbf{x}^{(k)}), \\ \mathbf{J}_{k+1} &:= \mathbf{J}_k + \frac{F(\mathbf{x}^{(k+1)})(\Delta \mathbf{x}^{(k)})^\top}{\|\Delta \mathbf{x}^{(k)}\|_2^2}. \end{aligned} \quad (8.4.5.6)$$

To start the iteration we have to initialize  $\mathbf{J}_0$ , e.g. with the exact Jacobi matrix  $\mathbf{D} F(\mathbf{x}^{(0)})$ .

**Remark 8.4.5.7 (Minimality property of Broyden’s rank-1-modification)** in another sense  $\mathbf{J}_l$  is closest to  $\mathbf{J}_{k-1}$  under the constraint of the secant condition (8.4.5.2):

Let  $\mathbf{x}^{(k)}$  and  $\mathbf{J}_k$  be the iterates and matrices, respectively, from Broyden's method (8.4.5.6), and let  $\mathbf{J} \in \mathbb{R}^{n,n}$  satisfy the same secant condition (8.4.5.2) as  $\mathbf{J}_{k+1}$ :

$$\mathbf{J}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = F(\mathbf{x}^{(k+1)}) - F(\mathbf{x}^{(k)}). \quad (8.4.5.8)$$

Then from  $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} - \mathbf{J}_k^{-1}F(\mathbf{x}^{(k)})$  we obtain

$$(\mathbf{I} - \mathbf{J}_k^{-1}\mathbf{J})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = -\mathbf{J}_k^{-1}F(\mathbf{x}^{(k)}) - \mathbf{J}_k^{-1}(F(\mathbf{x}^{(k+1)}) - F(\mathbf{x}^{(k)})) = -\mathbf{J}_k^{-1}F(\mathbf{x}^{(k+1)}). \quad (8.4.5.9)$$

From this we get the identity

$$\begin{aligned} \mathbf{I} - \mathbf{J}_k^{-1}\mathbf{J}_{k+1} &= \mathbf{I} - \mathbf{J}_k^{-1}\left(\mathbf{J}_k + \frac{F(\mathbf{x}^{(k+1)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})^\top}{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2^2}\right) \\ &= -\mathbf{J}_k^{-1}F(\mathbf{x}^{(k+1)})\frac{(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})^\top}{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2^2} = \\ &\stackrel{(8.4.5.9)}{=} (\mathbf{I} - \mathbf{J}_k^{-1}\mathbf{J})\frac{(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})^\top}{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2^2}. \end{aligned}$$

Using the submultiplicative property (1.5.5.11) of the Euclidean matrix norm, we conclude

$$\|\mathbf{I} - \mathbf{J}_k^{-1}\mathbf{J}_{k+1}\| \leq \|\mathbf{I} - \mathbf{J}_k^{-1}\mathbf{J}\|, \text{ because } \left\| \frac{(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})^\top}{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2^2} \right\|_2 \leq 1,$$

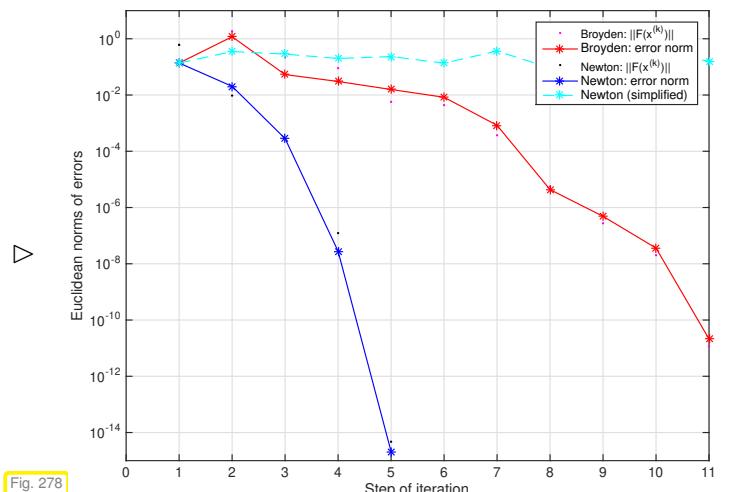
which we saw in Ex. 1.5.5.20. This estimate holds for *all* matrices  $\mathbf{J}$  satisfying (8.4.5.8).

We may read this as follows: (8.4.5.5) gives the  $\|\cdot\|_2$ -minimal relative correction of  $\mathbf{J}_{k-1}$ , such that the secant condition (8.4.5.2) holds.  $\square$

**EXPERIMENT 8.4.5.10 (Broydens quasi-Newton method: Convergence)** We revisit the  $2 \times 2$  non-linear system of the Exp. 8.4.2.1,

$$F(\mathbf{x}) = \begin{bmatrix} x_1^2 - x_2^4 \\ x_1 - x_2^3 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2 \quad \text{with solution} \quad F\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) = 0, \quad (8.4.2.2)$$

and take  $\mathbf{x}^{(0)} = [0.7, 0.7]^T$ . As starting value for the matrix iteration we use  $\mathbf{J}_0 = D F(\mathbf{x}^{(0)})$ .



The numerical example shows that, in terms of convergence, the method is:

- slower than Newton method (8.4.1.1),
- faster than the simplified Newton method (see Rem. 8.4.1.39)

**Remark 8.4.5.11 (Convergence monitors)** In general, any iterative methods for non-linear systems of equations convergence can fail, that is it may stall or even diverge.

Demand on good numerical software: Algorithms should warn users of impending failure. For iterative methods this is the task of **convergence monitors**, that is, conditions, *cheaply verifiable* a posteriori during the iteration, that indicate stalled convergence or divergence.

For the damped Newton's method this role can be played by the natural monotonicity test, see Code 8.4.4.5; if it fails repeatedly, then the iteration should terminate with an error status.

For Broyden's quasi-Newton method, a similar strategy can rely on the relative size of the “simplified Broyden correction”  $\mathbf{J}_k F(\mathbf{x}^{(k+1)})$ :

$$\text{Convergence monitor for (8.4.5.6)} : \quad \mu := \frac{\|\mathbf{J}_{k-1}^{-1} F(\mathbf{x}^{(k)})\|}{\|\Delta \mathbf{x}^{(k-1)}\|} < 1 ? \quad (8.4.5.12)$$

### EXPERIMENT 8.4.5.13 (Monitoring convergence for Broyden's quasi-Newton method)

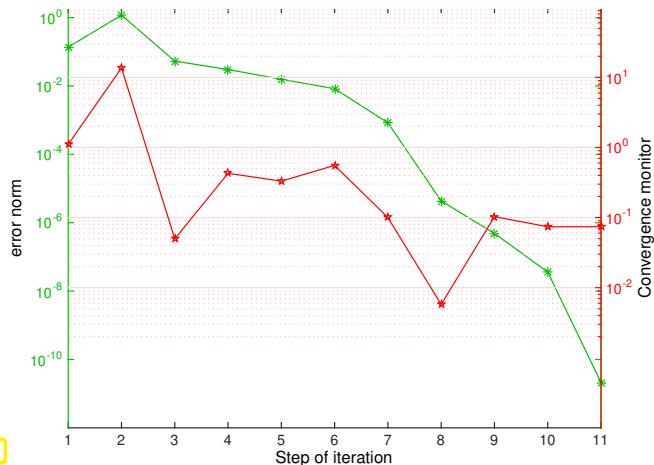


Fig. 279

We rely on the setting of Exp. 8.4.5.10.

We track

1. the Euclidean norm of the iteration error,
2. and the value of the convergence monitor from (8.4.5.12).

▷ Decay of (norm of) iteration error and  $\mu$  are well correlated.

**Remark 8.4.5.14 (Damped Broyden method)** Option to improve robustness (increase region of local convergence):

damped Broyden method (cf. same idea for Newton's method, Section 8.4.4)

**§8.4.5.15 (Implementation of Broyden's quasi-Newton method)** As remarked, (8.4.5.6) represents a **rank-1-update** as already discussed in § 2.6.0.12.

Idea: use Sherman-Morrison-Woodbury update-formula from Lemma 2.6.0.21,

$$(\mathbf{A} + \mathbf{U}\mathbf{V}^H)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{V}^H\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}^H\mathbf{A}^{-1},$$

for the special case  $k = 1$ ,  $\mathbb{K} = \mathbb{R}$ ,

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^T)^{-1} = \left( \mathbf{I} - \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}^T}{1 + \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u}} \right) \mathbf{A}^{-1}, \quad 1 + \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u} \neq 0,$$

which yields, with  $\mathbf{u} := F(\mathbf{x}^{(k+1)})$ ,  $\mathbf{v} := \Delta\mathbf{x}^{(k)}$ ,

$$\mathbf{J}_{k+1}^{-1} = \left( \mathbf{I} - \frac{\mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)})(\Delta\mathbf{x}^{(k)})^T}{\|\Delta\mathbf{x}^{(k)}\|_2^2 + \Delta\mathbf{x}^{(k)} \cdot \mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)})} \right) \mathbf{J}_k^{-1} = \left( \mathbf{I} + \frac{\Delta\mathbf{x}^{(k+1)}(\Delta\mathbf{x}^{(k)})^T}{\|\Delta\mathbf{x}^{(k)}\|_2^2} \right) \mathbf{J}_k^{-1}. \quad (8.4.5.16)$$

This gives a well defined  $\mathbf{J}_{k+1}$ , if

$$\left\| \mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)}) \right\|_2 < \left\| \Delta\mathbf{x}^{(k)} \right\|_2. \quad (8.4.5.17)$$

"simplified Quasi-Newton correction"

Note that the condition (8.4.5.17) is checked by the convergence monitor (8.4.5.12).

Iterated application of (8.4.5.16) pays off, if iteration terminates after only a few steps. For large  $n \gg 1$  it is not advisable to form the matrices  $\mathbf{J}_k^{-1}$  (which will usually be dense in contrast to  $\mathbf{J}_k$ ), but we employ fast successive multiplications with rank-1-matrices ( $\rightarrow$  Ex. 1.4.3.1) to apply  $\mathbf{J}_k^{-1}$  to a vector. This is implemented in the following code.

MATLAB-code : Broyden method (8.4.5.6)

```

1  function x = broyden(F,x,J,tol)
2  k = 1;
3  [L,U] = lu(J);
4  s = U\ (L\ F(x));
5  sn = dot(s,s);
6  dx = [s];
7  x = x - s;
8  f = F(x);
9
10 while (sqrt(sn) > tol)
11    w = U\ (L\ f);
12    for l=2:k-1
13      w = w+dx(:,l)* (dx(:,l-1)' *w) ...
14        /dxn(l-1);
15    end
16    if (norm(w)>=sn)
17      warning('Dubious step %d!',k);
18    end
19    z = s'*w;
20    s = (1+z/(sn-z))*w;
21    sn=s'*s;
22    dx = [dx,s];
23    dxn = [dxn,sn];
24    x = x - s;
25    f = F(x);
26  end
```

Computational cost :  
 $N$  steps

- ◆  $O(N^2 \cdot n)$  operations with vectors, (Level I)
- ◆ 1 LU-decomposition of  $J$ ,  $N \times$  solutions of SLEs, see Section 2.3.2
- ◆  $N$  evaluations of  $F$  !

Memory cost :

- $N$  steps
- ◆ LU-factors of  $J$  + auxiliary vectors  $\in \mathbb{R}^n$
  - ◆  $N$  vectors  $\mathbf{x}^{(k)} \in \mathbb{R}^n$

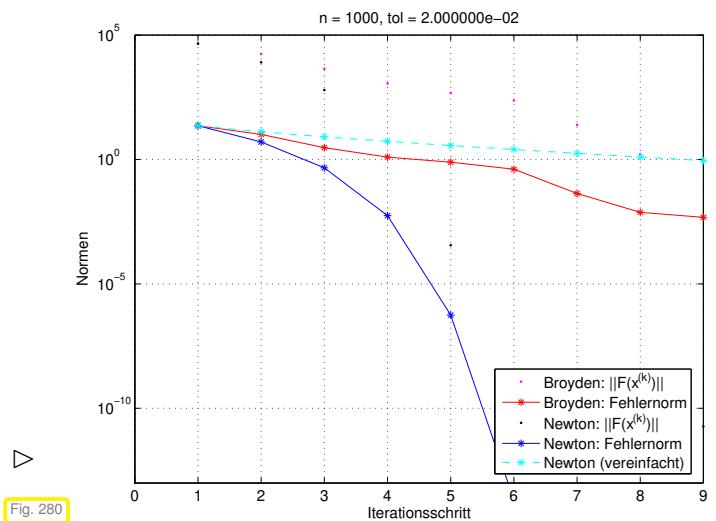
### EXPERIMENT 8.4.5.18 (Broyden method for a large non-linear system)

$$F(\mathbf{x}) = \begin{cases} \mathbb{R}^n \rightarrow \mathbb{R}^n \\ \mathbf{x} \mapsto \text{diag}(\mathbf{x})\mathbf{A}\mathbf{x} - \mathbf{b}, \end{cases}$$

$\mathbf{b} = [1, 2, \dots, n] \in \mathbb{R}^n,$   
 $\mathbf{A} = \mathbf{I} + \mathbf{a}\mathbf{a}^T \in \mathbb{R}^{n,n},$   
 $\mathbf{a} = \frac{1}{\sqrt{1 \cdot \mathbf{b} - 1}}(\mathbf{b} - \mathbf{1}).$

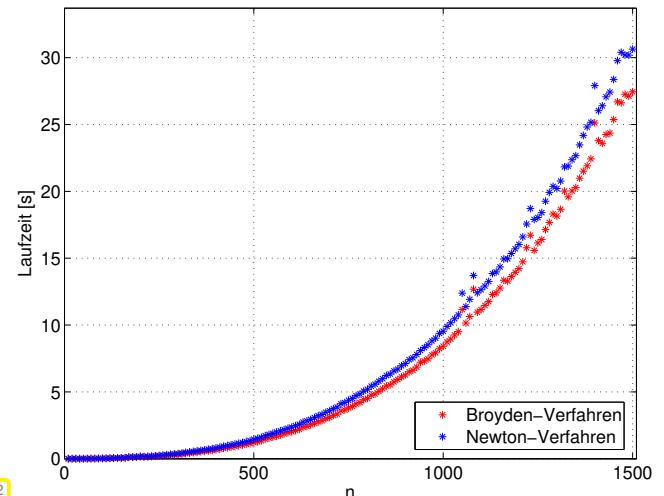
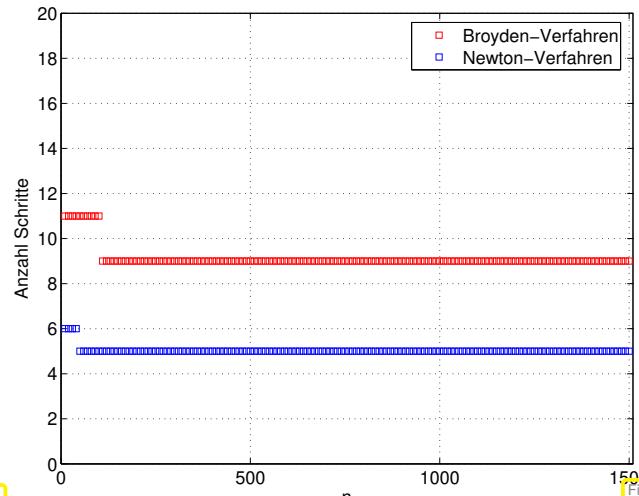
Initial guess:  $\mathbf{h} = 2/n; \mathbf{x}_0 = (2:h:4-\mathbf{h})'$ ;

The results resemble those of Exp. 8.4.5.10



Efficiency comparison:

Broyden method  $\longleftrightarrow$  Newton method:  
(in case of dimension  $n$  use tolerance  $\text{tol} = 2n \cdot 10^{-5}$ ,  $\mathbf{h} = 2/n; \mathbf{x}_0 = (2:h:4-\mathbf{h})'$ ;



In conclusion,

the Broyden method is worthwhile for dimensions  $n \gg 1$  and low accuracy requirements.



*Supplementary literature.* A comprehensive monograph about all aspects of Newton's methods and generalizations in [Deu11]. The multi-dimensional Newton method is also presented in [Han02, Sect. 19], [DR08, Sect. 5.6], [AG11, Sect. 9.1].

## 8.5 Non-linear Least Squares [DR08, Ch. 6]

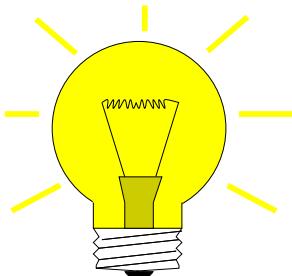
So far we have studied non-linear systems of equations  $F(\mathbf{x}) = 0$  with the same number  $n \in \mathbb{N}$  of unknowns and equations:  $F : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ . This generalizes square linear systems of equations whose numerical treatment was the subjects of Chapter 2. Then, in Chapter 3, we turned our attention to *overdetermined* linear systems of equations  $\mathbf{Ax} = \mathbf{b}$  with  $\mathbf{A} \in \mathbb{R}^{m,n}$ ,  $m > n$ . Now we take the same step for non-linear systems of equations  $F(\mathbf{x}) = 0$  and admit non-linear functions  $F : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m > n$ .

For overdetermined linear systems of equations we had to introduce the concept of a *least-squares solution*.

tion in Section 3.1, Def. 3.1.1.1. The same concept will apply in the non-linear case.

**EXAMPLE 8.5.0.1 (Least squares data fitting)** In Section 5.1 we discussed the reconstruction of a *parameterized function*  $f(x_1, \dots, x_n; \cdot) : D \subset \mathbb{R} \mapsto \mathbb{R}$  from data points  $(t_i, y_i), i = 1, \dots, n$ , by imposing interpolation conditions. We demanded that the number  $n$  of parameters agreed with number of data points. Thus, in the case of a general dependence of  $f$  on the parameters, the interpolation conditions (5.1.0.2) yield a non-linear system of equations.

The interpolation approach is justified in the case of *highly accurate data*. However, we frequently encountered *inaccurate data*, for instance, due to *measurement errors*. As we discussed in Section 5.7 this renders the interpolation approach dubious, also in light of the impact of “outliers”.



Idea: Mitigate the impact of data uncertainty by choosing *fewer parameters than data points*.  
➤ Thus measurement errors can “average out”.

Guided by this idea we arrive at a particular version of the least squares data fitting problem from Section 5.7, *cf.* (5.7.0.2).

### Non-linear least squares fitting problem

Given: ♦ data points  $(t_i, y_i), i = 1, \dots, m$   
♦ (symbolic formula) for parameterized function  
 $f(x_1, \dots, x_n; \cdot) : D \subset \mathbb{R} \mapsto \mathbb{R}, n < m$

Sought: parameter values  $x_1^*, \dots, x_n^* \in \mathbb{R}$  such that

$$(x_1^*, \dots, x_n^*) = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \sum_{i=1}^m |f(x_1, \dots, x_n; t_i) - y_i|^2. \quad (8.5.0.3)$$

As we did in § 5.7.0.13 for the linear case, we can rewrite (8.5.0.3) by introducing

$$F(\mathbf{x}) := \begin{bmatrix} f(x_1, \dots, x_n, t_1) - y_1 \\ \vdots \\ f(x_1, \dots, x_n, t_m) - y_m \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}. \quad (8.5.0.4)$$

With this notation we have the equivalence

$$(8.5.0.3) \Leftrightarrow \mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \sum_{i=1}^m |f(x_1, \dots, x_n, t_i) - y_i|^2 = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \|F(\mathbf{x})\|_2^2. \quad (8.5.0.5)$$

The previous example motivates the following definition generalizing Def. 3.1.1.1.

### Definition 8.5.0.6. Non-linear least-squares solution

Given  $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^m, m, n \in \mathbb{N}, m > n$ , we call  $\mathbf{x}^*$  a **non-linear least squares solution** of  $F(\mathbf{x}) = \mathbf{0}$ , if

$$\mathbf{x}^* = \underset{\mathbf{x} \in D}{\operatorname{argmin}} \|F(\mathbf{x})\|^2.$$

The search for such non-linear least-squares solutions is our current concern.

### Non-linear least squares problem

$$\begin{aligned} \text{Given: } & F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^m, \quad m, n \in \mathbb{N}, \quad m > n. \\ \text{Find: } & \mathbf{x}^* \in D: \quad \mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in D} \Phi(\mathbf{x}), \quad \Phi(\mathbf{x}) := \frac{1}{2} \|F(\mathbf{x})\|_2^2. \end{aligned} \quad (8.5.0.7)$$

Terminology:  $D \hat{=} \text{parameter space}$ ,  $x_1, \dots, x_n \hat{=} \text{parameter}$ .

As in the case of linear least squares problems ( $\rightarrow$  Section 3.1.1): a non-linear least squares problem is related to an overdetermined non-linear system of equations  $F(\mathbf{x}) = \mathbf{0}$ .

As for non-linear systems of equations discussed in the beginning of this chapter, existence and uniqueness of  $\mathbf{x}^*$  in (8.5.0.7) has to be established in each concrete case!

**Remark 8.5.0.8 (“Full-rank condition”)** Recall from Rem. 3.1.2.15, Ex. 3.1.2.17, Rem. 3.1.2.18 that for a linear least-squares problem  $\|\mathbf{Ax} - \mathbf{b}\| \rightarrow \min$  full rank of the matrix  $\mathbf{A} \in \mathbb{R}^{m,n}$  was linked to a “good model”, in which every parameter had an influence independently of the others.

Also in the non-linear setting we require “independence for each parameter”:

$$\exists \text{ neighbourhood } \mathcal{U}(\mathbf{x}^*) \text{ such that } D F(\mathbf{x}) \text{ has full rank } n \quad \forall \mathbf{x} \in \mathcal{U}(\mathbf{x}^*). \quad (8.5.0.9)$$

This means that the columns of the Jacobi matrix  $D F(\mathbf{x})$  must be linearly independent.

If (8.5.0.9) is not satisfied, then the parameters are redundant in the sense that fewer parameters would be enough to model the same dependence (locally at  $\mathbf{x}^*$ ), cf. Rem. 3.1.2.18.  $\square$

### 8.5.1 (Damped) Newton Method

We examine a first, natural approach to solving the non-linear least squares problem

$$\mathbf{x}^* \in D: \quad \mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in D} \Phi(\mathbf{x}), \quad \Phi(\mathbf{x}) := \frac{1}{2} \|F(\mathbf{x})\|_2^2. \quad (8.5.0.7)$$

We assume that  $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^m$  is *twice continuously differentiable*. Then the non-linear least-squares solution  $\mathbf{x}^*$  has to be a zero of the derivative of  $\mathbf{x} \mapsto \Phi(\mathbf{x})$ :

$$\Phi(\mathbf{x}^*) = \min \Rightarrow \mathbf{grad} \Phi(\mathbf{x}) = \mathbf{0}, \quad \text{where} \quad \mathbf{grad} \Phi(\mathbf{x}) := \left[ \frac{\partial \Phi}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial \Phi}{\partial x_n}(\mathbf{x}) \right]^\top \in \mathbb{R}^n.$$

Note that  $\mathbf{grad} \Phi : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ . The simple idea is to use Newton's method ( $\rightarrow$  Section 8.4) to solve the non-linear  $n \times n$  system of equations  $\mathbf{grad} \Phi(\mathbf{x}) = \mathbf{0}$ .

The Newton iteration (8.4.1.1) for non-linear system of equations  $\mathbf{grad} \Phi(\mathbf{x}) = \mathbf{0}$  is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{H} \Phi(\mathbf{x}^{(k)})^{-1} \mathbf{grad} \Phi(\mathbf{x}^{(k)}), \quad (8.5.1.1)$$

where  $\mathbf{H} \Phi(\mathbf{x}) \in \mathbb{R}^{n,n}$  is the **Hessian matrix** of  $\Phi$ , see Def. 8.4.1.11:

$$\mathbf{H} \Phi(\mathbf{x}) = \left[ \frac{\partial^2 \Phi}{\partial x_i \partial x_j} \right]_{i,j=1}^n \in \mathbb{R}^{n,n}.$$

Using the definition  $\Phi(\mathbf{x}) := \frac{1}{2} \|F(\mathbf{x})\|^2$  we can express  $\mathbf{grad} \Phi$  and  $\mathbf{H} \Phi$  in terms of  $F : \mathbb{R}^n \mapsto \mathbb{R}^n$ . First, since  $\Phi(\mathbf{x}) = (G \circ F)(\mathbf{x})$  with  $G : \mathbb{R}^m \rightarrow \mathbb{R}$ ,  $G(\mathbf{z}) := \frac{1}{2} \|\mathbf{z}\|^2$ , the chain rule (8.4.1.9) gives

$$\mathbf{D} \Phi(\mathbf{x}) \mathbf{h} = \mathbf{D} G(F(\mathbf{x})) \mathbf{D} F(\mathbf{x}) \mathbf{h} = F(\mathbf{x})^\top \mathbf{D} F(\mathbf{x}) \mathbf{h} \quad \mathbf{grad} \Phi(\mathbf{x}) = \mathbf{D} F(\mathbf{x})^\top F(\mathbf{x}) . \quad (8.5.1.2)$$

Second, we can apply the product rule (8.4.1.10) to  $\mathbf{x} \mapsto \mathbf{D} F(\mathbf{x})^\top F(\mathbf{x})$  and get

$$\begin{aligned} \mathbf{H} \Phi(\mathbf{x}) &:= \mathbf{D}(\mathbf{grad} \Phi)(\mathbf{x}) = \mathbf{D} F(\mathbf{x})^\top \mathbf{D} F(\mathbf{x}) + \sum_{j=1}^m F_j(\mathbf{x}) \mathbf{D}^2 F_j(\mathbf{x}) , \\ &\Updownarrow \\ (\mathbf{H} \Phi(\mathbf{x}))_{i,k} &= \sum_{j=1}^n \left\{ \frac{\partial^2 F_j}{\partial x_i \partial x_k}(\mathbf{x}) F_j(\mathbf{x}) + \frac{\partial F_j}{\partial x_k}(\mathbf{x}) \frac{\partial F_j}{\partial x_i}(\mathbf{x}) \right\} . \end{aligned} \quad (8.5.1.3)$$

We make the recommendation, *cf.* § 8.4.1.8, that when in doubt, the reader should differentiate components of matrices and vectors!

The above derivative formulas permit us to rewrite (8.5.1.1) in concrete terms. We obtain the Newton correction  $\mathbf{s} \in \mathbb{R}^n$  to the current Newton iterate  $\mathbf{x}^{(k)}$  by solving the  $n \times n$  linear system of equations

$$\underbrace{\left( \mathbf{D} F(\mathbf{x}^{(k)})^\top \mathbf{D} F(\mathbf{x}^{(k)}) + \sum_{j=1}^m F_j(\mathbf{x}^{(k)}) \mathbf{D}^2 F_j(\mathbf{x}^{(k)}) \right)}_{=\mathbf{H} \Phi(\mathbf{x}^{(k)})} \mathbf{s} = -\underbrace{\mathbf{D} F(\mathbf{x}^{(k)})^\top F(\mathbf{x}^{(k)})}_{=\mathbf{grad} \Phi(\mathbf{x}^{(k)})} . \quad (8.5.1.4)$$

All the techniques presented in Section 8.4 (damping, termination) can now be applied to the particular Newton iteration for  $\mathbf{grad} \Phi = \mathbf{0}$ . We refer to that section.

**Remark 8.5.1.5 (Newton method and minimization of quadratic functional)** Newton's method (8.5.1.1) for (8.5.0.7) can be read as *successive minimization* of a local **quadratic approximation** of  $\Phi$ :

$$\Phi(\mathbf{x}) \approx Q(\mathbf{s}) := \Phi(\mathbf{x}^{(k)}) + \mathbf{grad} \Phi(\mathbf{x}^{(k)})^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \mathbf{H} \Phi(\mathbf{x}^{(k)}) \mathbf{s} , \quad (8.5.1.6)$$

$$\mathbf{grad} Q(\mathbf{s}) = 0 \Leftrightarrow \mathbf{H} \Phi(\mathbf{x}^{(k)}) \mathbf{s} + \mathbf{grad} \Phi(\mathbf{x}^{(k)}) = 0 \Leftrightarrow (8.5.1.4) .$$

- So we deal with yet another model function method (→ Section 8.3.2) with quadratic model function  $Q$  for  $\Phi$ .

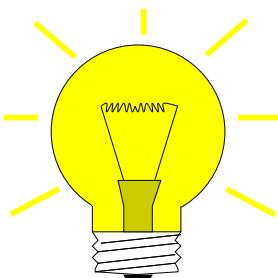
□

## 8.5.2 Gauss-Newton Method

The Newton method derived in Section 8.5.1 hinges on the availability of *second derivatives* of  $F$ . This compounds difficulties of implementation, in particular, if  $F$  is given only implicitly or in procedural form. Now we will learn about a method for the non-linear least-squares problem

$$\mathbf{x}^* \in D: \quad \mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in D} \Phi(\mathbf{x}) , \quad \Phi(\mathbf{x}) := \frac{1}{2} \|F(\mathbf{x})\|_2^2 , \quad (8.5.0.7)$$

that relies on first derivatives of  $F$  only.



Idea:

*Local linearization* of  $F$ :  $F(\mathbf{x}) \approx F(\mathbf{y}) + \mathbf{D} F(\mathbf{y})(\mathbf{x} - \mathbf{y})$

➤ Leads to a sequence of *linear* least squares problems

Employing local linearization we find that the minimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \|F(\mathbf{x})\|_2 \quad \text{is approximated by} \quad \underbrace{\underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \|F(\mathbf{x}_0) + D F(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)\|_2}_{(\spadesuit)} ,$$

where  $\mathbf{x}_0$  is an approximation of the solution  $\mathbf{x}^*$  of (8.5.0.7).

$$(\spadesuit) \Leftrightarrow \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \|\mathbf{Ax} - \mathbf{b}\| \quad \text{with} \quad \begin{aligned} \mathbf{A} &:= D F(\mathbf{x}_0) \in \mathbb{R}^{m,n}, \\ \mathbf{b} &:= -F(\mathbf{x}_0) + D F(\mathbf{x}_0)\mathbf{x}_0 \in \mathbb{R}^m . \end{aligned}$$

This is a linear least squares problem in the standard form given in Def. 3.1.1.1.

Note that by condition (8.5.0.9)  $\mathbf{A}$  has full rank, if  $\mathbf{x}_0$  sufficiently close to  $\mathbf{x}^*$ .

Also be aware that this approach is different from local quadratic approximation of  $\Phi$  underlying Newton's method for (8.5.0.7), see Section 8.5.1, Rem. 8.5.1.5.

The idea of local linearization leads to the **Gauss-Newton iteration** (under full-rank assumption (8.5.0.9))

$$\begin{aligned} \text{Initial guess } \mathbf{x}^{(0)} &\in D \\ \mathbf{x}^{(k+1)} &:= \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \left\| F(\mathbf{x}^{(k)}) + D F(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}) \right\|_2 . \end{aligned} \tag{8.5.2.1}$$

↑  
a linear least squares problem!

### C++ code 8.5.2.2: Template for Gauss-Newton method

```

2 template <class Function, class Jacobian>
3 VectorXd gn(const Eigen::VectorXd &init, Function &&F, Jacobian &&J,
4             double rtol = 1.0E-6, double atol = 1.0E-8) {
5     Eigen::VectorXd x = init; // Vector for iterates  $\mathbf{x}^{(k)}$ 
6     // Vector for Gauss-Newton correction  $\mathbf{s}$ 
7     Eigen::VectorXd s = J(x).householderQr().solve(F(x)); //
8     x = x - s;
9     // A posteriori termination based on absolute and relative tolerances
10    while ((s.norm() > rtol * x.norm()) && (s.norm() > atol)) { //
11        s = J(x).householderQr().solve(F(x)); //
12        x = x - s;
13    }
14    return x;
15 }
```

Comments on Code 8.5.2.2:

- ☞ Recall EIGEN's built-in functions for solving linear least-squares problems used in Line 7 and Line 11, see Code 3.3.4.2.
- ☞ Argument  $\mathbf{x}$  passes initial guess  $\mathbf{x}^{(0)} \in \mathbb{R}^n$ , argument  $F$  must be a *handle* to a function  $F : \mathbb{R}^n \mapsto \mathbb{R}^m$ , argument  $J$  provides the Jacobian of  $F$ , namely  $D F : \mathbb{R}^n \mapsto \mathbb{R}^{m,n}$ , arguments  $rtol$  and  $atol$  specify the tolerances for termination
- ☞ Line 10: iteration terminates if relative norm of correction is below threshold specified in  $tol$ .

Note that the function of Code 8.5.2.2 also implements Newton's method (→ Section 8.4.1) in the case  $m = n$ !

**Remark 8.5.2.3 (Gauss-Newton versus Newton)** Let us summarize the pros and cons of using the Gauss-Newton approach:

Advantage of the Gauss-Newton method : second derivative of  $F$  not needed.

Drawback of the Gauss-Newton method : no local quadratic convergence.

**EXAMPLE 8.5.2.4 (Non-linear data fitting (II))** → **Ex. 8.5.0.1**) We consider the non-linear data fitting problem (8.5.0.5) for the parameterized function  $f(x_1, x_2, x_3; t) := x_1 + x_2 \exp(-x_3 t)$ . In this case

$$F(\mathbf{x}) = \begin{bmatrix} x_1 + x_2 \exp(-x_3 t_1) - y_1 \\ \vdots \\ x_1 + x_2 \exp(-x_3 t_m) - y_m \end{bmatrix} : \mathbb{R}^3 \mapsto \mathbb{R}^m, \quad D F(\mathbf{x}) = \begin{bmatrix} 1 & e^{-x_3 t_1} & -x_2 t_1 e^{-x_3 t_1} \\ \vdots & \vdots & \vdots \\ 1 & e^{-x_3 t_m} & -x_2 t_m e^{-x_3 t_m} \end{bmatrix}$$

We use data points  $(t_j, y_j)$ ,  $j = 1, \dots, m$ ,  $m = 21$ ,  $t_j = 1 + 0.3j$ ,  $j = 1, \dots, 21$  and the  $y_j$  generated by the following C++ code

```
Eigen::VectorXd gnrndinit(const Eigen::VectorXd &x) {
    std::srand((unsigned int)time(0));
    auto t = Eigen::VectorXd::LinSpaced((7.0 - 1.0) / 0.3 - 1, 1.0,
                                         7.0);
    auto y = x(0) + x(1) * ((-x(2) * t).array().exp());
    return y + 0.1 * (Eigen::VectorXd::Random(y.size()).array() - 0.5);
}
```

We study the convergence of the Newton method, damped Newton method (→ Section 8.4.4) and Gauss-Newton method for different initial guesses.

- ◆ initial value  $(1.8, 1.8, 0.1)^T$  (red curves, blue curves)
- ◆ initial value  $(1.5, 1.5, 0.1)^T$  (cyan curves, green curves)

First experiment (→ Section 8.5.1): iterative solution of non-linear least squares data fitting problem by means of the Newton method (8.5.1.4) and the damped Newton method from Code 8.4.4.5.

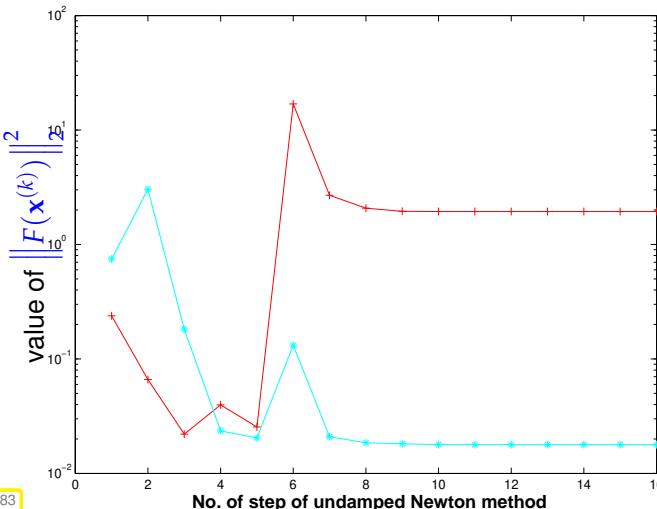


Fig. 283

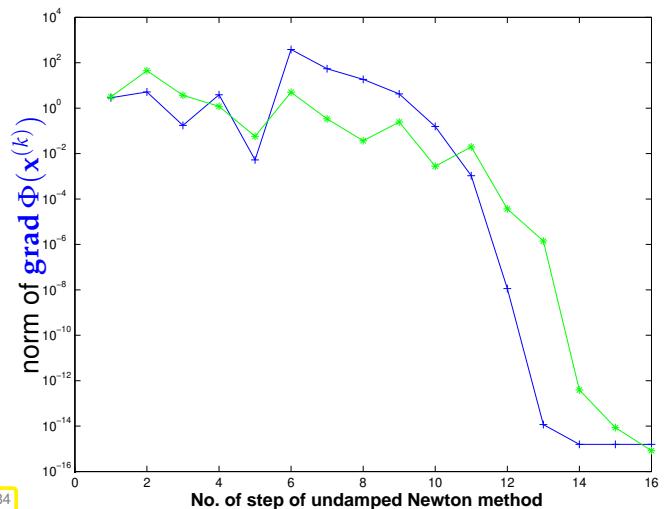
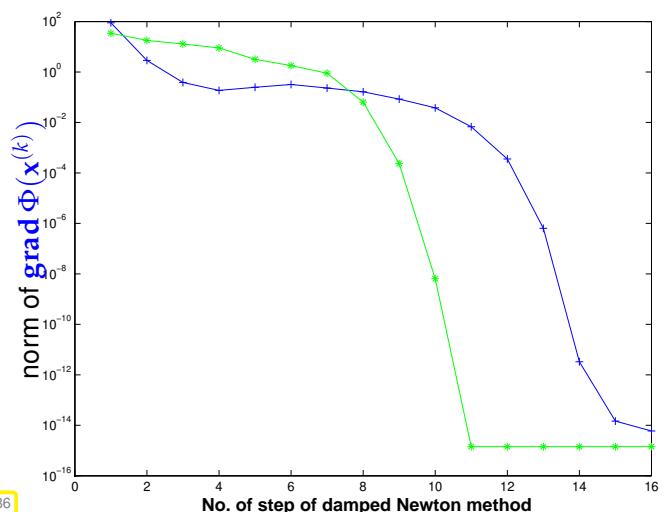
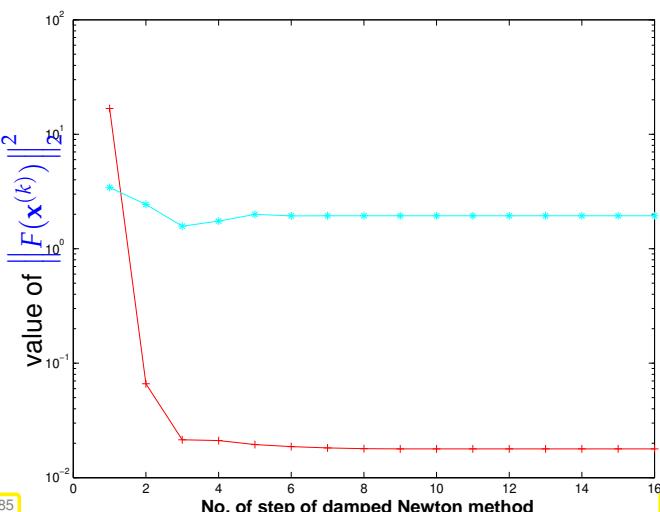


Fig. 284

Convergence behaviour of plain Newton method:

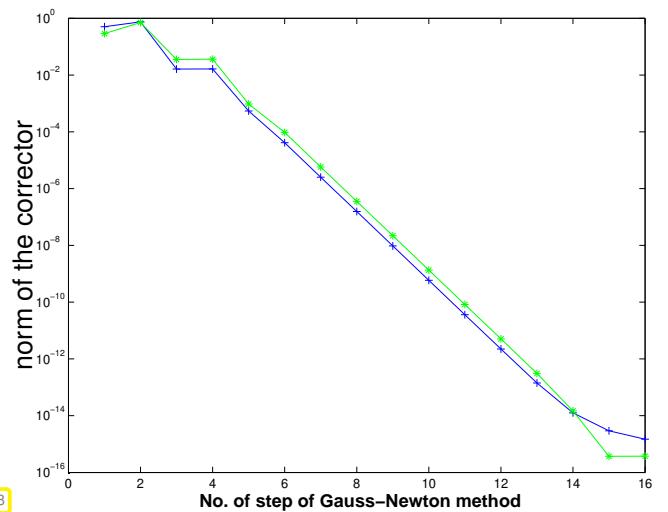
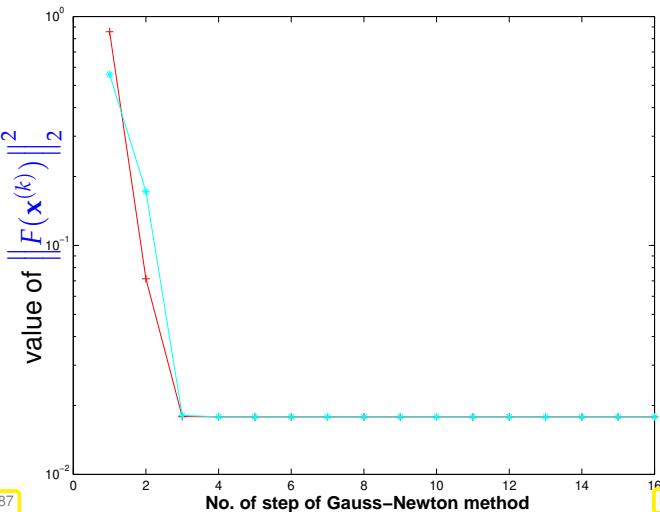
- initial value  $(1.8, 1.8, 0.1)^T$  (red curve) ► Newton method caught in **local minimum**,  
 initial value  $(1.5, 1.5, 0.1)^T$  (cyan curve) ► fast (locally quadratic) convergence.



Convergence behavior of damped Newton method:

initial value  $(1.8, 1.8, 0.1)^T$  (red curve) ➤ fast (locally quadratic) convergence,  
initial value  $(1.5, 1.5, 0.1)^T$  (cyan curve) ➤ Newton method caught in local minimum.

Second experiment: iterative solution of non-linear least squares data fitting problem by means of the Gauss-Newton method (8.5.2.1), see Code 8.5.2.2.



We observe: linear convergence for all initial values, cf. Def. 8.1.1.1, Rem. 8.1.1.6. □

### 8.5.3 Trust Region Method (Levenberg-Marquardt Method)

As in the case of Newton's method for non-linear systems of equations, see Section 8.4.4, often overshooting of Gauss-Newton corrections occurs. We can resort to a similar remedy as in the case of Newton's method: **damping** of Gauss-Newton corrections.

Idea: damping of the Gauss-Newton correction in (8.5.2.1) using a **penalty term**

$$\text{instead of } \left\| F(\mathbf{x}^{(k)}) + D F(\mathbf{x}^{(k)}) \mathbf{s} \right\|^2 \text{ minimize } \left\| F(\mathbf{x}^{(k)}) + D F(\mathbf{x}^{(k)}) \mathbf{s} \right\|^2 + \lambda \|\mathbf{s}\|_2^2,$$

where  $\lambda > 0$  is a penalty parameter.

A central issue is the choice of the penalty parameter  $\lambda > 0$ . Since rigorous rules are elusive, heuristic

strategies are used in practice, for instance

$$\lambda = \gamma \|F(\mathbf{x}^{(k)})\|_2, \quad \gamma := \begin{cases} 10 & , \text{ if } \|F(\mathbf{x}^{(k)})\|_2 \geq 10, \\ 1 & , \text{ if } 1 < \|F(\mathbf{x}^{(k)})\|_2 < 10, \\ 0.01 & , \text{ if } \|F(\mathbf{x}^{(k)})\|_2 \leq 1. \end{cases} \quad (8.5.3.1)$$

The minimization problem

$$\mathbf{s} = \underset{\mathbf{z} \in \mathbb{R}^n}{\operatorname{argmin}} \|F(\mathbf{x}^{(k)}) + D F(\mathbf{x}^{(k)}) \mathbf{z}\|^2 + \lambda \|\mathbf{z}\|_2^2$$

is a standard linear least-squares problem and leads to the following  $n \times n$  linear system of **normal equations** for the damped Gauss-Newton correction  $\mathbf{s}$ , see Thm. 3.1.2.1:

$$(D F(\mathbf{x}^{(k)})^T D F(\mathbf{x}^{(k)}) + \lambda I) \mathbf{s} = -D F(\mathbf{x}^{(k)}) F(\mathbf{x}^{(k)}). \quad (8.5.3.2)$$

## Learning Outcomes

- Knowledge about concepts related to the speed of convergence of an iteration for solving a non-linear system of equations.
- Ability to estimate type and orders of convergence from empiric data.
- Ability to predict asymptotic linear, quadratic and cubic convergence by inspection of the iteration function.
- Familiarity with (damped) Newton's method for general non-linear systems of equations and with the secant method in 1D.
- Ability to derive the Newton iteration for an (implicitly) given non-linear system of equations.
- Knowledge about quasi-Newton method as multi-dimensional generalizations of the secant method.

# Bibliography

- [Ale12] A. Alexanderian. *A basic note on iterative matrix inversion*. Online document. 2012 (cit. on p. 530).
- [AG11] Uri M. Ascher and Chen Greif. *A first course in numerical methods*. Vol. 7. Computational Science & Engineering. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2011, pp. xxii+552. DOI: [10.1137/1.9780898719987](https://doi.org/10.1137/1.9780898719987) (cit. on pp. 495, 498, 507, 509, 516, 544).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 493, 497, 498, 502, 503, 507, 509, 516, 531, 544–551).
- [DH03] P. Deuflhard and A. Hohmann. *Numerical Analysis in Modern Scientific Computing*. Vol. 43. Texts in Applied Mathematics. Springer, 2003 (cit. on p. 533).
- [Deu11] Peter Deuflhard. *Newton methods for nonlinear problems*. Vol. 35. Springer Series in Computational Mathematics. Heidelberg: Springer, 2011, pp. xii+424. DOI: [10.1007/978-3-642-23899-4](https://doi.org/10.1007/978-3-642-23899-4) (cit. on pp. 524, 533, 537, 544).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 490, 493, 503, 509, 512, 516, 517, 544).
- [Mol04] C. Moler. *Numerical Computing with MATLAB*. Philadelphia, PA: SIAM, 2004 (cit. on p. 519).
- [PS91] Victor Pan and Robert Schreiber. “An Improved Newton Iteration for the Generalized Inverse of a Matrix, with Applications”. In: *SIAM Journal on Scientific and Statistical Computing* 12.5 (1991), pp. 1109–1130. DOI: [10.1137/0912058](https://doi.org/10.1137/0912058) (cit. on p. 530).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 493, 498, 508, 539).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on pp. 494, 499, 502, 503, 505, 507, 511, 518, 526, 530).
- [Wer92] J. Werner. *Numerische Mathematik I. Lineare und nichtlineare Gleichungssysteme, Interpolation, numerische Integration*. vieweg studium. Aufbaukurs Mathematik. Braunschweig: Vieweg, 1992 (cit. on p. 539).

# Chapter 9

## Computation of Eigenvalues and Eigenvectors

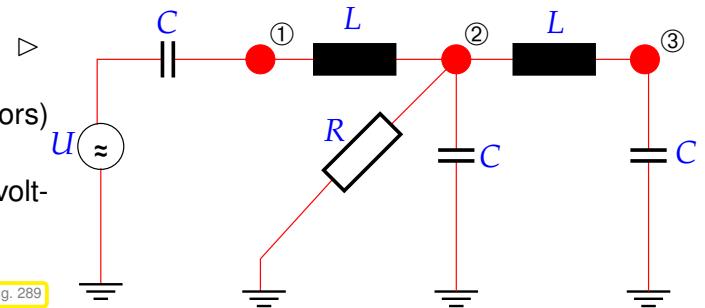


*Supplementary literature.* [Bai+00] offers comprehensive presentation of numerical methods

for the solution of eigenvalue problems from an algorithmic point of view.

### EXAMPLE 9.0.0.1 (Resonances of linear electric circuits)

Simple electric circuit, cf. Ex. 2.1.0.3



- ◆ linear components (resistors, coils, capacitors) only,
- ◆ time-harmonic excitation (alternating voltage/current)
- ◆ “frequency domain” circuit model

Ex. 2.1.0.3: nodal analysis of linear ( $\leftrightarrow$  composed of resistors, inductors, capacitors) electric circuit **in frequency domain (at angular frequency  $\omega > 0$ )**, see (2.1.0.6)

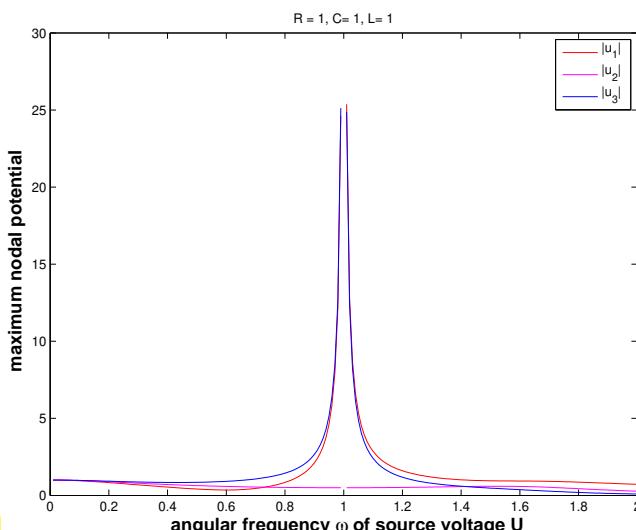
➢ linear system of equations for nodal potentials with complex system matrix  $\mathbf{A}$

For circuit of Code 9.0.0.3: three unknown nodal potentials

➢ system matrix from nodal analysis at angular frequency  $\omega > 0$ :

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} i\omega C + \frac{1}{i\omega L} & -\frac{1}{i\omega L} & 0 \\ -\frac{1}{i\omega L} & i\omega C + \frac{1}{R} + \frac{2}{i\omega L} & -\frac{1}{i\omega L} \\ 0 & -\frac{1}{i\omega L} & i\omega C + \frac{1}{i\omega L} \end{bmatrix} \\ &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{1}{R} & 0 \\ 0 & 0 & 0 \end{bmatrix} + i\omega \begin{bmatrix} C & 0 & 0 \\ 0 & C & 0 \\ 0 & 0 & C \end{bmatrix} + 1/i\omega \begin{bmatrix} \frac{1}{L} & -\frac{1}{L} & 0 \\ -\frac{1}{L} & \frac{2}{L} & -\frac{1}{L} \\ 0 & -\frac{1}{L} & \frac{1}{L} \end{bmatrix}. \end{aligned}$$

$$\mathbf{A}(\omega) := \mathbf{W} + i\omega \mathbf{C} - i\omega^{-1} \mathbf{S} \quad , \quad \mathbf{W}, \mathbf{C}, \mathbf{S} \in \mathbb{R}^{n,n} \text{ symmetric .} \quad (9.0.0.2)$$



$\triangleleft$  plot of  $|u_i(U)|$ ,  $i = 1, 2, 3$  for  $R = L = C = 1$  (scaled model)

Blow-up of some nodal potentials for certain  $\omega$  !

### MATLAB-code 9.0.0.3: Computation of nodal potential for circuit of Code 9.0.0.3

```

1 function rescirc(R,L,C)
2 % Ex. ???: Numerical nodal analysis of the resonant circuit
3 % R, L, C ≈ network component parameters
4
5 Z = 1/R; K = 1/L;
6
7 % Matrices W, C, S for nodal analysis of circuit
8 Wmat = [0 0 0; 0 Z 0; 0 0 0];
9 Cmat = [C 0 0; 0 C 0; 0 0 C];
10 Smat = [K -K 0; -K 2*K -K; 0 -K K];
11 % System matrix from nodal analysis
12 Amat = @(w) (Wmat+i*w*Cmat+Smat/(i*w));
13
14 % Scanning source currents
15 res = [];
16 for w=0.01:0.01:2
17 res = [res; w, abs(Amat(w)\[C;0;0])'];
18 end
19
20 figure('name','resonant circuit');
21 plot(res(:,1),res(:,2),'r-',res(:,1),res(:,3),'m-',res(:,1),res(:,4),'b-')
22 xlabel('{\bf angular frequency \omega of source voltage U}', 'fontsize', 14);
23 ylabel('{\bf maximum nodal potential}', 'fontsize', 14);
24 title(sprintf('R = %d, C= %d, L= %d', R,L,C));
25 legend('|u_1|', '|u_2|', '|u_3|');
26
27 print -depsc2 '../PICTURES/rescircpot.eps'
28
29 % Solving generalized eigenvalue problem (9.0.0.5)
30 Zmat = zeros(3,3); Imat = eye(3,3);
31 % Assemble 6×6-matrices M and B
32 Mmat = [Wmat,Smat; Imat, Zmat];
33 Bmat = [-i*Cmat, Zmat; Zmat , i*Imat];

```

```

34 % Solve generalized eigenvalue problem, cf. (9.0.0.6)
35 omega = eig(Mmat,Bmat);
36
37 figure('name','resonances');
38 plot(real(omega),imag(omega),'r*'); hold on;
39 ax = axis;
40 plot([ax(1) ax(2)], [0 0], 'k-');
41 plot([ 0 0], [ax(3) ax(4)], 'k-');
42 grid on;
43 xlabel('{\bf Re}(\omega)', 'fontsize',14);
44 ylabel('{\bf Im}(\omega)', 'fontsize',14);
45 title(sprintf('R = %d, C= %d, L= %d', R,L,C));
46 legend('omega');

47
48 print -depsc2 '../PICTURES/rescirmomega.eps'

```

$$\text{resonant frequencies} = \omega \in \{\omega \in \mathbb{R}: \mathbf{A}(\omega) \text{ singular}\}$$

If the circuit is operated at a real resonant frequency, the circuit equations will not possess a solution. Of course, the real circuit will always behave in a well-defined way, but the linear model will break down due to extremely large currents and voltages. In an experiment this breakdown manifests itself as a rather explosive meltdown of circuits components. Hence, it is vital to determine resonant frequencies of circuits in order to avoid their destruction.

→ relevance of numerical methods for solving:

$$\text{Find } \omega \in \mathbb{C} \setminus \{0\}: \mathbf{W} + \iota\omega\mathbf{C} + \frac{1}{\iota\omega}\mathbf{S} \text{ singular.}$$

This is a **quadratic eigenvalue problem**: find  $\mathbf{x} \neq 0, \omega \in \mathbb{C} \setminus \{0\}$ ,

$$\mathbf{A}(\omega)\mathbf{x} = (\mathbf{W} + \iota\omega\mathbf{C} + \frac{1}{\iota\omega}\mathbf{S})\mathbf{x} = 0. \quad (9.0.0.4)$$

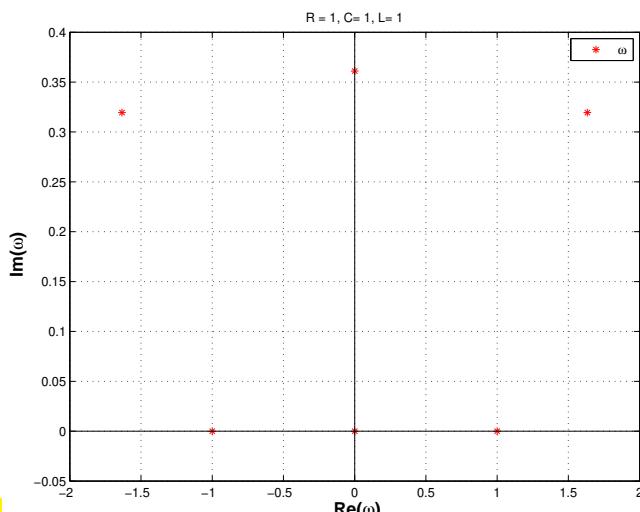
Substitution:  $\mathbf{y} = \frac{1}{\iota\omega}\mathbf{x} \leftrightarrow \mathbf{x} = \iota\omega\mathbf{y}$  [TM01, Sect. 3.4]:

$$(9.0.0.4) \Leftrightarrow \underbrace{\begin{bmatrix} \mathbf{W} & \mathbf{S} \\ \mathbf{I} & 0 \end{bmatrix}}_{:=\mathbf{M}} \underbrace{\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}}_{:=\mathbf{z}} = \omega \underbrace{\begin{bmatrix} -\iota\mathbf{C} & 0 \\ 0 & -\iota\mathbf{I} \end{bmatrix}}_{:=\mathbf{B}} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \quad (9.0.0.5)$$

➤ **generalized linear eigenvalue problem** of the form: find  $\omega \in \mathbb{C}, \mathbf{z} \in \mathbb{C}^{2n} \setminus \{0\}$  such that

$$\mathbf{M}\mathbf{z} = \omega\mathbf{B}\mathbf{z}. \quad (9.0.0.6)$$

In this example one is mainly interested in the **eigenvalues**  $\omega$ , whereas the **eigenvectors**  $\mathbf{z}$  usually need not be computed.



▷ resonant frequencies for circuit from Code 9.0.0.3  
(including decaying modes with  $\text{Im}(\omega) > 0$ )

### EXAMPLE 9.0.0.7 (Analytic solution of homogeneous linear ordinary differential equations → [Str09, Remark 5.6.1], [Gut09, Sect. 10.1], [NS02, Sect. 8.1], [DR08, Ex. 7.3])

Autonomous homogeneous linear ordinary differential equation (ODE):

$$\dot{\mathbf{y}} = \mathbf{A}\mathbf{y} , \quad \mathbf{A} \in \mathbb{C}^{n,n} . \quad (9.0.0.8)$$

$$\mathbf{A} = \mathbf{S} \underbrace{\begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}}_{=: \mathbf{D}} \mathbf{S}^{-1} , \quad \mathbf{S} \in \mathbb{C}^{n,n} \text{ regular} \implies (\dot{\mathbf{y}} = \mathbf{A}\mathbf{y} \xrightarrow{\mathbf{z} = \mathbf{S}^{-1}\mathbf{y}} \dot{\mathbf{z}} = \mathbf{D}\mathbf{z}) .$$

➤ solution of initial value problem:

$$\dot{\mathbf{y}} = \mathbf{A}\mathbf{y} , \quad \mathbf{y}(0) = \mathbf{y}_0 \in \mathbb{C}^n \Rightarrow \mathbf{y}(t) = \mathbf{S}\mathbf{z}(t) , \quad \dot{\mathbf{z}} = \mathbf{D}\mathbf{z} , \quad \mathbf{z}(0) = \mathbf{S}^{-1}\mathbf{y}_0 .$$

The initial value problem for the *decoupled* homogeneous linear ODE  $\dot{\mathbf{z}} = \mathbf{D}\mathbf{z}$  has a simple analytic solution

$$\mathbf{z}_i(t) = \exp(\lambda_i t)(\mathbf{z}_0)_i = \exp(\lambda_i t)((\mathbf{S}^{-1})_{i,:}^T \mathbf{y}_0) .$$

In light of Rem. 1.3.1.3:

$$\mathbf{A} = \mathbf{S} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \mathbf{S}^{-1} \Leftrightarrow \mathbf{A}((\mathbf{S})_{:,i}) = \lambda_i((\mathbf{S})_{:,i}) \quad i = 1, \dots, n . \quad (9.0.0.9)$$

In order to find the transformation matrix  $\mathbf{S}$  all non-zero solution vectors (= **eigenvectors**)  $\mathbf{x} \in \mathbb{C}^n$  of the **linear eigenvalue problem**

$$\mathbf{Ax} = \lambda \mathbf{x}$$

have to be found.

## Contents

9.1	Theory of eigenvalue problems	557
9.2	"Direct" Eigensolvers	559
9.3	Power Methods	563
9.3.1	Direct power method	563
9.3.2	Inverse Iteration [DR08, Sect. 7.6], [QSS00, Sect. 5.3.2]	572
9.3.3	Preconditioned inverse iteration (PINVIT)	584
9.3.4	Subspace iterations	587
9.4	Krylov Subspace Methods	601

## 9.1 Theory of eigenvalue problems



*Supplementary literature.* [NS02, Ch. 7], [Gut09, Ch. 9], [QSS00, Sect. 1.7]

**Definition 9.1.0.1. Eigenvalues and eigenvectors** → [NS02, Sects. 7.1, 7.2], [Gut09, Sect. 9.1]

- $\lambda \in \mathbb{C}$  eigenvalue (ger.: Eigenwert) of  $\mathbf{A} \in \mathbb{K}^{n,n}$  : $\Leftrightarrow$   $\underbrace{\det(\lambda\mathbf{I} - \mathbf{A})}_\text{characteristic polynomial } \chi(\lambda) = 0$
- spectrum of  $\mathbf{A} \in \mathbb{K}^{n,n}$ :  $\sigma(\mathbf{A}) := \{\lambda \in \mathbb{C}: \lambda \text{ eigenvalue of } \mathbf{A}\}$
- eigenspace (ger.: Eigenraum) associated with eigenvalue  $\lambda \in \sigma(\mathbf{A})$ :  
 $\mathbf{Eig}\mathbf{A}\lambda := \mathcal{N}\lambda\mathbf{I} - \mathbf{A}$
- $\mathbf{x} \in \mathbf{Eig}\mathbf{A}\lambda \setminus \{0\} \Rightarrow \mathbf{x}$  is eigenvector
- Geometric multiplicity (ger.: Vielfachheit) of an eigenvalue  $\lambda \in \sigma(\mathbf{A})$ :  
 $m(\lambda) := \dim \mathbf{Eig}\mathbf{A}\lambda$

Two simple facts:

$$\lambda \in \sigma(\mathbf{A}) \Rightarrow \dim \mathbf{Eig}\mathbf{A}\lambda > 0, \quad (9.1.0.2)$$

$$\det(\mathbf{A}) = \det(\mathbf{A}^T) \quad \forall \mathbf{A} \in \mathbb{K}^{n,n} \Rightarrow \sigma(\mathbf{A}) = \sigma(\mathbf{A}^T). \quad (9.1.0.3)$$

notation:  $\rho(\mathbf{A}) := \max\{|\lambda|: \lambda \in \sigma(\mathbf{A})\} \triangleq$  spectral radius of  $\mathbf{A} \in \mathbb{K}^{n,n}$

**Theorem 9.1.0.4. Bound for spectral radius**

For any matrix norm  $\|\cdot\|$  induced by a vector norm (→ Def. 1.5.5.10)

$$\rho(\mathbf{A}) \leq \|\mathbf{A}\|.$$

*Proof.* Let  $\mathbf{z} \in \mathbb{C}^n \setminus \{0\}$  be an eigenvector to the largest (in modulus) eigenvalue  $\lambda$  of  $\mathbf{A} \in \mathbb{C}^{n,n}$ . Then

$$\|\mathbf{A}\| := \sup_{\mathbf{x} \in \mathbb{C}^{n,n} \setminus \{0\}} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} \geq \frac{\|\mathbf{Az}\|}{\|\mathbf{z}\|} = |\lambda| = \rho(\mathbf{A}).$$

□

**Lemma 9.1.0.5. Gershgorin circle theorem** → [DR08, Thm. 7.13], [Han02, Thm. 32.1], [QSS00, Sect. 5.1]

For any  $\mathbf{A} \in \mathbb{K}^{n,n}$  holds true

$$\sigma(\mathbf{A}) \subset \bigcup_{j=1}^n \left\{ z \in \mathbb{C} : |z - a_{jj}| \leq \sum_{i \neq j} |a_{ji}| \right\}.$$

**Lemma 9.1.0.6. Similarity and spectrum** → [Gut09, Thm. 9.7], [DR08, Lemma 7.6], [NS02, Thm. 7.2]

The spectrum of a matrix is invariant with respect to *similarity transformations*:

$$\forall \mathbf{A} \in \mathbb{K}^{n,n}: \quad \sigma(\mathbf{S}^{-1}\mathbf{AS}) = \sigma(\mathbf{A}) \quad \forall \text{ regular } \mathbf{S} \in \mathbb{K}^{n,n}.$$

**Lemma 9.1.0.7.**

Existence of a one-dimensional invariant subspace

$$\forall \mathbf{C} \in \mathbb{C}^{n,n}: \quad \exists \mathbf{u} \in \mathbb{C}^n: \quad \mathbf{C}(\text{Span}\{\mathbf{u}\}) \subset \text{Span}\{\mathbf{u}\}.$$

**Theorem 9.1.0.8. Schur normal form** → [Hac94, Thm .2.8.1]

$$\forall \mathbf{A} \in \mathbb{K}^{n,n}: \quad \exists \mathbf{U} \in \mathbb{C}^{n,n} \text{ unitary: } \mathbf{U}^H \mathbf{A} \mathbf{U} = \mathbf{T} \quad \text{with } \mathbf{T} \in \mathbb{C}^{n,n} \text{ upper triangular}.$$

**Corollary 9.1.0.9. Principal axis transformation**

$$\mathbf{A} \in \mathbb{K}^{n,n}, \mathbf{A} \mathbf{A}^H = \mathbf{A}^H \mathbf{A}: \quad \exists \mathbf{U} \in \mathbb{C}^{n,n} \text{ unitary: } \mathbf{U}^H \mathbf{A} \mathbf{U} = \text{diag}(\lambda_1, \dots, \lambda_n), \quad \lambda_i \in \mathbb{C}.$$

A matrix  $\mathbf{A} \in \mathbb{K}^{n,n}$  with  $\mathbf{A} \mathbf{A}^H = \mathbf{A}^H \mathbf{A}$  is called *normal*.

- Examples of normal matrices are
- Hermitian matrices:  $\mathbf{A}^H = \mathbf{A}$  ➤  $\sigma(\mathbf{A}) \subset \mathbb{R}$
  - unitary matrices:  $\mathbf{A}^H = \mathbf{A}^{-1}$  ➤  $|\sigma(\mathbf{A})| = 1$
  - skew-Hermitian matrices:  $\mathbf{A} = -\mathbf{A}^H$  ➤  $\sigma(\mathbf{A}) \subset i\mathbb{R}$



Normal matrices can be diagonalized by *unitary* similarity transformations

Symmetric real matrices can be diagonalized by *orthogonal* similarity transformations

In Cor. 9.1.0.9:

- $\lambda_1, \dots, \lambda_n$  = eigenvalues of  $\mathbf{A}$
- Columns of  $\mathbf{U}$  = orthonormal basis of eigenvectors of  $\mathbf{A}$

Classes of relevant **eigenvalue problems (EVP)**:

- ① Given  $\mathbf{A} \in \mathbb{K}^{n,n}$  find **all eigenvalues** (= spectrum of  $\mathbf{A}$ ).
- ② Given  $\mathbf{A} \in \mathbb{K}^{n,n}$  find  $\sigma(\mathbf{A})$  plus **all eigenvectors**.
- ③ Given  $\mathbf{A} \in \mathbb{K}^{n,n}$  find **a few** eigenvalues and associated eigenvectors.

(Linear) **generalized eigenvalue problem**:

Given  $\mathbf{A} \in \mathbb{C}^{n,n}$ , regular  $\mathbf{B} \in \mathbb{C}^{n,n}$ , seek  $\mathbf{x} \neq 0, \lambda \in \mathbb{C}$

$$\mathbf{Ax} = \lambda \mathbf{Bx} \Leftrightarrow \mathbf{B}^{-1} \mathbf{Ax} = \lambda \mathbf{x}. \quad (9.1.0.10)$$

$\mathbf{x} \triangleq$  generalized eigenvector,  $\lambda \triangleq$  generalized eigenvalue

Obviously every generalized eigenvalue problem is equivalent to a standard eigenvalue problem

$$\mathbf{Ax} = \lambda \mathbf{Bx} \Leftrightarrow \mathbf{B}^{-1} \mathbf{Ax} = \lambda \mathbf{x}.$$

However, usually it is not advisable to use this equivalence for numerical purposes!

**Remark 9.1.0.11 (Generalized eigenvalue problems and Cholesky factorization)**

If  $\mathbf{B} = \mathbf{B}^H$  s.p.d. ( $\rightarrow$  Def. 1.1.2.6) with Cholesky factorization  $\mathbf{B} = \mathbf{R}^H \mathbf{R}$

$$\mathbf{Ax} = \lambda \mathbf{Bx} \Leftrightarrow \tilde{\mathbf{A}}\mathbf{y} = \lambda \mathbf{y} \text{ where } \tilde{\mathbf{A}} := \mathbf{R}^{-H} \mathbf{A} \mathbf{R}^{-1}, \mathbf{y} := \mathbf{Rx}.$$

→ This transformation can be used for efficient computations. □

## 9.2 “Direct” Eigensolvers

Purpose: solution of eigenvalue problems ①, ② for **dense** matrices “up to machine precision”

MATLAB-function:

`eig`

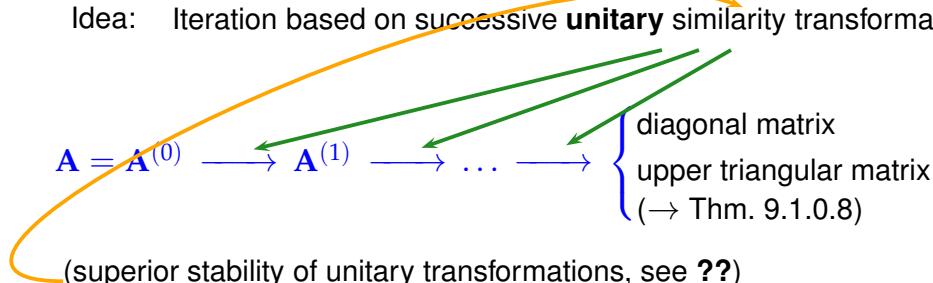
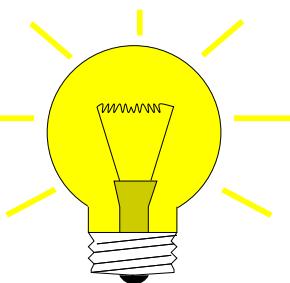
$\mathbf{d} = \text{eig}(\mathbf{A})$  : computes spectrum  $\sigma(\mathbf{A}) = \{d_1, \dots, d_n\}$  of  $\mathbf{A} \in \mathbb{C}^{n,n}$   
 $[\mathbf{V}, \mathbf{D}] = \text{eig}(\mathbf{A})$  : computes  $\mathbf{V} \in \mathbb{C}^{n,n}$ , *diagonal*  $\mathbf{D} \in \mathbb{C}^{n,n}$  such that  $\mathbf{AV} = \mathbf{VD}$

**Remark 9.2.0.1 (QR-Algorithm** → [GV89, Sect. 7.5], [NS02, Sect. 10.3], [Han02, Ch. 26], [QSS00, Sect. 5.5-5.7])

Note:

All “direct” eigensolvers are iterative methods

Idea: Iteration based on successive **unitary** similarity transformations



► QR-algorithm (with shift)

- ◆ in general: quadratic convergence
  - ◆ cubic convergence for norm matrices
- (→ [GV89, Sect. 7.5, 8.2])

### MATLAB-code 9.2.0.2: QR-algorithm with shift

```

1 function d = eigqr(A,tol)
2 n = size(A,1);
3 while (norm(tril(A,-1)) >
4 tol*norm(A))
5 % shift by eigenvalue of lower right 2x2
6 % block closest to (A)n,n
7 sc = eig(A(n-1:n,n-1:n));
8 [dummy,si] = min(abs(sc-A(n,n)));
9 shift = sc(si);
10 [Q,R] = qr( A - shift * eye(n));
11 A = Q' * A * Q;
12 end
13 d = diag(A);

```

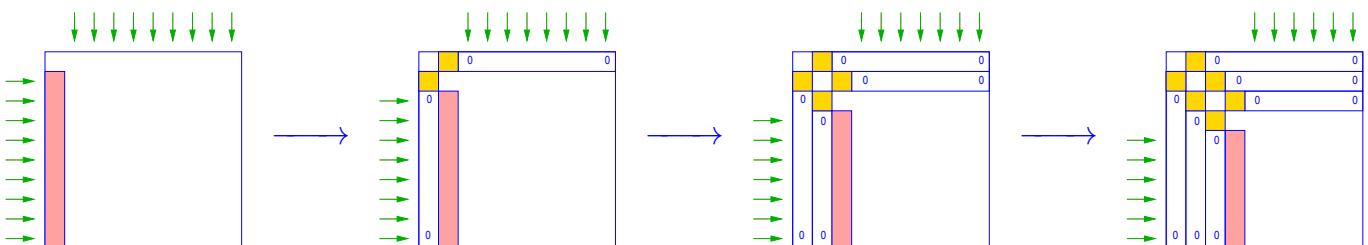
Computational cost:  $O(n^3)$  operations per step of the QR-algorithm

Library implementations of the QR-algorithm provide *numerically stable* eigensolvers (→ Def.1.5.5.19)

### Remark 9.2.0.3 (Unitary similarity transformation to tridiagonal form)

Successive Householder similarity transformations of  $\mathbf{A} = \mathbf{A}^H$ :

(→  $\hat{\wedge}$  affected rows/columns,  $\hat{\square}$   $\hat{\wedge}$  targeted vector)



► transformation to tridiagonal form ! (for general matrices a similar strategy can achieve a similarity transformation to upper Hessenberg form)

► this transformation is used as a preprocessing step for QR-algorithm ➤ **eig**.

Similar functionality for generalized EVP  $\mathbf{Ax} = \lambda \mathbf{Bx}$ ,  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$

$\mathbf{d} = \mathbf{eig}(\mathbf{A}, \mathbf{B})$  : computes all generalized eigenvalues  
 $[\mathbf{V}, \mathbf{D}] = \mathbf{eig}(\mathbf{A}, \mathbf{B})$  : computes  $\mathbf{V} \in \mathbb{C}^{n,n}$ , diagonal  $\mathbf{D} \in \mathbb{C}^{n,n}$  such that  $\mathbf{AV} = \mathbf{BVD}$

Note: (Generalized) eigenvectors can be recovered as columns of  $\mathbf{V}$ :

$$\mathbf{AV} = \mathbf{VD} \Leftrightarrow \mathbf{A}(\mathbf{V})_{:,i} = (\mathbf{D})_{i,i} \mathbf{V}_{:,i},$$

if  $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$ .

**Remark 9.2.0.4 (Computational effort for eigenvalue computations)**Computational effort (#elementary operations) for `eig()`:

eigenvalues & eigenvectors of $\mathbf{A} \in \mathbb{K}^{n,n}$	$\sim 25n^3 + O(n^2)$	$\left. \begin{array}{l} \sim 25n^3 + O(n^2) \\ \sim 10n^3 + O(n^2) \\ \sim 9n^3 + O(n^2) \\ \sim \frac{4}{3}n^3 + O(n^2) \\ \sim 30n^2 + O(n) \end{array} \right\}$
only eigenvalues of $\mathbf{A} \in \mathbb{K}^{n,n}$	$\sim 10n^3 + O(n^2)$	
eigenvalues and eigenvectors $\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n}$	$\sim 9n^3 + O(n^2)$	
only eigenvalues of $\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n}$	$\sim \frac{4}{3}n^3 + O(n^2)$	
only eigenvalues of tridiagonal $\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n}$	$\sim 30n^2 + O(n)$	

$O(n^3)!$

Note:

`eig` not available for sparse matrix arguments

Exception:

`d=eig(A)` for sparse *Hermitian* matrices**EXAMPLE 9.2.0.5 (Runtimes of `eig`)****MATLAB-code 9.2.0.6:**

```
1 A = rand(500,500); B = A'*A; C = gallery('tridiag',500,1,3,1);
```

- ♦ **A** generic dense matrix
- ♦ **B** symmetric (s.p.d. → Def. 1.1.2.6) matrix
- ♦ **C** s.p.d. *tridiagonal* matrix

**MATLAB-code 9.2.0.7: measuring runtimes of `eig`**

```

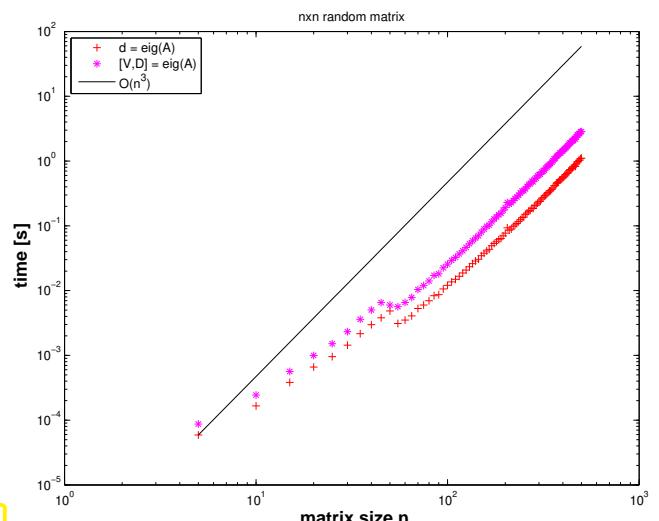
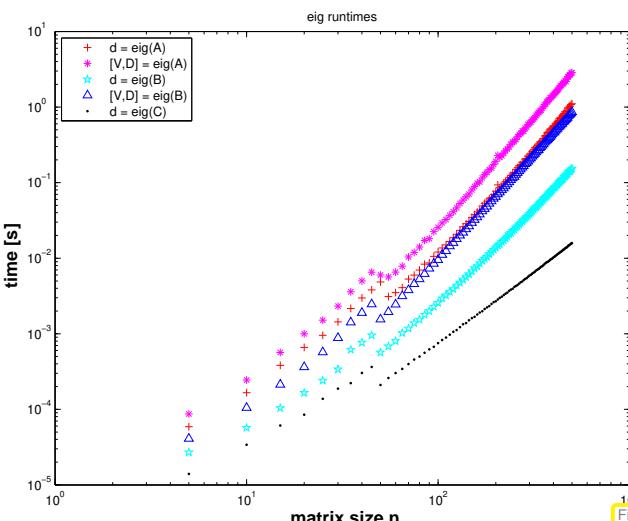
1 function eigtiming
2
3 A = rand(500,500); B = A'*A;
4 C = gallery('tridiag',500,1,3,1);
5 times = [];
6 for n=5:5:500
7     An = A(1:n,1:n); Bn = B(1:n,1:n); Cn = C(1:n,1:n);
8     t1 = 1000; for k=1:3, tic; d = eig(An); t1 = min(t1,toc); end
9     t2 = 1000; for k=1:3, tic; [V,D] = eig(An); t2 = min(t2,toc);
10    end
11    t3 = 1000; for k=1:3, tic; d = eig(Bn); t3 = min(t3,toc); end
12    t4 = 1000; for k=1:3, tic; [V,D] = eig(Bn); t4 = min(t4,toc);
13    end
14    t5 = 1000; for k=1:3, tic; d = eig(Cn); t5 = min(t5,toc); end
15    times = [times; n t1 t2 t3 t4 t5];
16
17 figure;
18 loglog(times(:,1),times(:,2),'r+', times(:,1),times(:,3),'m*',...
19         times(:,1),times(:,4),'cp', times(:,1),times(:,5),'b^',...
20         times(:,1),times(:,6),'k.');
21 xlabel('{\bf matrix size n}', 'fontsize', 14);
22 ylabel('{\bf time [s]}', 'fontsize', 14);
23 title('eig runtimes');
24 legend('d = eig(A)', '[V,D] = eig(A)', 'd = eig(B)', '[V,D] = ...
25         eig(B)', 'd = eig(C)', ...

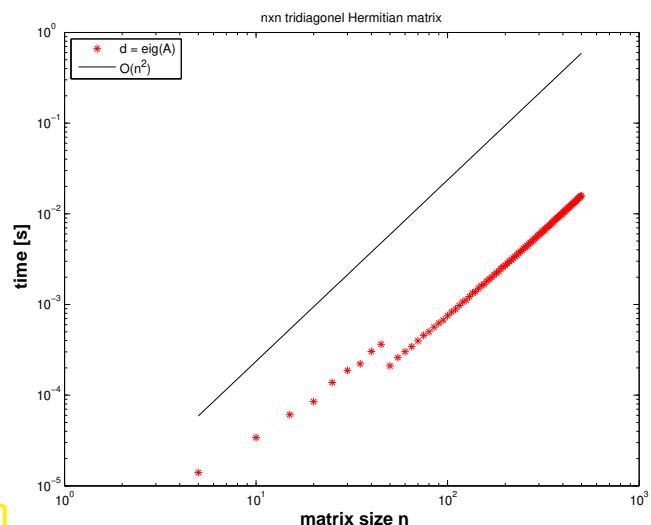
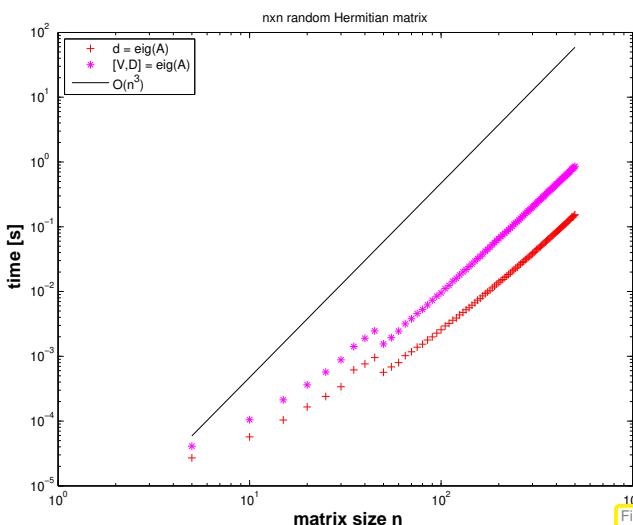
```

```

24         'location','northwest');
25
26 print -depsc2 '../PICTURES/eigtimingall.eps'
27
28 figure;
29 loglog(times(:,1),times(:,2),'r+', times(:,1),times(:,3),'m*',...
30         times(:,1),(times(:,1).^3)/(times(1,1)^3)*times(1,2),'k-');
31 xlabel('{\bf matrix size n}','fontsize',14);
32 ylabel('{\bf time [s]}','fontsize',14);
33 title('nxn random matrix');
34 legend('d = eig(A)', '[V,D] = eig(A)', 'O(n^3)', 'location','northwest');
35
36 print -depsc2 '../PICTURES/eigtimingA.eps'
37
38 figure;
39 loglog(times(:,1),times(:,4),'r+', times(:,1),times(:,5),'m*',...
40         times(:,1),(times(:,1).^3)/(times(1,1)^3)*times(1,2),'k-');
41 xlabel('{\bf matrix size n}','fontsize',14);
42 ylabel('{\bf time [s]}','fontsize',14);
43 title('nxn random Hermitian matrix');
44 legend('d = eig(A)', '[V,D] = eig(A)', 'O(n^3)', 'location','northwest');
45
46 print -depsc2 '../PICTURES/eigtimingB.eps'
47
48 figure;
49 loglog(times(:,1),times(:,6),'r*',...
50         times(:,1),(times(:,1).^2)/(times(1,1)^2)*times(1,2),'k-');
51 xlabel('{\bf matrix size n}','fontsize',14);
52 ylabel('{\bf time [s]}','fontsize',14);
53 title('nxn tridiagonel Hermitian matrix');
54 legend('d = eig(A)', 'O(n^2)', 'location','northwest');
55
56 print -depsc2 '../PICTURES/eigtimingC.eps'

```





- For the sake of efficiency: think which information you really need when computing eigenvalues/eigenvectors of dense matrices

Potentially more efficient methods for *sparse matrices* will be introduced below in Section 9.3, 9.4.

## 9.3 Power Methods

### 9.3.1 Direct power method



*Supplementary literature.* [DR08, Sect. 7.5], [QSS00, Sect. 5.3.1], [QSS00, Sect. 5.3]

#### EXAMPLE 9.3.1.1 ((Simplified) Page rank algorithm → [Gle15; LM06])

Model: Random surfer visits a web page, stays there for fixed time  $\Delta t$ , and then

- ① either follows each of  $\ell$  links on a page with probability  $1/\ell$ .
- ② or resumes surfing at a randomly (with equal probability) selected page

Option ② is chosen with probability  $d$ ,  $0 \leq d \leq 1$ , option ① with probability  $1 - d$ .

- Stationary **Markov chain**, state space  $\hat{=}$  set of all web pages

Question: Fraction of time spent by random surfer on  $i$ -th page (= **page rank**  $x_i \in [0, 1]$ )

This number  $\in ]0, 1[$  can be used to gauge the “importance” of a web page, which, in turns, offers a way to sort the hits resulting from a keyword query: the GOOGLE idea.

Method: Stochastic simulation



#### MATLAB-code 9.3.1.2: stochastic page rank simulation

```

1 function prstochsim(Nhops)
2 % Load web graph data stored in N × N-matrix G
3 load harvard500.mat;
4 N = size(G,1); d = 0.15;
5 count = zeros(1,N); cp = 1;
6 figure('position',[0 0 1200 1000]); pause;

```

```

7 | for n=1:Nhops
8 |   % Find links from current page cp
9 |   idx = find(G(:,cp)); l = size(idx,1); rn = rand(); %
10 |   % If no links, jump to any other pages with equal probability
11 |   if (isempty(idx)), cp = floor(rn*N)+1;
12 |   % With probability d jump to any other page
13 |   elseif (rn < d), cp = floor(rn/d*N)+1;
14 |   % Follow outgoing links with equal probability
15 |   else cp = idx(floor((rn-d)/(1-d)*l)+1,1);
16 |   end
17 |   count(cp) = count(cp) + 1;
18 |   plot(1:N, count/n, 'r.');
19 |   axis([0 N+1 0 0.1]);
20 |   xlabel('{\bf harvard500: no. of page}', 'fontsize', 14);
21 |   ylabel('{\bf page rank}', 'fontsize', 14);
22 |   title(sprintf('{\bf page rank, harvard500: %d
23 |   hops}', n), 'fontsize', 14);
24 |   drawnow;
25 | end

```

Explanations Code 9.3.1.2:

- ◆ Line 9: `rand` generates uniformly distributed pseudo-random numbers  $\in [0,1]$
- ◆ Web graph encoded in  $\mathbf{G} \in \{0,1\}^{N,N}$ :

$$(\mathbf{G})_{ij} = 1 \Rightarrow \text{link } j \rightarrow i,$$

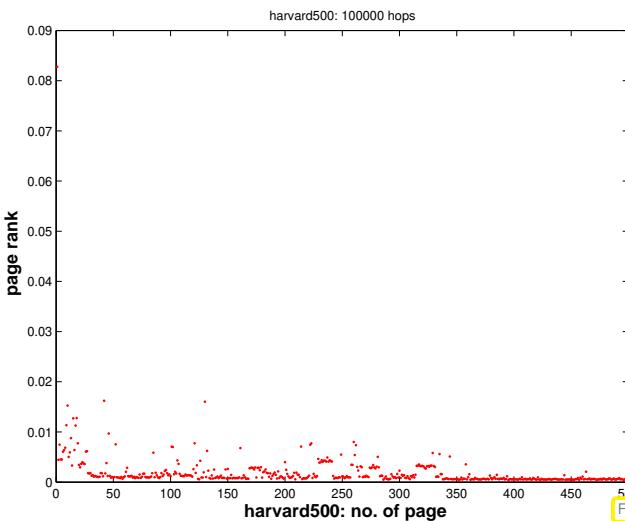


Fig. 296

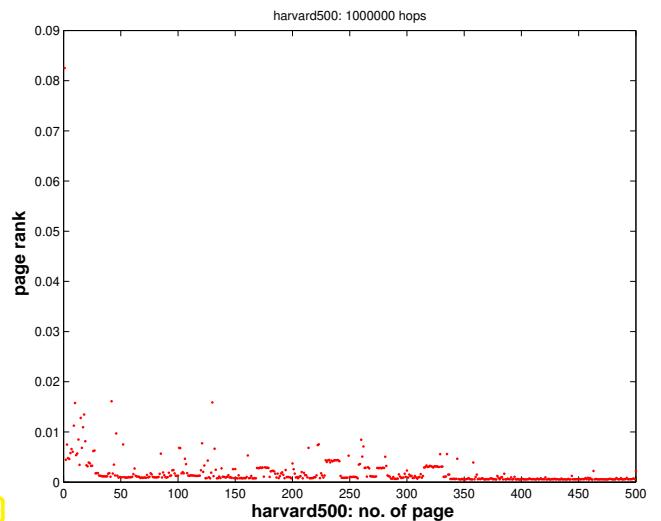


Fig. 297

Observation: relative visit times stabilize as the number of hops in the stochastic simulation  $\rightarrow \infty$ .

The limit distribution is called **stationary distribution**/invariant measure of the Markov chain. This is what we seek.

- ◆ Numbering of pages  $1, \dots, N$ ,  $\ell_i \hat{=} \text{number of links from page } i$
- ◆  $N \times N$ -matrix of transition probabilities page  $j \rightarrow$  page  $i$ :  $\mathbf{A} = (a_{ij})_{i,j=1}^N \in \mathbb{R}^{N,N}$

$$a_{ij} \in [0, 1] \hat{=} \text{probability to jump from page } j \text{ to page } i.$$

$$\Rightarrow \sum_{i=1}^N a_{ij} = 1 . \quad (9.3.1.3)$$

A matrix  $\mathbf{A} \in [0,1]^{N,N}$  with the property (9.3.1.3) is called a (column) **stochastic matrix**.

“Meaning” of  $\mathbf{A}$ : given  $\mathbf{x} \in [0, 1]^N$ ,  $\|\mathbf{x}\|_1 = 1$ , where  $x_i$  is the probability of the surfer to visit page  $i$ ,  $i = 1, \dots, N$ , at an instance  $t$  in time,  $\mathbf{y} = \mathbf{Ax}$  satisfies

$$y_j \geq 0 \quad , \quad \sum_{j=1}^N y_j = \sum_{j=1}^N \sum_{i=1}^N a_{ji} x_i = \sum_{i=1}^N x_i \underbrace{\sum_{j=1}^N a_{ij}}_{=1} = \sum_{i=1}^N x_i = 1 \; .$$

  $y_j \doteq$  probability for visiting page  $j$  at time  $t + \Delta t$ .

Transition probability matrix for page rank computation

$$(\mathbf{A})_{ij} = \begin{cases} \frac{1}{N} & , \text{if } (\mathbf{G})_{ij} = 0 \forall i = 1, \dots, N , \\ d/N + (1-d) \frac{(\mathbf{G})_{ij}}{\ell_j} & \text{else.} \end{cases} \quad (9.3.1.4)$$

random jump to any other page                            follow link

### MATLAB-code 9.3.1.5: transition probability matrix for page rank

```

1 function A = prbbuildA(G,d)
2 N = size(G,1);
3 l = full(sum(G)); idx = find(l>0);
4 s = zeros(N,1); s(idx) = 1./l(idx);
5 ds = ones(N,1)/N; ds(idx) = d/N;
6 A = ones(N,1)*ones(1,N)*diag(ds) +
    (1-d)*G*diag(s);

```

Note: special treatment of zero columns of  $\mathbf{G}$ , cf. (9.3.1.4)!



Stochastic simulation based on a single surfer is *slow*. Alternatives?

Thought experiment: Instead of a single random surfer we may consider  $m \in \mathbb{N}$ ,  $m \gg 1$ , of them who visit pages independently. The fraction of time  $m \cdot T$  they all together spend on page  $i$  will obviously be the same for  $T \rightarrow \infty$  as that for a single random surfer.

Instead of counting the surfers we watch the proportions of them visiting particular web pages at an instance of time. Thus, after the  $k$ -th hop we can assign a number  $x_i^{(k)} \in [0, 1]$  to web page  $i$ , which gives the proportion of surfers currently on that page:  $x_i^{(k)} := \frac{n_i^{(k)}}{m}$ , where  $n_i^{(k)} \in \mathbb{N}_0$  designates the number of surfers on page  $i$  after the  $k$ -th hop.

Now consider  $m \rightarrow \infty$ . The law of large numbers suggests that the ("infinitely many") surfers visiting page  $j$  will move on to other pages proportional to the transition probabilities  $a_{ij}$ : in terms of proportions, for  $m \rightarrow \infty$  the stochastic evolution becomes a deterministic discrete dynamical system and we find

$$x_i^{(k+1)} = \sum_{j=1}^N a_{ij} x_j^{(k)}, \quad (9.3.1.6)$$

that is, the proportion of surfers ending up on page  $i$  equals the sum of the proportions on the “source pages” weighted with the transition probabilities.

Notice that (9.3.1.6) amounts to matrix  $\times$  vector. Thus, writing  $\mathbf{x}^{(0)} \in [0, 1]^N$ ,  $\|\mathbf{x}^{(0)}\| = 1$  for the initial distribution of the surfers on the net we find

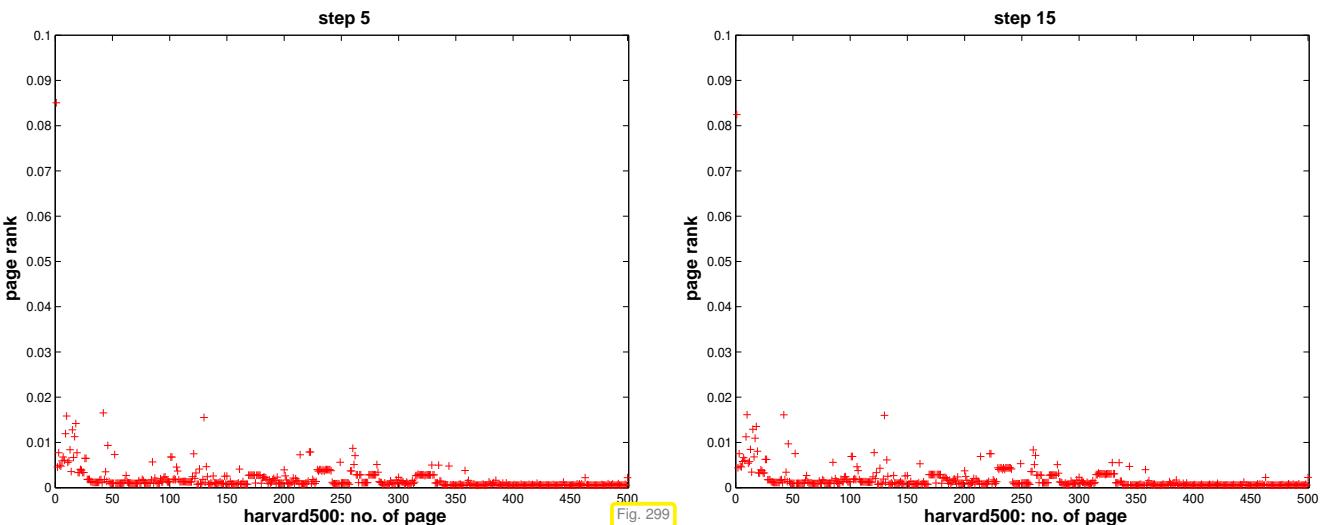
$$\mathbf{x}^{(k)} = \mathbf{A}^k \mathbf{x}^{(0)}$$

will be their mass distribution after  $k$  hops. If the limit exists, the  $i$ -th component of  $\mathbf{x}^* := \lim_{k \rightarrow \infty} \mathbf{x}^{(k)}$  tells us which fraction of the (infinitely many) surfers will be visiting page  $i$  most of the time. Thus,  $\mathbf{x}^*$  yields the stationary distribution of the Markov chain.

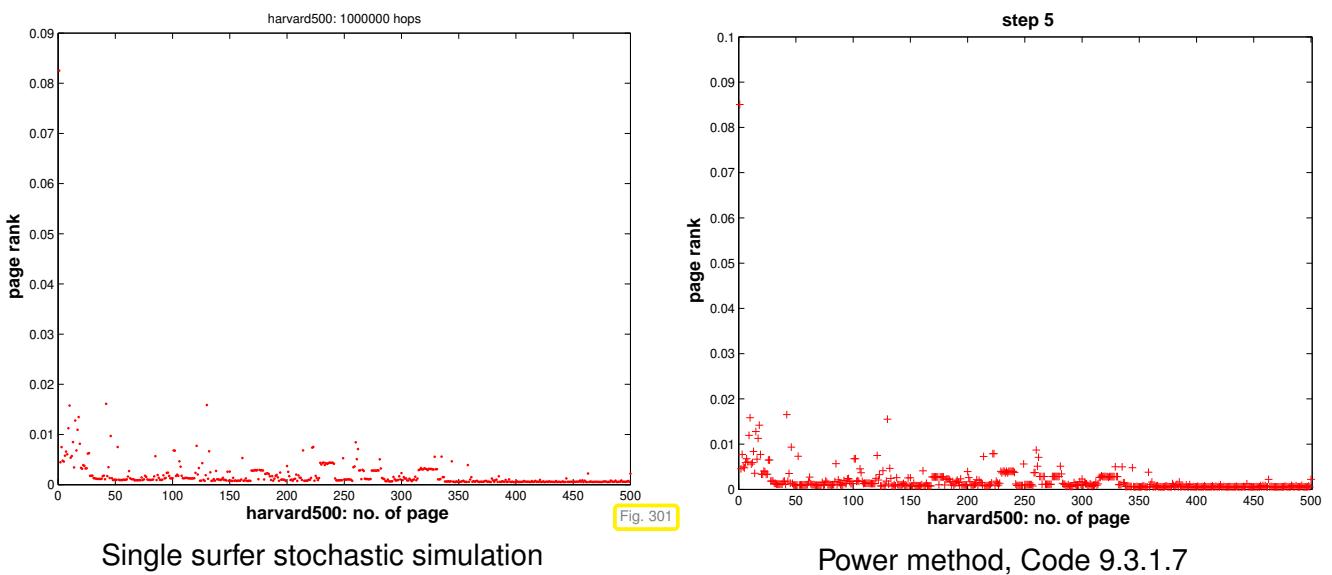
#### MATLAB-code 9.3.1.7: tracking fractions of many surfers

```

1 function prpowitsim(d,Nsteps)
2 % MATLAB way of specifying Default arguments
3 if (nargin < 2), Nsteps = 5; end
4 if (nargin < 1), d = 0.15; end
5 % load connectivity matrix and build transition matrix
6 load harvard500.mat; A = prbuildA(G,d);
7 N = size(A,1); x = ones(N,1)/N;
8
9 figure('position',[0 0 1200 1000]);
10 plot(1:N,x,'r+'); axis([0 N+1 0 0.1]);
11 % Plain power iteration for stochastic matrix A
12 for l=1:Nsteps
13   pause; x = A*x; plot(1:N,x,'r+'); axis([0 N+1 0 0.1]);
14   title(sprintf('{\bf step %d}',l),'fontsize',14);
15   xlabel('{\bf harvard500: no. of page}','fontsize',14);
16   ylabel('{\bf page rank}','fontsize',14); drawnow;
17 end
```



Comparison:



Observation: Convergence of the  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ , and the limit must be a fixed point of the iteration function:

$$\Rightarrow \quad \mathbf{A}\mathbf{x}^* = \mathbf{x}^* \quad \Rightarrow \quad \mathbf{x}^* \in \mathbf{EigA1} .$$

Does  $\mathbf{A}$  possess an eigenvalue  $= 1$ ? Does the associated eigenvector really provide a probability distribution (after scaling), that is, are all of its entries non-negative? Is this probability distribution unique? To answer these questions we have to study the matrix  $\mathbf{A}$ :

For every stochastic matrix  $\mathbf{A}$ , by definition (9.3.1.3)

$$\begin{aligned} \mathbf{A}^T \mathbf{1} &= \mathbf{1} \stackrel{(9.1.0.3)}{\Rightarrow} \mathbf{1} \in \sigma(\mathbf{A}) , \\ (1.5.5.14) \quad \Rightarrow \quad \|\mathbf{A}\|_1 &= 1 \stackrel{\text{Thm. 9.1.0.4}}{\Rightarrow} \rho(\mathbf{A}) = 1 , \end{aligned}$$

where  $\rho(\mathbf{A})$  is the spectral radius of the matrix  $\mathbf{A}$ , see Section 9.1.

For  $\mathbf{r} \in \mathbf{EigA1}$ , that is,  $\mathbf{Ar} = \mathbf{r}$ , denote by  $|\mathbf{r}|$  the vector  $(|r_i|)_{i=1}^N$ . Since all entries of  $\mathbf{A}$  are non-negative, we conclude by the triangle inequality that  $\|\mathbf{Ar}\|_1 \leq \|\mathbf{A}|\mathbf{r}|\|_1$

$$\begin{aligned} \Rightarrow \quad 1 &= \|\mathbf{A}\|_1 = \sup_{\mathbf{x} \in \mathbb{R}^N} \frac{\|\mathbf{Ax}\|_1}{\|\mathbf{x}\|_1} \geq \frac{\|\mathbf{A}|\mathbf{r}|\|_1}{\|\mathbf{r}\|_1} \geq \frac{\|\mathbf{Ar}\|_1}{\|\mathbf{r}\|_1} = 1 . \\ \Rightarrow \quad \|\mathbf{A}|\mathbf{r}|\|_1 &= \|\mathbf{Ar}\|_1 \quad \text{if } a_{ij} > 0 \quad |\mathbf{r}| = \pm \mathbf{r} . \end{aligned}$$

Hence, different components of  $\mathbf{r}$  cannot have opposite sign, which means, that  $\mathbf{r}$  can be chosen to have non-negative entries, if the entries of  $\mathbf{A}$  are strictly positive, which is the case for  $\mathbf{A}$  from (9.3.1.4). After normalization  $\|\mathbf{r}\|_1 = 1$  the eigenvector can be regarded as a probability distribution on  $\{1, \dots, N\}$ .

If  $\mathbf{Ar} = \mathbf{r}$  and  $\mathbf{As} = \mathbf{s}$  with  $(\mathbf{r})_i \geq 0$ ,  $(\mathbf{s})_i \geq 0$ ,  $\|\mathbf{r}\|_1 = \|\mathbf{s}\|_1 = 1$ , then  $\mathbf{A}(\mathbf{r} - \mathbf{s}) = \mathbf{r} - \mathbf{s}$ . Hence, by the above considerations, also all the entries of  $\mathbf{r} - \mathbf{s}$  are either non-negative or non-positive. By the assumptions on  $\mathbf{r}$  and  $\mathbf{s}$  this is only possible, if  $\mathbf{r} - \mathbf{s} = 0$ . We conclude that

$$\mathbf{A} \in ]0, 1]^{N,N} \text{ stochastic} \quad \Rightarrow \quad \dim \mathbf{EigA1} = 1 . \quad (9.3.1.8)$$

Sorting the pages according to the size of the corresponding entries in  $\mathbf{r}$  yields the famous “page rank”.

**MATLAB-code 9.3.1.9: computing page rank vector  $\mathbf{r}$  via `eig`**

```
function prevp
load harvard500.mat; d = 0.15;
[V,D] = eig(prbuildA(G,d));

figure; r = V(:,1); N = length(r);
plot(1:N,r/sum(r), 'm.' ); axis([0 N+1 0 0.1]);
xlabel('{\bf harvard500: no. of page}', 'fontsize',14);
ylabel('{\bf entry of r-vector}', 'fontsize',14);
title('harvard_500: Perron–Frobenius vector');
print -depsc2 '../PICTURES/prevp.eps';
```

Plot of entries of unique vector  $\mathbf{r} \in \mathbb{R}^N$  with

$$0 \leq (\mathbf{r})_i \leq 1, \quad \|\mathbf{r}\|_1 = 1,$$

$$\boxed{\mathbf{Ar} = \mathbf{r}}.$$

Inefficient implementation!

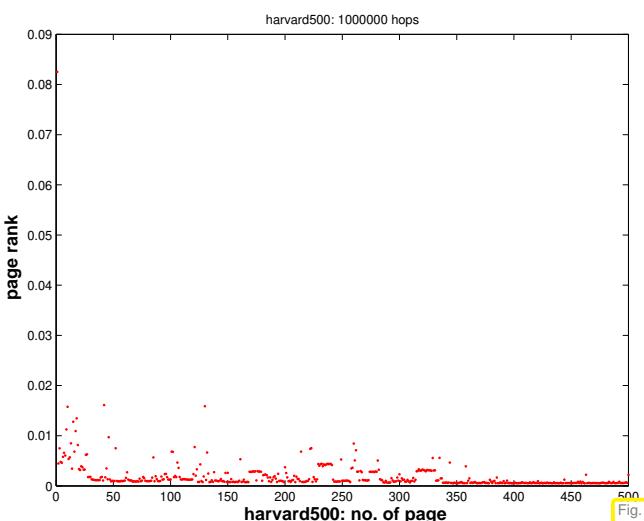
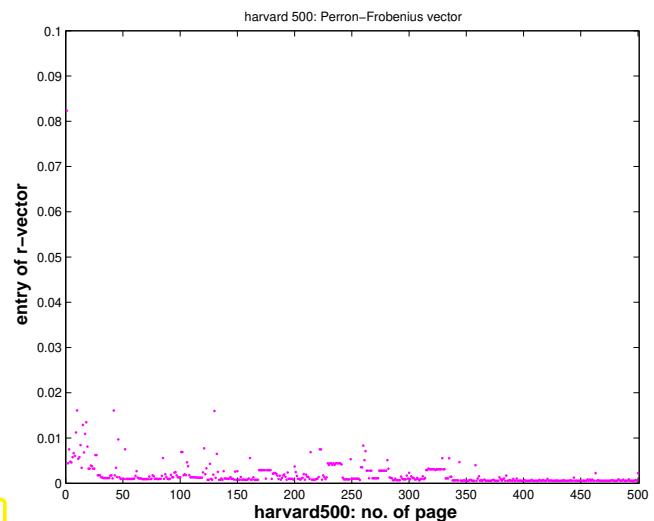


Fig. 302

stochastic simulation



eigenvector computation

The possibility to compute the stationary probability distribution of a Markov chain through an eigenvector of the transition probability matrix is due to a property of stationary Markov chains called **ergodicity**.

$\mathbf{A} \triangleq$  page rank transition probability matrix, see Code 9.3.1.5,  $d = 0.15$ , harvard500 example.

Errors:

$$\left\| \mathbf{A}^k \mathbf{x}_0 - \mathbf{r} \right\|_1,$$

with  $\mathbf{x}_0 = \mathbf{1}/N$ ,  $N = 500$ .

We observe **linear convergence!** ( $\rightarrow$  Def. 8.1.1.1, iteration error vs. iteration count  $\approx$  straight line lin-log plot)

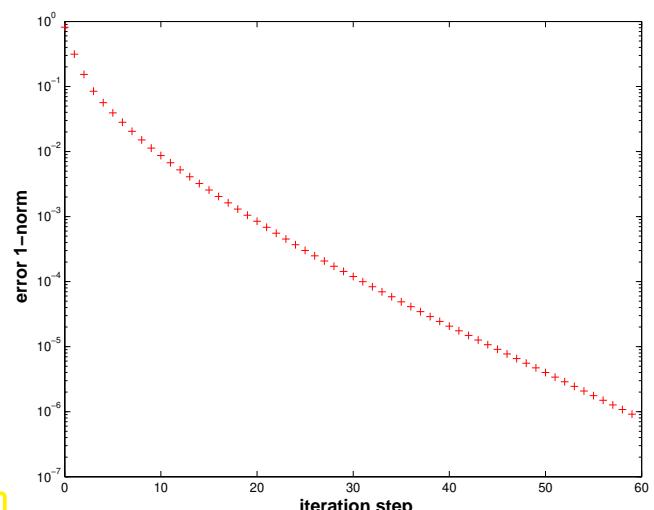


Fig. 304

The computation of page rank amounts to finding the eigenvector of the matrix  $\mathbf{A}$  of transition probabilities that belongs to its *largest* eigenvalue 1. This is addressed by an important class of practical eigenvalue problems:

Task:

given  $\mathbf{A} \in \mathbb{K}^{n,n}$ , find largest (in modulus) eigenvalue of  $\mathbf{A}$   
and (an) associated eigenvector.

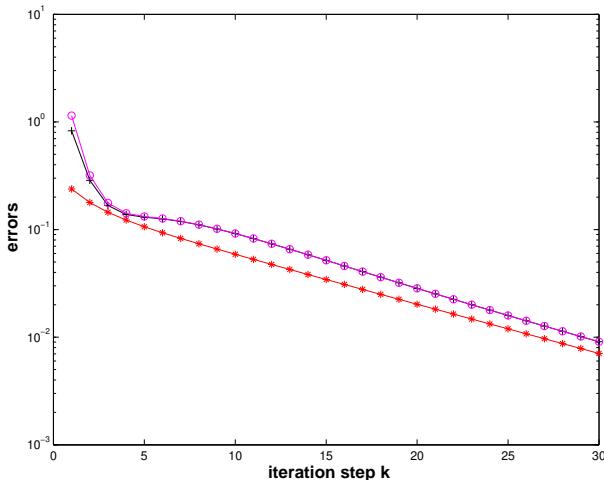


Idea: (suggested by page rank computation, Code 9.3.1.7)

Iteration:  $\mathbf{z}^{(k+1)} = \mathbf{A}\mathbf{z}^{(k)}$ ,  $\mathbf{z}^{(0)}$  arbitrary

### EXAMPLE 9.3.1.10 (Power iteration → Ex. 9.3.1.1)

Try the above iteration for general  $10 \times 10$ -matrix, largest eigenvalue 10, algebraic multiplicity 1.



#### MATLAB-code 9.3.1.11:

```
d = (1:10)'; n = length(d);
S = triu(diag(n:-1:1,0)+...
ones(n,n));
A = S*diag(d,0)*inv(S);
```

▷ error norm  $\left\| \frac{\mathbf{z}^{(k)}}{\|\mathbf{z}^{(k)}\|} - (\mathbf{S})_{:,10} \right\|$

(Note:  $(\mathbf{S})_{:,10} \hat{=} \text{eigenvector for eigenvalue 10}$ )

$\mathbf{z}^{(0)}$  = random vector

Fig. 305

Observation: linear convergence of (normalized) eigenvectors!

► Suggests direct power method (ger.: Potenzmethode): iterative method (→ Section 8.1)

$$\begin{aligned} \text{initial guess: } & \mathbf{z}^{(0)} \text{ "arbitrary"}, \\ \text{next iterate: } & \mathbf{w} := \mathbf{A}\mathbf{z}^{(k-1)}, \quad \mathbf{z}^{(k)} := \frac{\mathbf{w}}{\|\mathbf{w}\|_2}, \quad k = 1, 2, \dots . \end{aligned} \quad (9.3.1.12)$$

Note: the “normalization” of the iterates in (9.3.1.12) does not change anything (in exact arithmetic) and helps avoid overflow in floating point arithmetic.

Computational effort:  $1 \times \text{matrix} \times \text{vector}$  per step ➤ inexpensive for sparse matrices

A persuasive theoretical justification for the direct power method:

Assume  $\mathbf{A} \in \mathbb{K}^{n,n}$  to be diagonalizable:

$$\Leftrightarrow \exists \text{ basis } \{\mathbf{v}_1, \dots, \mathbf{v}_n\} \text{ of eigenvectors of } \mathbf{A}: \mathbf{A}\mathbf{v}_j = \lambda_j \mathbf{v}_j, \lambda_j \in \mathbb{C}.$$

Assume

$$|\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_{n-1}| < |\lambda_n|, \quad \|\mathbf{v}_j\|_2 = 1. \quad (9.3.1.13)$$

Key observations for power iteration (9.3.1.12)

$$\mathbf{z}^{(k)} = \mathbf{A}^k \mathbf{z}^{(0)} \quad (\rightarrow \text{name "power method"}) \quad (9.3.1.14)$$

$$\mathbf{z}^{(0)} = \sum_{j=1}^n \zeta_j \mathbf{v}_j \Rightarrow \mathbf{z}^{(k)} = \sum_{j=1}^n \zeta_j \lambda_j^k \mathbf{v}_j. \quad (9.3.1.15)$$

Due to (9.3.1.13) for large  $k \gg 1$  ( $\Rightarrow |\lambda_n^k| \gg |\lambda_j^k|$  for  $j \neq n$ ) the contribution of  $\mathbf{v}_n$  (size  $\zeta_n \lambda_n^k$ ) in the *eigenvector expansion* (9.3.1.15) will be much larger than the contribution (size  $\zeta_n \lambda_j^k$ ) of any other eigenvector (, if  $\zeta_n \neq 0$ ): the eigenvector for  $\lambda_n$  will swamp all other for  $k \rightarrow \infty$ .

Further (9.3.1.15) nutures expectation:  $\mathbf{v}_n$  will become dominant in  $\mathbf{z}^{(k)}$  the faster, the better  $|\lambda_n|$  is separated from  $|\lambda_{n-1}|$ , see Thm. 9.3.1.21 for rigorous statement.

$\mathbf{z}^{(k)} \rightarrow$  eigenvector, but how do we get the associated eigenvalue  $\lambda_n$  ?

When (9.3.1.12) has converged, two common ways to recover  $\lambda_{\max}$   $\rightarrow$  [DR08, Alg. 7.20]

$$\textcircled{1} \quad \mathbf{A}\mathbf{z}^{(k)} \approx \lambda_{\max} \mathbf{z}^{(k)} \quad \Rightarrow \quad |\lambda_n| \approx \frac{\|\mathbf{A}\mathbf{z}^{(k)}\|}{\|\mathbf{z}^{(k)}\|} \quad (\text{modulus only!})$$

$$\textcircled{2} \quad \lambda_{\max} \approx \underset{\theta \in \mathbb{R}}{\operatorname{argmin}} \left\| \mathbf{A}\mathbf{z}^{(k)} - \theta \mathbf{z}^{(k)} \right\|_2^2 \quad \Rightarrow \quad \lambda_{\max} \approx \frac{(\mathbf{z}^{(k)})^H \mathbf{A} \mathbf{z}^{(k)}}{\|\mathbf{z}^{(k)}\|_2^2}.$$

This latter formula is extremely useful, which has earned it a special name:

### Definition 9.3.1.16.

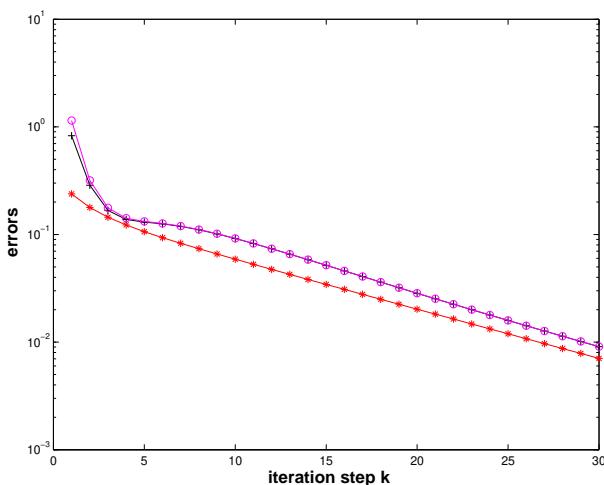
For  $\mathbf{A} \in \mathbb{K}^{n,n}$ ,  $\mathbf{u} \in \mathbb{K}^n$  the **Rayleigh quotient** is defined by

$$\rho_{\mathbf{A}}(\mathbf{u}) := \frac{\mathbf{u}^H \mathbf{A} \mathbf{u}}{\mathbf{u}^H \mathbf{u}}.$$

An immediate consequence of the definitions:

$$\lambda \in \sigma(\mathbf{A}) \quad , \quad \mathbf{z} \in \operatorname{Eig}_{\lambda} \mathbf{A} \quad \Rightarrow \quad \rho_{\mathbf{A}}(\mathbf{z}) = \lambda. \quad (9.3.1.17)$$

### EXAMPLE 9.3.1.18 (Direct power method $\rightarrow$ Ex. 9.3.1.18 cnt'd)



#### MATLAB-code 9.3.1.19:

```
n = length(d); S = triu(diag(n:-1:1,0)+...  
ones(n,n)); A = S*diag(d,0)*inv(S);
```

$\mathbf{d} = (1:10)'$ ;

$\circ$  : error  $|\lambda_n - \rho_{\mathbf{A}}(\mathbf{z}^{(k)})|$

$*$  : error norm  $\|\mathbf{z}^{(k)} - \mathbf{s}_{.,n}\|$

$+$  :  $\left| \lambda_n - \frac{\|\mathbf{A}\mathbf{z}^{(k-1)}\|_2}{\|\mathbf{z}^{(k-1)}\|_2} \right|$

$\mathbf{z}^{(0)}$  = random vector

Fig. 306

Test matrices:

- ①  $d = (1:10)'$ ;  $\Rightarrow |\lambda_{n-1}| : |\lambda_n| = 0.9$
- ②  $d = [\text{ones}(9,1); 2]'$ ;  $\Rightarrow |\lambda_{n-1}| : |\lambda_n| = 0.5$
- ③  $d = 1-2.^{(-(1:0.5:5)'')}$ ;  $\Rightarrow |\lambda_{n-1}| : |\lambda_n| = 0.9866$

### MATLAB-code 9.3.1.20: Investigating convergence of direct power method

```

1 % Demonstration of direct power method for Ex. 9.3.1.18
2 maxit = 30; d = [1:10]'; n = length(d);
3 % Initialize the matrix A
4 S = triu(diag(n:-1:1,0)+ones(n,n));
5 A = S*diag(d,0)*inv(S);
6 % This calculates the exact eigenvalues (for error calculation)
7 [V,D] = eig(A);

8
9 k = find(d == max(abs(d)));
10 if (length(k) > 1), error('No single largest EV'); end;
11 ev = X(:,k(1)); ev = ev/norm(ev); ev
12 ew = d(k(1)); ew
13 if (ew < 0), sgn = -1; else sgn = 1; end

14
15 z = rand(n,1); z = z/norm(z);
16 s = 1;
17 res = [];

18
19 % Actual direct power iteration
20 for i=1:maxit
21     w = A*z; l = norm(w); rq = real(dot(w,z)); z = w/l;
22     res = [res;i,l,norm(s*z-ev),abs(l-abs(ew)),abs(sgn*rq-ew)];
23     s = s * sgn;
24 end

25
26 % Plot the result
27 semilogy(res(:,1),res(:,3),'r-*',res(:,1),res(:,4),'k-+',res(:,1),res(:,5),'m-o');
28 xlabel({'\bf iteration step k'},'FontSize',14);
29 ylabel({'\bf errors'},'FontSize',14);
30 print -deps2c '../PICTURES/pml.eps';

```

$$\rho_{\text{EV}}^{(k)} := \frac{\|\mathbf{z}^{(k)} - \mathbf{s}_{\cdot,n}\|}{\|\mathbf{z}^{(k-1)} - \mathbf{s}_{\cdot,n}\|},$$

$$\rho_{\text{EW}}^{(k)} := \frac{|\rho_A(\mathbf{z}^{(k)}) - \lambda_n|}{|\rho_A(\mathbf{z}^{(k-1)}) - \lambda_n|}.$$

$k$	①		②		③	
	$\rho_{\text{EV}}^{(k)}$	$\rho_{\text{EW}}^{(k)}$	$\rho_{\text{EV}}^{(k)}$	$\rho_{\text{EW}}^{(k)}$	$\rho_{\text{EV}}^{(k)}$	$\rho_{\text{EW}}^{(k)}$
22	0.9102	0.9007	0.5000	0.5000	0.9900	0.9781
23	0.9092	0.9004	0.5000	0.5000	0.9900	0.9791
24	0.9083	0.9001	0.5000	0.5000	0.9901	0.9800
25	0.9075	0.9000	0.5000	0.5000	0.9901	0.9809
26	0.9068	0.8998	0.5000	0.5000	0.9901	0.9817
27	0.9061	0.8997	0.5000	0.5000	0.9901	0.9825
28	0.9055	0.8997	0.5000	0.5000	0.9901	0.9832
29	0.9049	0.8996	0.5000	0.5000	0.9901	0.9839
30	0.9045	0.8996	0.5000	0.5000	0.9901	0.9844

Observation:

linear convergence ( $\rightarrow ??$ )

»

**Theorem 9.3.1.21. Convergence of direct power method** → [DR08, Thm. 25.1]

Let  $\lambda_n > 0$  be the largest (in modulus) eigenvalue of  $\mathbf{A} \in \mathbb{K}^{n,n}$  and have (algebraic) multiplicity 1. Let  $\mathbf{v}, \mathbf{y}$  be the left and right eigenvectors of  $\mathbf{A}$  for  $\lambda_n$  normalized according to  $\|\mathbf{y}\|_2 = \|\mathbf{v}\|_2 = 1$ . Then there is convergence

$$\left\| \mathbf{A} \mathbf{z}^{(k)} \right\|_2 \rightarrow \lambda_n \quad , \quad \mathbf{z}^{(k)} \rightarrow \pm \mathbf{v} \quad \text{linearly with rate} \quad \frac{|\lambda_{n-1}|}{|\lambda_n|} \text{,}$$

where  $\mathbf{z}^{(k)}$  are the iterates of the direct power iteration and  $\mathbf{y}^H \mathbf{z}^{(0)} \neq 0$  is assumed.

**Remark 9.3.1.22 (Initial guess for power iteration)**

roundoff errors ➤  $\mathbf{y}^H \mathbf{z}^{(0)} \neq 0$  always satisfied in practical computations

Usual (not the best!) choice for  $\mathbf{x}^{(0)}$  = random vector

**Remark 9.3.1.23 (Termination criterion for direct power iteration)** (→ Section 8.1.2)

Adaptation of a posteriori termination criterion (8.2.2.20)

“relative change”  $\leq \text{tol}$ :

$$\begin{cases} \min \left\| \mathbf{z}^{(k)} \pm \mathbf{z}^{(k-1)} \right\| \leq (1/L - 1)\text{tol} , \\ \left| \frac{\|\mathbf{A} \mathbf{z}^{(k)}\|}{\|\mathbf{z}^{(k)}\|} - \frac{\|\mathbf{A} \mathbf{z}^{(k-1)}\|}{\|\mathbf{z}^{(k-1)}\|} \right| \leq (1/L - 1)\text{tol} \end{cases} \text{ see (8.1.2.7).}$$

Estimated rate of convergence

**9.3.2 Inverse Iteration [DR08, Sect. 7.6], [QSS00, Sect. 5.3.2]****EXAMPLE 9.3.2.1 (Image segmentation)**

Given: gray-scale image: intensity matrix  $\mathbf{P} \in \{0, \dots, 255\}^{m,n}$ ,  $m, n \in \mathbb{N}$   
 $((\mathbf{P})_{ij}) \leftrightarrow \text{pixel, } 0 \hat{=} \text{black, } 255 \hat{=} \text{white}$

Loading and displaying images  
in MATLAB ➤

**MATLAB-code 9.3.2.2: loading and displaying an image**

```

1 M = imread('eth.pbm');
2 [m, n] = size(M);
3 fprintf('%dx%d grey scale pixel\n', m, n);
4 figure; image(M); title('ETH view');
5 col = [0:1/215:1]*[1, 1, 1]; colormap(col);

```

(Fuzzy) task: Local segmentation

Find connected patches of image of the same shade/color

More general segmentation problem (non-local): identify parts of the image, not necessarily connected, with the same texture.

Next: Statement of (rigorously defined) problem, cf. ??:

Preparation: Numbering of pixels  $1 \dots, mn$ , e.g., lexicographic numbering:

- ◆ pixel set  $\mathcal{V} := \{1, \dots, nm\}$
- ◆ indexing:  $\text{index}(\text{pixel}_{i,j}) = (i-1)n + j$

☞ notation:  $p_k := (\mathbf{P})_{ij}$ , if  $k = \text{index}(\text{pixel}_{i,j}) = (i-1)n + j$ ,  $k = 1, \dots, N := mn$

Local similarity matrix:

$$\mathbf{W} \in \mathbb{R}^{N,N}, \quad N := mn, \quad (9.3.2.3)$$

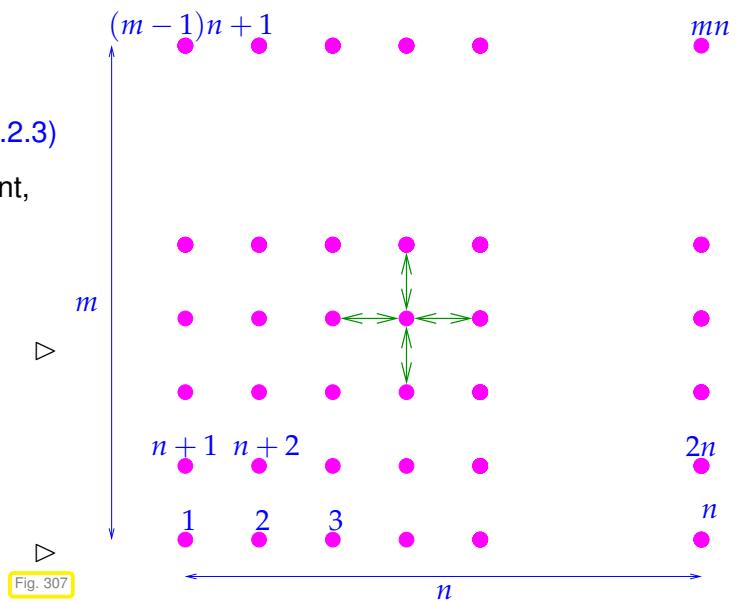
$$(\mathbf{W})_{ij} = \begin{cases} 0 & , \text{if pixels } i, j \text{ not adjacent,} \\ 0 & , \text{if } i = j, \\ \sigma(p_i, p_j) & , \text{if pixels } i, j \text{ adjacent.} \end{cases}$$

↔  $\hat{=}$  adjacent pixels

Similarity function, e.g., with  $\alpha > 0$

$$\sigma(x, y) := \exp(-\alpha(x - y)^2), \quad x, y \in \mathbb{R}.$$

Lexicographic numbering



The entries of the matrix  $\mathbf{W}$  measure the “similarity” of neighboring pixels: if  $(\mathbf{W})_{ij}$  is large, they encode (almost) the same intensity, if  $(\mathbf{W})_{ij}$  is close to zero, then they belong to parts of the picture with very different brightness. In the latter case, the boundary of the segment may separate the two pixels.

#### Definition 9.3.2.4. Normalized cut ( $\rightarrow$ [SM00, Sect. 2])

For  $\mathcal{X} \subset \mathcal{V}$  define the normalized cut as

$$\text{Ncut}(\mathcal{X}) := \frac{\text{cut}(\mathcal{X})}{\text{weight}(\mathcal{X})} + \frac{\text{cut}(\mathcal{X})}{\text{weight}(\mathcal{V} \setminus \mathcal{X})},$$

with  $\text{cut}(\mathcal{X}) := \sum_{i \in \mathcal{X}, j \notin \mathcal{X}} w_{ij}$ ,  $\text{weight}(\mathcal{X}) := \sum_{i \in \mathcal{X}, j \in \mathcal{X}} w_{ij}$ .

In light of local similarity relationship:

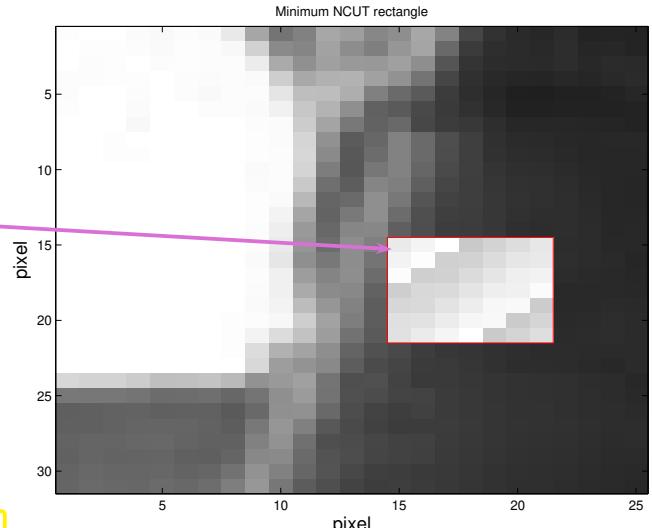
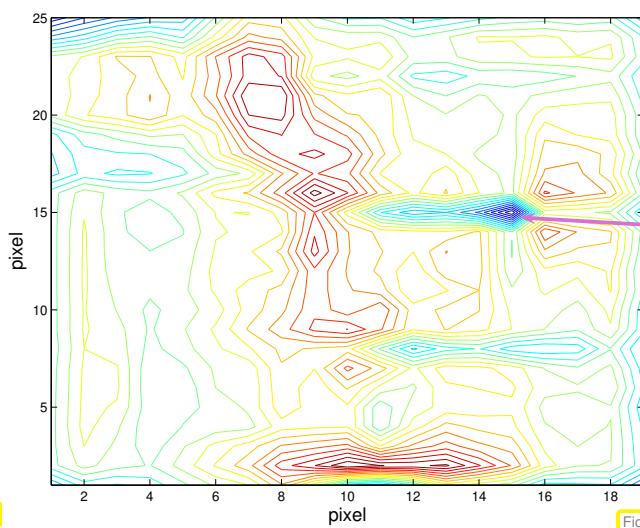
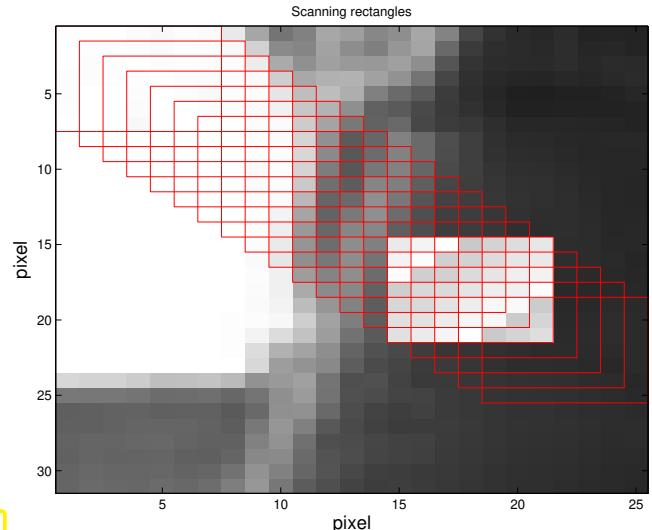
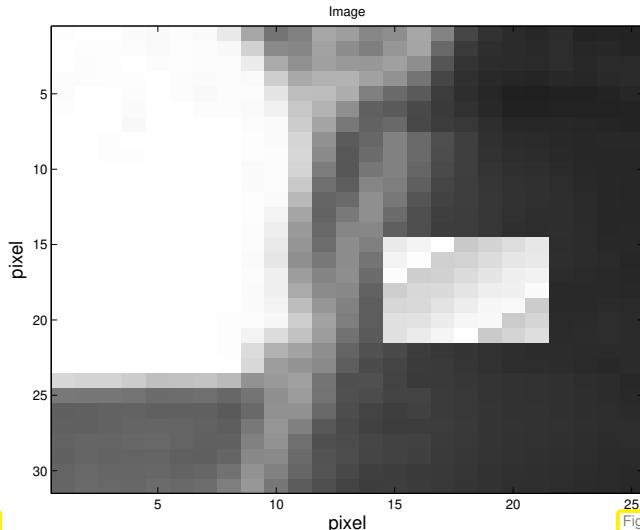
- $\text{cut}(\mathcal{X})$  big  $\Rightarrow$  substantial similarity of pixels across interface between  $\mathcal{X}$  and  $\mathcal{V} \setminus \mathcal{X}$ .
- $\text{weight}(\mathcal{X})$  big  $\Rightarrow$  a lot of similarity of adjacent pixels inside  $\mathcal{X}$ .

► Segmentation problem (rigorous statement):

$$\text{find } \mathcal{X}^* \subset \mathcal{V}: \quad \mathcal{X}^* = \underset{\mathcal{X} \subset \mathcal{V}}{\operatorname{argmin}} \text{Ncut}(\mathcal{X}). \quad (9.3.2.5)$$



NP-hard combinatorial optimization problem !



△  $\text{Ncut}(\mathcal{X})$  for pixel subsets  $\mathcal{X}$  defined by sliding rectangles, see Fig. 309.

*Equivalent reformulation:*

indicator function:  $z : \{1, \dots, N\} \mapsto \{-1, 1\}, \quad z_i := z(i) = \begin{cases} 1 & , \text{if } i \in \mathcal{X} , \\ -1 & , \text{if } i \notin \mathcal{X} . \end{cases} \quad (9.3.2.6)$

►  $\text{Ncut}(\mathcal{X}) = \frac{\sum_{z_i>0, z_j<0} -w_{ij}z_i z_j}{\sum_{z_i>0} d_i} + \frac{\sum_{z_i>0, z_j<0} -w_{ij}z_i z_j}{\sum_{z_i<0} d_i} , \quad (9.3.2.7)$

$$d_i = \sum_{j \in \mathcal{V}} w_{ij} = \text{weight}(\{i\}) . \quad (9.3.2.8)$$

Sparse matrices:

$$\mathbf{D} := \text{diag}(d_1, \dots, d_N) \in \mathbb{R}^{N,N}, \quad \mathbf{A} := \mathbf{D} - \mathbf{W} = \mathbf{A}^\top. \quad (9.3.2.9)$$

Summary: (obvious) properties of these matrices

- ◆  $\mathbf{A}$  has positive diagonal and non-positive off-diagonal entries.
- ◆  $\mathbf{A}$  is **diagonally dominant** ( $\rightarrow$  Def. 2.8.0.8)  $\Rightarrow$   $\mathbf{A}$  is positive semidefinite by Lemma 2.8.0.12.
- ◆  $\mathbf{A}$  has row sums = 0:

$$\mathbf{1}^\top \mathbf{A} = \mathbf{A} \mathbf{1} = 0. \quad (9.3.2.10)$$

#### MATLAB-code 9.3.2.11: assembly of $\mathbf{A}$ , $\mathbf{D}$

```

1 function [A,D] = imgsegmat(P)
2 P = double(P); [n,m] = size(P);
3 spdata = zeros(4*n*m,1); spidx = zeros(4*n*m,2);
4 k = 1;
5 for ni=1:n
6   for mi=1:m
7     mni = (mi-1)*n+ni;
8     if (ni-1>0), spidx(k,:) = [mni,mni-1];
9     spdata(k) = Sfun(P(ni,mi),P(ni-1,mi));
10    k = k + 1;
11  end
12  if (ni+1<=n), spidx(k,:) = [mni,mni+1];
13  spdata(k) = Sfun(P(ni,mi),P(ni+1,mi));
14  k = k + 1;
15 end
16  if (mi-1>0), spidx(k,:) = [mni,mni-n];
17  spdata(k) = Sfun(P(ni,mi),P(ni,mi-1));
18  k = k + 1;
19 end
20  if (mi+1<=m), spidx(k,:) = [mni,mni+n];
21  spdata(k) = Sfun(P(ni,mi),P(ni,mi+1));
22  k = k + 1;
23 end
24 end
25 end
26 % Efficient initialization, see Sect. ??, Ex. ???
27 W = sparse(spidx(1:k-1,1),spidx(1:k-1,2),spdata(1:k-1),n*m,n*m);
28 D = spdiags(full(sum(W')'),[0],n*m,n*m);
29 A = D-W;
```

**Lemma 9.3.2.12. Ncut and Rayleigh quotient** ( $\rightarrow$  [SM00, Sect. 2])

With  $\mathbf{z} \in \{-1, 1\}^N$  according to (9.3.2.6) there holds

$$\text{Ncut}(\mathcal{X}) = \frac{\mathbf{y}^\top \mathbf{A} \mathbf{y}}{\mathbf{y}^\top \mathbf{D} \mathbf{y}} , \quad \mathbf{y} := (\mathbf{1} + \mathbf{z}) - \beta(\mathbf{1} - \mathbf{z}) , \quad \beta := \frac{\sum_{z_i > 0} d_i}{\sum_{z_i < 0} d_i} .$$

generalized Rayleigh quotient  $\rho_{\mathbf{A}, \mathbf{D}}(\mathbf{y})$

*Proof.* Note that by (9.3.2.6)  $(\mathbf{1} - \mathbf{z})_i = 0 \Leftrightarrow i \in \mathcal{X}$ ,  $(\mathbf{1} + \mathbf{z})_i = 0 \Leftrightarrow i \notin \mathcal{X}$ . Hence, since  $(\mathbf{1} + \mathbf{z})^\top \mathbf{D} (\mathbf{1} - \mathbf{z}) = 0$ ,

$$4 \text{Ncut}(\mathcal{X}) = (\mathbf{1} + \mathbf{z})^\top \mathbf{A} (\mathbf{1} + \mathbf{z}) \left( \frac{1}{\kappa \mathbf{1}^\top \mathbf{D} \mathbf{1}} + \frac{1}{(1 - \kappa) \mathbf{1}^\top \mathbf{D} \mathbf{1}} \right) = \frac{\mathbf{y}^\top \mathbf{A} \mathbf{y}}{\beta \mathbf{1}^\top \mathbf{D} \mathbf{1}} ,$$

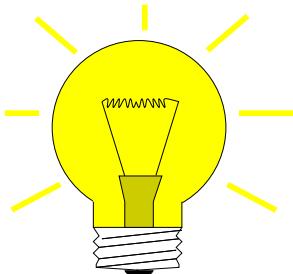
where  $\kappa := \sum_{z_i > 0} d_i / \sum_i d_i = \frac{\beta}{1 + \beta}$ . Also observe

$$\begin{aligned} \mathbf{y}^\top \mathbf{D} \mathbf{y} &= (\mathbf{1} + \mathbf{z})^\top \mathbf{D} (\mathbf{1} + \mathbf{z}) + \beta^2 (\mathbf{1} - \mathbf{z})^\top \mathbf{D} (\mathbf{1} - \mathbf{z}) = \\ &4 \left( \sum_{z_i > 0} d_i + \beta^2 \sum_{z_i < 0} d_i \right) = 4\beta \sum_i d_i = 4\beta \mathbf{1}^\top \mathbf{D} \mathbf{1} . \end{aligned}$$

This finishes the proof.  $\square$

- ◆ (9.3.2.10)  $\Rightarrow \mathbf{1} \in \text{Eig} \mathbf{A} 0$
- ◆ Lemma 2.8.0.12:  $\mathbf{A}$  diagonally dominant  $\implies \mathbf{A}$  is positive semidefinite ( $\rightarrow$  Def. 1.1.2.6)
- $\text{Ncut}(\mathcal{X}) \geq 0$  and 0 is the smallest eigenvalue of  $\mathbf{A}$ .

However, we are by no means interested in a minimizer  $\mathbf{y} \in \text{Span}\{\mathbf{1}\}$  (with constant entries) that does not provide a meaningful segmentation.



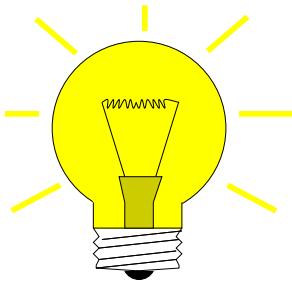
Idea: weed out undesirable constant minimizers by imposing **orthogonality constraint** (orthogonality w.r.t. inner product induced by  $\mathbf{D}$ , cf. Section 10.1)

$$\mathbf{y} \perp \mathbf{D} \mathbf{1} \Leftrightarrow \mathbf{1}^\top \mathbf{D} \mathbf{y} = 0 . \quad (9.3.2.13)$$

► segmentation problem (9.3.2.5)  $\Leftrightarrow \underset{\mathbf{y} \in \{2, -2\beta\}^N, \mathbf{1}^\top \mathbf{D} \mathbf{y} = 0}{\text{argmin}} \rho_{\mathbf{A}, \mathbf{D}}(\mathbf{y}) . \quad (9.3.2.14)$

still NP-hard

➤ Minimizing  $\text{Ncut}(\mathcal{X})$  amounts to minimizing a (generalized) Rayleigh quotient ( $\rightarrow$  Def. 9.3.1.16) over a discrete set of vectors, which is still an NP-hard problem.



Idea:

Relaxation

Discrete optimization problem → continuous optimization problem

$$(9.3.2.14) \rightarrow \underset{\mathbf{y} \in \mathbb{R}^N, \mathbf{y} \neq 0, \mathbf{1}^\top \mathbf{D}\mathbf{y} = 0}{\operatorname{argmin}} \rho_{\mathbf{A}, \mathbf{D}}(\mathbf{y}) . \quad (9.3.2.15)$$

Task: (9.3.2.15) ⇔ Find minimizer of (generalized) Rayleigh quotient under linear constraint

Here: linear constraint on  $\mathbf{y}$ :  $\mathbf{1}^\top \mathbf{D}\mathbf{y} = 0$ 

The next theorem establishes a link between argument vectors that render the Rayleigh quotient extremal and eigenspace for extremal eigenvalues.

### Theorem 9.3.2.16. Rayleigh quotient theorem

Let  $\lambda_1 < \lambda_2 < \dots < \lambda_m$ ,  $m \leq n$ , be the sorted sequence of all (real!) eigenvalues of  $\mathbf{A} = \mathbf{A}^H \in \mathbb{C}^{n,n}$ . Then

$$\operatorname{Eig}\mathbf{A}\lambda_1 = \underset{\mathbf{y} \in \mathbb{C}^{n,n} \setminus \{0\}}{\operatorname{argmin}} \rho_{\mathbf{A}}(\mathbf{y}) \quad \text{and} \quad \operatorname{Eig}\mathbf{A}\lambda_m = \underset{\mathbf{y} \in \mathbb{C}^{n,n} \setminus \{0\}}{\operatorname{argmax}} \rho_{\mathbf{A}}(\mathbf{y}) .$$

### Remark 9.3.2.17 (Min-max theorem)

Thm. 9.3.2.16 is a an immediate consequence of the following more general and fundamentally important result.

### Theorem 9.3.2.18. Courant-Fischer min-max theorem → [GV89, Thm. 8.1.2]

Let  $\lambda_1 < \lambda_2 < \dots < \lambda_m$ ,  $m \leq n$ , be the sorted sequence of the (real!) eigenvalues of  $\mathbf{A} = \mathbf{A}^H \in \mathbb{C}^{n,n}$ . Write

$$U_0 = \{0\} , \quad U_\ell := \sum_{j=1}^{\ell} \operatorname{Eig}\mathbf{A}\lambda_j , \quad \ell = 1, \dots, m \quad \text{and} \quad U_\ell^\perp := \{\mathbf{x} \in \mathbb{C}^n : \mathbf{u}^H \mathbf{x} = 0 \ \forall \mathbf{u} \in U_\ell\} .$$

Then

$$\min_{\mathbf{y} \in U_{\ell-1}^\perp \setminus \{0\}} \rho_{\mathbf{A}}(\mathbf{y}) = \lambda_\ell , \quad 1 \leq \ell \leq m , \quad \underset{\mathbf{y} \in U_{\ell-1}^\perp \setminus \{0\}}{\operatorname{argmin}} \rho_{\mathbf{A}}(\mathbf{y}) \subset \operatorname{Eig}\mathbf{A}\lambda_\ell .$$

*Proof.* For diagonal  $\mathbf{A} \in \mathbb{R}^{n,n}$  the assertion of the theorem is obvious. Thus, Cor. 9.1.0.9 settles everything.  $\square$

A simple conclusion from Thm. 9.3.2.18: If  $\mathbf{A} = \mathbf{A}^\top \in \mathbb{R}^{n,n}$  with eigenvalues  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ , then

$$\lambda_1 = \min_{\mathbf{z} \in \mathbb{R}^n} \rho_{\mathbf{A}}(\mathbf{z}) , \quad \lambda_2 = \min_{\mathbf{z} \in \mathbb{R}^n, \mathbf{z} \perp \mathbf{v}_1} \rho_{\mathbf{A}}(\mathbf{z}) , \quad (9.3.2.19)$$

where  $\mathbf{v}_1 \in \text{Eig}_{\mathbf{A}} \lambda_1 \setminus \{0\}$ . □

Well, in Lemma 9.3.2.12 we encounter a generalized Rayleigh quotient  $\rho_{\mathbf{A}, \mathbf{D}}(\mathbf{y})$ ! How can Thm. 9.3.2.16 be applied to it?

► Transformation idea:  $\rho_{\mathbf{A}, \mathbf{D}}(\mathbf{D}^{-1/2}\mathbf{z}) = \rho_{\mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}}(\mathbf{z})$ ,  $\mathbf{y} \in \mathbb{R}^n$ . (9.3.2.20)

► Apply Thm. 9.3.2.18 to transformed matrix  $\tilde{\mathbf{A}} := \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ . Elementary manipulations show

$$(9.3.2.15) \Leftrightarrow \begin{aligned} & \underset{\mathbf{1}^\top \mathbf{D}\mathbf{y}=0}{\operatorname{argmin}} \rho_{\mathbf{A}, \mathbf{D}}(\mathbf{y}) \stackrel{\mathbf{z}=\mathbf{D}^{1/2}\mathbf{y}}{=} \underset{\mathbf{1}^\top \mathbf{D}^{1/2}\mathbf{z}=0}{\operatorname{argmin}} \rho_{\mathbf{A}, \mathbf{D}}(\mathbf{D}^{-1/2}\mathbf{z}) \\ &= \underset{\mathbf{1}^\top \mathbf{D}^{1/2}\mathbf{z}=0}{\operatorname{argmin}} \rho_{\tilde{\mathbf{A}}}(\mathbf{z}) \quad \text{with } \tilde{\mathbf{A}} := \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}. \end{aligned} \quad (9.3.2.21)$$

Related: transformation of a generalized eigenvalue problem into a standard eigenvalue problem according to

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{B}\mathbf{x} \quad \xrightarrow{\mathbf{z}=\mathbf{B}^{1/2}\mathbf{x}} \quad \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}\mathbf{z} = \lambda \mathbf{z}. \quad (9.3.2.22)$$

$\mathbf{B}^{1/2}$  ≈ square root of s.p.d. matrix  $\mathbf{B}$  → Rem. 10.3.0.2.

For segmentation problem:  $\mathbf{B} = \mathbf{D}$  diagonal with positive diagonal entries, see (9.3.2.9)

►  $\mathbf{D}^{-1/2} = \text{diag}(d_1^{-1/2}, \dots, d_N^{-1/2})$  and  $\tilde{\mathbf{A}} := \mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$  can easily be computed.

► In the sequel consider minimization problem/related eigenvalue problem

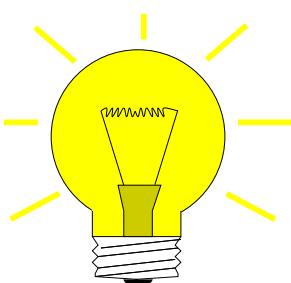
$$\boxed{\mathbf{z}^* = \underset{\mathbf{1}^\top \mathbf{D}^{1/2}\mathbf{z}=0}{\operatorname{argmin}} \rho_{\tilde{\mathbf{A}}}(\mathbf{z})} \quad \longleftrightarrow \quad \tilde{\mathbf{A}}\mathbf{z} = \lambda \mathbf{z}. \quad (9.3.2.23)$$

Recover solution  $\mathbf{y}^*$  of (9.3.2.15) as  $\mathbf{y}^* = \mathbf{D}^{-1/2}\mathbf{z}^*$ .



Still, Thm. 9.3.2.16 cannot be applied to (9.3.2.23):

How to deal with constraint  $\mathbf{1}^\top \mathbf{D}^{1/2}\mathbf{z} = 0$  ?



Idea:

Penalization

Add term  $P(\mathbf{z})$  to  $\rho_{\tilde{\mathbf{A}}}(\mathbf{z})$  that becomes “sufficiently large” in case the constraint is violated.

►  $\mathbf{z}^*$  can only be a minimizer of  $\rho_{\tilde{\mathbf{A}}}(\mathbf{z}) + P(\mathbf{z})$ , if  $P(\mathbf{z}) = 0$ .

How to choose the **penalty function**  $P(\mathbf{z})$  for the segmentation problem ?

$$\left\{ \mathbf{1}^\top \mathbf{D}^{1/2} \mathbf{z} \neq 0 \Rightarrow P(\mathbf{z}) > 0 \right\} \text{ satisfied for}$$

$$P(\mathbf{z}) = \mu \frac{|\mathbf{1}^\top \mathbf{D}^{1/2} \mathbf{z}|^2}{\|\mathbf{z}\|_2^2},$$

with **penalty parameter**  $\mu > 0$ .

► **penalized minimization problem**

$$\begin{aligned} \mathbf{z}^* &= \underset{\mathbf{z} \in \mathbb{R}^N \setminus \{0\}}{\operatorname{argmin}} \rho_{\tilde{\mathbf{A}}}(\mathbf{z}) + P(\mathbf{z}) = \underset{\mathbf{z} \in \mathbb{R}^N \setminus \{0\}}{\operatorname{argmin}} \rho_{\tilde{\mathbf{A}}}(\mathbf{z}) + \frac{\mathbf{z}^\top (\mathbf{D}^{1/2} \mathbf{1} \mathbf{1}^\top \mathbf{D}^{1/2}) \mathbf{z}}{\mathbf{z}^\top \mathbf{z}} \\ &= \underset{\mathbf{z} \in \mathbb{R}^N \setminus \{0\}}{\operatorname{argmin}} \rho_{\hat{\mathbf{A}}}(\mathbf{z}) \quad \text{with } \hat{\mathbf{A}} := \tilde{\mathbf{A}} + \mathbf{D}^{1/2} \mathbf{1} \mathbf{1}^\top \mathbf{D}^{1/2}. \end{aligned} \tag{9.3.2.24}$$

dense rank-1 matrix

How to choose the penalty parameter  $\mu$  ?

In general: finding a “suitable” value for  $\mu$  may be difficult or even impossible!

Here we are lucky:

$$(9.3.2.10) \Rightarrow \mathbf{A} \mathbf{1} = 0 \Rightarrow \tilde{\mathbf{A}}(\mathbf{D}^{1/2} \mathbf{1}) = 0 \Leftrightarrow \mathbf{D}^{1/2} \mathbf{1} \in \operatorname{Eig} \tilde{\mathbf{A}} 0.$$

► Constraint in (9.3.2.23) means

Minimize over the **orthogonal complement** of an *eigenvector*. (9.3.2.25)

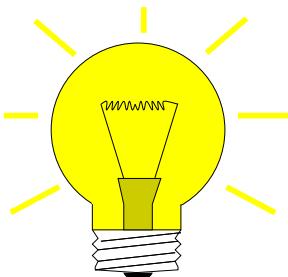
Cor. 9.1.0.9 ► The orthogonal complement of an eigenvector of a *symmetric* matrix is spanned by the other eigenvectors (orthonormalization of eigenvectors belonging to the same eigenvalue is assumed).

(9.3.2.25) ► The minimizer of (9.3.2.23) will be one of the other eigenvectors of  $\tilde{\mathbf{A}}$  that belongs to the smallest eigenvalue.

Note: This eigenvector  $\mathbf{z}^*$  will be orthogonal to  $\mathbf{D}^{1/2} \mathbf{1}$ , it satisfies the constraint, and, thus,  $P(\mathbf{z}^*) = 0$ !

Note: eigenspaces of  $\tilde{\mathbf{A}}$  and  $\hat{\mathbf{A}}$  agree.

Note: Lemma 2.8.0.12  $\implies \tilde{\mathbf{A}}$  is *positive semidefinite* ( $\rightarrow$  Def. 1.1.2.6) with smallest eigenvalue 0.n



Idea: Choose penalization parameter  $\mu$  in (9.3.2.24) such that  $\mathbf{D}^{1/2} \mathbf{1}$  is guaranteed not to be an eigenvector belonging to the smallest eigenvalue of  $\hat{\mathbf{A}}$ .

Safe choice: choose  $\mu$  such that  $\mathbf{D}^{1/2} \mathbf{1}$  will belong to the largest eigenvalue of  $\hat{\mathbf{A}}$ : Thm. 9.1.0.4 suggests

$$\mu = \|\tilde{\mathbf{A}}\|_\infty \stackrel{(1.5.5.13)}{=} 2. \tag{9.3.2.26}$$



$$\mathbf{z}^* = \underset{\mathbf{1}^\top \mathbf{D}^{1/2} \mathbf{z} = 0}{\operatorname{argmin}} \rho_{\tilde{\mathbf{A}}}(\mathbf{z}) = \underset{\mathbf{z} \neq 0}{\operatorname{argmin}} \rho_{\tilde{\mathbf{A}}}(\mathbf{z}). \quad (9.3.2.27)$$

By Thm. 9.3.2.16:

$$\begin{aligned} \mathbf{z}^* &= \text{eigenvector belonging to minimal eigenvalue of } \hat{\mathbf{A}}, \\ &\Updownarrow \\ \mathbf{z}^* &= \text{eigenvector } \perp \mathbf{D}^{1/2}\mathbf{1} \text{ belonging to minimal eigenvalue of } \tilde{\mathbf{A}}, \\ &\Updownarrow \\ \mathbf{D}^{-1/2}\mathbf{z}^* &= \text{minimizer for (9.3.2.15)}. \end{aligned}$$

### §9.3.2.28 (Algorithm outline: Binary grayscale image segmentation)

- ① Given similarity function  $\sigma$  compute (sparse!) matrices  $\mathbf{W}, \mathbf{D}, \mathbf{A} \in \mathbb{R}^{N,N}$ , see (9.3.2.3), (9.3.2.9).
- ② Compute  $\mathbf{y}^*$ ,  $\|\mathbf{y}^*\|_2 = 1$ , as eigenvector belonging to the *smallest eigenvalue* of  $\hat{\mathbf{A}} := \mathbf{D}^{-1/2}\mathbf{AD}^{-1/2} + 2(\mathbf{D}^{1/2}\mathbf{1})(\mathbf{D}^{1/2}\mathbf{1})^\top$ .
- ③ Set  $\mathbf{x}^* = \mathbf{D}^{-1/2}\mathbf{y}^*$  and define the image segment as pixel set

$$\mathcal{X} := \left\{ i \in \{1, \dots, N\} : x_i^* > \frac{1}{N} \sum_{i=1}^N x_i^* \right\}. \quad (9.3.2.29)$$

↑  
mean value of entries of  $\mathbf{x}^*$

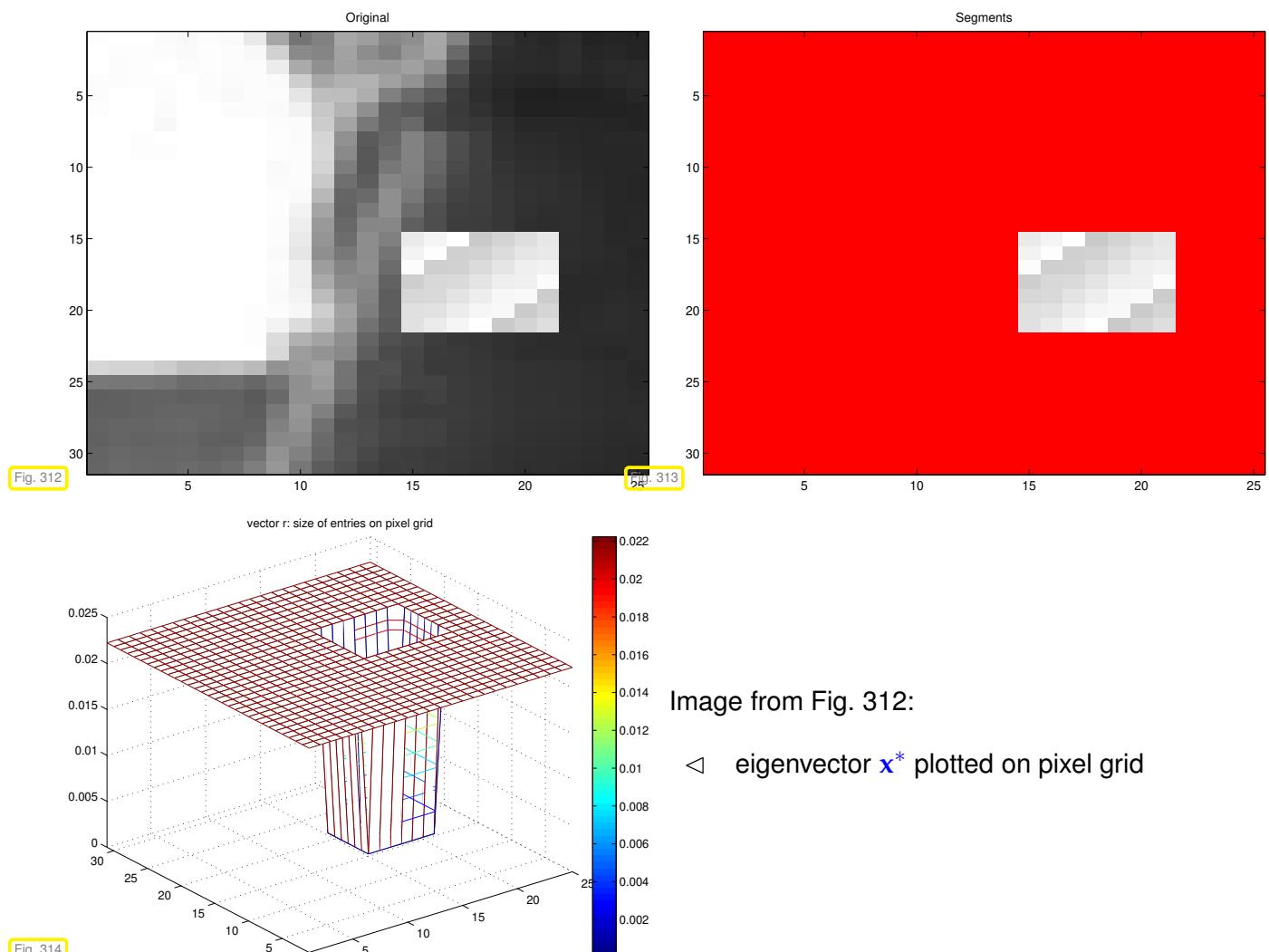
#### MATLAB-code 9.3.2.30: 1st stage of segmentation of grayscale image

```

1 % Read image and build matrices, see Code 9.3.2.11 and (9.3.2.9)
2 P = imread('image.pbm'); [m,n] = size(P); [A,D] = imgsegmat(P);
3 % Build scaling matrices
4 N = size(A,1); dv = sqrt(spdiags(A,0));
5 Dm = spdiags(1./dv,[0],N,N); % D^{-1/2}
6 Dp = spdiags(dv,[0],N,N); % D^{-1/2}
7 % Build (densely populated !) matrix \hat{A}
8 C = Dp*ones(N,1); Ah = Dm*A*Dm + 2*c*c';
9 % Compute and sort eigenvalues; grossly inefficient !
10 [W,E] = eig(full(Ah)); [ev,idx] = sort(diag(E)); W(:,idx) = W;
11 % Obtain eigenvector x^* belonging to 2nd smallest generalized
12 % eigenvalue of A and D
13 x = W(:,1); x = Dm*x;
14 % Extract segmented image
15 xs = reshape(x,m,n); Xidx = find(xs > (sum(sum(xs)) / (n*m)));

```

1st-stage of segmentation of  $31 \times 25$  grayscale pixel image (`root.pbm`, red pixels  $\hat{=} \mathcal{X}$ ,  $\sigma(x,y) = \exp(-(x-y/10)^2)$ )



To identify more segments, the same algorithm is *recursively applied* to segment parts of the image already determined.

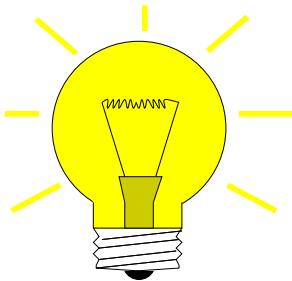
Practical segmentation algorithms rely on many more steps of which the above algorithm is only one, preceded by substantial preprocessing. Moreover, they dispense with the strictly local perspective adopted above and take into account more distant connections between image parts, often in a randomized fashion [SM00].

The image segmentation problem falls into the wider class of **graph partitioning problems**. Methods based on (a few of) the eigenvector of the connectivity matrix belonging to the smallest eigenvalues are known as **spectral partitioning methods**. The eigenvector belonging to the smallest non-zero eigenvalue that we computed above is usually called the **Fiedler vector** of the graph, see [AKY99; ST96].

The solution of the image segmentation problem by means of `eig` in Code 9.3.2.30 amounts a tremendous waste of computational resources: we compute *all* eigenvalues/eigenvectors of dense matrices, though only a single eigenvector associated with the smallest eigenvalue is of interest.

This motivates the quest to find *efficient* numerical methods for the following task.

Task: given  $\mathbf{A} \in \mathbb{K}^{n,n}$ , find **smallest** (in modulus) eigenvalue of regular  $\mathbf{A} \in \mathbb{K}^{n,n}$  and (an) associated eigenvector.



If  $\mathbf{A} \in \mathbb{K}^{n,n}$  regular:

$$\text{Smallest (in modulus) EV of } \mathbf{A} = \left( \text{Largest (in modulus) EV of } \mathbf{A}^{-1} \right)^{-1}$$

► Direct power method ( $\rightarrow$  Section 9.3.1) for  $\mathbf{A}^{-1} = \text{inverse iteration}$

#### MATLAB-code 9.3.2.31: inverse iteration for computing $\lambda_{\min}(\mathbf{A})$ and associated eigenvector

```

1 function [lmin,y] = invit(A,tol)
2 [L,U] = lu(A); % single intial LU-factorization, see Rem. 2.5.0.10
3 n = size(A,1); x = rand(n,1); x = x/norm(x); % random initial guess
4 y = U\ (L\x); lmin = 1/norm(y); y = y*lmin; lold = 0;
5 while (abs(lmin-lold) > tol*lmin) % termination, if small relative
   change
6   lold = lmin; x = y;
7   y = U\ (L\x); % core iteration: y = A-1x,
8   lmin = 1/norm(y); % new approximation of λmin(A)
9   y = y*lmin; % normalization y := y / ||y||2
10 end

```

Note: reuse of LU-factorization, see Rem. 2.5.0.10

#### Remark 9.3.2.32 (Shifted inverse iteration)

More general task:

For  $\alpha \in \mathbb{C}$  find  $\lambda \in \sigma(\mathbf{A})$  such that  $|\alpha - \lambda| = \min\{|\alpha - \mu|, \mu \in \sigma(\mathbf{A})\}$

► Shifted inverse iteration: [DR08, Alg .7.24]

$$\mathbf{z}^{(0)} \text{ arbitrary}, \quad \mathbf{w} = (\mathbf{A} - \alpha \mathbf{I})^{-1} \mathbf{z}^{(k-1)}, \quad \mathbf{z}^{(k)} := \frac{\mathbf{w}}{\|\mathbf{w}\|_2}, \quad k = 1, 2, \dots, \quad (9.3.2.33)$$

where:  $(\mathbf{A} - \alpha \mathbf{I})^{-1} \mathbf{z}^{(k-1)} \triangleq \text{solve } (\mathbf{A} - \alpha \mathbf{I})\mathbf{w} = \mathbf{z}^{(k-1)}$  based on Gaussian elimination ( $\leftrightarrow$  a single LU-factorization of  $\mathbf{A} - \alpha \mathbf{I}$  as in Code 9.3.2.31).

#### Remark 9.3.2.34 ((Nearly) singular LSE in shifted inverse iteration)

What if “by accident”  $\alpha \in \sigma(\mathbf{A})$  ( $\Leftrightarrow \mathbf{A} - \alpha \mathbf{I}$  singular)?

Stability of Gaussian elimination/LU-factorization ( $\rightarrow ??$ ) will ensure that “ $\mathbf{w}$  from (9.3.2.33) points in the right direction”

In other words, roundoff errors may badly affect the length of the solution  $\mathbf{w}$ , but not its direction.

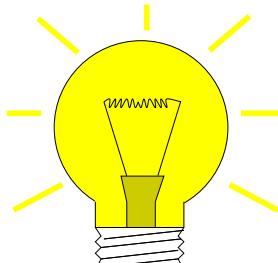
Practice [GT08]: If, in the course of Gaussian elimination/LU-factorization a zero pivot element is really encountered, then we just replace it with  $\text{eps}$ , in order to avoid  $\text{inf}$  values!

Thm. 9.3.1.21 ➤ Convergence of shifted inverse iteration for  $\mathbf{A}^H = \mathbf{A}$ :

Asymptotic linear convergence, Rayleigh quotient  $\rightarrow \lambda_j$  with rate

$$\frac{|\lambda_j - \alpha|}{\min\{|\lambda_i - \alpha|, i \neq j\}} \quad \text{with} \quad \lambda_j \in \sigma(\mathbf{A}), \quad |\alpha - \lambda_j| \leq |\alpha - \lambda| \quad \forall \lambda \in \sigma(\mathbf{A}).$$

► Extremely fast for  $\alpha \approx \lambda_j$ !



Idea:

A posteriori adaptation of shift

Use  $\alpha := \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})$  in  $k$ -th step of inverse iteration.

### §9.3.2.35 (Rayleigh quotient iteration) → [Han02, Alg. 25.2])

#### MATLAB-code 9.3.2.36: Rayleigh quotient iteration (for normal $\mathbf{A} \in \mathbb{R}^{n,n}$ )

```

1 function [z, lmin] = rqui(A, tol, maxit)
2 alpha = 0; n = size(A,1);
3 z = rand(size(A,1),1); z = z/norm(z); % z^(0)
4 for i=1:maxit
5     z = (A-alpha*speye(n))\z; % z^(k+1) = (A - rho_A(z^(k))I)^{-1}x^(k)
6     z = z/norm(z); lmin=dot(A*z,z); % Computation of rho_A(z^(k+1))
7     if (abs(alpha-lmin) < tol*lmin) % Desired relative accuracy
        reached ?
        break; end;
8     alpha = lmin;
9 end
10

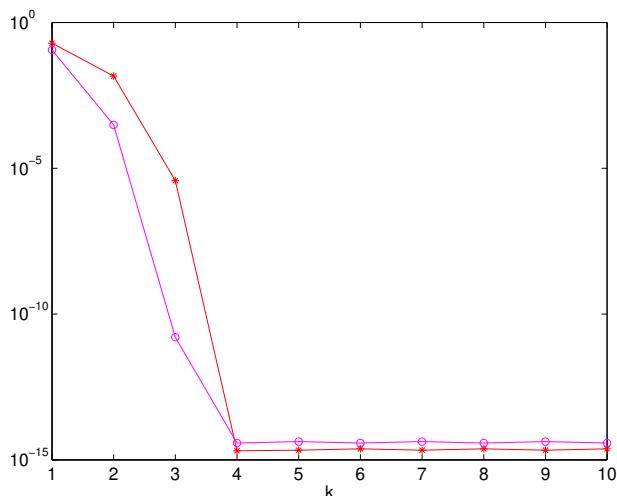
```

Line 5: note use of `speye` to preserve sparse matrix data format!

- ◆ Drawback compared with Code 9.3.2.31: reuse of LU-factorization no longer possible.
- ◆ Even if LSE nearly singular, stability of Gaussian elimination guarantees correct direction of  $\mathbf{z}$ , see discussion in Rem. 9.3.2.34.

#### EXAMPLE 9.3.2.37 (Rayleigh quotient iteration)

Monitored: iterates of Rayleigh quotient iteration (9.3.2.36) for s.p.d.  $\mathbf{A} \in \mathbb{R}^{n,n}$



$k$	$ \lambda_{\min} - \rho_{\mathbf{A}}(\mathbf{z}^{(k)}) $	$\ \mathbf{z}^{(k)} - \mathbf{x}_j\ $
1	0.09381702342056	0.20748822490698
2	0.00029035607981	0.01530829569530
3	0.000000000001783	0.000000411928759
4	0.000000000000000	0.000000000000000
5	0.000000000000000	0.000000000000000

**MATLAB-code 9.3.2.38:**

```
d = (1:10)';
n = length(d);
Z = diag(sqrt(1:n),0)+ones(n,n);
[Q,R] = qr(Z);
```

```
A = Q*diag(d,0)*Q';
```

o :  $|\lambda_{\min} - \rho_{\mathbf{A}}(\mathbf{z}^{(k)})|$   
\* :  $\|\mathbf{z}^{(k)} - \mathbf{x}_j\|$ ,  $\lambda_{\min} = \lambda_j$ ,  $\mathbf{x}_j \in \text{Eig}\mathbf{A}\lambda_j$ ,  
:  $\|\mathbf{x}_j\|_2 = 1$

**Theorem 9.3.2.39.** → [Han02, Thm. 25.4]

If  $\mathbf{A} = \mathbf{A}^H$ , then  $\rho_{\mathbf{A}}(\mathbf{z}^{(k)})$  converges locally of order 3 (→ Def. 8.1.1.10) to the smallest eigenvalue (in modulus), when  $\mathbf{z}^{(k)}$  are generated by the Rayleigh quotient iteration (9.3.2.36).

### 9.3.3 Preconditioned inverse iteration (PINVIT)

Task: given  $\mathbf{A} \in \mathbb{K}^{n,n}$ , find **smallest** (in modulus) eigenvalue of regular  $\mathbf{A} \in \mathbb{K}^{n,n}$  and (an) associated eigenvector.

► Options: inverse iteration (→ Code 9.3.2.31) and Rayleigh quotient iteration (9.3.2.36).



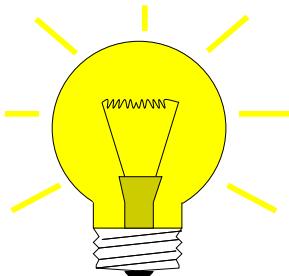
What if direct solution of  $\mathbf{Ax} = \mathbf{b}$  not feasible ?

This can happen, in case

- for large sparse  $\mathbf{A}$  the amount of fill-in exhausts memory, despite sparse elimination techniques (→ Section 2.7.4),
- $\mathbf{A}$  is available only through a routine `evalA(x)` providing  $\mathbf{A} \times$ vector.

We expect that an approximate solution of the linear systems of equations encountered during inverse iteration should be sufficient, because we are dealing with approximate eigenvectors anyway.

Thus, iterative solvers for solving  $\mathbf{Aw} = \mathbf{z}^{(k-1)}$  may be considered, see Chapter 10. However, the required accuracy is not clear a priori. Here we examine an approach that completely dispenses with an iterative solver and uses a *preconditioner* (→ Notion 10.3.0.3) instead.



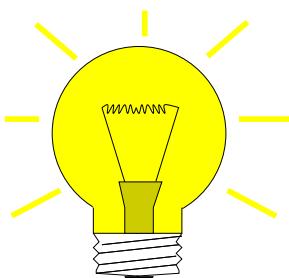
Idea: (for inverse iteration without shift,  $\mathbf{A} = \mathbf{A}^H$  s.p.d.)  
Instead of solving  $\mathbf{Aw} = \mathbf{z}^{(k-1)}$  compute  $\mathbf{w} = \mathbf{B}^{-1}\mathbf{z}^{(k-1)}$  with  
“inexpensive” s.p.d. approximate inverse  $\mathbf{B}^{-1} \approx \mathbf{A}^{-1}$

>  $\mathbf{B} \triangleq$  Preconditioner for  $\mathbf{A}$ , see Notion 10.3.0.3



Possible to replace  $\mathbf{A}^{-1}$  with  $\mathbf{B}^{-1}$  in inverse iteration ?

**NO**, because we are not interested in smallest eigenvalue of  $\mathbf{B}$  !



Replacement  $\mathbf{A}^{-1} \rightarrow \mathbf{B}^{-1}$  possible only when applied to residual quantity  
residual quantity = quantity that  $\rightarrow 0$  in the case of convergence to exact solution

Natural residual quantity for eigenvalue problem  $\mathbf{Ax} = \lambda\mathbf{x}$ :

$$\mathbf{r} := \mathbf{Az} - \rho_{\mathbf{A}}(\mathbf{z})\mathbf{z}, \quad \rho_{\mathbf{A}}(\mathbf{z}) = \text{Rayleigh quotient} \rightarrow \text{Def. 9.3.1.16}.$$

Note: only *direction* of  $\mathbf{A}^{-1}\mathbf{z}$  matters in inverse iteration (9.3.2.33)

$$(\mathbf{A}^{-1}\mathbf{z}) \parallel (\mathbf{z} - \mathbf{A}^{-1}(\mathbf{Az} - \rho_{\mathbf{A}}(\mathbf{z})\mathbf{z})) \Rightarrow \text{defines same next iterate!}$$

[Preconditioned inverse iteration (PINVIT) for s.p.d.  $\mathbf{A}$ ]

$$\begin{aligned} \mathbf{z}^{(0)} \text{ arbitrary, } \mathbf{w} &= \mathbf{z}^{(k-1)} - \mathbf{B}^{-1}(\mathbf{Az}^{(k-1)} - \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})\mathbf{z}^{(k-1)}), \\ \mathbf{z}^{(k)} &= \frac{\mathbf{w}}{\|\mathbf{w}\|_2}, \end{aligned} \quad k = 1, 2, \dots. \quad (9.3.3.1)$$

**MATLAB-code 9.3.3.2: preconditioned inverse iteration (9.3.3.1)**

```

1 function [lmin,z,res] =
2     pinvit(evalA,n,invB,tol,maxit)
3 % invB ≈ handle to function implementing
4 % preconditioner  $B^{-1}$ 
5 z = (1:n)'; z = z/norm(z); % initial guess
6 res = []; rho = 0;
7 for i=1:maxit
8     v = evalA(z); rhon = dot(v,z); % Rayleigh
9         quotient
10    r = v - rhon*z; % residual
11    z = z - invB(r); % iteration according to
12        (9.3.3.1)
13    z = z/norm(z); % normalization
14    res = [res; rhon]; % tracking iteration
15    if abs(rho-rhon) < tol*abs(rhon)), break;
16    else rho = rhon; end
17 end
18 lmin = dot(evalA(z),z); res = [res; lmin],

```

Computational effort:

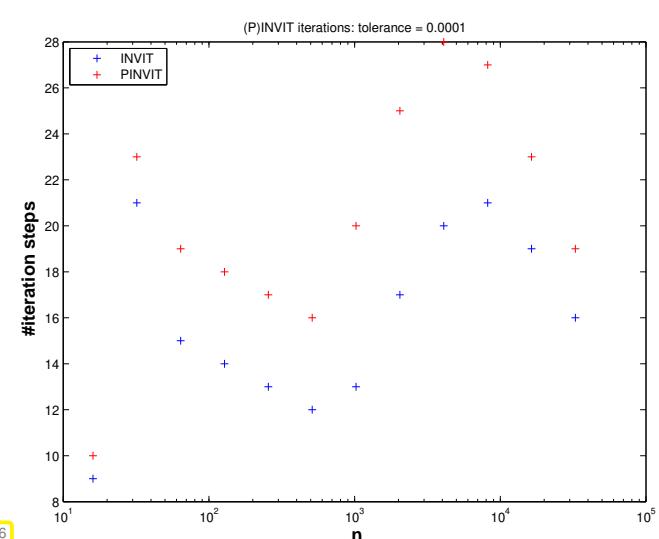
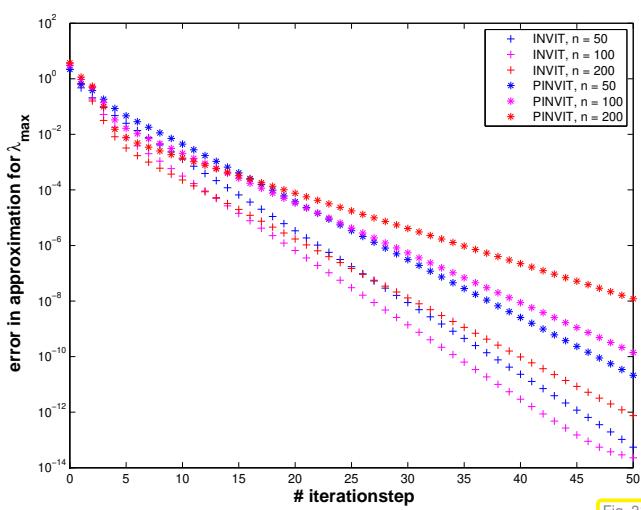
1 matrix  $\times$  vector  
 1 evaluation of pre-  
 conditioner  
 $A$  few  
 AXPY-operations

**EXAMPLE 9.3.3.3 (Convergence of PINVIT)**S.p.d. matrix  $A \in \mathbb{R}^{n,n}$ , tridiagonal preconditioner, see Ex. 10.3.0.11**MATLAB-code 9.3.3.4:**

```

1 A = spdiags(repmat([1/n,-1,2*(1+1/n),-1,1/n],n,1),
2 [-n/2,-1,0,1,n/2],n,n);
3 evalA = @(x) A*x;
4 % inverse iteration
5 invB = @(x) A\b;
6 % tridiagonal preconditioning
7 B = spdiags(spdiags(A,[-1,0,1]),[-1,0,1],n,n); invB = @(x) B\b;

```

Monitored: error decay during iteration of Code 9.3.3.2:  $|\rho_A(z^{(k)}) - \lambda_{\min}(A)|$ 

Observation: linear convergence of eigenvectors also for PINVIT.

↓

Theory [Ney99a; Ney99b]:

- ◆ linear convergence of (9.3.3.1)
- ◆ fast convergence, if spectral condition number  $\kappa(\mathbf{B}^{-1}\mathbf{A})$  small

The theory of PINVIT [Ney99a; Ney99b] is based on the identity

$$\mathbf{w} = \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})\mathbf{A}^{-1}\mathbf{z}^{(k-1)} + (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})(\mathbf{z}^{(k-1)} - \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})\mathbf{A}^{-1}\mathbf{z}^{(k-1)}). \quad (9.3.3.5)$$

For small residual  $\mathbf{A}\mathbf{z}^{(k-1)} - \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})\mathbf{z}^{(k-1)}$  PINVIT almost agrees with the regular inverse iteration.

### 9.3.4 Subspace iterations

#### Remark 9.3.4.1 (Excited resonances)

Consider the non-autonomous ODE (excited harmonic oscillator)

$$\ddot{y} + \lambda^2 y = \cos(\omega t), \quad (9.3.4.2)$$

with general solution

$$y(t) = \begin{cases} \frac{1}{\lambda^2 - \omega^2} \cos(\omega t) + A \cos(\lambda t) + B \sin(\lambda t) & , \text{ if } \lambda \neq \omega , \\ \frac{t}{2\omega} \sin(\omega t) + A \cos(\lambda t) + B \sin(\lambda t) & , \text{ if } \lambda = \omega . \end{cases} \quad (9.3.4.3)$$

► growing solutions possible in **resonance case**  $\lambda = \omega$  !

Now consider harmonically excited vibration modelled by ODE

$$\ddot{\mathbf{y}} + \mathbf{A}\mathbf{y} = \mathbf{b} \cos(\omega t), \quad (9.3.4.4)$$

with *symmetric, positive (semi)definite* matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ . By Cor. 9.1.0.9 there is an *orthogonal* matrix  $\mathbf{Q} \in \mathbb{R}^{n,n}$  such that

$$\mathbf{Q}^\top \mathbf{A} \mathbf{Q} = \mathbf{D} := \text{diag}(\lambda_1, \dots, \lambda_n).$$

where the  $0 \leq \lambda_1 < \lambda_2 < \dots < \lambda_n$  are the eigenvalues of  $\mathbf{A}$ .

► Transform ODE as in Ex. 9.0.0.7: with  $\mathbf{z} = \mathbf{Q}^\top \mathbf{y}$

$$(9.3.4.4) \quad \ddot{\mathbf{z}} + \mathbf{D}\mathbf{z} = \mathbf{Q}^\top \mathbf{b} \cos(\omega t).$$

We have obtained decoupled linear 2nd-order scalar ODEs of the type (9.3.4.2).

► (9.3.4.4) can have growing (with time) solutions, if  $\omega = \sqrt{\lambda_i}$  for some  $i = 1, \dots, n$

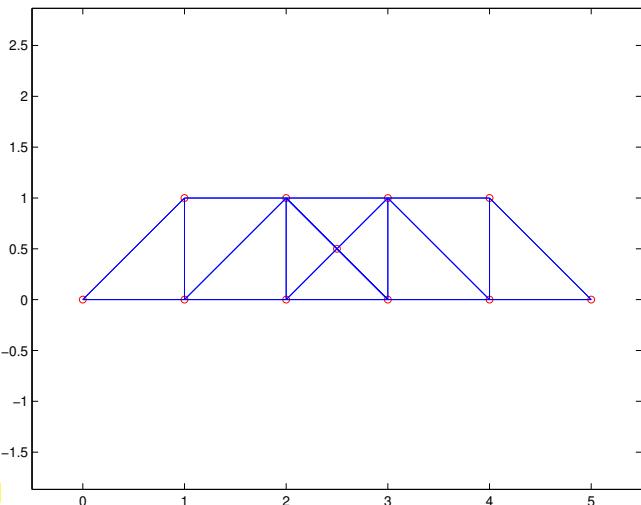
If  $\omega = \sqrt{\lambda_j}$  for one  $j \in \{1, \dots, n\}$ , then the solution for the initial value problem for (9.3.4.4) with  $\mathbf{y}(0) = \dot{\mathbf{y}}(0) = 0$  ( $\leftrightarrow \mathbf{z}(0) = \dot{\mathbf{z}}(0) = 0$ ) is

$$\begin{aligned} \mathbf{z}(t) &\sim \frac{t}{2\omega} \sin(\omega t) \mathbf{e}_j + \text{bounded oscillations} \\ &\Updownarrow \\ \mathbf{y}(t) &\sim \frac{t}{2\omega} \sin(\omega t) (\mathbf{Q})_{:,j} + \text{bounded oscillations}. \end{aligned}$$

$\downarrow$   
*j*-th eigenvector of  $\mathbf{A}$

► Eigenvectors of  $\mathbf{A} \leftrightarrow$  excitable states

### EXAMPLE 9.3.4.5 (Vibrations of a truss structure) cf. [Han02, Sect. 3], MATLAB's truss demo)



▷ a “bridge” truss]

A **truss** is a structure composed of (massless) rods and point masses; we consider in-plane (2D) trusses.

Encoding: positions of masses + (sparse) connectivity matrix

Fig. 317

#### MATLAB-code 9.3.4.6: Data for “bridge truss”

```

1 % Data for truss structure "bridge"
2 pos = [ 0 0; 1 0; 2 0; 3 0; 4 0; 5 0; 1 1; 2 1; 3 1; 4 1; 2.5 0.5];
3 con = [1 2; 2 3; 3 4; 4 5; 5 6; 1 7; 2 7; 3 8; 2 8; 4 8; 5 9; 5 10; 6 10; 7
8; 8 9; 9 10; 8 11; ...
; 9 11; 3 11; 4 11; 4 9];
4 n = size(pos,1);
5 top = sparse(con(:,1),con(:,2),ones(size(con,1),1),n,n);
6 top = sign(top+top');
7

```

Assumptions:

- ◆ Truss in static equilibrium (perfect balance of forces at each point mass).

- ◆ Rods are *perfectly elastic* (i.e., frictionless).

► Hook's law holds for force in the direction of a rod:

$$F = \alpha \frac{\Delta l}{l}, \quad (9.3.4.7)$$

where ◆  $l$  is the equilibrium length of the rod,

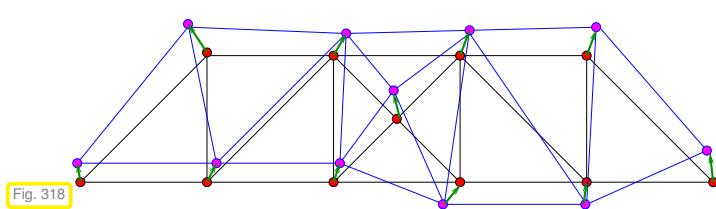
◆  $\Delta l$  is the elongation of the rod effected by the force  $F$  in the direction of the rod

◆  $\alpha$  is a material coefficient (Young's modulus).

$n$  point masses are numbered  $1, \dots, n$ :  $\mathbf{p}^i \in \mathbb{R}^2 \hat{=} \text{position of } i\text{-th mass}$

We consider a swaying truss: description by time-dependent **displacements**  $\mathbf{u}^i(t) \in \mathbb{R}^2$  of point masses:

$$\text{position of } i\text{-th mass at time } t = \mathbf{p}^i + \mathbf{u}^i(t);.$$



$\triangleleft$  deformed truss:

- $\hat{=}$  point masses at positions  $\mathbf{p}^i$
- $\rightarrow \hat{=}$  displacement vectors  $\mathbf{u}^i$
- $\hat{=}$  shifted masses at  $\mathbf{p}^i + \mathbf{u}^i$

Equilibrium length and (time-dependent) elongation of rod connecting point masses  $i$  and  $j$ ,  $i \neq j$ :

$$l_{ij} := \|\Delta \mathbf{p}^{ji}\|_2, \quad \Delta \mathbf{p}^{ji} := \mathbf{p}^j - \mathbf{p}^i, \quad (9.3.4.8)$$

$$\Delta l_{ij}(t) := \|\Delta \mathbf{p}^{ji} + \Delta \mathbf{u}^{ji}(t)\|_2 - l_{ij}, \quad \Delta \mathbf{u}^{ji}(t) := \mathbf{u}^j(t) - \mathbf{u}^i(t). \quad (9.3.4.9)$$

$\blacktriangleright$  Extra (reaction) force on masses  $i$  and  $j$ :

$$F_{ij}(t) = -\alpha_{ij} \frac{\Delta l_{ij}}{l_{ij}} \cdot \frac{\Delta \mathbf{p}^{ji} + \Delta \mathbf{u}^{ji}(t)}{\|\Delta \mathbf{p}^{ji} + \Delta \mathbf{u}^{ji}(t)\|_2}. \quad (9.3.4.10)$$

Assumption:

Small displacements

$\blacktriangleright$  Possibility of **linearization** by neglecting terms of order  $\|\mathbf{u}^i\|_2^2$

$$(9.3.4.8) \quad \blacktriangleright \quad F_{ij}(t) = \alpha_{ij} \left( \frac{1}{\|\Delta \mathbf{p}^{ji} + \Delta \mathbf{u}^{ji}(t)\|_2} - \frac{1}{\|\Delta \mathbf{p}^{ji}\|} \right) \cdot (\Delta \mathbf{p}^{ji} + \Delta \mathbf{u}^{ji}(t)). \quad (9.3.4.11)$$

### Lemma 9.3.4.12. Taylor expansion of inverse distance function

For  $\mathbf{x} \in \mathbb{R}^d \setminus \{0\}$ ,  $\mathbf{y} \in \mathbb{R}^d$ ,  $\|\mathbf{y}\|_2 < \|\mathbf{x}\|_2$  holds for  $\mathbf{y} \rightarrow 0$

$$\frac{1}{\|\mathbf{x} + \mathbf{y}\|_2} = \frac{1}{\|\mathbf{x}\|_2} - \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|_2^3} + O(\|\mathbf{y}\|_2^2).$$

*Proof.* Simple Taylor expansion up to linear term for  $f(\mathbf{x}) = (x_1^2 + \dots + x_d^2)^{-1/2}$  and  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + \mathbf{grad} f(\mathbf{x}) \cdot \mathbf{y} + O(\|\mathbf{y}\|_2^2)$ .  $\square$

Linearization of force: apply Lemma 9.3.4.12 to (9.3.4.11) and drop terms  $O(\|\Delta \mathbf{u}^{ji}\|_2^2)$ :

$$\begin{aligned} \blacktriangleright \quad F_{ij}(t) &\approx -\alpha_{ij} \frac{\Delta \mathbf{p}^{ji} \cdot \Delta \mathbf{u}^{ji}(t)}{l_{ij}^3} \cdot (\Delta \mathbf{p}^{ji} + \Delta \mathbf{u}^{ji}(t)) \\ &\approx -\alpha_{ij} \frac{\Delta \mathbf{p}^{ji} \cdot \Delta \mathbf{u}^{ji}(t)}{l_{ij}^3} \cdot \Delta \mathbf{p}^{ji}. \end{aligned} \quad (9.3.4.13)$$

Newton's second law of motion: ( $F_i \hat{=} \text{total force acting on } i\text{-th mass}$ )

$$m_i \frac{d^2}{dt^2} \mathbf{u}^i(t) = F_i = \sum_{\substack{j=1 \\ j \neq i}}^n -F_{ij}(t), \quad (9.3.4.14)$$

$m_i \hat{=} \text{mass of point mass } i$ .

$$\blacktriangleright m_i \frac{d^2}{dt^2} \mathbf{u}^i(t) = \sum_{\substack{j=1 \\ j \neq i}}^n \alpha_{ij} \frac{1}{l_{ij}^3} (\Delta \mathbf{p}^{ji} (\Delta \mathbf{p}^{ji})^\top) (\mathbf{u}^j(t) - \mathbf{u}^i(t)). \quad (9.3.4.15)$$

Compact notation: collect all displacements into one vector  $\mathbf{u}(t) = (\mathbf{u}^i(t))_{i=1}^n \in \mathbb{R}^{2n}$

$$(9.3.4.15) \quad \blacktriangleright \quad \mathbf{M} \frac{d\mathbf{u}}{dt^2}(t) + \mathbf{A}\mathbf{u}(t) = \mathbf{f}(t). \quad (9.3.4.16)$$

with **mass matrix**  $\mathbf{M} = \text{diag}(m_1, m_1, \dots, m_n, m_n)$

and **stiffness matrix**  $\mathbf{A} \in \mathbb{R}^{2n,2n}$  with  $2 \times 2$ -blocks

$$\begin{aligned} (\mathbf{A})_{2i-1:2i, 2i-1:2i} &= \sum_{\substack{j=1 \\ j \neq i}}^n \alpha_{ij} \frac{1}{l_{ij}^3} (\Delta \mathbf{p}^{ji} (\Delta \mathbf{p}^{ji})^\top), \quad i = 1, \dots, n, \\ (\mathbf{A})_{2i-1:2i, 2j-1:2j} &= -\alpha_{ij} \frac{1}{l_{ij}^3} (\Delta \mathbf{p}^{ji} (\Delta \mathbf{p}^{ji})^\top), \quad i \neq j. \end{aligned} \quad (9.3.4.17)$$

and external forces  $\mathbf{f}(t) = (\mathbf{f}^i(t))_{i=1}^n$ .

Note: stiffness matrix  $\mathbf{A}$  is *symmetric, positive semidefinite* ( $\rightarrow$  Def. 1.1.2.6).

Rem. 9.3.4.1: if periodic external forces  $\mathbf{f}(t) = \cos(\omega t)\mathbf{f}$ ,  $\mathbf{f} \in \mathbb{R}^{2n}$ , (wind, earthquake) act on the truss they can excite vibrations of (linearly in time) growing amplitude, if  $\omega$  coincides with  $\sqrt{\lambda_j}$  for an eigenvalue  $\lambda_j$  of  $\mathbf{A}$ .

Excited vibrations can lead to the collapse of a truss structure, cf. the notorious Tacoma-Narrows bridge disaster.

$\blacktriangleright$  It is essential to know whether eigenvalues of a truss structure fall into a range that can be excited by external forces.

These will typically<sup>(\*)</sup> be the **low modes**  $\leftrightarrow$  a few of the smallest eigenvalues.

(<sup>(\*)</sup>) Reason: fast oscillations will quickly be damped due to friction, which was neglected in our model.)

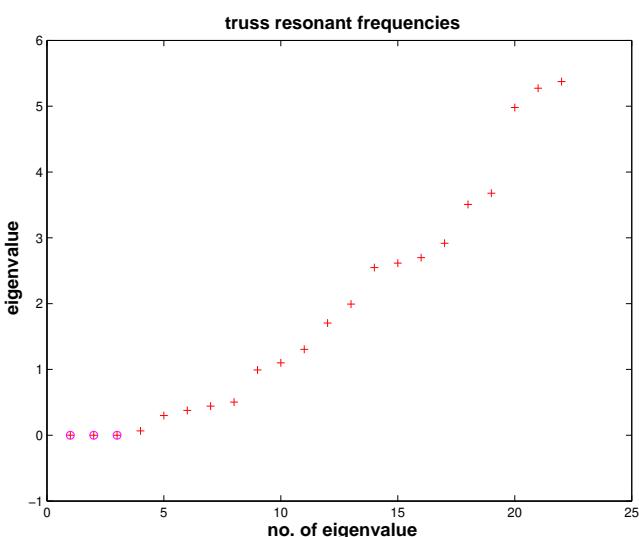
#### MATLAB-code 9.3.4.18: Computing resonant frequencies and modes of elastic truss

```
1 function [lambda,V] = trussvib(pos,top)
2 % Computes vibration modes of a truss structure, see Ex. 9.3.4.5. Mass
3 % point
4 % positions passed in the n×2-matrix pos and the connectivity encoded |
```

```

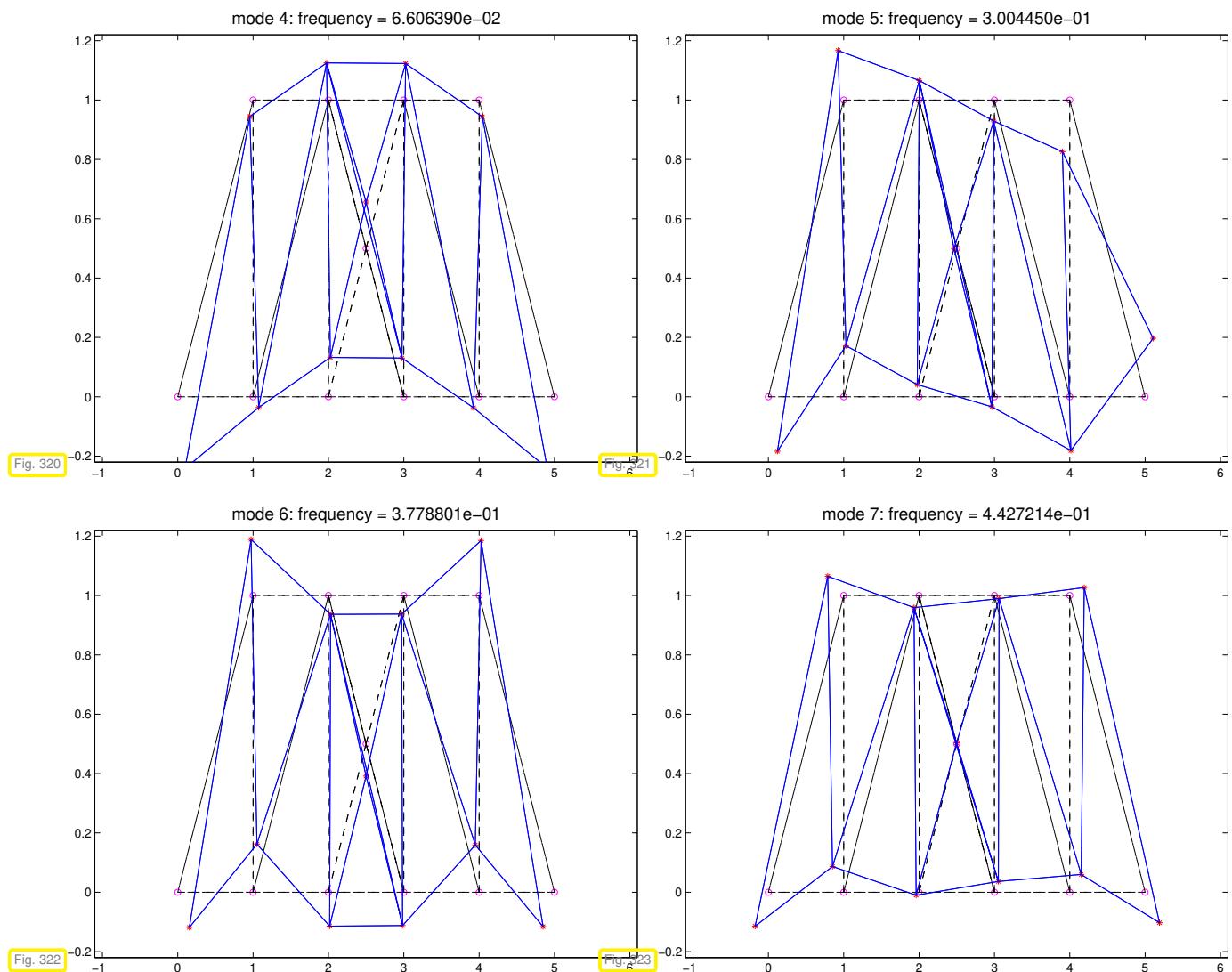
4 |   in
4 | % the sparse symmetric matrix top. In addition top(i,j) also stores
5 | % the
5 | % Young's moduli  $\alpha_{ij}$ .
6 | % The  $2n$  resonant frequencies are returned in the vector lambda, the
7 | % eigenmodes in the column of V, where entries at odd positions contain
8 | %  $x_1$ -coordinates, entries at even positions the  $x_2$ -coordinates
9 | n = size(pos,1); % no. of point masses
10| % Assembly of stiffness matrix according to (9.3.4.17)
11| A = zeros(2*n,2*n);
12| [Iidx,Jidx] = find(top); idx = [Iidx,Jidx]; % Find connected masses
13| for ij = idx'
14|   i = ij(1); j = ij(2);
15|   dp = [pos(j,1);pos(j,2)] - [pos(i,1);pos(i,2)]; %  $\Delta p^{ji}$ 
16|   lij = norm(dp); %  $l_{ij}$ 
17|   A(2*i-1:2*i,2*j-1:2*j) = -(dp*dp')/(lij^3);
18| end
19| % Set Young's moduli  $\alpha_{ij}$  (stored in top matrix)
20| A = A.*full(kron(top,[1 1;1 1]));
21| % Set  $2 \times 2$  diagonal blocks
22| for i=1:n
23|   A(2*i-1:2*i,2*i-1) = -sum(A(2*i-1:2*i,1:2:end))';
24|   A(2*i-1:2*i,2*i) = -sum(A(2*i-1:2*i,2:2:end))';
25| end
26| % Compute eigenvalues and eigenmodes
27| [V,D] = eig(A); lambda = diag(D);

```



▷ resonant frequencies of bridge truss from Fig. 317.  
 The stiffness matrix will always possess three zero eigenvalues corresponding to **rigid body modes** (= displacements without change of length of the rods)

Fig. 319



To compute *a few* of a truss's lowest resonant frequencies and excitable mode, we need efficient numerical methods for the following tasks. Obviously, Code 9.3.4.18 cannot be used for large trusses, because `eig` invariable operates on dense matrices and will be prohibitively slow and gobble up huge amounts of memory, also recall the discussion of Code 9.3.2.30.

Task:	Compute $m$ , $m \ll n$ , of the smallest/largest (in modulus) eigenvalues of $\mathbf{A} = \mathbf{A}^H \in \mathbb{C}^{n,n}$ and associated eigenvectors.
-------	---

Of course, we aim to tackle this task by iterative methods generalizing power iteration ( $\rightarrow$  Section 9.3.1) and inverse iteration ( $\rightarrow$  Section 9.3.2).

### 9.3.4.1 Orthogonalization

Preliminary considerations (in  $\mathbb{R}$ ,  $m = 2$ ):

According to Cor. 9.1.0.9: For  $\mathbf{A} = \mathbf{A}^\top \in \mathbb{R}^{n,n}$  there is a factorization  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^\top$  with  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$ ,  $\lambda_j \in \mathbb{R}$ ,  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ , and  $\mathbf{U}$  orthogonal. Thus,  $\mathbf{u}_j := (\mathbf{U})_{:,j}$ ,  $j = 1, \dots, n$ , are (mutually orthogonal) eigenvectors of  $\mathbf{A}$ .

Assume  $0 \leq \lambda_1 \leq \dots \leq \lambda_{n-2} < \lambda_{n-1} < \lambda_n$  (largest eigenvalues are simple).

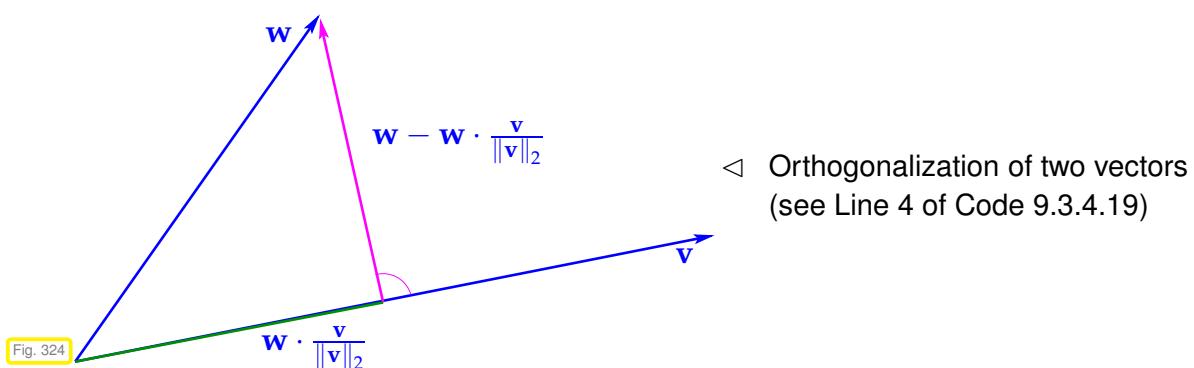
If we just carry out the direct power iteration (9.3.1.12) for two vectors both sequences will converge to the largest (in modulus) eigenvector. However, we recall that all eigenvectors are mutually orthogonal. This suggests that we orthogonalize the iterates of the second power iteration (that is to yield the eigenvector for the second largest eigenvalue) with respect to those of the first. This idea spawns the following iteration, cf. Gram-Schmidt orthogonalization in (10.2.2.4):

**MATLAB-code 9.3.4.19: one step of subspace power iteration,  $m = 2$** 

```

1 function [v,w] = sspowitstep(A,v,w)
2 v = A*v; w = A*w; % "power iteration", cf. (9.3.1.12)
3 % orthogonalization, cf. Gram-Schmidt orthogonalization (10.2.2.4)
4 v = v/norm(v); w = w - dot(v,w)*v; w = w/norm(w); % now w ⊥ v

```



Analysis through eigenvector expansions ( $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ ,  $\|\mathbf{v}\|_2 = \|\mathbf{w}\|_2 = 1$ )

$$\begin{aligned}
& \mathbf{v} = \sum_{i=1}^n \alpha_j \mathbf{u}_j , \quad \mathbf{w} = \sum_{i=1}^n \beta_j \mathbf{u}_j , \\
\Rightarrow & \mathbf{Av} = \sum_{i=1}^n \lambda_j \alpha_j \mathbf{u}_j , \quad \mathbf{Aw} = \sum_{i=1}^n \lambda_j \beta_j \mathbf{u}_j , \\
& \mathbf{v}_0 := \frac{\mathbf{v}}{\|\mathbf{v}\|_2} = \left( \sum_{i=1}^n \lambda_j^2 \alpha_j^2 \right)^{-1/2} \sum_{i=1}^n \lambda_j \alpha_j \mathbf{u}_j , \\
& \mathbf{Aw} - (\mathbf{v}_0^\top \mathbf{Aw}) \mathbf{v}_0 = \sum_{i=1}^n \left( \beta_j - \left( \sum_{i=1}^n \lambda_j^2 \alpha_j \beta_j / \sum_{i=1}^n \lambda_j^2 \alpha_j^2 \right) \alpha_j \right) \lambda_j \mathbf{u}_j .
\end{aligned}$$

We notice that  $\mathbf{v}$  is just mapped to the next iterate in the regular direct power iteration (9.3.1.12). After many steps, it will be very close to  $\mathbf{u}_n$ , and, therefore, we may now assume  $\mathbf{v} = \mathbf{u}_n \Leftrightarrow \alpha_j = \delta_{j,n}$  (Kronecker symbol).

$$\begin{aligned}
\mathbf{z} &:= \mathbf{Aw} - (\mathbf{v}_0^\top \mathbf{Aw}) \mathbf{v}_0 = 0 \cdot \mathbf{u}_n + \sum_{i=1}^{n-1} \lambda_i \beta_i \mathbf{u}_i , \\
\mathbf{w}^{(\text{new})} &:= \frac{\mathbf{z}}{\|\mathbf{z}\|_2} = \left( \sum_{i=1}^{n-1} \lambda_i^2 \beta_i^2 \right)^{-1/2} \sum_{i=1}^{n-1} \lambda_i \beta_i \mathbf{u}_i .
\end{aligned}$$

The sequence  $\mathbf{w}^{(k)}$  produced by repeated application of the mapping given by Code 9.3.4.19 asymptotically (that is, when  $\mathbf{v}^{(k)}$  has already converged to  $\mathbf{u}_n$ ) agrees with the sequence produced by the direct power method for  $\tilde{\mathbf{A}} := \mathbf{U} \text{diag}(\lambda_1, \dots, \lambda_{n-1}, 0)$ . Its convergence will be governed by the relative gap  $\lambda_{n-2}/\lambda_{n-1}$ , see Thm. 9.3.1.21.

However: if  $\mathbf{v}^{(k)}$  itself converges slowly, this reasoning does not apply.

**EXAMPLE 9.3.4.20 (Subspace power iteration with orthogonal projection)**

☞ construction of matrix  $\mathbf{A} = \mathbf{A}^T$  as in Ex. 9.3.2.37

**MATLAB-code 9.3.4.21: power iteration with orthogonal projection for two vectors**

```

1 function sppowitdriver(d,maxit)
2 % monitor power iteration with orthogonal projection for finding
3 % the two largest (in modulus) eigenvalues and associated eigenvectors
4 % of a symmetric matrix with prescribed eigenvalues passed in d
5 if (nargin < 10), maxit = 20; end
6 if (nargin < 1), d = (1:10)'; end
7 % Generate matrix
8 n = length(d);
9 Z = diag(sqrt(1:n),0) + ones(n,n);
10 [Q,R] = qr(Z); % generate orthogonal matrix
11 A = Q*diag(d,0)*Q'; % "synthetic"  $\mathbf{A} = \mathbf{A}^T$  with spectrum  $\sigma(\mathbf{A}) = \{d_1, \dots, d_n\}$ 
12 % Compute "exact" eigenvectors and eigenvalues
13 [V,D] = eig(A); [d,idx] = sort(diag(D)),
14 v_ex = V(:,idx(n)); w_ex = V(:,idx(n-1));
15 lv_ex = d(n); lw_ex = d(n-1);

16
17 v = ones(n,1); w = (-1).^v; % (Arbitrary) initial guess for
18 % eigenvectors
19 v = v/norm(v); w = w/norm(w);
20 result = [];
21 for k=1:maxit
22     v_new = A*v; w_new = A*w; % "power iteration", cf. (9.3.1.12)
23     % Rayleigh quotients provide approximate eigenvalues
24     lv = dot(v_new,v); lw = dot(w_new,w);
25     % orthogonalization, cf. Gram-Schmidt orthogonalization (10.2.2.4):
26     %  $w \perp v$ 
27     v = v_new/norm(v_new); w = w_new - dot(v,w_new)*v; w =
28         w/norm(w);
29     % Record errors in eigenvalue and eigenvector approximations. Note
30     % that the
31     % direction of the eigenvectors is not specified.
32     result = [result; k, abs(lv-lv_ex), abs(lw-lw_ex), ...
33             min(norm(v-v_ex),norm(v+v_ex)), ...
34             min(norm(w-w_ex),norm(w+w_ex))];
35
36 end

37
38 figure('name','sspowit');
39 semilogy(result(:,1),result(:,2),'m-+',...
40             result(:,1),result(:,3),'r-*',...
41             result(:,1),result(:,4),'k^-',...
42             result(:,1),result(:,5),'b-p');
43 title('d = [0.5*(1:8),9.5,10]');
44 xlabel('{\bf power iteration step}', 'fontsize',14);
45 ylabel('{\bf error}', 'fontsize',14);
46 legend('error in \lambda_n', 'error in \lambda_{n-1}', 'error in
47           v', 'error in w', 'location', 'northeast');
48 print -depsc2 '../PICTURES/sspowitcvgl.eps';

```

```

42
43 rates = result(2:end,2:end)./result(1:end-1,2:end);
44 figure('name','rates');
45 plot(result(2:end,1),rates(:,1),'m-+',...
46      result(2:end,1),rates(:,2),'r-*',...
47      result(2:end,1),rates(:,3),'k^-',...
48      result(2:end,1),rates(:,4),'b-p');
49 axis([0 maxit 0.5 1]);
50 title('d = [0.5*(1:8),9.5,10]');
51 xlabel('{\bf power iteration step}', 'fontsize', 14);
52 ylabel('{\bf error quotient}', 'fontsize', 14);
53 legend('error in \lambda_n', 'error in \lambda_{n-1}', 'error in v',...
54         'error in w', 'location', 'southeast');
55 print -depsc2 '../PICTURES/sspowitcvgrates1.eps';

```

$\sigma(\mathbf{A}) = \{1, 2, \dots, 10\}$ :

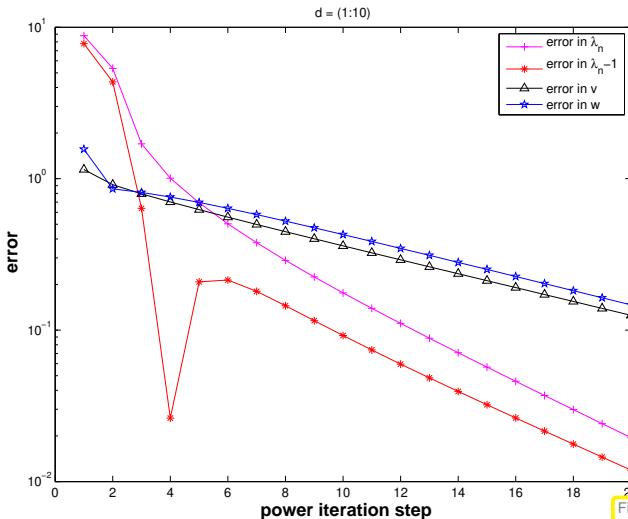


Fig. 325

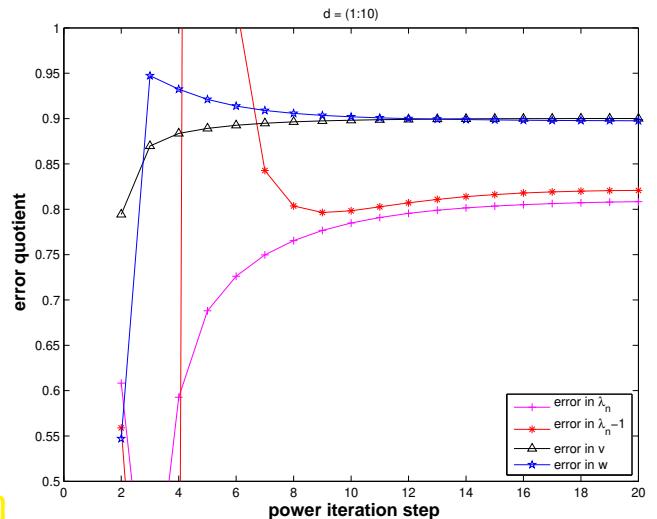


Fig. 326

$\sigma(\mathbf{A}) = \{0.5, 1, \dots, 4, 9.5, 10\}$ :

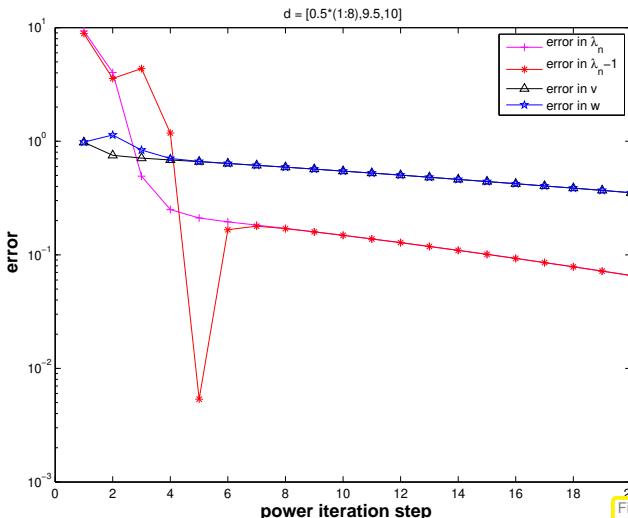


Fig. 327

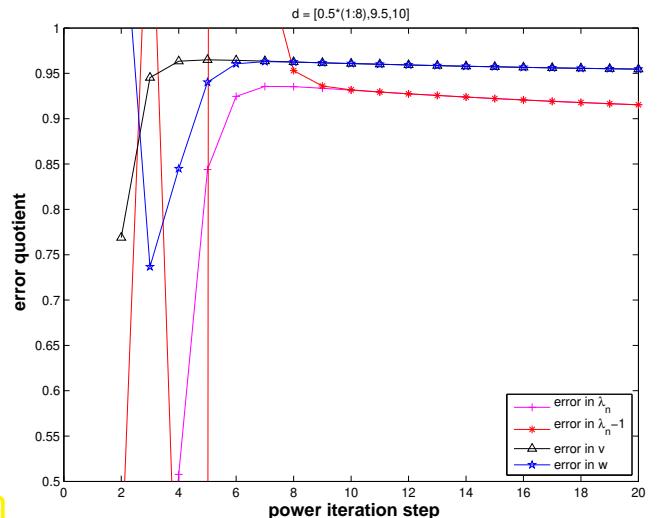


Fig. 328

Issue: generalization of orthogonalization idea to subspaces of dimension  $> 2$

Nothing new:

### Gram-Schmidt orthonormalization

(→ [NS02, Thm. 4.8], [Gut09, Alg. 6.1], [QSS00, Sect. 3.4.3])

Given: linearly independent vectors  $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^n$ ,  $m \in \mathbb{N}$

Sought: vectors  $\mathbf{q}_1, \dots, \mathbf{q}_m \in \mathbb{R}^n$  such that

$$\textcircled{2} \quad \mathbf{q}_l^\top \mathbf{q}_k = \delta_{lk} \quad (\text{orthonormality}), \quad (9.3.4.22)$$

$$\textcircled{2} \quad \text{Span}\{\mathbf{q}_1, \dots, \mathbf{q}_m\} = \text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_k\} \quad \text{for all } k = 1, \dots, m. \quad (9.3.4.23)$$

Constructive proof & algorithm for Gram-Schmidt orthonormalization:

$$\begin{aligned} \mathbf{z}_1 &= \mathbf{v}_1, \\ \mathbf{z}_2 &= \mathbf{v}_2 - \frac{\mathbf{v}_2^\top \mathbf{z}_1}{\mathbf{z}_1^\top \mathbf{z}_1} \mathbf{z}_1, \\ \mathbf{z}_3 &= \mathbf{v}_3 - \frac{\mathbf{v}_3^\top \mathbf{z}_1}{\mathbf{z}_1^\top \mathbf{z}_1} \mathbf{z}_1 - \frac{\mathbf{v}_3^\top \mathbf{z}_2}{\mathbf{z}_2^\top \mathbf{z}_2} \mathbf{z}_2, \\ &\vdots \end{aligned} \quad (9.3.4.24)$$

$$+ \text{normalization} \quad \mathbf{q}_k = \frac{\mathbf{z}_k}{\|\mathbf{z}_k\|_2}, \quad k = 1, \dots, m. \quad (9.3.4.25)$$

Easy computation: the vectors  $\mathbf{q}_1, \dots, \mathbf{q}_m$  produced by (9.3.4.24) satisfy (9.3.4.22) and (9.3.4.23).

#### MATLAB-code 9.3.4.26: Gram-Schmidt orthonormalization (do not use, unstable algorithm!)

```

1 function Q = gso(V)
2 % Gram-Schmidt orthonormalization of the columns of V ∈ ℝn,m, see
3 % (9.3.4.24). The vectors q1, ..., qm are returned as the columns of
4 % the orthogonal matrix Q.
5 m = size(V, 2);
6 Q = V(:, 1) / norm(V(:, 1)); % normalization
7 for l=2:m
8   q = V(:, l);
9   % orthogonalization
10  for k=1:l-1
11    q = q - dot(Q(:, k), V(:, l)) * Q(:, k);
12  end
13  Q = [Q, q / norm(q)]; % normalization
14 end

```



Warning! Code 9.3.4.26 provides an unstable implementation of Gram-Schmidt orthonormalization: for large  $n, m$  impact of round-off will destroy the orthogonality of the columns of  $Q$ .

A stable implementation of Gram-Schmidt orthogonalization of the columns of a matrix  $V \in \mathbb{K}^{n,m}$ ,  $m \leq n$ , is provided by the following MATLAB command:

$[Q, R] = qr(V, 0)$  (Asymptotic computational cost:  $O(m^2n)$ )  
dummy return value (for our purposes) dummy argument

Detailed description of the algorithm behind  $qr$  and meaning of the return value  $R$  → Section 3.3.3.

#### EXAMPLE 9.3.4.27 (qr based orthogonalization, $m = 2$ )

The following two MATLAB code snippets perform the same function, cf. Code 9.3.4.19:

**MATLAB-code 9.3.4.28:**

```

1 v = v/norm(v);
2 w = w - dot(v,w)*v;
3 w = w/norm(w);

```

**MATLAB-code 9.3.4.29:**

```

1 [Q,R] = qr([v,w],0);
2 v = Q(:,1);
3 w = Q(:,2);

```

Explanation ➤ discussion of Gram-Schmidt orthonormalization.


**MATLAB-code 9.3.4.30: General subspace power iteration step with qr based orthonormalization**

```

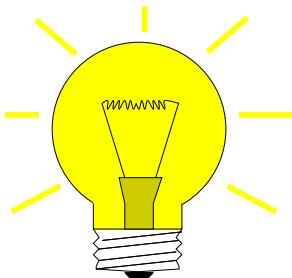
1 function V = sspowitstep(A,V)
2 % power iteration with orthonormalization for A = A^T.
3 % columns of matrix V span subspace for power iteration.
4 V = A*V; % actual power iteration on individual columns
5 [V,R] = qr(V,0); % Gram-Schmidt orthonormalization (9.3.4.24)

```

### 9.3.4.2 Ritz projection

Observations on Code 9.3.4.19:

- ◆ the first column of  $\mathbf{V}$ ,  $(\mathbf{V})_{:,1}$ , is a sequence of vectors created by the standard direct power method (9.3.1.12).
- ◆ reasoning: the other columns of  $\mathbf{V}$ , after each multiplication with  $\mathbf{A}$  can be expected to contain a significant component in the direction of the eigenvector associated with the eigenvalue of largest modulus.



Idea: use information in  $(\mathbf{V})_{:,2}, \dots, (\mathbf{V})_{:,end}$  to accelerate convergence of  $(\mathbf{V})_{:,1}$ .

Since the columns of  $\mathbf{V}$  span a *subspace* of  $\mathbb{R}^n$ , this idea can be recast as the following task:

Task: given  $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{K}^n$ ,  $k \ll n$ , extract (good approximations of) eigenvectors of  $\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n}$  contained in  $\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ .

We take for granted that  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  is linearly independent.

Assumption:

$$\text{Eig}\mathbf{A}\lambda \cap V := \text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_m\} \neq \{0\}$$

$\Leftrightarrow$

$V$  contains an eigenvector of  $\mathbf{A}$

$$\begin{aligned} &\Leftrightarrow \exists \mathbf{w} \in V \setminus \{0\}: \mathbf{Aw} = \lambda \mathbf{w} \\ &\Leftrightarrow \exists \mathbf{u} \in \mathbb{K}^m \setminus \{0\}: \mathbf{Av}_u = \lambda \mathbf{v}_u \end{aligned}$$

$$\Rightarrow \exists \mathbf{u} \in \mathbb{K}^m \setminus \{0\}: \mathbf{V}^H \mathbf{A} \mathbf{V} \mathbf{u} = \lambda \mathbf{V}^H \mathbf{V} \mathbf{u}, \quad (9.3.4.31)$$

where  $\mathbf{V} := (\mathbf{v}_1, \dots, \mathbf{v}_m) \in \mathbb{K}^{n,m}$  and we used

$$\mathbf{V} = \{\mathbf{V}\mathbf{u} : \mathbf{u} \in \mathbb{K}^m\} \quad (\text{linear combinations of the } \mathbf{v}_i).$$

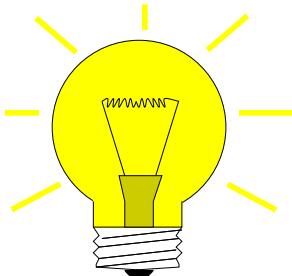
(9.3.4.31)  $\Rightarrow \mathbf{u} \in \mathbb{K}^k \setminus \{0\}$  solves  $m \times m$  generalized eigenvalue problem

$$(\mathbf{V}^H \mathbf{A} \mathbf{V}) \mathbf{u} = \lambda (\mathbf{V}^H \mathbf{V}) \mathbf{u}. \quad (9.3.4.32)$$

Note:

$$\begin{aligned} \{\mathbf{v}_1, \dots, \mathbf{v}_m\} \text{ linearly independent} \\ \Updownarrow \\ \mathbf{V} \text{ has full rank } m \quad (\rightarrow \text{Def. 2.2.1.3}) \\ \Updownarrow \\ \mathbf{V}^H \mathbf{V} \text{ is symmetric positive definite} \quad (\rightarrow \text{Def. 1.1.2.6}) \end{aligned}$$

If our initial assumption holds true and  $\mathbf{u}$  solves (9.3.4.32) and is a simple eigenvalue, a corresponding  $\mathbf{x} \in \text{Eig} \mathbf{A} \lambda$  can be recovered as  $\mathbf{x} = \mathbf{V}\mathbf{u}$ .



Idea: Given a subspace  $V = \text{Im}(\mathbf{V}) \subset \mathbb{K}^n$ ,  $\mathbf{V} \in \mathbb{K}^{n,m}$ ,  $\dim(V) = m$ , obtain *approximations* of (a few) eigenvalues and eigenvectors  $\mathbf{x}_1, \dots, \mathbf{x}_m$  of  $\mathbf{A}$  by

① solving the generalized eigenvalue problem (9.3.4.32)

$\rightarrow$  eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{K}^m$  (linearly independent),

② and transforming them back according to  $\mathbf{x}_k = \mathbf{V}\mathbf{u}_k, k = 1, \dots, m$ .

Terminology: (9.3.4.32) is called the **Ritz projection** of EVP  $\mathbf{Ax} = \lambda \mathbf{x}$  onto  $V$

Terminology:  $\sigma(\mathbf{V}^H \mathbf{A} \mathbf{V}) \hat{=} \text{Ritz values}$ ,  
eigenvectors of  $\mathbf{V}^H \mathbf{A} \mathbf{V} \hat{=} \text{Ritz vectors}$

Example: Ritz projection of  $\mathbf{Ax} = \lambda \mathbf{x}$  onto  $\text{Span}\{\mathbf{v}, \mathbf{w}\}$ :

$$(\mathbf{v}, \mathbf{w})^H \mathbf{A} (\mathbf{v}, \mathbf{w}) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \lambda (\mathbf{v}, \mathbf{w})^H (\mathbf{v}, \mathbf{w}) \begin{pmatrix} \alpha \\ \beta \end{pmatrix}.$$

Note: If  $\mathbf{V}$  is *unitary* ( $\rightarrow$  Def. 6.2.1.2), then the generalized eigenvalue problem (9.3.4.32) will become a standard linear eigenvalue problem.

### Remark 9.3.4.33 (Justification of Ritz projection by min-max theorem)

We revisit  $m = 2$ , see Code 9.3.4.19. Recall that by the min-max theorem Thm. 9.3.2.18

$$\mathbf{u}_n = \operatorname{argmax}_{\mathbf{x} \in \mathbb{R}^n} \rho_{\mathbf{A}}(\mathbf{x}), \quad \mathbf{u}_{n-1} = \operatorname{argmax}_{\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \perp \mathbf{u}_n} \rho_{\mathbf{A}}(\mathbf{x}). \quad (9.3.4.34)$$

Idea: maximize Rayleigh quotient over  $\text{Span}\{\mathbf{v}, \mathbf{w}\}$ , where  $\mathbf{v}, \mathbf{w}$  are output by Code 9.3.4.19. This leads to the optimization problem

$$(\alpha^*, \beta^*) := \operatorname{argmax}_{\alpha, \beta \in \mathbb{R}, \alpha^2 + \beta^2 = 1} \rho_{\mathbf{A}}(\alpha \mathbf{v} + \beta \mathbf{w}) = \operatorname{argmax}_{\alpha, \beta \in \mathbb{R}, \alpha^2 + \beta^2 = 1} \rho_{(\mathbf{v}, \mathbf{w})^T \mathbf{A} (\mathbf{v}, \mathbf{w})} \left( \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \right). \quad (9.3.4.35)$$

Then a better approximation for the eigenvector to the largest eigenvalue is

$$\mathbf{v}^* := \alpha^* \mathbf{v} + \beta^* \mathbf{w}.$$

Note that  $\|\mathbf{v}^*\|_2 = 1$ , if both  $\mathbf{v}$  and  $\mathbf{w}$  are normalized, which is guaranteed in Code 9.3.4.19.

Then, orthogonalizing  $\mathbf{w}$  w.r.t  $\mathbf{v}^*$  will produce a new iterate  $\mathbf{w}^*$ .

Again the min-max theorem Thm. 9.3.2.18 tells us that we can find  $(\alpha^*, \beta^*)^\top$  as eigenvector to the largest eigenvalue of

$$(\mathbf{v}, \mathbf{w})^\top \mathbf{A}(\mathbf{v}, \mathbf{w}) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \lambda \begin{pmatrix} \alpha \\ \beta \end{pmatrix}. \quad (9.3.4.36)$$

Since eigenvectors of symmetric matrices are mutually orthogonal, we find  $\mathbf{w}^* = \alpha_2 \mathbf{v} + \beta_2 \mathbf{w}$ , where  $(\alpha_2, \beta_2)^\top$  is the eigenvector of (9.3.4.36) belonging to the smallest eigenvalue. This assumes orthonormal vectors  $\mathbf{v}, \mathbf{w}$ .  $\square$

#### MATLAB-code 9.3.4.37: one step of subspace power iteration with Ritz projection, matrix version

```

1 function V = sspowitsteprp(A,V)
2 V = A*V; % power iteration applied to columns of V
3 [Q,R] = qr(V,0); % orthonormalization, see Section 9.3.4.1
4 [U,D] = eig(Q'*A*Q); % Solve Ritz projected m x m eigenvalue problem
5 V = Q*U; % recover approximate eigenvectors
6 ev = diag(D); % approximate eigenvalues

```

Note that the orthogonalization step in Code 9.3.4.37 is actually redundant, if exact arithmetic could be employed, because the Ritz projection could also be realized by solving the generalized eigenvalue problem.

However, prior orthogonalization is essential for numerical stability ( $\rightarrow$  Def. 1.5.5.19), cf. the discussion in Section 3.3.3.

#### EXAMPLE 9.3.4.38 (Power iteration with Ritz projection)

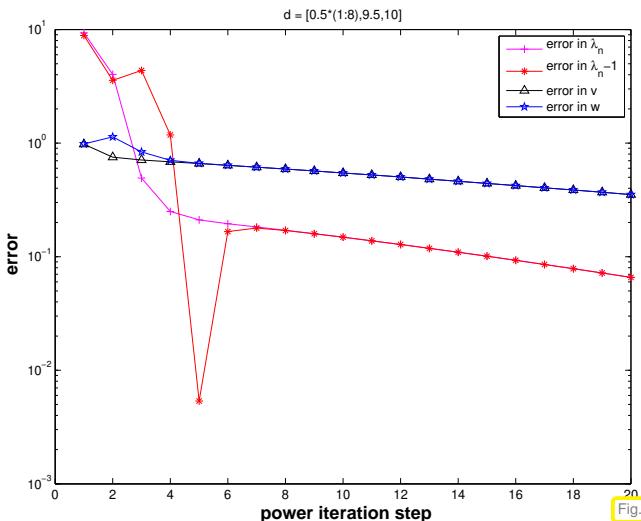
Listing 9.1: Main loop: power iteration with Ritz projection for two eigenvectors

```

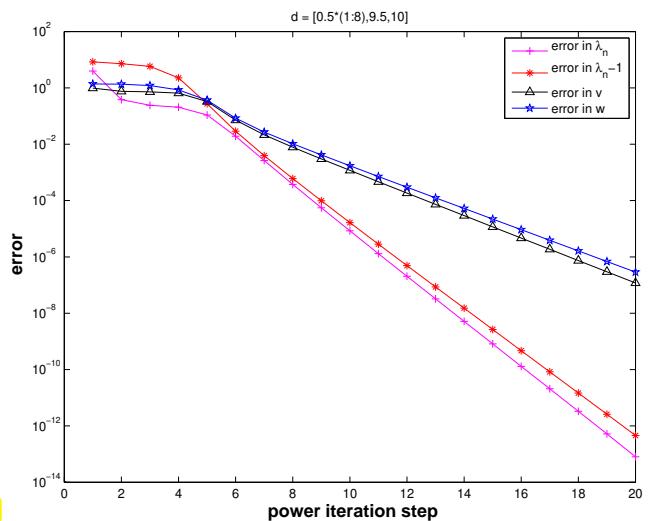
1 % See Code 9.3.4.21 for generation of matrix A and output
2 for k=1:maxit
3 v_new = A*v; w_new = A*w; % "power iteration", cf. (9.3.1.12)
4 [Q,R] = qr([v_new,w_new],0); % orthogonalization, see Sect. 9.3.4.1
5 [U,D] = eig(Q'*A*Q); % Solve Ritz projected eigenvalue problem
6 [ev,idx] = sort(abs(diag(D))), % Sort eigenvalues
7 w = Q*U(:,idx(1)); v = Q*U(:,idx(2)); % Recover approximate
     eigenvectors
8
9 % Record errors in eigenvalue and eigenvector approximations. Note that
10 % the
11 % direction of the eigenvectors is not specified.
12 result = [result; k, abs(ev(2)-lv_ex), abs(ev(1)-lw_ex), ...
13           min(norm(v-v_ex),norm(v+v_ex)),
14           min(norm(w-w_ex),norm(w+w_ex))];
end

```

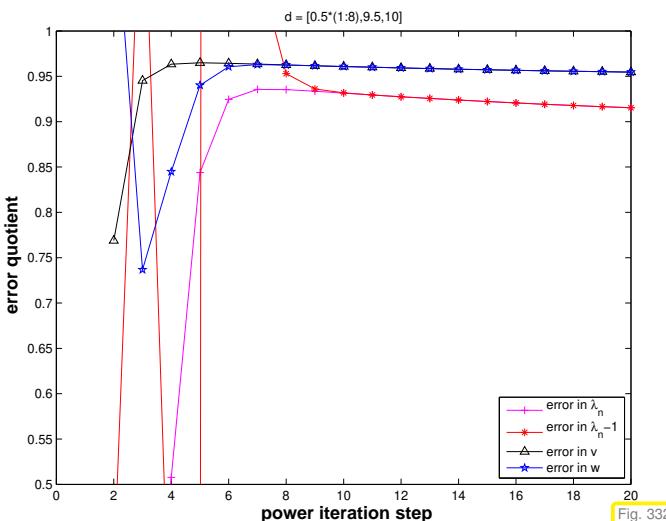
Matrix as in Ex. 9.3.4.20,  $\sigma(\mathbf{A}) = \{0.5, 1, \dots, 4, 9.5, 10\}$ :



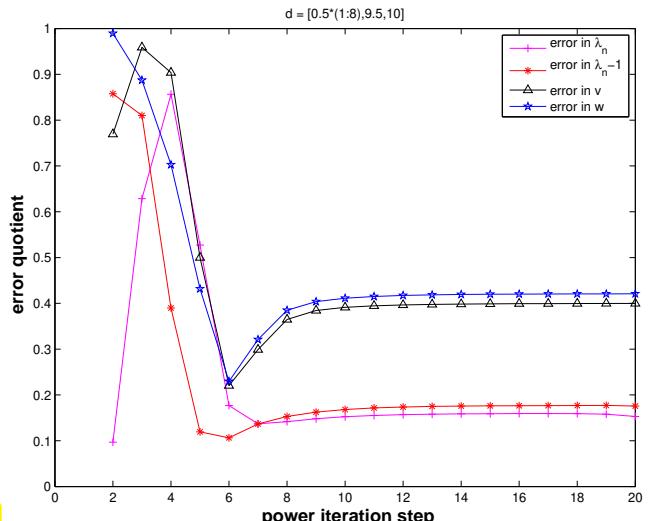
simple orthonormalization, Ex. 9.3.4.20



with Ritz projection, Code 9.3.4.37



simple orthonormalization, Ex. 9.3.4.20



with Ritz projection, Code 9.3.4.37

Observation: tremendous acceleration of power iteration through Ritz projection, convergence still linear but with much better rates.

In Code 9.3.4.37: diagonal entries of  $\mathbf{D}$  provide approximations of eigenvalues. Their (relative) changes can be used as a termination criterion.

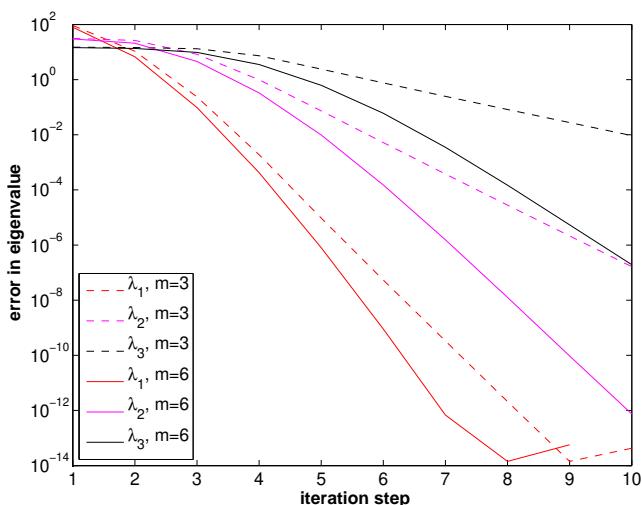
### §9.3.4.39 (Subspace variant of direct power method with Ritz projection)

**MATLAB-code 9.3.4.40: Subspace power iteration with Ritz projection**

```

1 function [ev,V] = sspowitrp(A,k,m,tol,maxit)
2 % Power iteration with Ritz projection for matrix  $\mathbf{A} = \mathbf{A}^T \in \mathbb{R}^{n,n}$ :
3 % Subspace of dimension  $m \leq n$  is used to compute the  $k \leq m$  largest
4 % eigenvalues of  $\mathbf{A}$  and associated eigenvectors.
5 n = size(A,1); V = eye(n,m); d = zeros(m,1); % (Arbitrary) initial
6 % eigenvectors
7 for i=1:maxit
8     [Q,R] = qr(A*V,0); % Power iteration and orthonormalization
9     [U,D] = eig(Q'*A*Q); % Small  $m \times m$  eigenvalue problem for Ritz
10    [ev,idx] = sort(diag(D)); % eigenvalue approximations in diagonal
11    V = Q*U; % 2nd part of Ritz projection
12    if (abs(ev-d) < tol*max(abs(ev))), break; end
13    d = ev;
14 end
15 ev = ev(m-k+1:end);
16 V = V(:,idx(m-k+1:end));

```

**EXAMPLE 9.3.4.41 (Convergence of subspace variant of direct power method)**

S.p.d. test matrix:  $a_{ij} := \min\left\{\frac{i}{j}, \frac{j}{i}\right\}$   
 $n=200$ ;  $A = \text{gallery('lehmer', } n)$ ;  
“Initial eigenvector guesses”:  
 $V = \text{eye}(n, m)$ ;

- Observation:  
linear convergence of eigenvalues
- choice  $m > k$  boosts convergence of eigenvalues

**Remark 9.3.4.42 (Subspace power methods)**

Analogous to § 9.3.4.39: construction of subspace variants of inverse iteration ( $\rightarrow$  Code 9.3.2.31), PIN-VIT (9.3.3.1), and Rayleigh quotient iteration (9.3.2.36).

**9.4 Krylov Subspace Methods**

*Supplementary literature.* [Han02, Sect. 30]

All power methods ( $\rightarrow$  Section 9.3) for the eigenvalue problem (EVP)  $\mathbf{Ax} = \lambda \mathbf{x}$  only rely on the last iterate to determine the next one (1-point methods, cf. (8.1.0.4))

$\Rightarrow$  NO MEMORY, cf. discussion in the beginning of Section 10.2.

“Memory for power iterations”: pursue same idea that led from the gradient method, § 10.1.3.3, to the conjugate gradient method, § 10.2.2.10: use information from previous iterates to achieve efficient minimization over larger and larger subspaces.

Min-max theorem, Thm. 9.3.2.18	:	$\mathbf{A} = \mathbf{A}^H \Rightarrow$	EVPs $\Leftrightarrow$	Finding extrema/stationary points of Rayleigh quotient ( $\rightarrow$ Def. 9.3.1.16)
-----------------------------------	---	---	------------------------	--

Setting: EVP  $\mathbf{Ax} = \lambda \mathbf{x}$  for real s.p.d. ( $\rightarrow$  Def. 1.1.2.6) matrix  $\mathbf{A} = \mathbf{A}^T \in \mathbb{R}^{n,n}$

notations used below:  $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ : eigenvalues of  $\mathbf{A}$ , counted with multiplicity, see Def. 9.1.0.1,

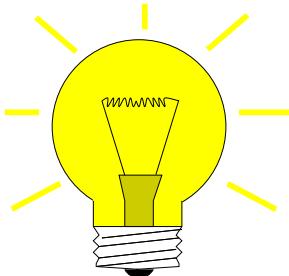
$\mathbf{u}_1, \dots, \mathbf{u}_n \hat{=}$  corresponding orthonormal eigenvectors, cf. Cor. 9.1.0.9.

$$\Rightarrow \mathbf{AU} = \mathbf{DU} , \quad \mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n) \in \mathbb{R}^{n,n} , \quad \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n) .$$

We recall

- ◆ the direct power method (9.3.1.12) from Section 9.3.1
- ◆ and the inverse iteration from Section 9.3.2

and how they produce sequences  $(\mathbf{z}^{(k)})_{k \in \mathbb{N}_0}$  of vectors that are supposed to converge to a vector  $\in \text{Eig}\mathbf{A}\lambda_1$  or  $\in \text{Eig}\mathbf{A}\lambda_n$ , respectively.



Idea: Better  $\mathbf{z}^{(k)}$  from Ritz projection onto  $\mathcal{V} := \text{Span}\{\mathbf{z}^{(0)}, \dots, \mathbf{z}^{(k)}\}$   
(= space spanned by previous iterates)

Recall ( $\rightarrow$  Code 9.3.4.37) Ritz projection of an EVP  $\mathbf{Ax} = \lambda \mathbf{x}$  onto a subspace  $\mathcal{V} := \text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ ,  $m < n \Rightarrow$  smaller  $m \times m$  generalized EVP

$$\underbrace{\mathbf{V}^T \mathbf{AV}}_{:= \mathbf{H}} \mathbf{x} = \lambda \mathbf{V}^T \mathbf{V} \mathbf{x} , \quad \mathbf{V} := (\mathbf{v}_1, \dots, \mathbf{v}_m) \in \mathbb{R}^{n,m} . \quad (9.4.0.1)$$

From Rayleigh quotient Thm. 9.3.2.16 and considerations in Section 9.3.4.2:

$$\begin{aligned} \mathbf{u}_n \in \mathcal{V} &\Rightarrow \text{largest eigenvalue of (9.4.0.1)} = \lambda_{\max}(\mathbf{A}) , \\ \mathbf{u}_1 \in \mathcal{V} &\Rightarrow \text{smallest eigenvalue of (9.4.0.1)} = \lambda_{\min}(\mathbf{A}) . \end{aligned}$$

Intuition: If  $\mathbf{u}_n$  ( $\mathbf{u}_1$ ) “well captured” by  $\mathcal{V}$  (that is, the angle between the vector and the space  $\mathcal{V}$  is small), then we can expect that the largest (smallest) eigenvalue of (9.4.0.1) is a good approximation for  $\lambda_{\max}(\mathbf{A})$  ( $\lambda_{\min}(\mathbf{A})$ ), and that, assuming normalization

$$\mathbf{V}\mathbf{w} \approx \mathbf{u}_1 \quad (\text{or } \mathbf{V}\mathbf{w} \approx \mathbf{u}_n) ,$$

where  $\mathbf{w}$  is the corresponding eigenvector of (9.4.0.1).

► For direct power method (9.3.1.12):

$$\mathbf{z}^{(k)} || \mathbf{A}^k \mathbf{z}^{(0)}$$

$$V = \text{Span}\{\mathbf{z}^{(0)}, \mathbf{A}\mathbf{z}^{(0)}, \dots, \mathbf{A}^{(k)}\mathbf{z}^{(0)}\} = \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{z}^{(0)}) \quad \text{a Krylov space, } \rightarrow \text{Def. 10.2.1.1 .} \quad (9.4.0.2)$$

### MATLAB-code 9.4.0.3: Ritz projections onto Krylov space (9.4.0.2)

```

1 function [V,D] = kryleig(A,m)
2 % Ritz projection onto Krylov subspace. An
3 % orthonormal basis of  $\mathcal{K}_m(\mathbf{A}, \mathbf{1})$  is assembled
4 % into the columns of V.
5 n = size(A,1); V = (1:n)'; V =
6 V/norm(V);
7 for l=1:m-1
8     V = [V, A*V(:,end)]; [Q,R] = qr(V,0);
9     [W,D] = eig(Q'*A*Q); V = Q*W;
10 end

```

► direct power method with Ritz projection onto Krylov space from (9.4.0.2), cf. § 9.3.4.39.

Note: implementation for demonstration purposes only (inefficient for sparse matrix  $\mathbf{A}$ !)

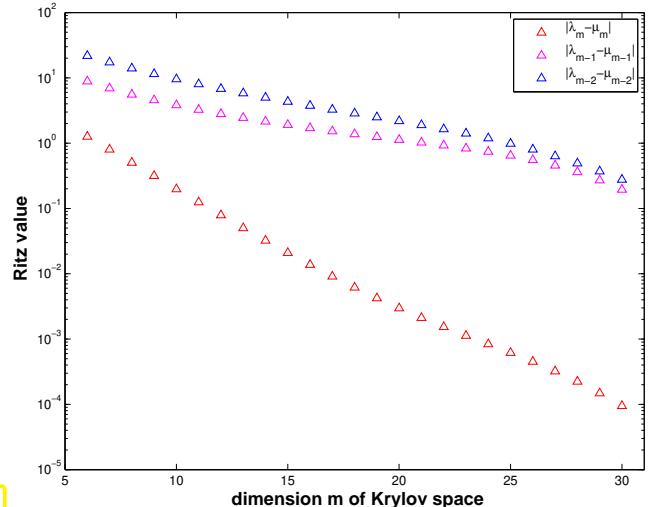
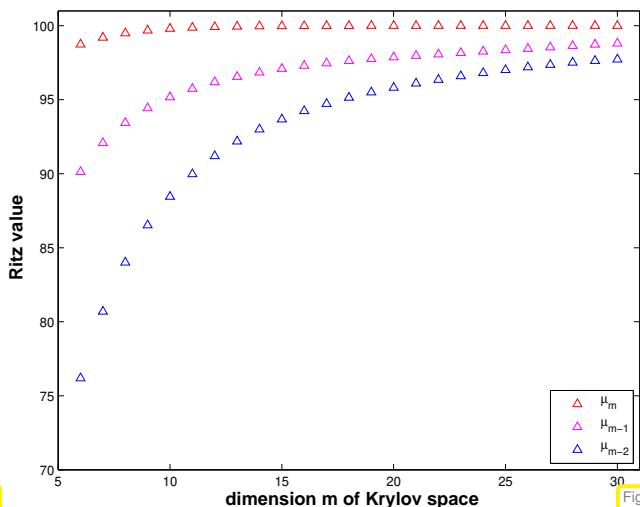
### EXAMPLE 9.4.0.4 (Ritz projections onto Krylov space)

#### MATLAB-code 9.4.0.5:

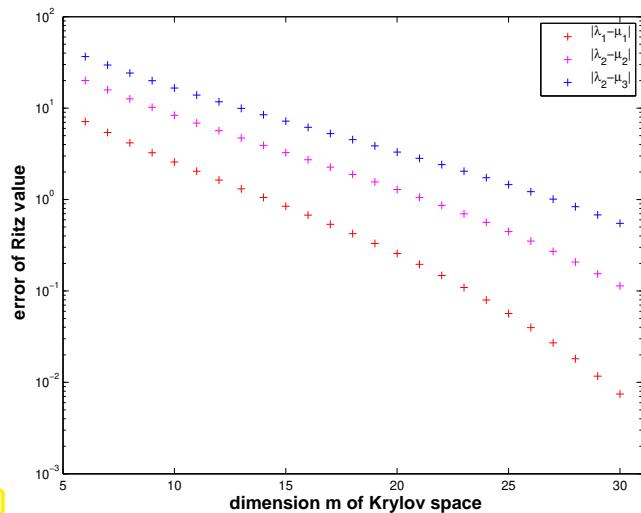
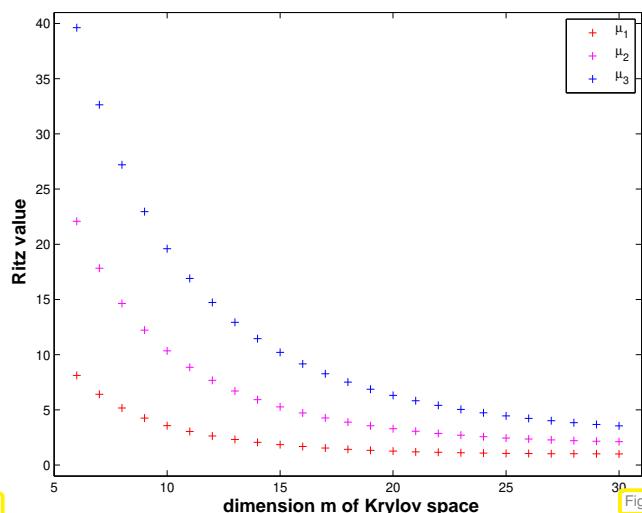
```

1 n=100;
2 M=gallery ('tridiag', -0.5*ones(n-1,1), 2*ones(n,1), -1.5*ones(n-1,1));
3 [Q,R]=qr(M); A=Q'*diag(1:n)*Q; % synthetic matrix,
4 sigma(A) = {1,2,3,...,100}

```



Observation: “vaguely linear” convergence of largest Ritz values (notation  $\mu_i$ ) to largest eigenvalues. Fastest convergence of largest Ritz value → largest eigenvalue of  $\mathbf{A}$



Observation: *Also the smallest Ritz values converge “vaguely linearly” to the smallest eigenvalues of  $\mathbf{A}$ .*  
Fastest convergence of smallest Ritz value  $\rightarrow$  smallest eigenvalue of  $\mathbf{A}$ . □



Why do smallest Ritz values converge to smallest eigenvalues of  $\mathbf{A}$ ?

Consider direct power method (9.3.1.12) for  $\tilde{\mathbf{A}} := \nu \mathbf{I} - \mathbf{A}$ ,  $\nu > \lambda_{\max}(\mathbf{A})$ :

$$\mathbf{z}^{(0)} \text{ arbitrary}, \quad \tilde{\mathbf{z}}^{(k+1)} = \frac{(\nu \mathbf{I} - \mathbf{A}) \tilde{\mathbf{z}}^{(k)}}{\|(\nu \mathbf{I} - \mathbf{A}) \tilde{\mathbf{z}}^{(k)}\|_2} \quad (9.4.0.6)$$

As  $\sigma(\nu \mathbf{I} - \mathbf{A}) = \nu - \sigma(\mathbf{A})$  and eigenspaces agree, we infer from Thm. 9.3.1.21

$$\lambda_1 < \lambda_2 \Rightarrow \mathbf{z}^{(k)} \xrightarrow{k \rightarrow \infty} \mathbf{u}_1 \quad \& \quad \rho_{\mathbf{A}}(\mathbf{z}^{(k)}) \xrightarrow{k \rightarrow \infty} \lambda_1 \quad \text{linearly}. \quad (9.4.0.7)$$

By the binomial theorem (also applies to matrices, if they commute)

$$(\nu \mathbf{I} - \mathbf{A})^k = \sum_{j=0}^k \binom{k}{j} \nu^{k-j} \mathbf{A}^j \Rightarrow (\nu \mathbf{I} - \mathbf{A})^k \tilde{\mathbf{z}}^{(0)} \in \mathcal{K}_k(\mathbf{A}, \mathbf{z}^{(0)}),$$

$\mathcal{K}_k(\nu \mathbf{I} - \mathbf{A}, \mathbf{x}) = \mathcal{K}_k(\mathbf{A}, \mathbf{x})$

(9.4.0.8)

➤  $\mathbf{u}_1$  can also be expected to be “well captured” by  $\mathcal{K}_k(\mathbf{A}, \mathbf{x})$  and the smallest Ritz value should provide a good approximation for  $\lambda_{\min}(\mathbf{A})$ .

Recall from Section 10.2.2 Lemma 10.2.2.5 :

Residuals  $\mathbf{r}_0, \dots, \mathbf{r}_{m-1}$  generated in CG iteration, § 10.2.2.10 applied to  $\mathbf{Ax} = \mathbf{z}$  with  $\mathbf{x}^{(0)} = 0$ , provide orthogonal basis for  $\mathcal{K}_m(\mathbf{A}, \mathbf{z})$  (, if  $\mathbf{r}_k \neq 0$ ).

► Inexpensive Ritz projection of  $\mathbf{Ax} = \lambda \mathbf{x}$  onto  $\mathcal{K}_m(\mathbf{A}, \mathbf{z})$ : orthogonal matrix

$$\mathbf{V}_m^T \mathbf{A} \mathbf{V}_m \mathbf{x} = \lambda \mathbf{x}, \quad \mathbf{V}_m := \left( \frac{\mathbf{r}_0}{\|\mathbf{r}_0\|}, \dots, \frac{\mathbf{r}_{m-1}}{\|\mathbf{r}_{m-1}\|} \right) \in \mathbb{R}^{n,m}. \quad (9.4.0.9)$$

recall: residuals generated by *short recursions*, see § 10.2.2.10

**Lemma 9.4.0.10. Tridiagonal Ritz projection from CG residuals**

$\mathbf{V}_m^T \mathbf{A} \mathbf{V}_m$  is a tridiagonal matrix.

*Proof.* Lemma 10.2.2.5:  $\{\mathbf{r}_0, \dots, \mathbf{r}_{\ell-1}\}$  is an orthogonal basis of  $\mathcal{K}_\ell(\mathbf{A}, \mathbf{r}_0)$ , if all the residuals are non-zero. As  $\mathbf{A}\mathcal{K}_{\ell-1}(\mathbf{A}, \mathbf{r}_0) \subset \mathcal{K}_\ell(\mathbf{A}, \mathbf{r}_0)$ , we conclude the orthogonality  $\mathbf{r}_m^T \mathbf{A} \mathbf{r}_j$  for all  $j = 0, \dots, m-2$ . Since

$$(\mathbf{V}_m^T \mathbf{A} \mathbf{V}_m)_{ij} = \mathbf{r}_{i-1}^T \mathbf{A} \mathbf{r}_{j-1}, \quad 1 \leq i, j \leq m,$$

the assertion of the theorem follows.  $\square$

$$\mathbf{V}_l^H \mathbf{A} \mathbf{V}_l = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \beta_2 & \alpha_3 & \ddots & \\ & & \ddots & \ddots & \\ & & & \ddots & \beta_{k-1} \\ & & & & \alpha_k \end{bmatrix} =: \mathbf{T}_l \in \mathbb{K}^{k,k} \quad [\text{tridiagonal matrix}] \quad (9.4.0.11)$$

Algorithm for computing  $\mathbf{V}_l$  and  $\mathbf{T}_l$ :  
**Lanczos process**

Computational effort/step:

- 1  $\times$   $\mathbf{A} \times$  vector
- 2 dot products
- 2 AXPY-operations
- 1 division

Closely related to CG iteration,  
§ 10.2.2.10, Code 10.2.2.11.

**MATLAB-code 9.4.0.12: Lanczos process, cf. Code 10.2.2.11**

```

1 function [V, alph, bet] = lanczos(A, k, z0)
2 % Note: this implementation of the Lanczos
3 % process also records the orthonormal CG
4 % residuals in the columns of the matrix V,
5 % which is not needed when only eigenvalue
6 % approximations are desired.
7 V = z0/norm(z0);
8 % Vectors storing entries of tridiagonal matrix
9 % (9.4.0.11)
10 alph=zeros(k,1); bet = zeros(k,1);
11 for j=1:k
12 q = A*V(:,j); alph(j) = dot(q,V(:,j));
13 w = q - alph(j)*V(:,j);
14 if (j > 1), w = w - bet(j-1)*V(:,j-1);
15 end;
16 bet(j) = norm(w); V = [V,w/bet(j)];
17 end
18 bet = bet(1:end-1);

```

Total computational effort for  $l$  steps of Lanczos process, if  $\mathbf{A}$  has at most  $k$  non-zero entries per row:  
 $O(nkl)$

Note: Code 9.4.0.12 assumes that no residual vanishes. This could happen, if  $\mathbf{z}_0$  exactly belonged to the span of a few eigenvectors. However, in practical computations inevitable round-off errors will always ensure that the iterates do not stay in an invariant subspace of  $\mathbf{A}$ , cf. Rem. 9.3.1.22.

Convergence (what we expect from the above considerations) → [DH03, Sect. 8.5])

$$\text{In } l\text{-th step: } \lambda_n \approx \mu_l^{(l)}, \lambda_{n-1} \approx \mu_{l-1}^{(l)}, \dots, \lambda_1 \approx \mu_1^{(l)}, \\ \sigma(\mathbf{T}_l) = \{\mu_1^{(l)}, \dots, \mu_l^{(l)}\}, \quad \mu_1^{(l)} \leq \mu_2^{(l)} \leq \dots \leq \mu_l^{(l)}.$$

### EXAMPLE 9.4.0.13 (Lanczos process for eigenvalue computation)

$\mathbf{A}$  from Ex. 9.4.0.4

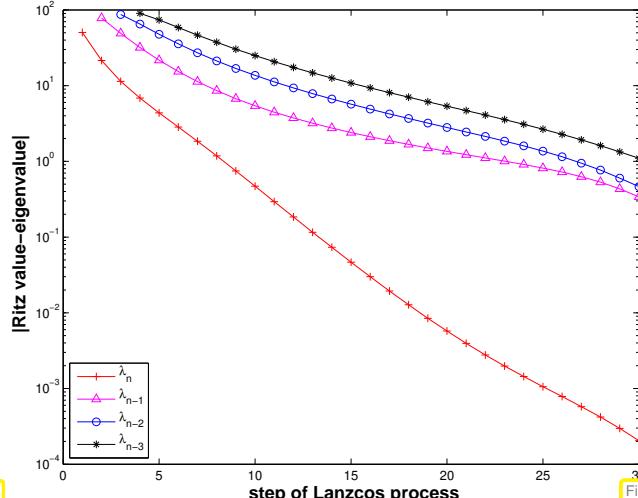


Fig. 338

$\mathbf{A} = \text{gallery('minij', 100);}$

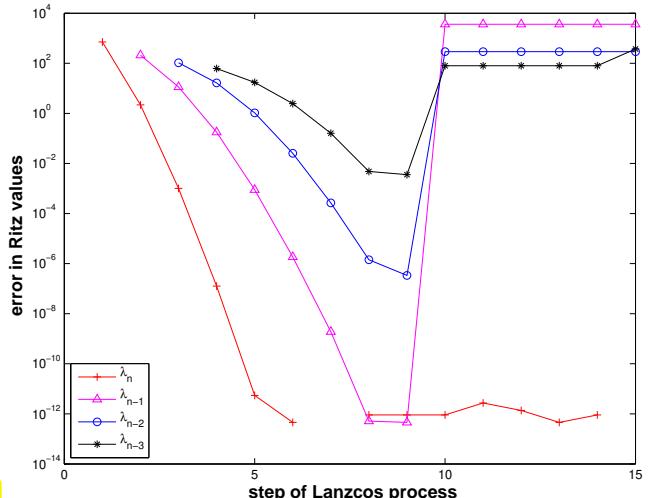


Fig. 339

Observation: same as in Ex. 9.4.0.4, linear convergence of Ritz values to eigenvalues.

However for  $\mathbf{A} \in \mathbb{R}^{10,10}$ ,  $a_{ij} = \min\{i, j\}$  good initial convergence, but sudden “jump” of Ritz values off eigenvalues!

Conjecture: Impact of roundoff errors, cf. Ex. 10.2.3.1

### EXAMPLE 9.4.0.14 (Impact of roundoff on Lanczos process)

$$\mathbf{A} \in \mathbb{R}^{10,10}, \quad a_{ij} = \min\{i, j\}.$$

$\mathbf{A} = \text{gallery('minij', 10);}$

Computed by  $[V, \alpha, \beta] = \text{lanczos}(\mathbf{A}, n, \text{ones}(n, 1));$ , see Code 9.4.0.12:

$$\mathbf{T} = \left( \begin{array}{ccccccccc} 38.500000 & 14.813845 & & & & & & & \\ 14.813845 & 9.642857 & 2.062955 & & & & & & \\ & 2.062955 & 2.720779 & 0.776284 & & & & & \\ & & 0.776284 & 1.336364 & 0.385013 & & & & \\ & & & 0.385013 & 0.826316 & 0.215431 & & & \\ & & & & 0.215431 & 0.582380 & 0.126781 & & \\ & & & & & 0.126781 & 0.446860 & 0.074650 & & \\ & & & & & & 0.074650 & 0.363803 & 0.043121 & & \\ & & & & & & & 0.043121 & 3.820888 & 11.991094 & \\ & & & & & & & & & 11.991094 & 41.254286 \end{array} \right)$$

$$\sigma(\mathbf{A}) = \{0.255680, 0.273787, 0.307979, 0.366209, 0.465233, 0.643104, 1.000000, 1.873023, 5.048917, 44.766069\}$$

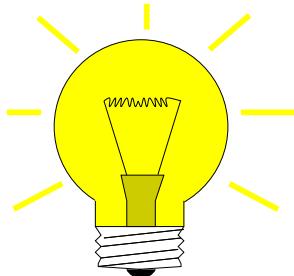
$$\sigma(\mathbf{T}) = \{0.263867, 0.303001, 0.365376, 0.465199, 0.643104, 1.000000, 1.873023, 5.048917, 44.765976, 44.766069\}$$

- Uncanny cluster of computed eigenvalues of  $\mathbf{T}$  (“ghost eigenvalues”, [GV89, Sect. 9.2.5])

$$\mathbf{V}^H \mathbf{V} = \begin{pmatrix} 1.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000251 & 0.258801 & \textcolor{red}{0.883711} \\ 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000106 & 0.109470 & 0.373799 \\ 0.000000 & -0.000000 & 1.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000005 & 0.005373 & 0.018347 \\ 0.000000 & 0.000000 & 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000096 & 0.0000328 & \\ 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & 0.000000 & 0.000000 & 0.000001 & 0.000003 & \\ 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & \\ 0.000000 & 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & \\ 0.000251 & 0.000106 & 0.000005 & 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 \\ 0.258801 & 0.109470 & 0.005373 & 0.000096 & 0.000001 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & 0.000000 \\ \textcolor{red}{0.883711} & \textcolor{red}{0.373799} & \textcolor{red}{0.018347} & \textcolor{blue}{0.000328} & \textcolor{blue}{0.000003} & \textcolor{blue}{0.000000} & \textcolor{blue}{0.000000} & \textcolor{blue}{0.000000} & \textcolor{blue}{0.000000} & 1.000000 \end{pmatrix}$$

- Loss of orthogonality of residual vectors due to roundoff  
(compare: impact of roundoff on CG iteration, Ex. 10.2.3.1)

$l$	$\sigma(\mathbf{T}_l)$								
1	38.500000								
2	3.392123 44.750734								
3	1.117692 4.979881 44.766064								
4	0.597664 1.788008 5.048259 44.766069								
5	0.415715 0.925441 1.870175 5.048916 44.766069								
6	0.336507 0.588906 0.995299 1.872997 5.048917 44.766069								
7	0.297303 0.431779 0.638542 0.999922 1.873023 5.048917 44.766069								
8	0.276160 0.349724 0.462449 0.643016 1.000000 1.873023 5.048917 44.766069								
9	0.276035 0.349451 0.462320 0.643006 1.000000 1.873023 3.821426 5.048917 44.766069								
10	0.263867 0.303001 0.365376 0.465199 0.643104 1.000000 1.873023 5.048917 44.765976 44.766069								



Idea:

- ◆ do not rely on orthogonality relations of Lemma 10.2.2.5
- ◆ use explicit **Gram-Schmidt orthogonalization** [NS02, Thm. 4.8], [Gut09, Alg .6.1]

Details: inductive approach: given  $\{\mathbf{v}_1, \dots, \mathbf{v}_l\}$  ONB of  $\mathcal{K}_l(\mathbf{A}, \mathbf{z})$

$$\blacktriangleright \tilde{\mathbf{v}}_{l+1} := \mathbf{A}\mathbf{v}_l - \sum_{j=1}^l (\mathbf{v}_j^H \mathbf{A}\mathbf{v}_l) \mathbf{v}_j, \quad \mathbf{v}_{l+1} := \frac{\tilde{\mathbf{v}}_{l+1}}{\|\tilde{\mathbf{v}}_{l+1}\|_2} \Rightarrow \mathbf{v}_{l+1} \perp \mathcal{K}_l(\mathbf{A}, \mathbf{z}). \quad (9.4.0.15)$$

(Gram-Schmidt, cf. (10.2.2.4))

orthogonal

- **Arnoldi process:** In step  $l$ :  $1 \times \mathbf{A} \times \text{vector}$   
 $l+1$  dot products  
 $l$  AXPY-operations  
 $n$  divisions

➤ Computational cost for  $l$  steps, if at most  $k$  non-zero entries in each row of  $\mathbf{A}$ :  $\mathcal{O}(nkl^2)$

**MATLAB-code 9.4.0.16: Arnoldi process**

```

1 function [V,H] = arnoldi(A,k,v0)
2 % Columns of V store orthonormal basis of Krylov spaces  $\mathcal{K}_l(\mathbf{A}, \mathbf{v}_0)$ .
3 % H returns Hessenberg matrix, see Lemma 9.4.0.17.
4 V = [v0/norm(v0)];
5 H = zeros(k+1,k);
6 for l=1:k
7   vt = A*V(:,l); % "power iteration", next basis vector
8   for j=1:l
9     % Gram-Schmidt orthogonalization, cf. Sect. 9.3.4.1
10    H(j,l) = dot(V(:,j),vt);
11    vt = vt - H(j,l)*V(:,j);
12  end
13  H(l+1,l) = norm(vt);
14  if (H(l+1,l) == 0), break; end % "theoretical" termination
15 V = [V, vt/H(l+1,l)];
16 end

```

If it does not stop prematurely, the Arnoldi process of Code 9.4.0.16 will yield an *orthonormal basis* (ONB) of  $\mathcal{K}_{k+1}(\mathbf{A}, \mathbf{v}_0)$  for a **general**  $\mathbf{A} \in \mathbb{C}^{n,n}$ .

Algebraic view of the Arnoldi process of Code 9.4.0.16, meaning of output H:

$$\mathbf{V}_l = [\mathbf{v}_1, \dots, \mathbf{v}_l] : \mathbf{AV}_l = \mathbf{V}_{l+1} \tilde{\mathbf{H}}_l , \quad \tilde{\mathbf{H}}_l \in \mathbb{K}^{l+1,l} \text{ mit } \tilde{h}_{ij} = \begin{cases} \mathbf{v}_i^H \mathbf{Av}_j & , \text{ if } i \leq j , \\ \|\tilde{\mathbf{v}}_i\|_2 & , \text{ if } i = j+1 , \\ 0 & \text{else.} \end{cases}$$

→  $\tilde{\mathbf{H}}_l$  = non-square upper Hessenberg matrices

$$\left( \mathbf{A} \right) \left( \begin{array}{c|c|c|c|c} \mathbf{v}_1 & \mathbf{v}_2 & \mathbf{v}_3 & \cdots & \mathbf{v}_l \end{array} \right) = \left( \begin{array}{c|c|c|c|c|c} \mathbf{v}_1 & \mathbf{v}_2 & \mathbf{v}_3 & \cdots & \mathbf{v}_l & \mathbf{1}^T \mathbf{v}_l \\ \hline & & & & & \end{array} \right) \left( \begin{array}{c|c|c|c|c} & & & & \mathbf{H}_l \\ \hline & & & & \\ \hline & & & & \\ \hline & & & & \end{array} \right)$$

Translate Code 9.4.0.16 to matrix calculus:

**Lemma 9.4.0.17. Theory of Arnoldi process**

For the matrices  $\mathbf{V}_l \in \mathbb{K}^{n,l}$ ,  $\tilde{\mathbf{H}}_l \in \mathbb{K}^{l+1,l}$  arising in the  $l$ -th step,  $l \leq n$ , of the Arnoldi process holds

- (i)  $\mathbf{V}_l^H \mathbf{V}_l = \mathbf{I}$  (unitary matrix),
- (ii)  $\mathbf{A}\mathbf{V}_l = \mathbf{V}_{l+1}\tilde{\mathbf{H}}_l$ ,  $\tilde{\mathbf{H}}_l$  is non-square upper Hessenberg matrix,
- (iii)  $\mathbf{V}_l^H \mathbf{A}\mathbf{V}_l = \mathbf{H}_l \in \mathbb{K}^{l,l}$ ,  $h_{ij} = \tilde{h}_{ij}$  for  $1 \leq i, j \leq l$ ,
- (iv) If  $\mathbf{A} = \mathbf{A}^H$  then  $\mathbf{H}_l$  is tridiagonal ( $\geq$  Lanczos process)

*Proof.* Direct from Gram-Schmidt orthogonalization and inspection of Code 9.4.0.16.  $\square$

**Remark 9.4.0.18 (Arnoldi process and Ritz projection)**

Interpretation of Lemma 9.4.0.17 (iii) & (i):

$\mathbf{H}_l \mathbf{x} = \lambda \mathbf{x}$  is a (generalized) Ritz projection of EVP  $\mathbf{Ax} = \lambda \mathbf{x}$ , cf. Section 9.3.4.2.

► Eigenvalue approximation for general EVP  $\mathbf{Ax} = \lambda \mathbf{x}$  by Arnoldi process:

**MATLAB-code 9.4.0.19: Arnoldi eigenvalue approximation**

```

1 function [dn,V,Ht] =
2     arnoldieig(A,v0,k,tol)
3 n = size(A,1); V = [v0/norm(v0)];
4 H = zeros(1,0); dn = zeros(k,1);
5 for l=1:n
6     d = dn;
7     Ht = [Ht, zeros(l,1); zeros(1,l)];
8     vt = A*V(:,l);
9     for j=1:l
10         Ht(j,l) = dot(V(:,j),vt);
11         vt = vt - Ht(j,l)*V(:,j);
12     end
13     ev = sort(eig(Ht(1:l,1:l)));
14     dn(1:min(l,k)) =
15         ev(end:-1:end-min(l,k)+1);
16     if (norm(d-dn) <
17         tol*norm(dn)) &\pnodenode{SPOWITX} &,
18         break; end;
19     Ht(l+1,l) = norm(vt);
20     V = [V, vt/Ht(l+1,l)];
21 end

```

**MATLAB-code 9.4.0.20: MATLAB-CODE Arnoldi eigenvalue approximation**

```

function [dn,V,Ht] = arnoldieig(A,v0,k,tol)
n = size(A,1); V = [v0/norm(v0)];
H = zeros(1,0); dn = zeros(k,1);
for l=1:n
    d = dn;
    Ht = [Ht, zeros(l,1); zeros(1,l)];
    vt = A*V(:,l);
    for j=1:l
        Ht(j,l) = dot(V(:,j),vt);
        vt = vt - Ht(j,l)*V(:,j);
    end
    ev = sort(eig(Ht(1:l,1:l)));
    dn(1:min(l,k)) = ev(end:-1:end-min(l,k)+1);
    if (norm(d-dn) < tol*norm(dn)), break; end;
    Ht(l+1,l) = norm(vt);
    V = [V, vt/Ht(l+1,l)];
end

```

Heuristic termination criterion

Arnoldi process for computing the  $k$  largest (in modulus) eigenvalues of  $\mathbf{A} \in \mathbb{C}^{n,n}$

1 A  $\times$  vector per step  
( $\gg$  attractive for sparse matrices)

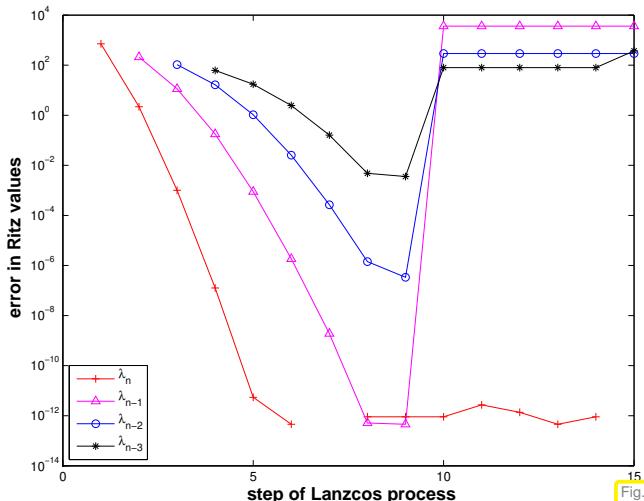
However: required storage increases with number of steps,  
*cf.* situation with GMRES, Section 10.4.1.

## Heuristic termination criterion

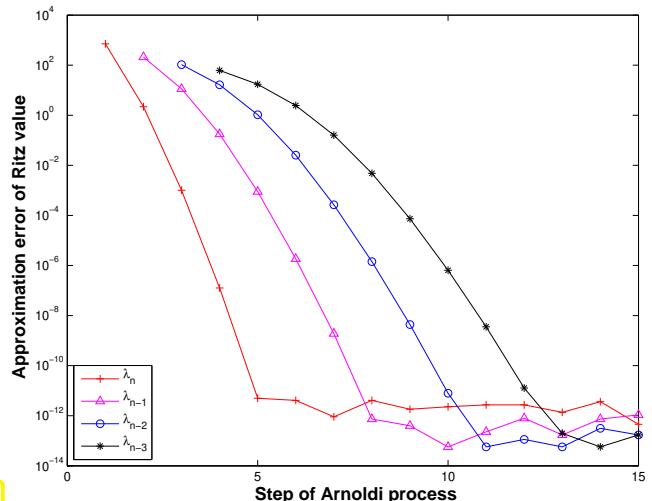
#### **EXAMPLE 9.4.0.21 (Stability of Arnoldi process)**

$$\mathbf{A} \in \mathbb{R}^{100,100} \quad , \quad a_{ij} = \min\{i,j\} \quad .$$

```
A = gallery('minij', 100);
```



## Lanczos process: Ritz values



## Arnoldi process: Ritz values

Ritz values during Arnoldi process for `A = gallery('minij', 10);`  $\leftrightarrow$  Ex. 9.4.0.13

Observation: (almost perfect approximation of spectrum of **A**)

For the above examples both the Arnoldi process and the Lanczos process are *algebraically equivalent*.

because they are applied to a symmetric matrix  $\mathbf{A} = \mathbf{A}^T$ . However, they behave strikingly differently, which indicates that they are *not numerically equivalent*.

The Arnoldi process is much less affected by roundoff than the Lanczos process, because it does not take for granted orthogonality of the “residual vector sequence”. Hence, the Arnoldi process enjoys superior numerical stability ( $\rightarrow ??$ , Def. 1.5.5.19) compared to the Lanczos process. ■

### EXAMPLE 9.4.0.22 (Eigenvalue computation with Arnoldi process)

Eigenvalue approximation from Arnoldi process for *non-symmetric*  $\mathbf{A}$ , initial vector `ones(100, 1)`;

#### MATLAB-code 9.4.0.23:

```
1 n=100;
2 M=full(gallery('tridiag', -0.5*ones(n-1, 1), 2*ones(n, 1), -1.5*ones(n-1, 1)));
3 A=M*diag(1:n)*inv(M);
```

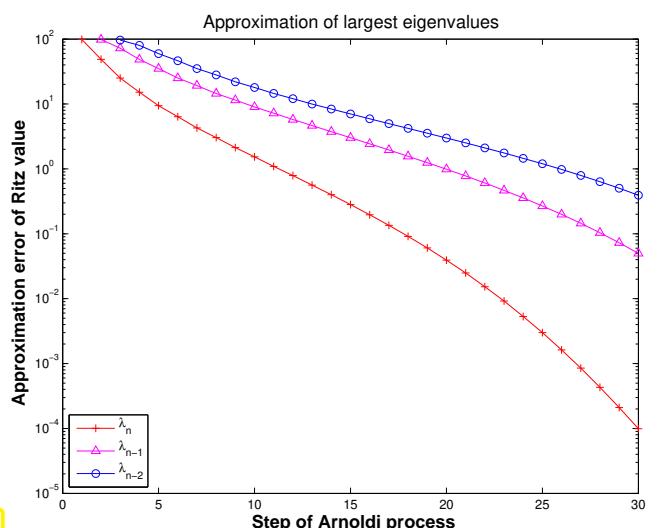
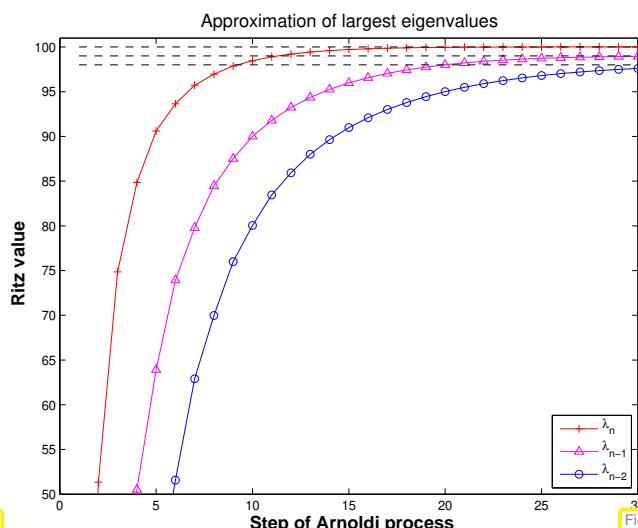


Fig. 342

Fig. 343

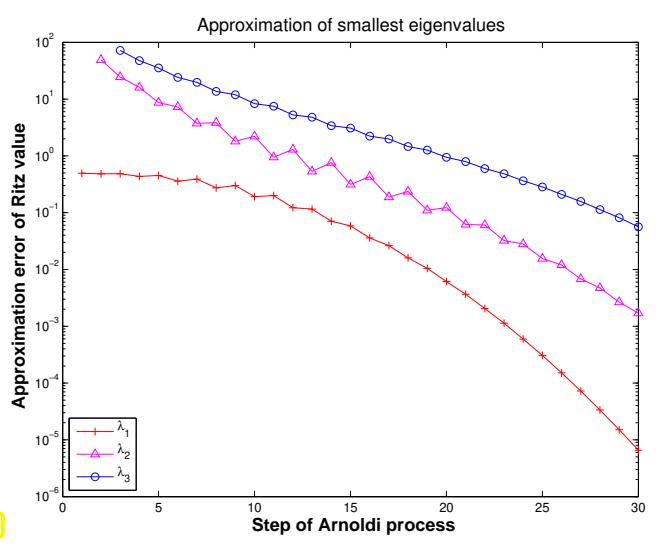
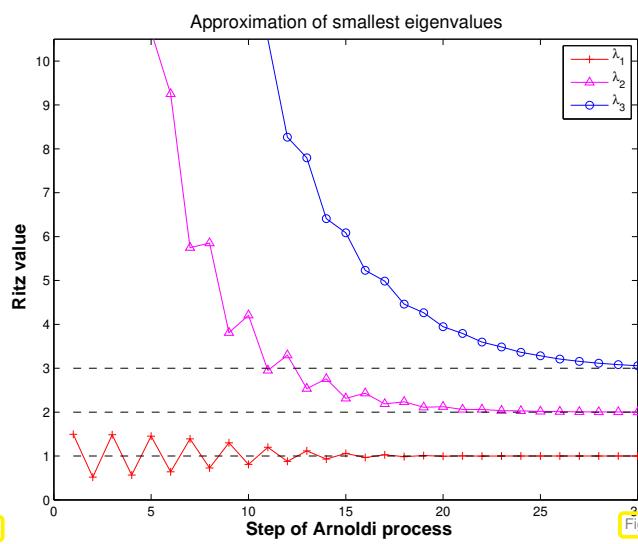


Fig. 345

Observation: “vaguely linear” convergence of largest and smallest eigenvalues, cf. Ex. 9.4.0.4. ■

Krylov subspace iteration methods (= Arnoldi process, Lanczos process) attractive for computing a few of the largest/smallest eigenvalues and associated eigenvectors of *large sparse matrices*.

**Remark 9.4.0.24 (Krylov subspace methods for generalized EVP)**

Adaptation of Krylov subspace iterative eigensolvers to generalized EVP:  $\mathbf{Ax} = \lambda \mathbf{Bx}$ ,  $\mathbf{B}$  s.p.d.: replace Euclidean inner product with “ $\mathbf{B}$ -inner product”  $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x}^H \mathbf{B} \mathbf{y}$ . □

MATLAB-functions:

```
d = eigs(A,k,sigma) : k largest/smallest eigenvalues of A
d = eigs(A,B,k,sigma): k largest/smallest eigenvalues for generalized EVP Ax = λBx, B
                           s.p.d.
d = eigs(Afun,n,k)   : Afun = handle to function providing matrix×vector for A/A⁻¹/A -
                           αI/(A - αB)⁻¹. (Use flags to tell eigs about special properties of
                           matrix behind Afun.)
```

`eigs` just calls routines of the open source **ARPACK** numerical library.

# Bibliography

- [AKY99] Charles J Alpert, Andrew B Kahng, and So-Zen Yao. “Spectral partitioning with multiple eigenvectors”. In: *Discrete Applied Mathematics* 90.1-3 (1999), pp. 3–26. DOI: [10.1016/S0166-218X\(98\)00083-3](https://doi.org/10.1016/S0166-218X(98)00083-3) (cit. on p. 581).
- [Bai+00] Z.-J. Bai, J. Demmel, J. Dongarra, A. Ruhe, and H. van der Vorst. *Templates for the Solution of Algebraic Eigenvalue Problems*. Philadelphia, PA: SIAM, 2000 (cit. on p. 553).
- [BF06] Yuri Boykov and Gareth Funka-Lea. “Graph Cuts and Efficient N-D Image Segmentation”. In: *International Journal of Computer Vision* 70.2 (Nov. 2006), pp. 109–131. DOI: [10.1007/s11263-006-7934-5](https://doi.org/10.1007/s11263-006-7934-5).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 556, 558, 563, 570, 572, 582).
- [DH03] P. Deuflhard and A. Hohmann. *Numerical Analysis in Modern Scientific Computing*. Vol. 43. Texts in Applied Mathematics. Springer, 2003 (cit. on p. 606).
- [Gle15] David F. Gleich. “PageRank Beyond the Web”. In: *SIAM Review* 57.3 (2015), pp. 321–363. DOI: [10.1137/140976649](https://doi.org/10.1137/140976649) (cit. on p. 563).
- [GV89] G.H. Golub and C.F. Van Loan. *Matrix computations*. 2nd. Baltimore, London: John Hopkins University Press, 1989 (cit. on pp. 559, 560, 577, 607).
- [GT08] Craig Gotsman and Sivan Toledo. “On the computation of null spaces of sparse rectangular matrices”. In: *SIAM J. Matrix Anal. Appl.* 30.2 (2008), pp. 445–463. DOI: [10.1137/050638369](https://doi.org/10.1137/050638369) (cit. on p. 582).
- [Gut09] M.H. Gutknecht. *Lineare Algebra*. Lecture Notes. SAM, ETH Zürich, 2009 (cit. on pp. 556–558, 596, 607).
- [Hac94] Wolfgang Hackbusch. *Iterative solution of large sparse systems of equations*. Vol. 95. Applied Mathematical Sciences. New York: Springer-Verlag, 1994, pp. xxii+429 (cit. on p. 558).
- [Hal70] K.M. Hall. “An r-dimensional quadratic placement algorithm”. In: *Management Science* 17.3 (1970), pp. 219–229.
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 558, 559, 583, 584, 588, 601).
- [LM06] A.N. Lengvile and C.D. Meyer. *Google’s PageRank and Beyond: The Science of Search Engine Rankings*. Princeton, NJ: Princeton University Press, 2006 (cit. on p. 563).
- [Ney99a] K. Neymeyr. *A geometric theory for preconditioned inverse iteration applied to a subspace*. Tech. rep. 130. Tübingen, Germany: SFB 382, Universität Tübingen, Nov. 1999 (cit. on p. 587).
- [Ney99b] K. Neymeyr. *A geometric theory for preconditioned inverse iteration: III. Sharp convergence estimates*. Tech. rep. 130. Tübingen, Germany: SFB 382, Universität Tübingen, Nov. 1999 (cit. on p. 587).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on pp. 556–559, 596, 607).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 557–559, 563, 572, 596).
- [SM00] J.-B. Shi and J. Malik. “Normalized cuts and image segmentation”. In: *IEEE Trans. Pattern Analysis and Machine Intelligence* 22.8 (2000), pp. 888–905 (cit. on pp. 573, 576, 581).

- [ST96] D.A. Spielman and Shang-Hua Teng. “Spectral partitioning works: planar graphs and finite element meshes”. In: *Foundations of Computer Science, 1996. Proceedings., 37th Annual Symposium on*. Oct. 1996, pp. 96–105. DOI: [10.1109/SFCS.1996.548468](https://doi.org/10.1109/SFCS.1996.548468) (cit. on p. 581).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on p. 556).
- [TM01] F. Tisseur and K. Meerbergen. “The quadratic eigenvalue problem”. In: *SIAM Review* 43.2 (2001), pp. 235–286 (cit. on p. 555).

# Chapter 10

## Krylov Methods for Linear Systems of Equations



*Supplementary literature.* There is a wealth of literature on iterative methods for the solution of

linear systems of equations: The two books [Hac94] and [Saa03] offer a comprehensive treatment of the topic (the latter is available online for ETH students and staff).

Concise presentations can be found in [QSS00, Ch. 4] and [DR08, Ch. 13].

Learning outcomes:

- Understanding when and why iterative solution of linear systems of equations may be preferred to direct solvers based on Gaussian elimination.
- 

= A class of **iterative methods** ( $\rightarrow$  Section 8.1) for approximate solution of large linear systems of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{K}^{n,n}$ .

BUT, we have reliable *direct* methods (Gauss elimination  $\rightarrow$  Section 2.3, LU-factorization  $\rightarrow$  § 2.3.2.15, QR-factorization  $\rightarrow$  ??) that provide an (apart from roundoff errors) exact solution with a *finite* number of elementary operations!

Alas, direct elimination may **not be feasible**, or may be grossly inefficient, because

- it may be too expensive (e.g. for  $\mathbf{A}$  too large, sparse),  $\rightarrow$  (2.3.2.10),
- inevitable fill-in may exhaust main memory,
- the system matrix may be available only as procedure  $y = \text{evalA}(x) \leftrightarrow \mathbf{y} = \mathbf{Ax}$

### Contents

10.1	Descent Methods [QSS00, Sect. 4.3.3]	617
10.1.1	Quadratic minimization context	617
10.1.2	Abstract steepest descent	618
10.1.3	Gradient method for s.p.d. linear system of equations	619
10.1.4	Convergence of the gradient method	620
10.2	Conjugate gradient method (CG) [Han02, Ch. 9], [DR08, Sect. 13.4], [QSS00, Sect. 4.3.4]	624
10.2.1	Krylov spaces	625
10.2.2	Implementation of CG	626
10.2.3	Convergence of CG	629
10.3	Preconditioning [DR08, Sect. 13.5], [Han02, Ch. 10], [QSS00, Sect. 4.3.5]	634
10.4	Survey of Krylov Subspace Methods	640

10.4.1 Minimal residual methods . . . . .	640
10.4.2 Iterations with short recursions [QSS00, Sect. 4.5] . . . . .	641

## 10.1 Descent Methods [QSS00, Sect. 4.3.3]

Focus:

Linear system of equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{n,n}$ ,  $\mathbf{b} \in \mathbb{R}^n$ ,  $n \in \mathbb{N}$  given,  
with **symmetric positive definite** (s.p.d.,  $\rightarrow$  Def. 1.1.2.6) system matrix  $\mathbf{A}$



$\mathbf{A}$ -inner product  $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x}^\top \mathbf{A} \mathbf{y}$   $\Rightarrow$  “ $\mathbf{A}$ -geometry”

**Definition 10.1.0.1. Energy norm**  $\rightarrow$  [Han02, Def. 9.1]

A s.p.d. matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  induces an **energy norm**

$$\|\mathbf{x}\|_A := (\mathbf{x}^\top \mathbf{A} \mathbf{x})^{1/2}, \quad \mathbf{x} \in \mathbb{R}^n.$$

**Remark 10.1.0.2 (Krylov methods for complex s.p.d. system matrices)** In this chapter, for the sake of simplicity, we restrict ourselves to  $\mathbb{K} = \mathbb{R}$ .

However, the (conjugate) gradient methods introduced below also work for LSE  $\mathbf{Ax} = \mathbf{b}$  with  $\mathbf{A} \in \mathbb{C}^{n,n}$ ,  $\mathbf{A} = \mathbf{A}^H$  s.p.d. when  $^\top$  is replaced with  $^H$  (Hermitian transposed). Then, all theoretical statements remain valid unaltered for  $\mathbb{K} = \mathbb{C}$ .  $\square$

### 10.1.1 Quadratic minimization context

**Lemma 10.1.1.1. S.p.d. LSE and quadratic minimization problem [DR08, (13.37)]**

A LSE with  $\mathbf{A} \in \mathbb{R}^{n,n}$  s.p.d. and  $\mathbf{b} \in \mathbb{R}^n$  is equivalent to a minimization problem:

$$\mathbf{Ax} = \mathbf{b} \Leftrightarrow \mathbf{x} = \arg \min_{\mathbf{y} \in \mathbb{R}^n} J(\mathbf{y}), \quad J(\mathbf{y}) := \frac{1}{2} \mathbf{y}^\top \mathbf{A} \mathbf{y} - \mathbf{b}^\top \mathbf{y}. \quad (10.1.1.2)$$

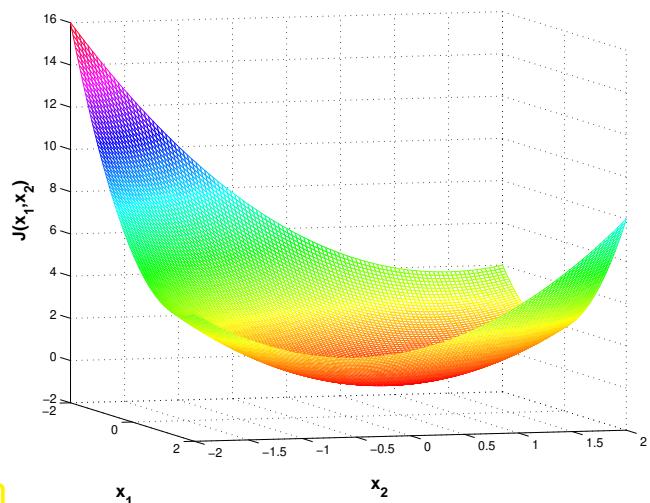
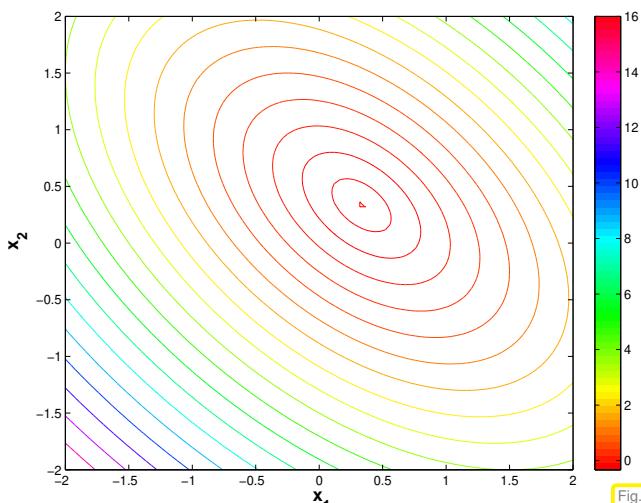
A **quadratic functional**

*Proof.* If  $\mathbf{x}^* := \mathbf{A}^{-1} \mathbf{b}$  a straightforward computation using  $\mathbf{A} = \mathbf{A}^\top$  shows

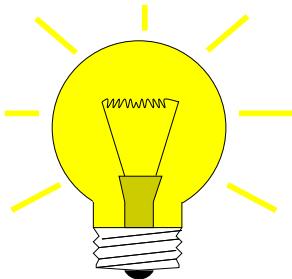
$$\begin{aligned} J(\mathbf{x}) - J(\mathbf{x}^*) &= \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x} - \frac{1}{2} (\mathbf{x}^*)^\top \mathbf{A} \mathbf{x}^* + \mathbf{b}^\top \mathbf{x}^* \\ &\stackrel{\mathbf{b}=\mathbf{Ax}^*}{=} \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - (\mathbf{x}^*)^\top \mathbf{A} \mathbf{x} + \frac{1}{2} (\mathbf{x}^*)^\top \mathbf{A} \mathbf{x}^* \\ &= \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|_A^2. \end{aligned} \quad (10.1.1.3)$$

Then the assertion follows from the properties of the energy norm.  $\square$

**EXAMPLE 10.1.1.4 (Quadratic functional in 2D)** Plot of  $J$  from (10.1.1.2) for  $\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$ ,  $\mathbf{b} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ .



Level lines of quadratic functionals with s.p.d.  $\mathbf{A}$  are (hyper)ellipses



Algorithmic idea: (Lemma 10.1.1.1  $\Rightarrow$ ) Solve  $\mathbf{Ax} = \mathbf{b}$  iteratively by successive solution of *simpler* minimization problems

## 10.1.2 Abstract steepest descent

Task:

Given continuously differentiable  $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}$ ,  
find minimizer  $\mathbf{x}^* \in D$ :  $\mathbf{x}^* = \underset{\mathbf{x} \in D}{\operatorname{argmin}} F(\mathbf{x})$

Note that a minimizer need not exist, if  $F$  is not bounded from below (e.g.,  $F(x) = x^3$ ,  $x \in \mathbb{R}$ , or  $F(x) = \log x$ ,  $x > 0$ ), or if  $D$  is open (e.g.,  $F(x) = \sqrt{x}$ ,  $x > 0$ ).

The existence of a minimizer is guaranteed if  $F$  is bounded from below and  $D$  is closed ( $\rightarrow$  Analysis).

The most natural iteration:

### §10.1.2.1 (Steepest descent (ger.: steilster Abstieg))

```

Initial guess  $\mathbf{x}^{(0)} \in D$ ,  $k = 0$ 
repeat
   $\mathbf{d}_k := -\operatorname{grad} F(\mathbf{x}^{(k)})$ 
   $t^* := \underset{t \in \mathbb{R}}{\operatorname{argmin}} F(\mathbf{x}^{(k)} + t\mathbf{d}_k)$  (line search)
   $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + t^*\mathbf{d}_k$ 
   $k := k + 1$ 
until  $\left( \begin{array}{l} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau_{\text{rel}} \|\mathbf{x}^{(k)}\| \\ \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau_{\text{abs}} \end{array} \right)$ 

```

- ◆  $\mathbf{d}_k \hat{=} \text{direction of steepest descent}$
- ◆ linear search  $\hat{=} 1\text{D}$  minimization:  
use Newton's method ( $\rightarrow$  Section 8.3.2.1) on derivative
- ◆ correction based a posteriori termination criterion, see Section 8.1.2 for a discussion.  
( $\tau \hat{=} \text{prescribed tolerance}$ )

The **gradient** ( $\rightarrow$  [Str09, Kapitel 7])

$$\mathbf{grad} F(\mathbf{x}) = \begin{bmatrix} \frac{\partial F}{\partial x_1} \\ \vdots \\ \frac{\partial F}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n \quad (10.1.2.2)$$

provides the direction of **local** steepest ascent/descent of  $F$

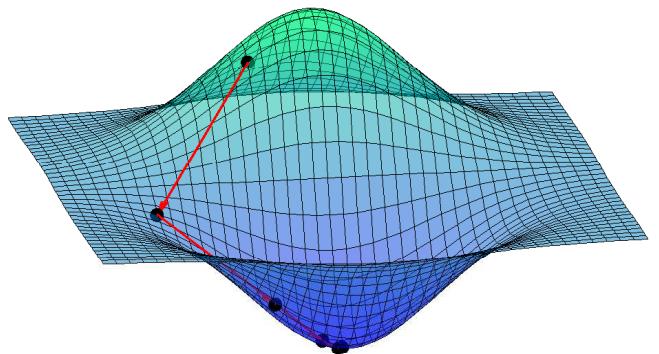


Fig. 348

Of course this very algorithm can encounter plenty of difficulties:

- iteration may get stuck in a *local minimum*,
- iteration may diverge or lead out of  $D$ ,
- line search may not be feasible.

### 10.1.3 Gradient method for s.p.d. linear system of equations

However, for the quadratic minimization problem (10.1.1.2) § 10.1.2.1 will converge:

(“Geometric intuition”, see Fig. 346: quadratic functional  $J$  with s.p.d.  $A$  has unique global minimum,  $\mathbf{grad} J \neq 0$  away from minimum, pointing towards it.)

Adaptation: steepest descent algorithm § 10.1.2.1 for quadratic minimization problem (10.1.1.2), see [QSS00, Sect. 7.2.4]:

$$F(\mathbf{x}) := J(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x} \Rightarrow \mathbf{grad} J(\mathbf{x}) = \mathbf{A} \mathbf{x} - \mathbf{b}. \quad (10.1.3.1)$$

This follows from  $\mathbf{A} = \mathbf{A}^\top$ , the componentwise expression

$$J(\mathbf{x}) = \frac{1}{2} \sum_{i,j=1}^n a_{ij} x_i x_j - \sum_{i=1}^n b_i x_i$$

and the definition (10.1.2.2) of the gradient.

➤ For the descent direction in § 10.1.2.1 applied to the minimization of  $J$  from (10.1.1.2) holds

$$\mathbf{d}_k = \mathbf{b} - \mathbf{A} \mathbf{x}^{(k)} =: \mathbf{r}_k \quad \text{the residual} (\rightarrow \text{Def. 2.4.0.1}) \text{ for } \mathbf{x}^{(k-1)}.$$

§ 10.1.2.1 for  $F = J$  from (10.1.1.2): function to be minimized in line search step:

$$\varphi(t) := J(\mathbf{x}^{(k)} + t \mathbf{d}_k) = J(\mathbf{x}^{(k)}) + t \mathbf{d}_k^\top (\mathbf{A} \mathbf{x}^{(k)} - \mathbf{b}) + \frac{1}{2} t^2 \mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k \rightarrow \text{a parabola!}.$$

$$\frac{d\varphi}{dt}(t^*) = 0 \Leftrightarrow \boxed{t^* = \frac{\mathbf{d}_k^\top \mathbf{d}_k}{\mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k}} \quad (\text{unique minimizer}). \quad (10.1.3.2)$$

Note:

$$\mathbf{d}_k = 0 \Leftrightarrow \mathbf{A} \mathbf{x}^{(k)} = \mathbf{b} \quad (\text{solution found!})$$

Note:  $\mathbf{A}$  s.p.d. ( $\rightarrow$  Def. 1.1.2.6)  $\Rightarrow \mathbf{d}_k^\top \mathbf{A} \mathbf{d}_k > 0$ , if  $\mathbf{d}_k \neq 0$

►  $\varphi(t)$  is a parabola that is bounded from below (upward opening)

Based on (10.1.3.1) and (10.1.3.2) we obtain the following steepest descent method for the minimization problem (10.1.1.2):

Steepest descent iteration = **gradient method** for LSE  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{n,n}$  s.p.d.,  $\mathbf{b} \in \mathbb{R}^n$ :

### §10.1.3.3 (Gradient method for s.p.d. LSE)

```

Initial guess  $\mathbf{x}^{(0)} \in \mathbb{R}^n$ ,  $k = 0$ 
 $\mathbf{r}_0 := \mathbf{b} - \mathbf{Ax}^{(0)}$ 
repeat
     $t^* := \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{Ar}_k}$ 
     $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + t^* \mathbf{r}_k$ 
     $\mathbf{r}_{k+1} := \mathbf{r}_k - t^* \mathbf{Ar}_k$ 
     $k := k + 1$ 
until  $\left( \begin{array}{l} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau_{\text{rel}} \|\mathbf{x}^{(k)}\| \\ \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau_{\text{abs}} \end{array} \right)$  or

```

**MATLAB-code 10.1.3.2**  
 $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A}$  s.p.d.

```

1 function x
2 gradit(A,b)
3 r = b-A*x;
4 Def. 2.4.
5 for k=1:max
6 p = A*r;
7 ts = (r'*r);
8 (10.1.3.2)
9 x = x + t;
10 cn = (abs(r));
11 correc(cn);
12 if ((cn < tol))
13 (cn < tol);
14 return;
15 r = r - t;
16 end

```

Recursion for residuals, see Line 11 of Code 10.1.3.4:

$$\mathbf{r}_{k+1} = \mathbf{b} - \mathbf{Ax}^{(k+1)} = \mathbf{b} - \mathbf{A}(\mathbf{x}^{(k)} + t^* \mathbf{r}_k) = \mathbf{r}_k - t^* \mathbf{Ar}_k. \quad (10.1.3.5)$$

One step of gradient method involves

- ◆ *A single matrix  $\times$  vector product with  $\mathbf{A}$ ,*
- ◆ *2 AXPY-operations ( $\rightarrow$  Section 1.3.2) on vectors of length  $n$ ,*
- ◆ *2 dot products in  $\mathbb{R}^n$ .*

Computational cost (per step) =  $\text{cost}(\text{matrix} \times \text{vector}) + O(n)$

- If  $\mathbf{A} \in \mathbb{R}^{n,n}$  is a sparse matrix ( $\rightarrow$  ??) with “ $O(n)$  nonzero entries”, and the data structures allow to perform the matrix  $\times$  vector product with a computational effort  $O(n)$ , then a single step of the gradient method costs  $O(n)$  elementary operations.
- Gradient method of § 10.1.3.3 only needs  $\mathbf{A} \times$  vector in procedural form  $y = \text{evalA}(x)$ .

## 10.1.4 Convergence of the gradient method

**EXAMPLE 10.1.4.1 (Gradient method in 2D)** S.p.d. matrices  $\in \mathbb{R}^{2,2}$ :

$$\mathbf{A}_1 = \begin{bmatrix} 1.9412 & -0.2353 \\ -0.2353 & 1.0588 \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} 7.5353 & -1.8588 \\ -1.8588 & 0.5647 \end{bmatrix}$$

Eigenvalues:  $\sigma(\mathbf{A}_1) = \{1, 2\}$ ,  $\sigma(\mathbf{A}_2) = \{0.1, 8\}$

notation: spectrum of a matrix  $\in \mathbb{K}^{n,n}$   $\sigma(\mathbf{M}) := \{\lambda \in \mathbb{C} : \lambda \text{ is eigenvalue of } \mathbf{M}\}$

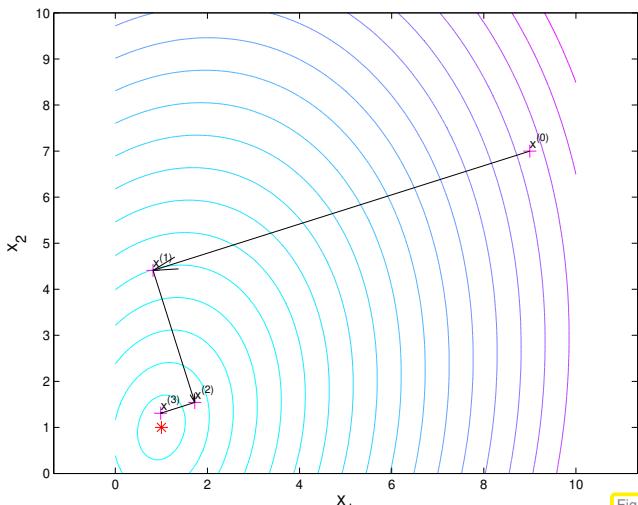


Fig. 349

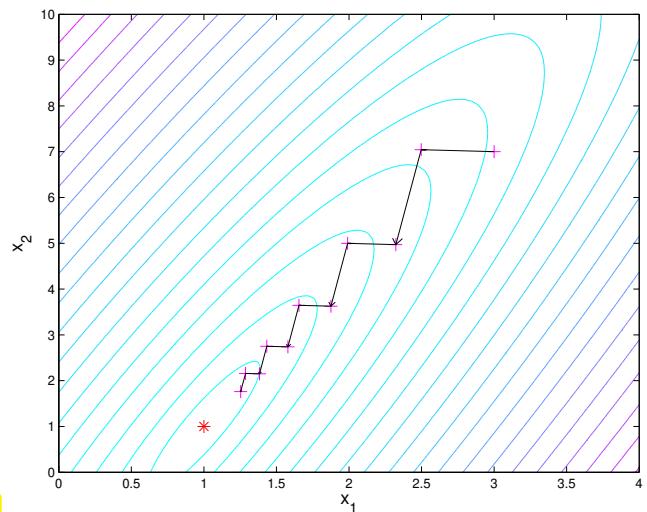
iterates of § 10.1.3.3 for  $\mathbf{A}_1$ 

Fig. 350

iterates of § 10.1.3.3 for  $\mathbf{A}_2$ 

Recall theorem on principal axis transformation: every real *symmetric* matrix can be diagonalized by *orthogonal* similarity transformations, see Cor. 9.1.0.9, [NS02, Thm. 7.8], [Gut09, Satz 9.15],

$$\mathbf{A} = \mathbf{A}^\top \in \mathbb{R}^{n,n} \Rightarrow \exists \mathbf{Q} \in \mathbb{R}^{n,n} \text{ orthogonal: } \mathbf{A} = \mathbf{Q} \mathbf{D} \mathbf{Q}^\top, \quad \mathbf{D} = \text{diag}(d_1, \dots, d_n) \in \mathbb{R}^{n,n} \text{ diagonal}. \quad (10.1.4.2)$$

$$J(\mathbf{Q}\hat{\mathbf{y}}) = \frac{1}{2}\hat{\mathbf{y}}^\top \mathbf{D}\hat{\mathbf{y}} - \underbrace{(\mathbf{Q}^\top \mathbf{b})^\top \hat{\mathbf{y}}}_{=: \hat{\mathbf{b}}^\top} = \frac{1}{2} \sum_{i=1}^n d_i \hat{y}_i^2 - \hat{b}_i \hat{y}_i.$$

Hence, a rigid transformation (rotation, reflection) maps the level surfaces of  $J$  from (10.1.1.2) to ellipses with principal axes  $d_i$ . As  $\mathbf{A}$  s.p.d.  $d_i > 0$  is guaranteed.

Observations:

- Larger spread of spectrum leads to more elongated ellipses as level lines  $\Rightarrow$  slower convergence of gradient method, see Fig. 350.
- Orthogonality of successive residuals  $\mathbf{r}_k, \mathbf{r}_{k+1}$ .

Clear from definition of § 10.1.3.3:

$$\mathbf{r}_k^\top \mathbf{r}_{k+1} = \mathbf{r}_k^\top \mathbf{r}_k - \mathbf{r}_k^\top \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{A} \mathbf{r}_k} \mathbf{A} \mathbf{r}_k = 0. \quad (10.1.4.3)$$

□

#### EXAMPLE 10.1.4.4 (Convergence of gradient method)

Convergence of gradient method for diagonal matrices,  $\mathbf{x}^* = [1, \dots, 1]^\top$ ,  $\mathbf{x}^{(0)} = 0$ :

```

1 d = 1:0.01:2; A1 = diag(d);
2 d = 1:0.1:11; A2 = diag(d);
3 d = 1:1:101; A3 = diag(d);

```

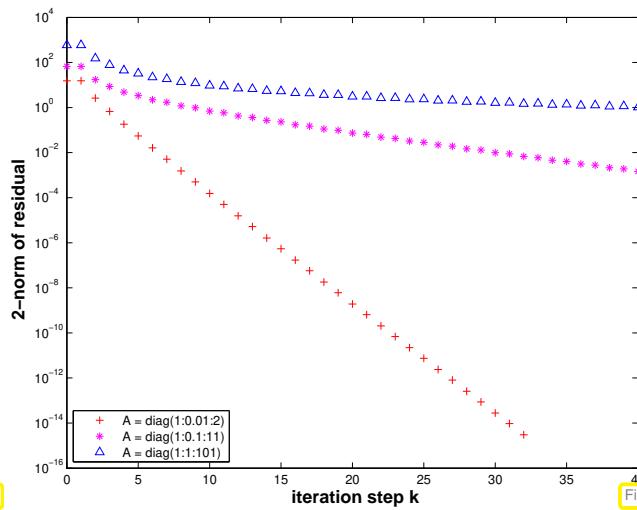


Fig. 351

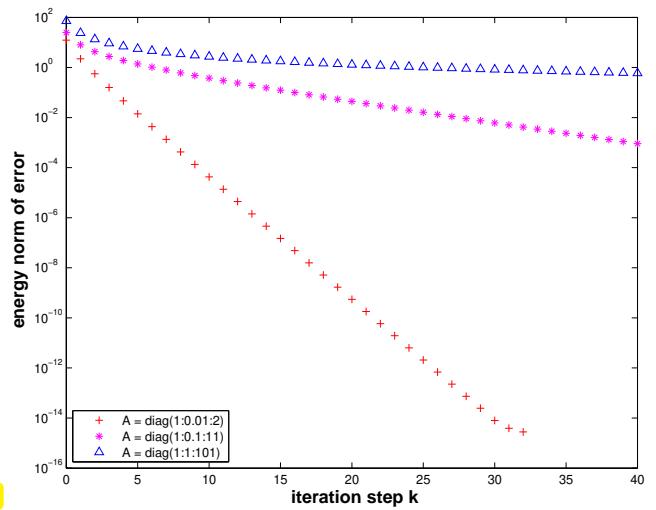


Fig. 352

Note: To study convergence it is *sufficient to consider diagonal matrices*, because

1. (10.1.4.2): for every  $\mathbf{A} \in \mathbb{R}^{n,n}$  with  $\mathbf{A}^\top = \mathbf{A}$  there is an orthogonal matrix  $\mathbf{Q} \in \mathbb{R}^{n,n}$  such that  $\mathbf{A} = \mathbf{Q}^\top \mathbf{D} \mathbf{Q}$  with a diagonal matrix  $\mathbf{D}$  (principal axis transformation),  $\rightarrow$  Cor. 9.1.0.9, [NS02, Thm. 7.8], [Gut09, Satz 9.15],
2. when applying the gradient method § 10.1.3.3 to both  $\mathbf{Ax} = \mathbf{b}$  and  $\mathbf{D}\tilde{\mathbf{x}} = \tilde{\mathbf{b}} := \mathbf{Q}\mathbf{b}$ , then the iterates  $\mathbf{x}^{(k)}$  and  $\tilde{\mathbf{x}}^{(k)}$  are related by  $\mathbf{Q}\mathbf{x}^{(k)} = \tilde{\mathbf{x}}^{(k)}$ .

With  $\tilde{\mathbf{r}}_k := \mathbf{Q}\mathbf{r}_k$ ,  $\tilde{\mathbf{x}}^{(k)} := \mathbf{Q}\mathbf{x}^{(k)}$ , using  $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ :

Initial guess  $\mathbf{x}^{(0)} \in \mathbb{R}^n$ ,  $k = 0$   
 $\mathbf{r}_0 := \mathbf{b} - \mathbf{Q}^\top \mathbf{D} \mathbf{Q} \mathbf{x}^{(0)}$   
**repeat**  
 $t^* := \frac{\mathbf{r}_k^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{Q}^\top \mathbf{D} \mathbf{Q} \mathbf{r}_k}$   
 $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + t^* \mathbf{r}_k$   
 $\mathbf{r}_{k+1} := \mathbf{r}_k - t^* \mathbf{Q}^\top \mathbf{D} \mathbf{Q} \mathbf{r}_k$   
 $k := k + 1$   
**until**  $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau \|\mathbf{x}^{(k)}\|$

Initial guess  $\tilde{\mathbf{x}}^{(0)} \in \mathbb{R}^n$ ,  $k = 0$   
 $\tilde{\mathbf{r}}_0 := \tilde{\mathbf{b}} - \mathbf{D}\tilde{\mathbf{x}}^{(0)}$   
**repeat**  
 $t^* := \frac{\tilde{\mathbf{r}}_k^\top \tilde{\mathbf{r}}_k}{\tilde{\mathbf{r}}_k^\top \mathbf{D}\tilde{\mathbf{r}}_k}$   
 $\tilde{\mathbf{x}}^{(k+1)} := \tilde{\mathbf{x}}^{(k)} + t^* \tilde{\mathbf{r}}_k$   
 $\tilde{\mathbf{r}}_{k+1} := \tilde{\mathbf{r}}_k - t^* \mathbf{D}\tilde{\mathbf{r}}_k$   
 $k := k + 1$   
**until**  $\|\tilde{\mathbf{x}}^{(k)} - \tilde{\mathbf{x}}^{(k-1)}\| \leq \tau \|\tilde{\mathbf{x}}^{(k)}\|$

Observation:

- ◆ linear convergence ( $\rightarrow$  Def. 8.1.1.1), see also Rem. 8.1.1.6
- ◆ rate of convergence increases ( $\leftrightarrow$  speed of convergence decreases) with spread of spectrum of  $\mathbf{A}$

Impact of distribution of diagonal entries ( $\leftrightarrow$  eigenvalues) of (diagonal matrix)  $\mathbf{A}$

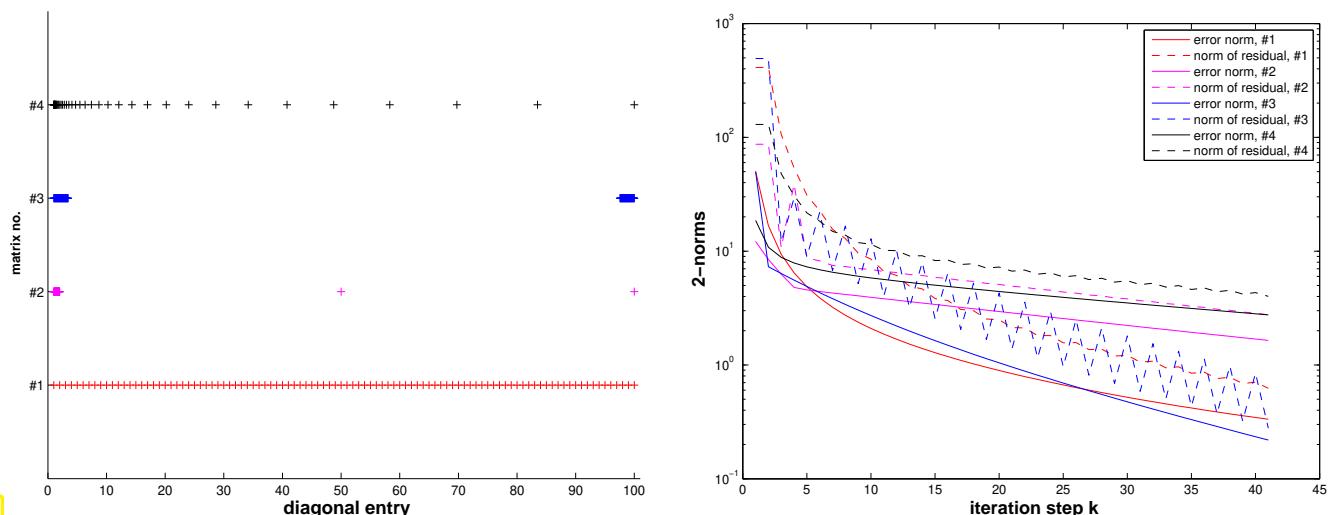
( $\mathbf{b} = \mathbf{x}^* = \mathbf{0}$ ,  $\mathbf{x} = \cos((1:n)')$ ;

Test matrix #1:  $\mathbf{A} = \text{diag}(\mathbf{d})$ ;  $\mathbf{d} = (1:100)$ ;

Test matrix #2:  $\mathbf{A} = \text{diag}(\mathbf{d})$ ;  $\mathbf{d} = [1+(0:97)/97, 50, 100]$ ;

Test matrix #3:  $\mathbf{A} = \text{diag}(\mathbf{d})$ ;  $\mathbf{d} = [1+(0:49)*0.05, 100-(0:49)*0.05]$ ;

Test matrix #4: eigenvalues exponentially dense at 1



Observation: Matrices #1, #2 & #4  $\Rightarrow$  little impact of distribution of eigenvalues on *asymptotic* convergence (exception: matrix #2)

Theory [Hac91, Sect. 9.2.2], [QSS00, Sect. 7.2.4]:

#### Theorem 10.1.4.5. Convergence of gradient method/steepest descent

The iterates of the gradient method of § 10.1.3.3 satisfy

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_A \leq L \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_A, \quad L := \frac{\text{cond}_2(\mathbf{A}) - 1}{\text{cond}_2(\mathbf{A}) + 1},$$

that is, the iteration converges at least linearly ( $\rightarrow$  Def. 8.1.1.1) w.r.t. energy norm ( $\rightarrow$  Def. 10.1.0.1).

notation:  $\text{cond}_2(\mathbf{A}) \hat{=} \text{condition number}$  ( $\rightarrow$  Def. 2.2.2.7) of  $\mathbf{A}$  induced by 2-norm

**Remark 10.1.4.6 (2-norm from eigenvalues  $\rightarrow$  [Gut09, Sect. 10.6], [NS02, Sect. 7.4])**

$$\begin{aligned} \mathbf{A} = \mathbf{A}^\top \Rightarrow \quad & \|\mathbf{A}\|_2 = \max(|\sigma(\mathbf{A})|), \\ & \|\mathbf{A}^{-1}\|_2 = \min(|\sigma(\mathbf{A})|)^{-1}, \text{ if } \mathbf{A} \text{ regular.} \end{aligned} \quad (10.1.4.7)$$

$$\mathbf{A} = \mathbf{A}^\top \Rightarrow \quad \text{cond}_2(\mathbf{A}) = \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})}, \quad \text{where} \quad \begin{aligned} \lambda_{\max}(\mathbf{A}) &:= \max(|\sigma(\mathbf{A})|), \\ \lambda_{\min}(\mathbf{A}) &:= \min(|\sigma(\mathbf{A})|). \end{aligned} \quad (10.1.4.8)$$

other notation  $\kappa(\mathbf{A}) := \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})} \hat{=} \text{spectral condition number}$  of  $\mathbf{A}$

(for general  $\mathbf{A}$ :  $\lambda_{\max}(\mathbf{A})/\lambda_{\min}(\mathbf{A})$  largest/smallest eigenvalue *in modulus*)

These results are an immediate consequence of the fact that

$$\forall \mathbf{A} \in \mathbb{R}^{n,n}, \quad \mathbf{A}^\top = \mathbf{A} \quad \exists \mathbf{U} \in \mathbb{R}^{n,n}, \quad \mathbf{U}^{-1} = \mathbf{U}^\top: \quad \mathbf{U}^\top \mathbf{A} \mathbf{U} \quad \text{is diagonal,}$$

see (10.1.4.2), Cor. 9.1.0.9, [NS02, Thm. 7.8], [Gut09, Satz 9.15].

Please note that for general regular  $\mathbf{M} \in \mathbb{R}^{n,n}$  we *cannot* expect  $\text{cond}_2(\mathbf{M}) = \kappa(\mathbf{M})$ .  $\square$

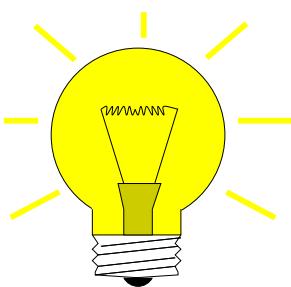
## 10.2 Conjugate gradient method (CG) [Han02, Ch. 9], [DR08, Sect. 13.4], [QSS00, Sect. 4.3.4]

Again we consider a linear system of equations  $\mathbf{Ax} = \mathbf{b}$  with s.p.d. ( $\rightarrow$  Def. 1.1.2.6) system matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$  and given  $\mathbf{b} \in \mathbb{R}^n$ .

Liability of gradient method of Section 10.1.3:

**NO MEMORY**

1D line search in § 10.1.3.3 is oblivious of former line searches, which rules out reuse of information gained in previous steps of the iteration. This is a typical drawback of 1-point iterative methods.



Idea:

Replace linear search with **subspace correction**

Given:

♦ initial guess  $\mathbf{x}^{(0)}$

♦ nested subspaces  $U_1 \subset U_2 \subset U_3 \subset \dots \subset U_n = \mathbb{R}^n$ ,  $\dim U_k = k$

$$\mathbf{x}^{(k)} := \underset{\mathbf{x} \in U_k + \mathbf{x}^{(0)}}{\operatorname{argmin}} J(\mathbf{x}), \quad (10.2.0.1)$$

quadratic functional from (10.1.1.2)

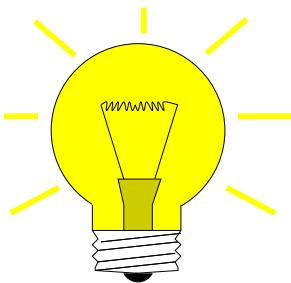
Note: Once the subspaces  $U_k$  and  $\mathbf{x}^{(0)}$  are fixed, the iteration (10.2.0.1) is well defined, because  $J|_{U_k + \mathbf{x}^{(0)}}$  always possess a unique minimizer.

Obvious (from Lemma 10.1.1.1):

$$\mathbf{x}^{(n)} = \mathbf{x}^* = \mathbf{A}^{-1}\mathbf{b}$$

Thanks to (10.1.1.3), definition (10.2.0.1) ensures:

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_A \leq \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_A$$



How to find suitable subspaces  $U_k$ ?

Idea:

$U_{k+1} \leftarrow U_k + \text{"local steepest descent direction"}$

given by  $-\operatorname{grad} J(\mathbf{x}^{(k)}) = \mathbf{b} - \mathbf{Ax}^{(k)} = \mathbf{r}_k$  (residual  $\rightarrow$  Def. 2.4.0.1)

$$U_{k+1} = \operatorname{Span}\{U_k, \mathbf{r}_k\}, \quad \mathbf{x}^{(k)} \text{ from (10.2.0.1)}. \quad (10.2.0.2)$$

Obvious:  $\mathbf{r}_k = 0 \Rightarrow \mathbf{x}^{(k)} = \mathbf{x}^* := \mathbf{A}^{-1}\mathbf{b}$  done ✓

**Lemma 10.2.0.3.**  $\mathbf{r}_k \perp U_k$

With  $\mathbf{x}^{(k)}$  according to (10.2.0.1),  $U_k$  from (10.2.0.2) the residual  $\mathbf{r}_k := \mathbf{b} - \mathbf{Ax}^{(k)}$  satisfies

$$\mathbf{r}_k^\top \mathbf{u} = 0 \quad \forall \mathbf{u} \in U_k \quad ("r_k \perp U_k").$$

Geometric consideration: since  $\mathbf{x}^{(k)}$  is the minimizer of  $J$  over the affine space  $U_k + \mathbf{x}^{(0)}$ , the projection of the steepest descent direction  $\mathbf{grad} J(\mathbf{x}^{(k)})$  onto  $U_k$  has to vanish:

$$\mathbf{x}^{(k)} := \underset{\mathbf{x} \in U_k + \mathbf{x}^{(0)}}{\operatorname{argmin}} J(x) \Rightarrow \mathbf{grad} J(\mathbf{x}^{(k)}) \perp U_k . \quad (10.2.0.4)$$

*Proof.* Consider

$$\psi(t) = J(\mathbf{x}^{(k)} + t\mathbf{u}) , \quad \mathbf{u} \in U_k, \quad t \in \mathbb{R} .$$

By (10.2.0.1),  $t \mapsto \psi(t)$  has a global minimum in  $t = 0$ , which implies

$$\frac{d\psi}{dt}(0) = \mathbf{grad} J(\mathbf{x}^{(k)})^\top \mathbf{u} = (\mathbf{A}\mathbf{x}^{(k)} - \mathbf{b})^\top \mathbf{u} = 0 .$$

Since  $\mathbf{u} \in U_k$  was arbitrary, the lemma is proved.  $\square$

### Corollary 10.2.0.5.

If  $\mathbf{r}_l \neq 0$  for  $l = 0, \dots, k$ ,  $k \leq n$ , then  $\{\mathbf{r}_0, \dots, \mathbf{r}_k\}$  is an *orthogonal basis* of  $U_k$ .

Lemma 10.2.0.3 also implies that, if  $U_0 = \{0\}$ , then  $\dim U_k = k$  as long as  $\mathbf{x}^{(k)} \neq \mathbf{x}^*$ , that is, before we have converged to the exact solution.

(10.2.0.1) and (10.2.0.2) define the **conjugate gradient method** (CG) for the iterative solution of  
 $\mathbf{Ax} = \mathbf{b}$

(hailed as a “**top ten algorithm**” of the 20th century, SIAM News, 33(4))

## 10.2.1 Krylov spaces

### Definition 10.2.1.1. Krylov space

For  $\mathbf{A} \in \mathbb{R}^{n,n}$ ,  $\mathbf{z} \in \mathbb{R}^n$ ,  $\mathbf{z} \neq 0$ , the  $l$ -th **Krylov space** is defined as

$$\mathcal{K}_l(\mathbf{A}, \mathbf{z}) := \operatorname{Span}\{\mathbf{z}, \mathbf{Az}, \dots, \mathbf{A}^{l-1}\mathbf{z}\} .$$

Equivalent definition:  $\mathcal{K}_l(\mathbf{A}, \mathbf{z}) = \{p(\mathbf{A})\mathbf{z}: p \text{ polynomial of degree } \leq l\}$

### Lemma 10.2.1.2.

The subspaces  $U_k \subset \mathbb{R}^n$ ,  $k \geq 1$ , defined by (10.2.0.1) and (10.2.0.2) satisfy

$$U_k = \operatorname{Span}\{\mathbf{r}_0, \mathbf{Ar}_0, \dots, \mathbf{A}^{k-1}\mathbf{r}_0\} = \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) ,$$

where  $\mathbf{r}_0 = \mathbf{b} - \mathbf{Ax}^{(0)}$  is the initial residual.

*Proof.* (by induction) Obviously  $\mathbf{A}\mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) \subset \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{r}_0)$ . In addition

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}(\mathbf{x}^{(0)} + \mathbf{z}) \quad \text{for some } \mathbf{z} \in U_k \Rightarrow \mathbf{r}_k = \underbrace{\mathbf{r}_0}_{\in \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{r}_0)} - \underbrace{\mathbf{Az}}_{\in \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{r}_0)} .$$

Since  $U_{k+1} = \text{Span}\{U_k, \mathbf{r}_k\}$ , we obtain  $U_{k+1} \subset \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{r}_0)$ . Dimensional considerations based on Lemma 10.2.0.3 finish the proof.  $\square$

## 10.2.2 Implementation of CG

Assume: basis  $\{\mathbf{p}_1, \dots, \mathbf{p}_l\}$ ,  $l = 1, \dots, n$ , of  $\mathcal{K}_l(\mathbf{A}, \mathbf{r})$  available

$$\mathbf{x}^{(l)} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0) \Rightarrow \text{set } \mathbf{x}^{(l)} = \mathbf{x}^{(0)} + \gamma_1 \mathbf{p}_1 + \dots + \gamma_l \mathbf{p}_l.$$

For  $\psi(\gamma_1, \dots, \gamma_l) := J(\mathbf{x}^{(0)} + \gamma_1 \mathbf{p}_1 + \dots + \gamma_l \mathbf{p}_l)$  holds

$$(10.2.0.1) \Leftrightarrow \frac{\partial \psi}{\partial \gamma_j} = 0, \quad j = 1, \dots, l.$$

This leads to a linear system of equations by which the coefficients  $\gamma_j$  can be computed:

$$\begin{bmatrix} \mathbf{p}_1^\top \mathbf{A} \mathbf{p}_1 & \cdots & \mathbf{p}_1^\top \mathbf{A} \mathbf{p}_l \\ \vdots & \ddots & \vdots \\ \mathbf{p}_l^\top \mathbf{A} \mathbf{p}_1 & \cdots & \mathbf{p}_l^\top \mathbf{A} \mathbf{p}_l \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_l \end{bmatrix} = \begin{bmatrix} \mathbf{p}_1^\top \mathbf{r} \\ \vdots \\ \mathbf{p}_l^\top \mathbf{r} \end{bmatrix}, \quad \mathbf{r} := \mathbf{b} - \mathbf{A} \mathbf{x}^{(0)}. \quad (10.2.2.1)$$

Great simplification, if  $\{\mathbf{p}_1, \dots, \mathbf{p}_l\}$  **A-orthogonal basis** of  $\mathcal{K}_l(\mathbf{A}, \mathbf{r})$ :  $\mathbf{p}_j^\top \mathbf{A} \mathbf{p}_i = 0$  for  $i \neq j$ .

Recall: s.p.d.  $\mathbf{A}$  induces an inner product  $\geq$  concept of orthogonality [NS02, Sect. 4.4], [Gut09, Sect. 6.2]. “ $\mathbf{A}$ -geometry” like standard Euclidean space.

Assume: **A-orthogonal basis**  $\{\mathbf{p}_1, \dots, \mathbf{p}_n\}$  of  $\mathbb{R}^n$  available, such that

$$\text{Span}\{\mathbf{p}_1, \dots, \mathbf{p}_l\} = \mathcal{K}_l(\mathbf{A}, \mathbf{r}).$$

► (Efficient) successive computation of  $\mathbf{x}^{(l)}$  becomes possible, see [DR08, Lemma 13.24]  
(LSE (10.2.2.1) becomes diagonal !)

**Input:** : initial guess  $\mathbf{x}^{(0)} \in \mathbb{R}^n$

**Given:** : **A-orthogonal bases**  $\{\mathbf{p}_1, \dots, \mathbf{p}_l\}$  of  $\mathcal{K}_l(\mathbf{A}, \mathbf{r}_0)$ ,  $l = 1, \dots, n$

**Output:** : approximate solution  $\mathbf{x}^{(l)} \in \mathbb{R}^n$  of  $\mathbf{Ax} = \mathbf{b}$

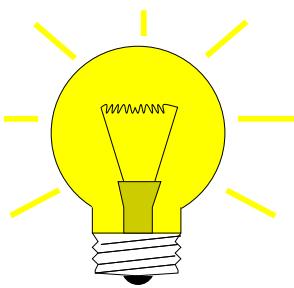
$$\begin{aligned} \mathbf{r}_0 &:= \mathbf{b} - \mathbf{Ax}^{(0)}; \\ \text{for } j = 1 \text{ to } l \text{ do } \{ &\quad \mathbf{x}^{(j)} := \mathbf{x}^{(j-1)} + \frac{\mathbf{p}_j^\top \mathbf{r}_0}{\mathbf{p}_j^\top \mathbf{A} \mathbf{p}_j} \mathbf{p}_j \quad \} \end{aligned} \quad (10.2.2.2)$$

**Task:** Efficient computation of **A-orthogonal vectors**  $\{\mathbf{p}_1, \dots, \mathbf{p}_l\}$  spanning  $\mathcal{K}_l(\mathbf{A}, \mathbf{r}_0)$   
during the CG iteration.  
**A-orthogonalities/orthogonalities**  $\Rightarrow$  short recursions

Lemma 10.2.0.3 implies orthogonality  $\mathbf{p}_j \perp \mathbf{r}_m := \mathbf{b} - \mathbf{Ax}^{(m)}$ ,  $1 \leq j \leq m \leq l$ . Also by **A-orthogonality** of the  $\mathbf{p}_k$

$$\mathbf{p}_j^\top (\mathbf{b} - \mathbf{Ax}^{(m)}) = \mathbf{p}_j^\top \left( \underbrace{\mathbf{b} - \mathbf{Ax}^{(0)}}_{= \mathbf{r}_0} - \sum_{k=1}^m \frac{\mathbf{p}_k^\top \mathbf{r}_0}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k} \mathbf{A} \mathbf{p}_k \right) = 0. \quad (10.2.2.3)$$

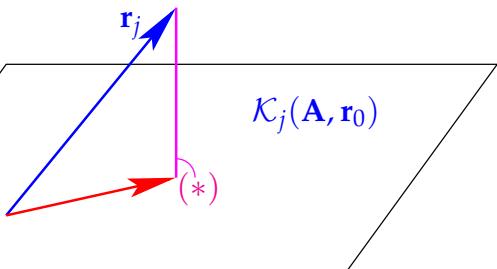
From linear algebra we already know a way to construct orthogonal basis vectors:

(10.2.2.3)  $\Rightarrow$  Idea:

Gram-Schmidt orthogonalization [NS02, Thm. 4.8], [Gut09, Alg. 6.1] of residuals  $\mathbf{r}_j := \mathbf{b} - \mathbf{Ax}^{(j)}$  w.r.t.  $\mathbf{A}$ -inner product:

$$\mathbf{p}_1 := \mathbf{r}_0, \mathbf{p}_{j+1} := (\underbrace{\mathbf{b} - \mathbf{Ax}^{(j)}}_{=:\mathbf{r}_j}) - \underbrace{\sum_{k=1}^j \frac{\mathbf{p}_k^\top \mathbf{A} \mathbf{r}_j}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k} \mathbf{p}_k}_{(*)}, \quad j = 1, \dots, l-1 .$$
(10.2.2.4)

Fig. 354



Geometric interpretation of (10.2.2.4):

$(*) \hat{=} \text{orthogonal projection of } \mathbf{r}_j \text{ on the subspace } \text{Span}\{\mathbf{p}_1, \dots, \mathbf{p}_j\}$

### Lemma 10.2.2.5. Bases for Krylov spaces in CG

If they do not vanish, the vectors  $\mathbf{p}_j$ ,  $1 \leq j \leq l$ , and  $\mathbf{r}_j := \mathbf{b} - \mathbf{Ax}^{(j)}$ ,  $0 \leq j \leq l$ , from (10.2.2.2), (10.2.2.4) satisfy

- (i)  $\{\mathbf{p}_1, \dots, \mathbf{p}_j\}$  is  $\mathbf{A}$ -orthogonal basis von  $\mathcal{K}_j(\mathbf{A}, \mathbf{r}_0)$ ,
- (ii)  $\{\mathbf{r}_0, \dots, \mathbf{r}_{j-1}\}$  is orthogonal basis of  $\mathcal{K}_j(\mathbf{A}, \mathbf{r}_0)$ , cf. Cor. 10.2.0.5

*Proof.*  $\mathbf{A}$ -orthogonality of  $\mathbf{p}_j$  by construction, study (10.2.2.4).

$$\begin{aligned} (10.2.2.2) \& (10.2.2.4) \Rightarrow \mathbf{p}_{j+1} = \mathbf{r}_0 - \sum_{k=1}^j \frac{\mathbf{p}_k^\top \mathbf{r}_0}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k} \mathbf{A} \mathbf{p}_k - \sum_{k=1}^j \frac{\mathbf{p}_k^\top \mathbf{A} \mathbf{r}_j}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k} \mathbf{p}_k . \\ & \Rightarrow \mathbf{p}_{j+1} \in \text{Span}\{\mathbf{r}_0, \mathbf{p}_1, \dots, \mathbf{p}_j, \mathbf{A} \mathbf{p}_1, \dots, \mathbf{A} \mathbf{p}_j\} . \end{aligned}$$

A simple induction argument confirms (i)

$$(10.2.2.4) \Rightarrow \mathbf{r}_j \in \text{Span}\{\mathbf{p}_1, \dots, \mathbf{p}_{j+1}\} \quad \& \quad \mathbf{p}_j \in \text{Span}\{\mathbf{r}_0, \dots, \mathbf{r}_{j-1}\} . \quad (10.2.2.6)$$

$$\blacktriangleright \boxed{\text{Span}\{\mathbf{p}_1, \dots, \mathbf{p}_j\} = \text{Span}\{\mathbf{r}_0, \dots, \mathbf{r}_{j-1}\} = \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0)} . \quad (10.2.2.7)$$

$$(10.2.2.3) \Rightarrow \mathbf{r}_j \perp \text{Span}\{\mathbf{p}_1, \dots, \mathbf{p}_j\} = \text{Span}\{\mathbf{r}_0, \dots, \mathbf{r}_{j-1}\} . \quad (10.2.2.8)$$

□

Orthogonalities from Lemma 10.2.2.5  $\blacktriangleright$  short recursions for  $\mathbf{p}_k$ ,  $\mathbf{r}_k$ ,  $\mathbf{x}^{(k)}$ !

$$(10.2.2.3) \Rightarrow (10.2.2.4) \text{ collapses to } \mathbf{p}_{j+1} := \mathbf{r}_j - \frac{\mathbf{p}_j^\top \mathbf{A} \mathbf{r}_j}{\mathbf{p}_j^\top \mathbf{A} \mathbf{p}_j} \mathbf{p}_j , \quad j = 1, \dots, l .$$

recursion for residuals:

$$(10.2.2.2) \blacktriangleright \mathbf{r}_j = \mathbf{r}_{j-1} - \frac{\mathbf{p}_j^\top \mathbf{r}_0}{\mathbf{p}_j^\top \mathbf{A} \mathbf{p}_j} \mathbf{A} \mathbf{p}_j .$$

$$\text{Lemma 10.2.2.5, (i)} \quad \Rightarrow \quad \mathbf{r}_{j-1}^H \mathbf{p}_j = \left( \mathbf{r}_0 + \sum_{k=1}^{m-1} \frac{\mathbf{r}_0^\top \mathbf{p}_k}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k} \mathbf{A} \mathbf{p}_k \right)^T \mathbf{p}_j = \mathbf{r}_0^\top \mathbf{p}_j. \quad (10.2.2.9)$$

The orthogonality (10.2.2.9) together with (10.2.2.8) permits us to replace  $\mathbf{r}_0$  with  $\mathbf{r}_{j-1}$  in the actual implementation.

### §10.2.2.10 (CG method for solving $\mathbf{Ax} = \mathbf{b}$ , A s.p.d. $\rightarrow$ [DR08, Alg. 13.27])

Input : initial guess  $\mathbf{x}^{(0)} \in \mathbb{R}^n$

Output : approximate solution  $\mathbf{x}^{(l)} \in \mathbb{R}^n$

```

 $\mathbf{p}_1 = \mathbf{r}_0 := \mathbf{b} - \mathbf{Ax}^{(0)}$ ;
 $\text{for } j = 1 \text{ to } l \text{ do } \{$ 
     $\mathbf{x}^{(j)} := \mathbf{x}^{(j-1)} + \frac{\mathbf{p}_j^\top \mathbf{r}_{j-1}}{\mathbf{p}_j^\top \mathbf{Ap}_j} \mathbf{p}_j;$ 
     $\mathbf{r}_j = \mathbf{r}_{j-1} - \frac{\mathbf{p}_j^\top \mathbf{r}_{j-1}}{\mathbf{p}_j^\top \mathbf{Ap}_j} \mathbf{Ap}_j;$ 
     $\mathbf{p}_{j+1} = \mathbf{r}_j - \frac{(\mathbf{Ap}_j)^\top \mathbf{r}_j}{\mathbf{p}_j^\top \mathbf{Ap}_j} \mathbf{p}_j;$ 
 $\}$ 

```

Input: initial guess  $\mathbf{x} \doteq \mathbf{x}^{(0)} \in \mathbb{R}^n$

tolerance  $\tau > 0$

Output: approximate solution  $\mathbf{x} \doteq \mathbf{x}^{(l)}$

```

 $\mathbf{p} := \mathbf{r}_0 := \mathbf{r} := \mathbf{b} - \mathbf{Ax};$ 
 $\text{for } j = 1 \text{ to } l_{\max} \text{ do } \{$ 
     $\beta := \mathbf{r}^\top \mathbf{r};$ 
     $\mathbf{h} := \mathbf{Ap};$ 
     $\alpha := \frac{\beta}{\mathbf{p}^\top \mathbf{h}};$ 
     $\mathbf{x} := \mathbf{x} + \alpha \mathbf{p};$ 
     $\mathbf{r} := \mathbf{r} - \alpha \mathbf{h};$ 
     $\text{if } \|\mathbf{r}\| \leq \tau \|\mathbf{r}_0\| \text{ then stop};$ 
     $\beta := \frac{\mathbf{r}^\top \mathbf{r}}{\beta};$ 
     $\mathbf{p} := \mathbf{r} + \beta \mathbf{p};$ 
 $\}$ 

```

In CG algorithm:  $\mathbf{r}_j = \mathbf{b} - \mathbf{Ax}^{(k)}$  agrees with the residual associated with the current iterate (in exact arithmetic, cf. Ex. 10.2.3.1), but computation through short recursion is more efficient.

➤ We find that the CG method possesses all the algorithmic advantages of the gradient method, cf. the discussion in Section 10.1.3.

1 matrix  $\times$  vector product, 3 dot products, 3 AXPY-operations per step:  
If  $\mathbf{A}$  sparse,  $\text{nnz}(\mathbf{A}) \sim n$  ➤ computational effort  $O(n)$  per step

### MATLAB-code 10.2.2.11: basic CG iteration for solving $\mathbf{Ax} = \mathbf{b}$ , § 10.2.2.10

```

1 function x = cg(evalA,b,x,tol,maxit)
2 % x supplies initial guess, maxit maximal number of CG steps
3 % evalA must pass a handle to a MATLAB function realizing A*x
4 r = b - evalA(x); rho = 1; n0 = norm(r);
5 for i = 1 : maxit
6     rho1 = rho; rho = r' * r;
7     if (i == 1), p = r;
8     else beta = rho/rho1; p = r + beta * p; end
9     q = evalA(p); alpha = rho/(p' * q);
10    x = x + alpha * p;      % update of approximate solution
11    if (norm(b-A*x) <= tol*n0) return; end % termination, see
12        Rem. 10.2.2.12
13    r = r - alpha * q;       % update of residual

```

13 | end

MATLAB-function:

`x=pcg(A,b,tol,maxit,[],[],x0)` : Solve  $\mathbf{Ax} = \mathbf{b}$  with at most maxit CG steps: stop, when  $\|\mathbf{r}_l\| : \|\mathbf{r}_0\| < \text{tol}$ .  
`x=pcg(Afun,b,tol,maxit,[],[],x0)` : Afun = handle to function for computing  $\mathbf{A} \times \text{vector}$ .  
`[x,flag,relr,it,resv] = pcg(...)` : diagnostic information about iteration

**Remark 10.2.2.12 (A posteriori termination criterion for plain CG)**

For any vector norm and associated matrix norm ( $\rightarrow$  Def. 1.5.5.10) hold (with residual  $\mathbf{r}_l := \mathbf{b} - \mathbf{Ax}^{(l)}$ )

$$\frac{1}{\text{cond}(\mathbf{A})} \frac{\|\mathbf{r}_l\|}{\|\mathbf{r}_0\|} \leq \frac{\|\mathbf{x}^{(l)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{r}_l\|}{\|\mathbf{r}_0\|}. \quad (10.2.2.13)$$

↑  
relative decrease of iteration error

(10.2.2.13) can easily be deduced from the error equation  $\mathbf{A}(\mathbf{x}^{(k)} - \mathbf{x}^*) = \mathbf{r}_k$ , see Def. 2.4.0.1 and 2.4.0.12.  $\square$

**10.2.3 Convergence of CG**

Note: CG is a *direct solver*, because (in exact arithmetic)  $\mathbf{x}^{(k)} = \mathbf{x}^*$  for some  $k \leq n$

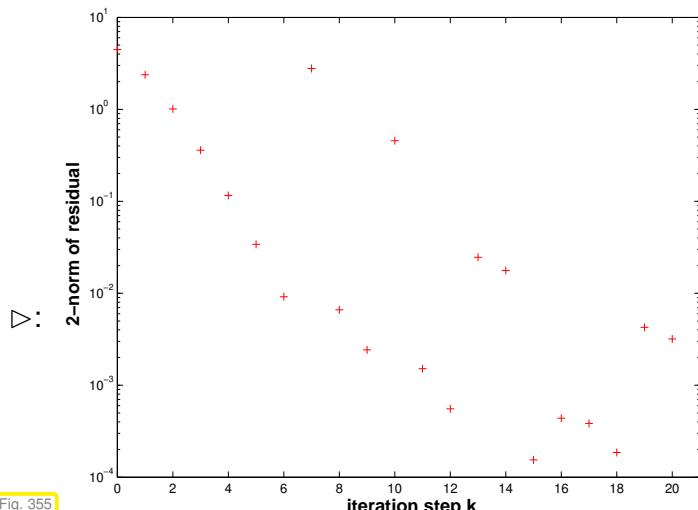
**EXAMPLE 10.2.3.1 (Impact of roundoff errors on CG  $\rightarrow$  [QSS00, Rem. 4.3])**

Numerical experiment:  $\mathbf{A} = \text{hilb}(20)$ ,  
 $\mathbf{x}^{(0)} = \mathbf{0}$ ,  $\mathbf{b} = [1, \dots, 1]^T$

Hilbert-Matrix: extremely ill-conditioned

residual norms during CG iteration

$$\mathbf{R} = [\mathbf{r}_0, \dots, \mathbf{r}^{(10)}]$$



$$\mathbf{R}^\top \mathbf{R} =$$

1.000000	-0.000000	0.000000	-0.000000	0.000000	-0.000000	0.016019	-0.795816	-0.430569	0.348133
-0.000000	1.000000	-0.000000	0.000000	-0.000000	0.000000	-0.012075	0.600068	-0.520610	0.420903
0.000000	-0.000000	1.000000	-0.000000	0.000000	-0.000000	0.001582	-0.078664	0.384453	-0.310577
-0.000000	0.000000	-0.000000	1.000000	-0.000000	0.000000	-0.000024	0.001218	-0.024115	0.019394
0.000000	-0.000000	0.000000	-0.000000	1.000000	-0.000000	0.000000	-0.000002	0.000151	-0.000118
-0.000000	0.000000	-0.000000	0.000000	-0.000000	1.000000	-0.000000	0.000000	-0.000000	0.000000
0.016019	-0.012075	0.001582	-0.000024	0.000000	-0.000000	1.000000	-0.000000	-0.000000	0.000000
-0.795816	0.600068	-0.078664	0.001218	-0.000002	0.000000	-0.000000	1.000000	0.000000	-0.000000
-0.430569	-0.520610	0.384453	-0.024115	0.000151	-0.000000	-0.000000	0.000000	1.000000	0.000000
0.348133	0.420903	-0.310577	0.019394	-0.000118	0.000000	0.000000	-0.000000	0.000000	1.000000

## ➤ Roundoff

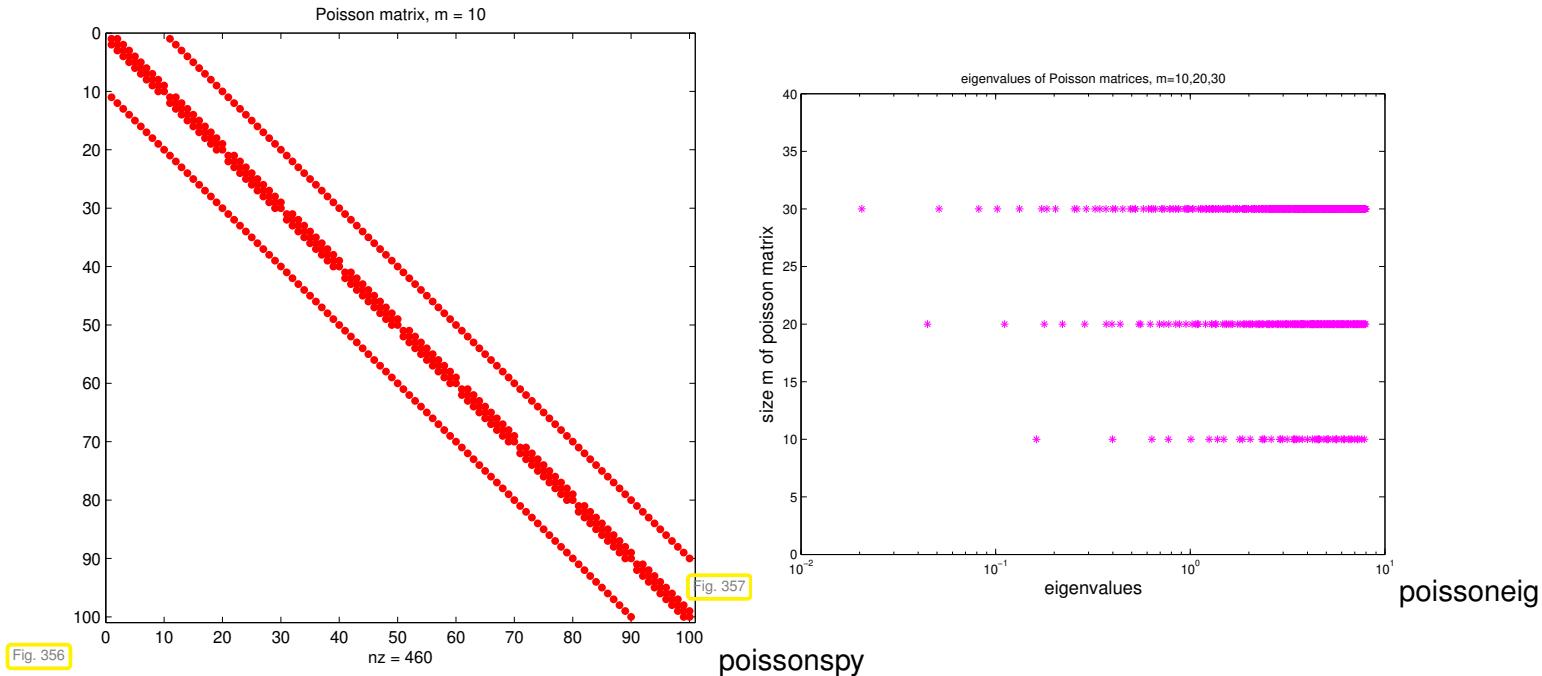
- ◆ destroys orthogonality of residuals
- ◆ prevents computation of exact solution after  $n$  steps.

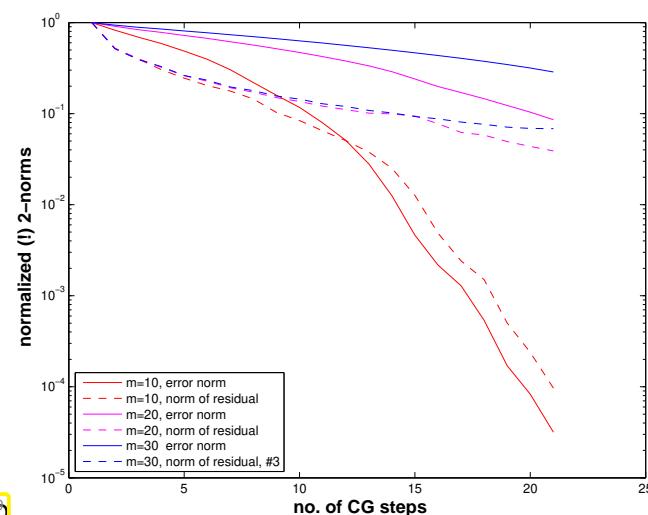
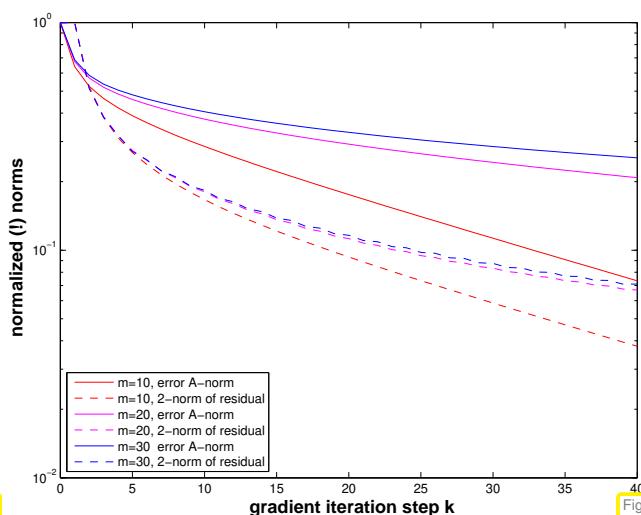
## ➤ Numerical instability (→ Def. 1.5.5.19) ➤ pointless to (try to) use CG as direct solver!

Practice: CG used for large  $n$  as *iterative solver*:  $\mathbf{x}^{(k)}$  for some  $k \ll n$  is expected to provide good approximation for  $\mathbf{x}^*$

**EXAMPLE 10.2.3.2 (Convergence of CG as iterative solver)** CG (Code 10.2.2.11) & gradient method (Code 10.1.3.4) for LSE with sparse s.p.d. “Poisson matrix”

```
A = gallery('poisson', m); x0 = (1:n)'; b = zeros(n, 1);
```

➤  $\mathbf{A} \in \mathbb{R}^{m^2, m^2}$ 



Observations:

- CG much faster than gradient method (as expected, because it has “memory”)
- Both, CG and gradient method converge more slowly for larger sizes of Poisson matrices.

**Convergence theory:** [Hac94, Sect. 9.4.3]

A simple consequence of (10.1.1.3) and (10.2.0.1):

**Corollary 10.2.3.3. “Optimality” of CG iterates**

Writing  $\mathbf{x}^* \in \mathbb{R}^n$  for the exact solution of  $\mathbf{Ax} = \mathbf{b}$  the CG iterates satisfy

$$\left\| \mathbf{x}^* - \mathbf{x}^{(l)} \right\|_A = \min \{ \|\mathbf{y} - \mathbf{x}^*\|_A : \mathbf{y} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0) \} , \quad \mathbf{r}_0 := \mathbf{b} - \mathbf{Ax}^{(0)} .$$

This paves the way for a quantitative convergence estimate:

$$\mathbf{y} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}) \Leftrightarrow \mathbf{y} = \mathbf{x}^{(0)} + \mathbf{A} p(\mathbf{A})(\mathbf{x} - \mathbf{x}^{(0)}) , \quad p = \text{polynomial of degree } \leq l-1 .$$

$$\blacktriangleright \quad \mathbf{x} - \mathbf{y} = q(\mathbf{A})(\mathbf{x} - \mathbf{x}^{(0)}) , \quad q = \text{polynomial of degree } \leq l , \quad q(0) = 1 .$$



$$\left\| \mathbf{x} - \mathbf{x}^{(l)} \right\|_A \leq \boxed{\min \left\{ \max_{\lambda \in \sigma(\mathbf{A})} |q(\lambda)| : q \text{ polynomial of degree } \leq l , \quad q(0) = 1 \right\}} \cdot \left\| \mathbf{x} - \mathbf{x}^{(0)} \right\|_A . \quad (10.2.3.4)$$

Bound this minimum for  $\lambda \in [\lambda_{\min}(\mathbf{A}), \lambda_{\max}(\mathbf{A})]$  by using suitable “polynomial candidates”

Tool: **Chebychev polynomials** ( $\rightarrow$  Section 6.1.3.1)  $\Rightarrow$  lead to the following estimate [Hac91, Satz 9.4.2], [DR08, Satz 13.29]

### Theorem 10.2.3.5. Convergence of CG method

The iterates of the CG method for solving  $\mathbf{Ax} = \mathbf{b}$  (see Code 10.2.2.11) with  $\mathbf{A} = \mathbf{A}^\top$  s.p.d. satisfy

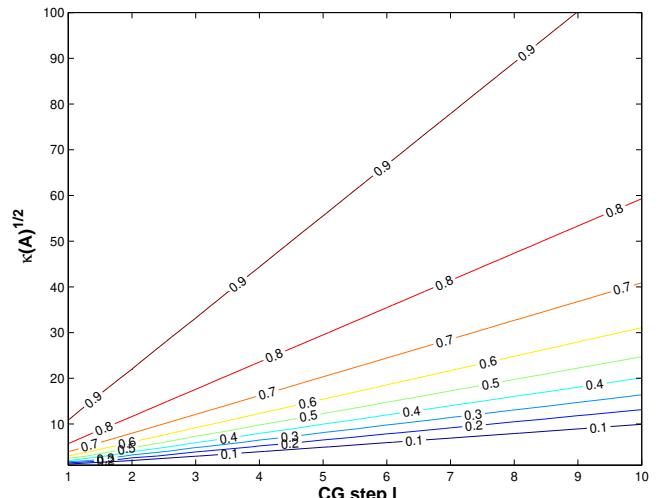
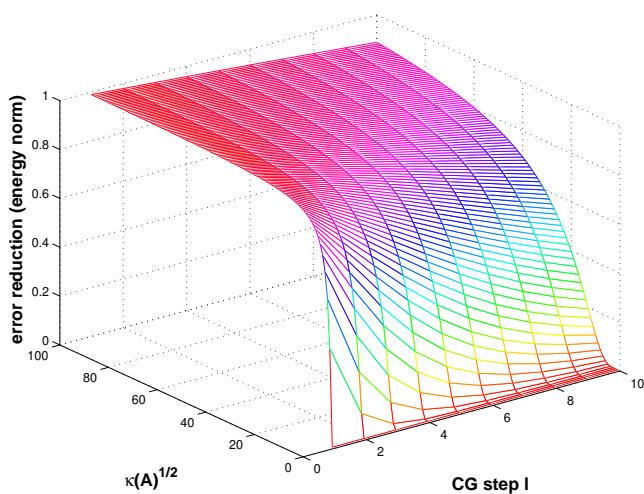
$$\begin{aligned}\|\mathbf{x} - \mathbf{x}^{(l)}\|_A &\leq \frac{2\left(1 - \frac{1}{\sqrt{\kappa(\mathbf{A})}}\right)^l}{\left(1 + \frac{1}{\sqrt{\kappa(\mathbf{A})}}\right)^{2l} + \left(1 - \frac{1}{\sqrt{\kappa(\mathbf{A})}}\right)^{2l}} \|\mathbf{x} - \mathbf{x}^{(0)}\|_A \\ &\leq 2\left(\frac{\sqrt{\kappa(\mathbf{A})} - 1}{\sqrt{\kappa(\mathbf{A})} + 1}\right)^l \|\mathbf{x} - \mathbf{x}^{(0)}\|_A.\end{aligned}$$

(recall:  $\kappa(\mathbf{A})$  = spectral condition number of  $\mathbf{A}$ ,  $\kappa(\mathbf{A}) = \text{cond}_2(\mathbf{A})$ )

The estimate of this theorem confirms *asymptotic linear convergence* of the CG method ( $\rightarrow$  Def. 8.1.1.1)

with a rate of  $\frac{\sqrt{\kappa(\mathbf{A})} - 1}{\sqrt{\kappa(\mathbf{A})} + 1}$

Plots of bounds for error reduction (in energy norm) during CG iteration from Thm. 10.2.3.5:



### MATLAB-code 10.2.3.6: plotting theoretical bounds for CG convergence rate

```

1 function plottheorate
2 % Python: ..../PYTHON/ConjugateGradient.py plottheorate()
3
4 [X,Y] = meshgrid(1:10,1:100); R = zeros(100,10);
5 for I=1:100
6 t = 1/I;
7 for j=1:10
8 R(I,j) = 2*(1-t)^j / ((1+t)^(2*j)+(1-t)^(2*j));
9 end
10 end
11
12 figure; view([-45,28]); mesh(X,Y,R); colormap hsv;
13 xlabel('{\bf CG step 1}', 'Fontsize', 14);
14 ylabel('{\bf \kappa(A)^{1/2}}', 'Fontsize', 14);

```

```

15 zlabel('{\bf error reduction (energy norm)}','Fontsize',14);
16
17 print -depsc2 '../PICTURES/theorate1.eps';
18
19 figure; [C,h] = contour(X,Y,R); clabel(C,h);
20 xlabel('{\bf CG step 1}','Fontsize',14);
21 ylabel('{\bf \kappa(A)^{1/2}}','Fontsize',14);
22
23 print -depsc2 '../PICTURES/theorate2.eps';

```

### EXAMPLE 10.2.3.7 (Convergence rates for CG method)

#### MATLAB-code 10.2.3.8: CG for Poisson matrix

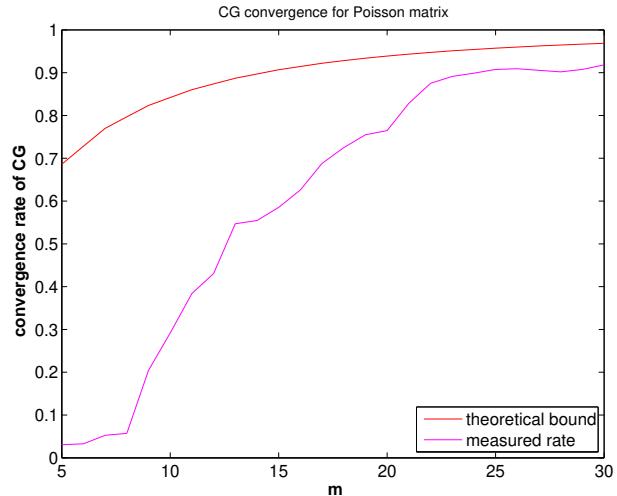
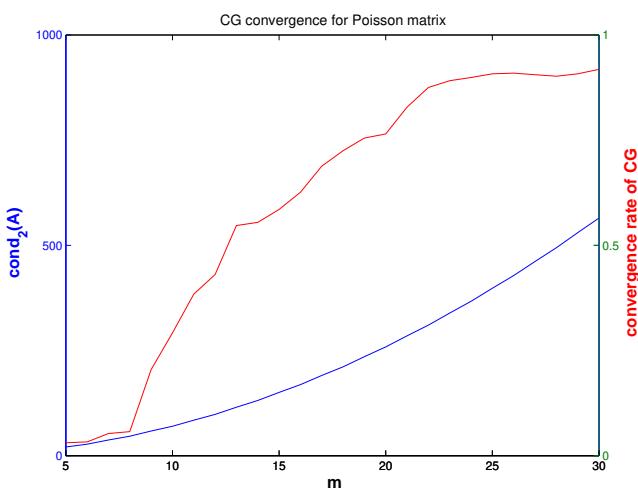
```

1 A = gallery ('poisson',m); n =
  size (A,1);
2 x0 = (1:n)'; b = ones(n,1); maxit
  30; tol =0;
3 [x,flag ,relres,iter,resvec] =
  pcg (A,b,tol,maxit,[],[],x0);

```

Measurement of rate of (linear) convergence:

$$\text{rate} \approx \sqrt[10]{\frac{\|\mathbf{r}_{30}\|_2}{\|\mathbf{r}_{20}\|_2}}. \quad (10.2.3.9)$$



Justification for estimating the rate of linear convergence ( $\rightarrow$  Def. 8.1.1.1) of  $\|\mathbf{r}_k\|_2 \rightarrow 0$ :

$$\|\mathbf{r}_{k+1}\|_2 \approx L\|\mathbf{r}_k\|_2 \Rightarrow \|\mathbf{r}_{k+m}\|_2 \approx L^m \|\mathbf{r}_k\|_2.$$

### EXAMPLE 10.2.3.10 (CG convergence and spectrum $\rightarrow$ Ex. 10.1.4.4)

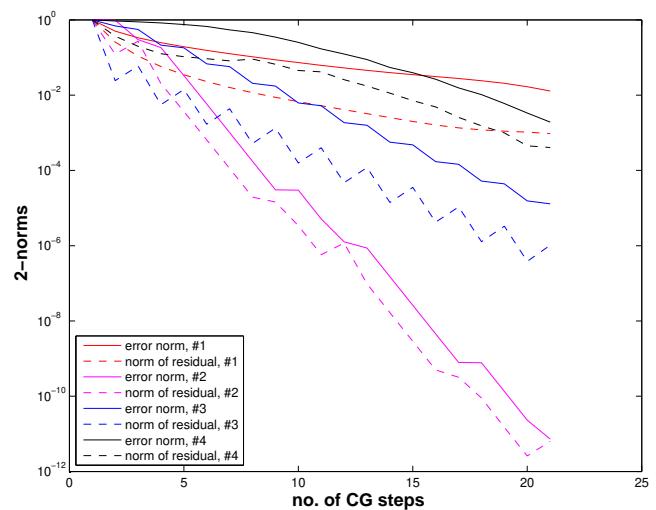
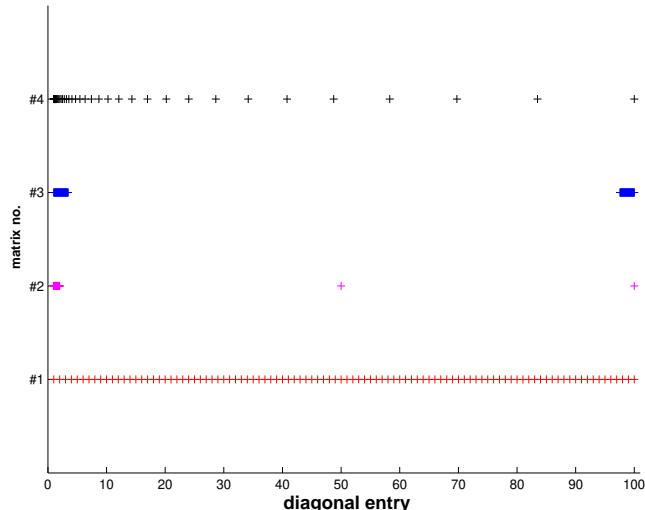
Test matrix #1:  $A=\text{diag}(d); d = (1:100);$

Test matrix #2:  $A=\text{diag}(d); d = [1+(0:97)/97, 50, 100];$

Test matrix #3:  $A=\text{diag}(d); d = [1+(0:49)*0.05, 100-(0:49)*0.05];$

Test matrix #4: eigenvalues exponentially dense at 1

```
x0 = cos((1:n)'); b = zeros(n,1);
```



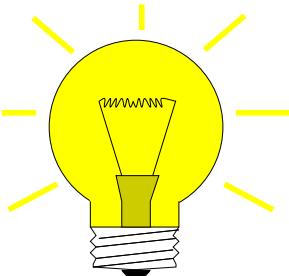
**Observations:** Distribution of eigenvalues has crucial impact on convergence of CG  
 (This is clear from the convergence theory, because detailed information about the spectrum allows a much better choice of “candidate polynomial” in (10.2.3.4) than merely using Chebychev polynomials)

- Clustering of eigenvalues leads to faster convergence of CG  
 (in stark contrast to the behavior of the gradient method, see Ex. 10.1.4.4)

CG convergence boosted by clustering of eigenvalues

## 10.3 Preconditioning [DR08, Sect. 13.5], [Han02, Ch. 10], [QSS00, Sect. 4.3.5]

Thm. 10.2.3.5 ➤ (Potentially) slow convergence of CG in case  $\kappa(\mathbf{A}) \gg 1$ .



Idea:

Preconditioning

Apply CG method to transformed linear system

$$\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}, \quad \tilde{\mathbf{A}} := \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}, \quad \tilde{\mathbf{x}} := \mathbf{B}^{1/2}\mathbf{x}, \quad \tilde{\mathbf{b}} := \mathbf{B}^{-1/2}\mathbf{b}, \quad (10.3.0.1)$$

with “small”  $\kappa(\tilde{\mathbf{A}})$ ,  $\mathbf{B} = \mathbf{B}^\top \in \mathbb{R}^{N,N}$  s.p.d.  $\hat{=}$  preconditioner.

### Remark 10.3.0.2 (Square root of a s.p.d. matrix)

What is meant by the “square root”  $\mathbf{B}^{1/2}$  of a s.p.d. matrix  $\mathbf{B}$ ?

Recall (10.1.4.2) : for every  $\mathbf{B} \in \mathbb{R}^{n,n}$  with  $\mathbf{B}^\top = \mathbf{B}$  there is an orthogonal matrix  $\mathbf{Q} \in \mathbb{R}^{n,n}$  such that  $\mathbf{B} = \mathbf{Q}^\top \mathbf{D} \mathbf{Q}$  with a diagonal matrix  $\mathbf{D}$  ( $\rightarrow$  Cor. 9.1.0.9, [NS02, Thm. 7.8], [LS09, Satz 9.15]). If  $\mathbf{B}$  is s.p.d. the (diagonal) entries of  $\mathbf{D}$  are strictly positive and we can define

$$\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n), \quad \lambda_i > 0 \quad \Rightarrow \quad \mathbf{D}^{1/2} := \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}).$$

This is generalized to

$$\mathbf{B}^{1/2} := \mathbf{Q}^\top \mathbf{D}^{1/2} \mathbf{Q},$$

and one easily verifies, using  $\mathbf{Q}^\top = \mathbf{Q}^{-1}$ , that  $(\mathbf{B}^{1/2})^2 = \mathbf{B}$  and that  $\mathbf{B}^{1/2}$  is s.p.d. In fact, these two requirements already determine  $\mathbf{B}^{1/2}$  uniquely:

$\mathbf{B}^{1/2}$  is the unique s.p.d. matrix such that  $(\mathbf{B}^{1/2})^2 = \mathbf{B}$ .

J

### Notion 10.3.0.3. Preconditioner

A s.p.d. matrix  $\mathbf{B} \in \mathbb{R}^{n,n}$  is called a **preconditioner** (ger.: Vorkonditionierer) for the s.p.d. matrix  $\mathbf{A} \in \mathbb{R}^{n,n}$ , if

1.  $\kappa(\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2})$  is “small” and
2. the evaluation of  $\mathbf{B}^{-1}\mathbf{x}$  is about as expensive (in terms of elementary operations) as the matrix  $\times$  vector multiplication  $\mathbf{A}\mathbf{x}$ ,  $\mathbf{x} \in \mathbb{R}^n$ .

Recall: spectral condition number  $\kappa(\mathbf{A}) := \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})}$ , see (10.1.4.8)

There are several equivalent ways to express that  $\kappa(\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2})$  is “small”:

- $\kappa(\mathbf{B}^{-1}\mathbf{A})$  is “small”,  
because spectra agree  $\sigma(\mathbf{B}^{-1}\mathbf{A}) = \sigma(\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2})$  due to similarity ( $\rightarrow$  Lemma 9.1.0.6)
- $\exists 0 < \gamma < \Gamma$ ,  $\Gamma/\gamma$  “small”:  $\gamma(\mathbf{x}^\top \mathbf{B} \mathbf{x}) \leq \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq \Gamma(\mathbf{x}^\top \mathbf{B} \mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{R}^n$ ,  
where equivalence is seen by transforming  $\mathbf{y} := \mathbf{B}^{-1/2}\mathbf{x}$  and appealing to the min-max Thm. 9.3.2.18.

“Reader’s digest” version of Notion 10.3.0.3:

S.p.d.  $\mathbf{B}$  preconditioner  $\Leftrightarrow \mathbf{B}^{-1}$  = cheap approximate inverse of  $\mathbf{A}$

Problem:  $\mathbf{B}^{1/2}$ , which occurs prominently in (10.3.0.1) is usually not available with acceptable computational costs.

However, if one formally applies § 10.2.2.10 to the transformed system

$$\tilde{\mathbf{A}}\tilde{\mathbf{x}} := (\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2})(\mathbf{B}^{1/2}\mathbf{x}) = \tilde{\mathbf{b}} := \mathbf{B}^{-1/2}\mathbf{b},$$

from (10.3.0.1), it becomes apparent that, after suitable transformation of the iteration variables  $\mathbf{p}_j$  and  $\mathbf{r}_j$ ,  $\mathbf{B}^{1/2}$  and  $\mathbf{B}^{-1/2}$  invariably occur in products  $\mathbf{B}^{-1/2}\mathbf{B}^{-1/2} = \mathbf{B}^{-1}$  and  $\mathbf{B}^{1/2}\mathbf{B}^{-1/2} = \mathbf{I}$ . Thus, thanks to this **intrinsic transformation** square roots of  $\mathbf{B}$  are not required for the implementation!

CG for  $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$ 

Input : initial guess  $\tilde{\mathbf{x}}^{(0)} \in \mathbb{R}^n$   
 Output : approximate solution  $\tilde{\mathbf{x}}^{(l)} \in \mathbb{R}^n$

```

 $\tilde{\mathbf{p}}_1 := \tilde{\mathbf{r}}_0 := \tilde{\mathbf{b}} - \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{x}}^{(0)}$ ;
for  $j = 1$  to  $l$  do {
   $\alpha := \frac{\tilde{\mathbf{p}}_j^T \tilde{\mathbf{r}}_{j-1}}{\tilde{\mathbf{p}}_j^T \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j}$ 
   $\tilde{\mathbf{x}}^{(j)} := \tilde{\mathbf{x}}^{(j-1)} + \alpha \tilde{\mathbf{p}}_j$ ;
   $\tilde{\mathbf{r}}_j = \tilde{\mathbf{r}}_{j-1} - \alpha \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{1/2}\tilde{\mathbf{p}}_j$ ;
   $\tilde{\mathbf{p}}_{j+1} = \tilde{\mathbf{r}}_j - \frac{(\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \tilde{\mathbf{r}}_j}{\tilde{\mathbf{p}}_j^T \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j} \tilde{\mathbf{p}}_j$ ;
}
  
```

## Equivalent CG with transformed variables

Input : initial guess  $\mathbf{x}^{(0)} \in \mathbb{R}^n$   
 Output : approximate solution  $\mathbf{x}^{(l)} \in \mathbb{R}^n$

```

 $\mathbf{B}^{1/2}\tilde{\mathbf{r}}_0 := \mathbf{B}^{1/2}\mathbf{b} - \mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{x}}^{(0)}$ ;
 $\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_1 := \mathbf{B}^{-1}(\mathbf{B}^{1/2}\tilde{\mathbf{r}}_0)$ ;
for  $j = 1$  to  $l$  do {
   $\alpha := \frac{(\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \mathbf{B}^{1/2}\tilde{\mathbf{r}}_{j-1}}{(\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j}$ 
   $\mathbf{B}^{-1/2}\tilde{\mathbf{x}}^{(j)} := \mathbf{B}^{-1/2}\tilde{\mathbf{x}}^{(j-1)} + \alpha \mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j$ ;
   $\mathbf{B}^{1/2}\tilde{\mathbf{r}}_j = \mathbf{B}^{1/2}\tilde{\mathbf{r}}_{j-1} - \alpha \mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j$ ;
   $\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_{j+1} = \mathbf{B}^{-1}(\mathbf{B}^{-1/2}\tilde{\mathbf{r}}_j) - \frac{(\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \mathbf{A}\mathbf{B}^{-1}(\mathbf{B}^{1/2}\tilde{\mathbf{r}}_j)}{(\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j} \mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j$ ;
}
  
```

with the transformations:

$$\tilde{\mathbf{x}}^{(k)} = \mathbf{B}^{1/2}\mathbf{x}^{(k)}, \quad \tilde{\mathbf{r}}_k = \mathbf{B}^{-1/2}\mathbf{r}_k, \quad \tilde{\mathbf{p}}_k = \mathbf{B}^{-1/2}\mathbf{p}_k. \quad (10.3.0.4)$$

**§10.3.0.5 (Preconditioned CG method (PCG)) [DR08, Alg. 13.32], [Han02, Alg. 10.1])**

Input: initial guess  $\mathbf{x} \in \mathbb{R}^n \doteq \mathbf{x}^{(0)} \in \mathbb{R}^n$ , tolerance  $\tau > 0$   
 Output: approximate solution  $\mathbf{x} \doteq \mathbf{x}^{(l)}$

```

 $\mathbf{p} := \mathbf{r} := \mathbf{b} - \mathbf{A}\mathbf{x}; \quad \mathbf{p} := \mathbf{B}^{-1}\mathbf{r}; \quad \mathbf{q} := \mathbf{p}; \quad \tau_0 := \mathbf{p}^T \mathbf{r}$ ;
for  $l = 1$  to  $l_{\max}$  do {
   $\beta := \mathbf{r}^T \mathbf{q}; \quad \mathbf{h} := \mathbf{A}\mathbf{p}; \quad \alpha := \frac{\beta}{\mathbf{p}^T \mathbf{h}}$ ;
   $\mathbf{x} := \mathbf{x} + \alpha \mathbf{p}$ ;
   $\mathbf{r} := \mathbf{r} - \alpha \mathbf{h}$ ;
   $\mathbf{q} := \mathbf{B}^{-1}\mathbf{r}; \quad \beta := \frac{\mathbf{r}^T \mathbf{q}}{\beta}$ ;
  if  $|\mathbf{q}^T \mathbf{r}| \leq \tau \cdot \tau_0$  then stop;
   $\mathbf{p} := \mathbf{q} + \beta \mathbf{p}$ ;
}
  
```

**MATLAB-code 10.3.0.7: simple implementation of PCG algorithm § 10.3.0.5**

```

1 function [x, rn, xk] = pcgbase(evalA, b, tol, maxit, invB, x)
2 % evalA must pass a handle to a function implementing A*x
3 % invB is to be a handle to a function providing the action of the
4 % preconditioner on a vector. The other arguments like for MATLAB's
5 % pcg.
6 r = b - evalA(x); rho = 1; rn = [];
7 if (nargout > 2), xk = x; end
  
```

```

7  for i = 1 : maxit
8      y = invB(r);
9      rho_old = rho; rho = r' * y; rn = [rn,rho];
10     if (i == 1), p = y; rho0 = rho;
11     elseif (rho < rho0*tol), return;
12     else beta = rho/rho_old; p = y+beta*p; end
13     q = evalA(p); alpha = rho / (p' * q);
14     x = x + alpha * p;
15     r = r - alpha * q;
16     if (nargout > 2), xk = [xk,x]; end
17 end

```

Computational effort per step: 1 evaluation  $\mathbf{A} \times \text{vector}$ ,  
 1 evaluation  $\mathbf{B}^{-1} \times \text{vector}$ ,  
 3 dot products, 3 AXPY-operations

**Remark 10.3.0.8 (Convergence theory for PCG)** Assertions of Thm. 10.2.3.5 remain valid with  $\kappa(\mathbf{A})$  replaced with  $\kappa(\mathbf{B}^{-1}\mathbf{A})$  and energy norm based on  $\tilde{\mathbf{A}}$  instead of  $\mathbf{A}$ .  $\square$

### EXAMPLE 10.3.0.9 (Simple preconditioners)

$\mathbf{B}$  = easily invertible “part” of  $\mathbf{A}$

♦  $\mathbf{B} = \text{diag}(\mathbf{A})$ : Jacobi preconditioner (diagonal scaling)

♦  $(\mathbf{B})_{ij} = \begin{cases} (\mathbf{A})_{ij} & , \text{if } |i-j| \leq k, \\ 0 & \text{else,} \end{cases} \quad \text{for some } k \ll n.$

♦ Symmetric Gauss-Seidel preconditioner

Idea: Solve  $\mathbf{Ax} = \mathbf{b}$  approximately in two stages:

- ① Approximation  $\mathbf{A}^{-1} \approx \text{tril}(\mathbf{A})$  (lower triangular part):  $\tilde{\mathbf{x}} = \text{tril}(\mathbf{A})^{-1}\mathbf{b}$
- ② Approximation  $\mathbf{A}^{-1} \approx \text{triu}(\mathbf{A})$  (upper triangular part) and use this to approximately “solve” the error equation  $\mathbf{A}(\mathbf{x} - \tilde{\mathbf{x}}) = \mathbf{r}$ , with residual  $\mathbf{r} := \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$ :

$$\mathbf{x} = \tilde{\mathbf{x}} + \text{triu}(\mathbf{A})^{-1}(\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}).$$

With  $\mathbf{L}_A := \text{tril}(\mathbf{A})$ ,  $\mathbf{U}_A := \text{triu}(\mathbf{A})$  one finds

$$\mathbf{x} = (\mathbf{L}_A^{-1} + \mathbf{U}_A^{-1} - \mathbf{U}_A^{-1}\mathbf{A}\mathbf{L}_A^{-1})\mathbf{b} \quad \blacktriangleright \quad \mathbf{B}^{-1} = \mathbf{L}_A^{-1} + \mathbf{U}_A^{-1} - \mathbf{U}_A^{-1}\mathbf{A}\mathbf{L}_A^{-1}. \quad (10.3.0.10)$$

For all these approaches the evaluation of  $\mathbf{B}^{-1}\mathbf{r}$  can be done with effort of  $\mathcal{O}(n)$  in the case of a sparse matrix  $\mathbf{A}$  (e.g. with  $\mathcal{O}(1)$  non-zero entries per row). However, there is absolutely no guarantee that  $\kappa(\mathbf{B}^{-1}\mathbf{A})$  will be reasonably small. It will crucially depend on  $\mathbf{A}$ , if this can be expected.  $\square$

More complicated preconditioning strategies:

- ♦ Incomplete Cholesky factorization, MATLAB-`ichol`, [DR08, Sect. 13.5]]

◆ Sparse approximate inverse preconditioner (SPAI)

### EXAMPLE 10.3.0.11 (Tridiagonal preconditioning)

Efficacy of preconditioning of sparse LSE with tridiagonal part:

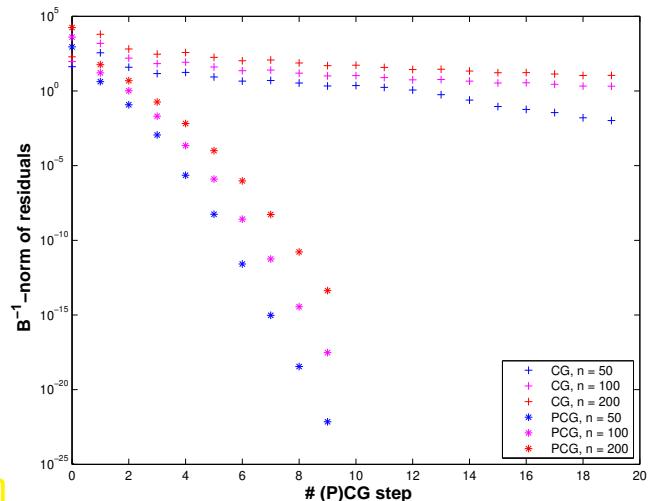
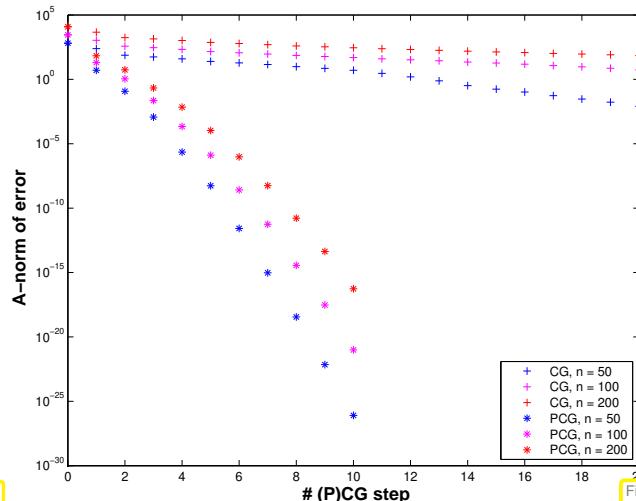
#### MATLAB-code 10.3.0.12: LSE for Ex. 10.3.0.11

```

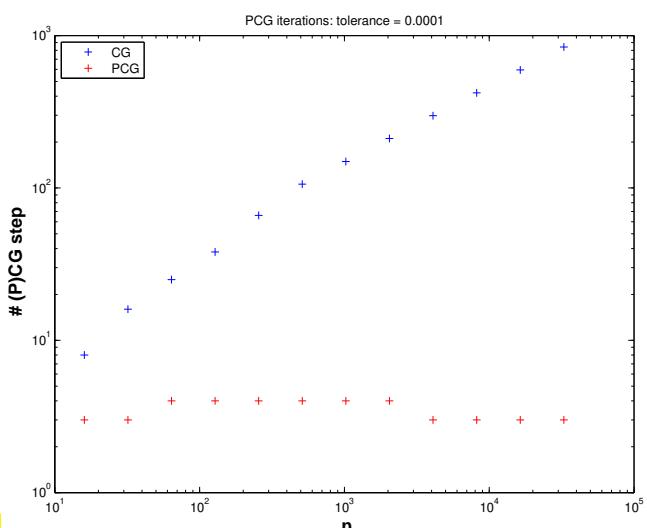
1 A =
2   spdiags (repmat ([1/n,-1,2+2/n,-1,1/n],n,1), [-n/2,-1,0,1,n/2],n,n);
3 b = ones(n,1); x0 = ones(n,1); tol = 1.0E-4; maxit = 1000;
4 evalA = @(x) A*x;
5 % no preconditioning, see Code 10.3.0.7
6 invB = @(x) x; [x,rn] = pcgbase(evalA,b,tol,maxit,invB,x0);
7
8 % tridiagonal preconditioning, see Code 10.3.0.7
9 B = spdiags (spdiags (A, [-1,0,1]), [-1,0,1],n,n);
10 invB = @(x) B\x; [x,rnpc] = pcgbase(evalA,b,tol,maxit,invB,x0);
11 %

```

The Code 10.3.0.12 highlights the use of a preconditioner in the context of the PCG method; it only takes a function that realizes the application of  $B^{-1}$  to a vector. In Line 10 of the code this function is passed as function handle `invB`.



$n$	# CG steps	# PCG steps
16	8	3
32	16	3
64	25	4
128	38	4
256	66	4
512	106	4
1024	149	4
2048	211	4
4096	298	3
8192	421	3
16384	595	3
32768	841	3



Clearly in this example the tridiagonal part of the matrix is dominant for large  $n$ . In addition, its condition number grows  $\sim n^2$  as is revealed by a closer inspection of the spectrum.

Preconditioning with the tridiagonal part manages to suppress this growth of the condition number of  $\mathbf{B}^{-1}\mathbf{A}$  and ensures fast convergence of the preconditioned CG method

**Remark 10.3.0.13 (Termination of PCG)** Recall Rem. 10.2.2.12, (10.2.2.13):

$$\frac{1}{\text{cond}(\mathbf{A})} \frac{\|\mathbf{r}_l\|}{\|\mathbf{r}_0\|} \leq \frac{\|\mathbf{x}^{(l)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{r}_l\|}{\|\mathbf{r}_0\|}. \quad (10.2.2.13)$$

$\mathbf{B}$  good preconditioner  $\Rightarrow \text{cond}_2(\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2})$  small ( $\rightarrow$  Notion 10.3.0.3)

Idea: consider (10.2.2.13) for

- ◆ Euclidean norm  $\|\cdot\| = \|\cdot\|_2 \leftrightarrow \text{cond}_2$
- ◆ transformed quantities  $\tilde{\mathbf{x}}, \tilde{\mathbf{r}}$ , see (10.3.0.1), (10.3.0.4)

Monitor 2-norm of transformed residual:

$$\tilde{\mathbf{r}} = \tilde{\mathbf{b}} - \tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{B}^{-1/2}\mathbf{r} \Rightarrow \|\tilde{\mathbf{r}}\|_2^2 = \mathbf{r}^\top \mathbf{B}^{-1} \mathbf{r}.$$

(10.2.2.13)

► estimate for 2-norm of transformed iteration errors:  $\|\tilde{\mathbf{e}}^{(l)}\|_2^2 = (\mathbf{e}^{(l)})^\top \mathbf{B} \mathbf{e}^{(l)}$

Analogous to (10.2.2.13), estimates for energy norm ( $\rightarrow$  Def. 10.1.0.1) of error  $\mathbf{e}^{(l)} := \mathbf{x} - \mathbf{x}^{(l)}$ ,  $\mathbf{x}^* := \mathbf{A}^{-1}\mathbf{b}$ :

Use error equation  $\mathbf{A}\mathbf{e}^{(l)} = \mathbf{r}_l$ :

$$\mathbf{r}_l^\top \mathbf{B}^{-1} \mathbf{r}_l = (\mathbf{B}^{-1} \mathbf{A} \mathbf{e}^{(l)})^\top \mathbf{A} \mathbf{e}^{(l)} \leq \lambda_{\max}(\mathbf{B}^{-1} \mathbf{A}) \|\mathbf{e}^{(l)}\|_A^2,$$

$$\|\mathbf{e}^{(l)}\|_A^2 = (\mathbf{A} \mathbf{e}^{(l)})^\top \mathbf{e}^{(l)} = \mathbf{r}_l^\top \mathbf{A}^{-1} \mathbf{r}_l = \mathbf{B}^{-1} \mathbf{r}_l^\top \mathbf{B} \mathbf{A}^{-1} \mathbf{r}_l \leq \lambda_{\max}(\mathbf{B} \mathbf{A}^{-1}) (\mathbf{B}^{-1} \mathbf{r}_l)^\top \mathbf{r}_l.$$

available during PCG iteration (10.3.0.6)

$$\frac{1}{\kappa(\mathbf{B}^{-1}\mathbf{A})} \frac{\|\mathbf{e}^{(l)}\|_A^2}{\|\mathbf{e}^{(0)}\|_A^2} \leq \frac{(\mathbf{B}^{-1}\mathbf{r}_l)^\top \mathbf{r}_l}{(\mathbf{B}^{-1}\mathbf{r}_0)^\top \mathbf{r}_0} \leq \kappa(\mathbf{B}^{-1}\mathbf{A}) \frac{\|\mathbf{e}^{(l)}\|_A^2}{\|\mathbf{e}^{(0)}\|_A^2} \quad (10.3.0.14)$$

$\kappa(\mathbf{B}^{-1}\mathbf{A})$  “small”  $\rightarrow \mathbf{B}^{-1}$ -energy norm of residual  $\approx \mathbf{A}$ -norm of error !  
 $(\mathbf{r}_l \cdot \mathbf{B}^{-1}\mathbf{r}_l = \mathbf{q}^\top \mathbf{r}$  in Algorithm (10.3.0.6))

MATLAB-function:  $[x, flag, relr, it, rv] = \text{pcg}(\mathbf{A}, \mathbf{b}, \text{tol}, \text{maxit}, \mathbf{B}, [], \mathbf{x}_0);$   
 $(\mathbf{A}, \mathbf{B}$  may be handles to functions providing  $\mathbf{A}\mathbf{x}$  and  $\mathbf{B}^{-1}\mathbf{x}$ , resp.)

**Remark 10.3.0.15 (Termination criterion in MATLAB-`pcg` → [QSS00, Sect. 4.6])**

Implementation (skeleton) of MATLAB built-in `pcg`:

Listing 10.1: MATLAB-code : PCG algorithm

```

1 function x = pcg(Afun,b,tol,maxit,Binvfun,x0)
2 x = x0; r = b - feval(Afun,x); rho = 1;
3 for i = 1 : maxit
4   y = feval(Binvfun,r);
5   rho1 = rho; rho = r' * y;
6   if (i == 1)
7     p = y;
8   else
9     beta = rho / rho1;
10    p = y + beta * p;
11  end
12  q = feval(Afun,p);
13  alpha = rho / (p' * q);
14  x = x + alpha * p;
15  if (norm(b - evalf(Afun,x)) <= tol * norm(b)) , return; end
16  r = r - alpha * q;
17 end

```

Dubious termination criterion !

## 10.4 Survey of Krylov Subspace Methods

### 10.4.1 Minimal residual methods

Idea: Replace Euclidean inner product in CG with  $\mathbf{A}$ -inner product

$$\|\mathbf{x}^{(l)} - \mathbf{x}\|_A \text{ replaced with } \|\mathbf{A}(\mathbf{x}^{(l)} - \mathbf{x})\|_2 = \|\mathbf{r}_l\|_2$$

► MINRES method [Hac91, Sect. 9.5.2] (for *any* symmetric matrix !)

**Theorem 10.4.1.1.**

For  $\mathbf{A} = \mathbf{A}^H \in \mathbb{R}^{n,n}$  the residuals  $\mathbf{r}_l$  generated in the MINRES iteration satisfy

$$\|\mathbf{r}_l\|_2 = \min\{\|\mathbf{Ay} - \mathbf{b}\|_2 : \mathbf{y} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0)\}$$

$$\Rightarrow \|\mathbf{r}_l\|_2 \leq \frac{2\left(1 - \frac{1}{\kappa(\mathbf{A})}\right)^l}{\left(1 + \frac{1}{\kappa(\mathbf{A})}\right)^{2l} + \left(1 - \frac{1}{\kappa(\mathbf{A})}\right)^{2l}} \|\mathbf{r}_0\|_2.$$

Note: similar formula for (linear) rate of convergence as for CG, see Thm. 10.2.3.5, but with  $\sqrt{\kappa(\mathbf{A})}$  replaced with  $\kappa(\mathbf{A})$ !

Iterative solver for  $\mathbf{Ax} = \mathbf{b}$  with *symmetric* system matrix  $\mathbf{A}$ :

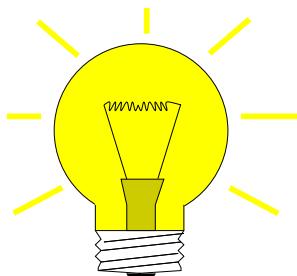
MATLAB-functions:

- `[x, flg, res, it, resv] = minres(A, b, tol, maxit, B, [], x0);`
- `[...] = minres(Afun, b, tol, maxit, Binvfun, [], x0);`

Computational costs: 1  $\mathbf{A} \times$ vector, 1  $\mathbf{B}^{-1} \times$ vector per step, a few dot products & SAXPYs

Memory requirement: a few vectors  $\in \mathbb{R}^n$

Extension to general regular  $\mathbf{A} \in \mathbb{R}^{n,n}$ :



Idea: Solver overdetermined linear system of equations

$$\mathbf{x}^{(l)} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0): \quad \mathbf{Ax}^{(l)} = \mathbf{b}$$

in *least squares sense*,  $\rightarrow$  Chapter 3.

$$\mathbf{x}^{(l)} = \operatorname{argmin}\{\|\mathbf{Ay} - \mathbf{b}\|_2 : \mathbf{y} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0)\}.$$

► **GMRES** method for general matrices  $\mathbf{A} \in \mathbb{R}^{n,n}$   $\rightarrow$  [Han02, Ch. 16], [QSS00, Sect. 4.4.2]

MATLAB-function:

- `[x, flag, relr, it, rv] = gmres(A, b, rs, tol, maxit, B, [], x0);`
- `[...] = gmres(Afun, b, rs, tol, maxit, Binvfun, [], x0);`

Computational costs : 1  $\mathbf{A} \times$ vector, 1  $\mathbf{B}^{-1} \times$ vector per step,

:  $O(l)$  dot products & SAXPYs in  $l$ -th step

Memory requirements:  $O(l)$  vectors  $\in \mathbb{K}^n$  in  $l$ -th step

**Remark 10.4.1.2 (Restarted GMRES)** After many steps of GMRES we face considerable computational costs and memory requirements for every further step. Thus, the iteration may be *restarted* with the current iterate  $\mathbf{x}^{(l)}$  as initial guess  $\rightarrow$  `rs`-parameter triggers restart after every `rs` steps (Danger: failure to converge).  $\square$

## 10.4.2 Iterations with short recursions [QSS00, Sect. 4.5]

Iterative methods for *general* regular system matrix  $\mathbf{A}$ :



Idea: Given  $\mathbf{x}^{(0)} \in \mathbb{R}^n$  determine (better) approximation  $\mathbf{x}^{(l)}$  through Petrov-Galerkin condition

$$\mathbf{x}^{(l)} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0): \quad \mathbf{p}^H(\mathbf{b} - \mathbf{A}\mathbf{x}^{(l)}) = 0 \quad \forall \mathbf{p} \in W_l,$$

with suitable test space  $W_l$ ,  $\dim W_l = l$ , e.g.  $W_l := \mathcal{K}_l(\mathbf{A}^H, \mathbf{r}_0)$  ( $\rightarrow$  bi-conjugate gradients, BiCG)

► Zoo of methods with short recursions (i.e. constant effort per step)

MATLAB-function:

- `[x, flag, r, it, rv] = bicgstab(A, b, tol, maxit, B, [], x0)`
- `[...] = bicgstab(Afun, b, tol, maxit, Binvfun, [], x0);`

Computational costs : 2  $\mathbf{A} \times$ vector, 2  $\mathbf{B}^{-1} \times$ vector, 4 dot products, 6 SAXPYs per step

Memory requirements: 8 vectors  $\in \mathbb{R}^n$

MATLAB-function:

- `[x, flag, r, it, rv] = qmr(A, b, tol, maxit, B, [], x0)`
- `[...] = qmr(Afun, b, tol, maxit, Binvfun, [], x0);`

Computational costs : 2  $\mathbf{A} \times$ vector, 2  $\mathbf{B}^{-1} \times$ vector, 2 dot products, 12 SAXPYs per step

Memory requirements: 10 vectors  $\in \mathbb{R}^n$



- ◆ little (useful) convergence theory available
- ◆ stagnation & “breakdowns” commonly occur

### EXAMPLE 10.4.2.1 (Failure of Krylov iterative solvers)

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & \cdots & & \cdots & 0 \\ 0 & 0 & 1 & 0 & & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & & \\ & & & & & \ddots & \vdots \\ \vdots & & & & & \ddots & 0 \\ 0 & & & & & 0 & 1 \\ 1 & 0 & \cdots & & \cdots & 0 & \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ \vdots \\ \vdots \\ 0 \\ 1 \end{bmatrix} \quad \Rightarrow \quad \mathbf{x} = \mathbf{e}_1.$$

$$\mathbf{x}^{(0)} = 0 \quad \Rightarrow \quad \mathbf{r}_0 = \mathbf{e}_n \quad \Rightarrow \quad \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0) = \text{Span}\{\mathbf{e}_n, \mathbf{e}_{n-1}, \dots, \mathbf{e}_{n-l+1}\}$$

$$\Rightarrow \min\{\|\mathbf{y} - \mathbf{x}\|_2: \mathbf{y} \in \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0)\} = \begin{cases} 1 & , \text{ if } l \leq n, \\ 0 & , \text{ for } l = n. \end{cases}$$

TRY & PRAY

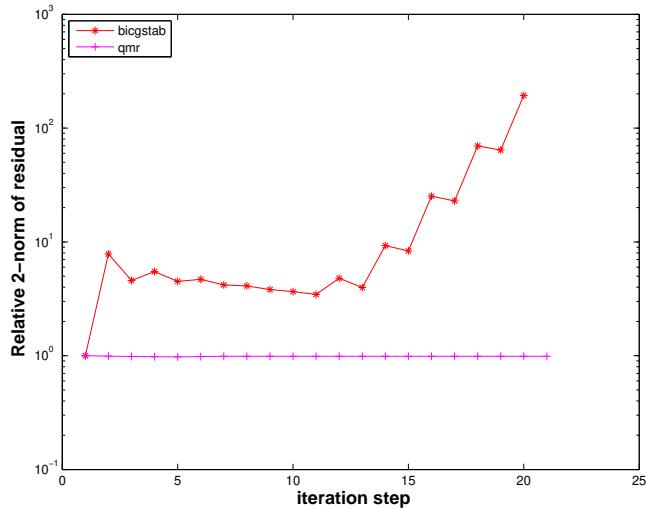
### EXAMPLE 10.4.2.2 (Convergence of Krylov subspace methods for non-symmetric system matrix)

```

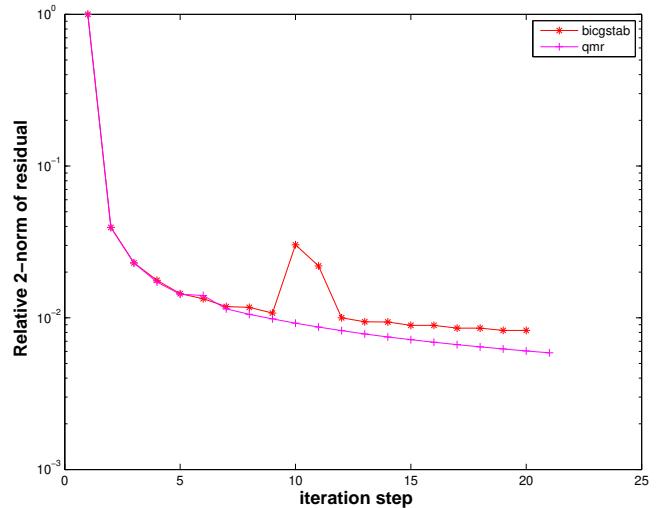
1 A = gallery('tridiag', -0.5*ones(n-1,1), 2*ones(n,1), -1.5*ones(n-1,1));
2 B = gallery('tridiag', 0.5*ones(n-1,1), 2*ones(n,1), 1.5*ones(n-1,1));

```

Plotted:  $\|\mathbf{r}_l\|_2 : \|\mathbf{r}_0\|_2$ :



tridiagonal matrix **A**



tridiagonal matrix **B**

## Summary:

Advantages of Krylov methods vs. direct elimination (**IF** they converge at all/sufficiently fast).

- They require system matrix **A** in procedural form  $\mathbf{y} = \text{evalA}(\mathbf{x}) \leftrightarrow \mathbf{y} = \mathbf{Ax}$  only.
- They can perfectly exploit sparsity of system matrix.
- They can cash in on low accuracy requirements (**IF** viable termination criterion available).
- They can benefit from a good initial guess.

# Bibliography

- [AK07] Owe Axelsson and János Karátson. “Mesh independent superlinear PCG rates via compact-equivalent operators”. In: *SIAM J. Numer. Anal.* 45.4 (2007), pp. 1495–1516. DOI: [10.1137/06066391X](https://doi.org/10.1137/06066391X).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 616, 617, 624–639).
- [Gut09] M.H. Gutknecht. *Lineare Algebra*. Lecture Notes. SAM, ETH Zürich, 2009 (cit. on pp. 621–623, 626, 627).
- [Hac91] W. Hackbusch. *Iterative Lösung großer linearer Gleichungssysteme*. B.G. Teubner–Verlag, Stuttgart, 1991 (cit. on pp. 623, 631, 640).
- [Hac94] Wolfgang Hackbusch. *Iterative solution of large sparse systems of equations*. Vol. 95. Applied Mathematical Sciences. New York: Springer-Verlag, 1994, pp. xxii+429 (cit. on pp. 616, 631).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 617, 624–639, 641).
- [IM97] I.C.F. Ipsen and C.D. Meyer. *The idea behind Krylov methods*. Technical Report 97-3. Raleigh, NC: Math. Dep., North Carolina State University, Jan. 1997.
- [LS09] A.R. Laliena and F.-J. Sayas. “Theoretical aspects of the application of convolution quadrature to scattering of acoustic waves”. In: *Numer. Math.* 112.4 (2009), pp. 637–678 (cit. on p. 634).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on pp. 621–623, 626, 627, 634).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 616–641).
- [Saa03] Yousef Saad. *Iterative methods for sparse linear systems*. Second. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2003, pp. xviii+528 (cit. on p. 616).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on p. 619).
- [Win80] Ragnar Winther. “Some superlinear convergence results for the conjugate gradient method”. In: *SIAM J. Numer. Anal.* 17.1 (1980), pp. 14–17.

# Chapter 11

## Numerical Integration – Single Step Methods

For historical reasons the approximate solution of initial value problems for ordinary differential equations is called “Numerical Integration”. This chapter will introduce the most important class of numerical methods for that purpose.

### Contents

---

<b>11.1 Initial-Value Problems (IVP) for ODEs</b> . . . . .	<b>645</b>
11.1.1 Modeling with ordinary differential equations: Examples . . . . .	647
11.1.2 Theory of Initial-Value-Problems (IVPs) . . . . .	650
11.1.3 Evolution Operators . . . . .	654
<b>11.2 Introduction: Polygonal Approximation Methods</b> . . . . .	<b>656</b>
11.2.1 Explicit Euler method . . . . .	657
11.2.2 Implicit Euler method . . . . .	659
11.2.3 Implicit midpoint method . . . . .	660
<b>11.3 General single step methods</b> . . . . .	<b>661</b>
11.3.1 Definition . . . . .	661
11.3.2 (Asymptotic) Convergence of Single-Step Methods . . . . .	664
<b>11.4 Explicit Runge-Kutta Methods</b> . . . . .	<b>670</b>
<b>11.5 Adaptive Stepsize Control</b> . . . . .	<b>677</b>

---

### 11.1 Initial-Value Problems (IVP) for ODEs

The title of this section contains a widely used acronym: **ODE**  $\triangleq$  ordinary differential equation . In this section we present a few examples of models involving ODEs and briefly review the relevant mathematical theory.

**§11.1.0.1 (Terminology and notations related to ordinary differential equations)** In our parlance, a **(first-order) ordinary differential equation** (ODE) is an equation of the form

$$\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}) , \quad (11.1.0.2)$$

with

- ☞ a (continuous) **right hand side function** (r.h.s)  $f : I \times D \rightarrow \mathbb{R}^d$  of **time**  $t \in \mathbb{R}$  and **state**  $\mathbf{y} \in \mathbb{R}^d$
- ☞ defined on a (finite) time interval  $I \subset \mathbb{R}$ , and **state space**  $D \subset \mathbb{R}^d$ .

In the context of mathematical modeling the state vector  $\mathbf{y} \in \mathbb{R}^d$  is supposed to provide a complete (in the sense of the model) description of a system. Then (11.1.0.2) models a **finite-dimensional dynamical system**. Examples will be provided below, see Ex. 11.1.1.1, Ex. 11.1.1.5, and Ex. 11.1.1.7.

For  $d > 1$   $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$  can be viewed as a **system of ordinary differential equations**:

$$\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}) \iff \begin{bmatrix} \dot{y}_1 \\ \vdots \\ \dot{y}_d \end{bmatrix} = \begin{bmatrix} f_1(t, y_1, \dots, y_d) \\ \vdots \\ f_d(t, y_1, \dots, y_d) \end{bmatrix}.$$

☞ Notation (Newton): dot  $\cdot$   $\hat{=}$  (total) derivative with respect to time  $t$

### Definition 11.1.0.3. Solution of an ordinary differential equation

A **solution** of the ODE  $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$  with continuous right hand side function  $\mathbf{f}$  is a continuously differentiable **function** "of time  $t$ "  $\mathbf{y} : J \subset I \rightarrow D$ , defined on an **open** interval  $J$ , for which  $\dot{\mathbf{y}}(t) = \mathbf{f}(t, \mathbf{y}(t))$  holds for all  $t \in J$ .

A solution describes a continuous **trajectory** in state space, a one-parameter family of states, parameterized by time.

It goes without saying that smoothness of the right hand side function  $\mathbf{f}$  is inherited by solutions of the ODE:

### Lemma 11.1.0.4. Smoothness of solutions of ODEs

Let  $\mathbf{y} : I \subset \mathbb{R} \rightarrow D$  be a solution of the ODE  $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$  on the time interval  $I$ .

If  $\mathbf{f} : I \times D \rightarrow \mathbb{R}^d$  is  $r$ -times continuously differentiable with respect to both arguments,  $r \in \mathbb{N}_0$ , then the trajectory  $t \mapsto \mathbf{y}(t)$  is  $r+1$ -times continuously differentiable in the interior of  $I$ .



**Supplementary literature.** Some grasp of the meaning and theory of ordinary differential equa-

tions (ODEs) is indispensable for understanding the construction and properties of numerical methods. Relevant information can be found in [Str09, Sect. 5.6, 5.7, 6.5].

Books dedicated to numerical methods for ordinary differential equations:

- [DB02] excellent textbook, but geared to the needs of students of mathematics.
- [HNW93] and [HW11] : the standard reference.
- [HLW06]: wonderful book conveying deep insight, with emphasis on mathematical concepts.

### 11.1.1 Modeling with ordinary differential equations: Examples

Most models of physical systems and phenomena that are *continuously changing with time* involve ordinary differential equations.

**EXAMPLE 11.1.1.1 (Growth with limited resources [Ama83, Sect. 1.1], [Han02, Ch. 60])** This is an example from **population dynamics** with a one-dimensional state space  $D = \mathbb{R}_0^+, d = 1$ :

$y : [0, T] \mapsto \mathbb{R}$ : bacterial population density as a function of time

ODE-based model: autonomous **logistic differential equations** [Str09, Ex. 5.7.2]

$$\dot{y} = f(y) := (\alpha - \beta y) y \quad (11.1.1.2)$$

- ◆  $y \hat{=} \text{population density, } [y] = \frac{1}{\text{m}^2}$ 
  - $\dot{y} \hat{=} \text{instantaneous change (growth/decay) of population density}$
- ◆ growth rate  $\alpha - \beta y$  with growth coefficients  $\alpha, \beta > 0$ ,  $[\alpha] = \frac{1}{\text{s}}$ ,  $[\beta] = \frac{\text{m}^2}{\text{s}}$ : decreases due to more fierce competition as population density increases.

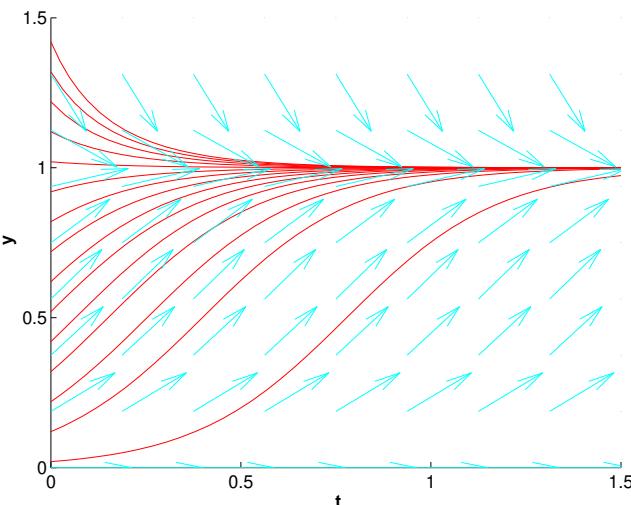


Fig. 363

Solution for different  $y(0)$  ( $\alpha, \beta = 5$ )

Note that by fixing the initial value  $y(0)$  we can single out a *unique* representative from the family of solutions. This will turn out to be a general principle, see Section 11.1.2. □

#### Definition 11.1.1.4. Autonomous ODE

An ODE of the form  $\dot{y} = f(y)$ , that is, with a right hand side function that does not depend on time, but only on state, is called **autonomous**.

For an autonomous ODE the right hand side function defines a **vector field** (“velocity field”)  $\mathbf{y} \mapsto \mathbf{f}(\mathbf{y})$  on state space.

**EXAMPLE 11.1.1.5 (Predator-prey model [Ama83, Sect. 1.1],[HLW06, Sect. 1.1.1],[Han02, Ch. 60],[DR08, Ex. 11.3])** A model from population dynamics:

Predators and prey coexist in an ecosystem. Without predators the population of prey would be governed by a simple exponential growth law. However, the growth rate of prey will decrease with increasing

numbers of predators and, eventually, become negative. Similar considerations apply to the predator population and lead to an ODE model.

ODE-based model: autonomous **Lotka-Volterra ODE**:

$$\begin{aligned}\dot{u} &= (\alpha - \beta v)u \\ \dot{v} &= (\delta u - \gamma)v\end{aligned} \leftrightarrow \dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) \quad \text{with} \quad \mathbf{y} = \begin{bmatrix} u \\ v \end{bmatrix}, \quad \mathbf{f}(\mathbf{y}) = \begin{bmatrix} (\alpha - \beta v)u \\ (\delta u - \gamma)v \end{bmatrix}, \quad (11.1.1.6)$$

with positive model parameters  $\alpha, \beta, \gamma, \delta > 0$ .

population densities:

$u(t) \rightarrow$  density of prey at time  $t$ ,  
 $v(t) \rightarrow$  density of predators at time  $t$

Right hand side vector field  $\mathbf{f}$  for Lotka-Volterra ODE

▷

Solution curves are trajectories of particles carried along by velocity field  $\mathbf{f}$ .

(Parameter values for Fig. 364:  $\alpha = 2, \beta = 1, \delta = 1, \gamma = 1$ .)

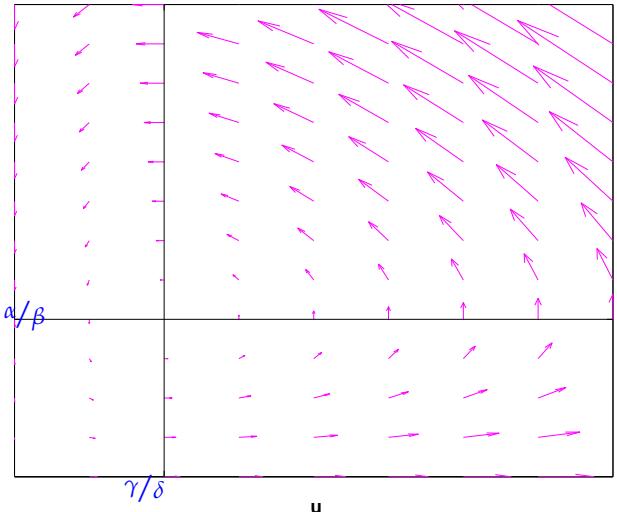


Fig. 364

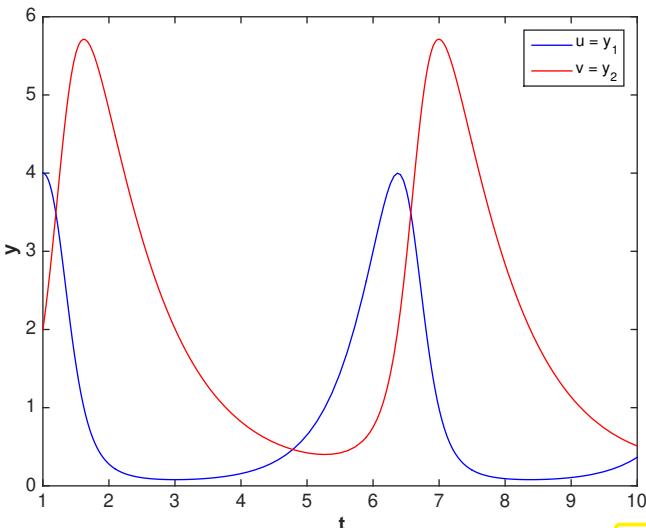
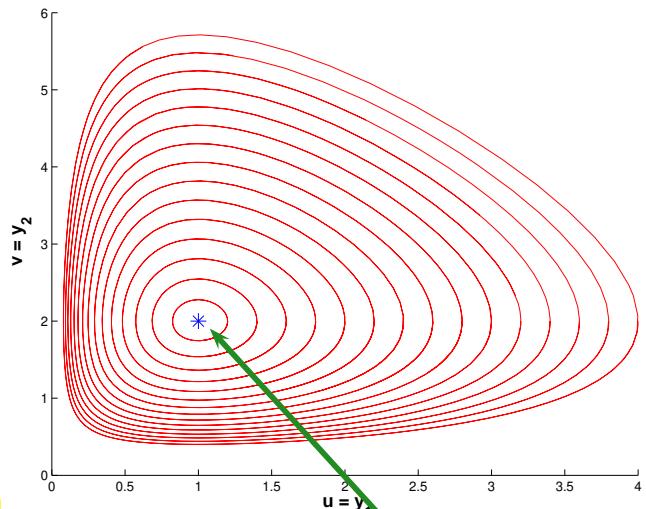


Fig. 366



Solution curves for (11.1.1.6)

stationary point

Parameter values for Fig. 366, 365:  $\alpha = 2, \beta = 1, \delta = 1, \gamma = 1$

**EXAMPLE 11.1.1.7 (Heartbeat model → [Dea80, p. 655])** This example deals with a *phenomenological model* from physiology. A model is called phenomenological, if it is entirely motivated by observations without appealing to underlying mechanisms or first principles.

State of heart described by quantities:  $l = l(t) \hat{=} \text{length of muscle fiber}$   
 $p = p(t) \hat{=} \text{electro-chemical potential}$

$$\begin{aligned}\text{Phenomenological model:} \quad \dot{l} &= -(l^3 - \alpha l + p), \\ \dot{p} &= \beta l,\end{aligned} \quad (11.1.1.8)$$

with parameters:  $\alpha \hat{=} \text{pre-tension of muscle fiber}$   
 $\beta \hat{=} \text{(phenomenological) feedback parameter}$

This is the so-called Zeeman model: it is a phenomenological model entirely based on macroscopic observations without relying on knowledge about the underlying molecular mechanisms.

Plots of vector fields for (11.1.1.8) and solutions for different choices of parameters are given next:

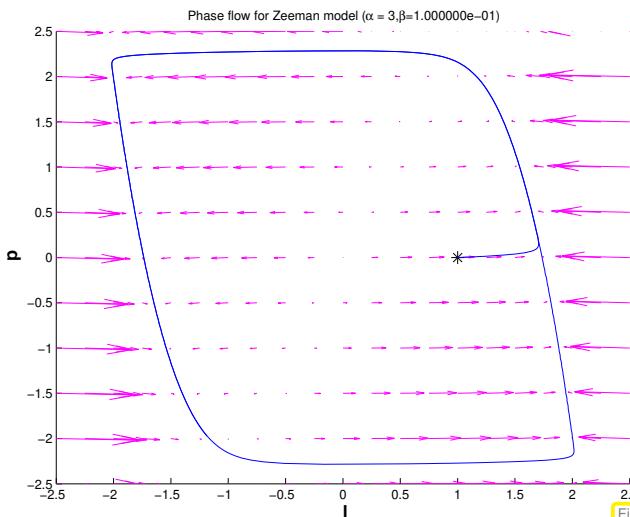


Fig. 367

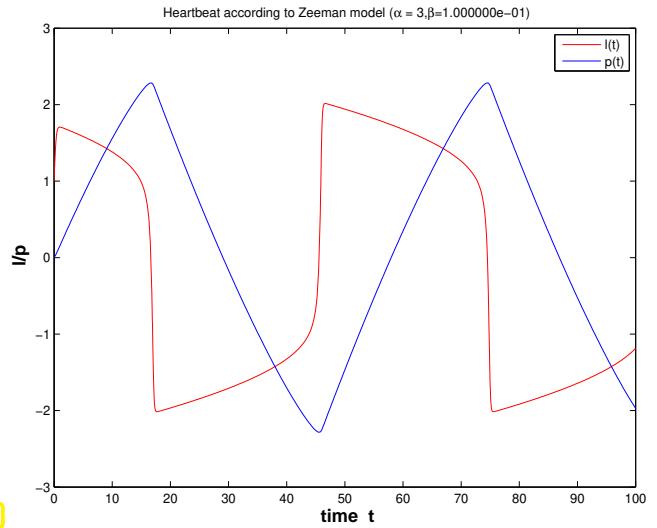


Fig. 368

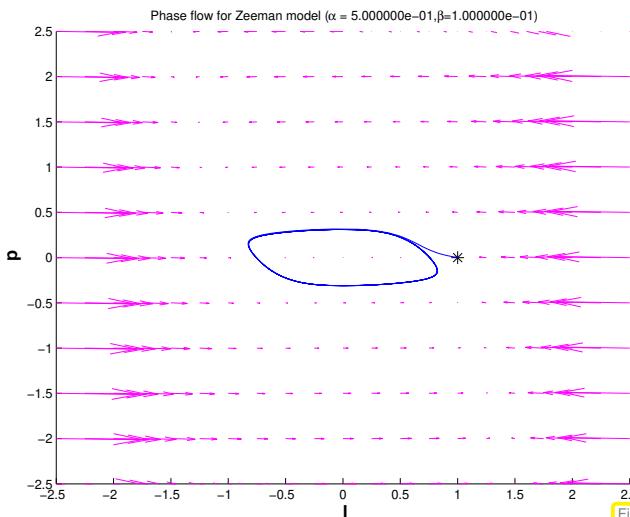


Fig. 369

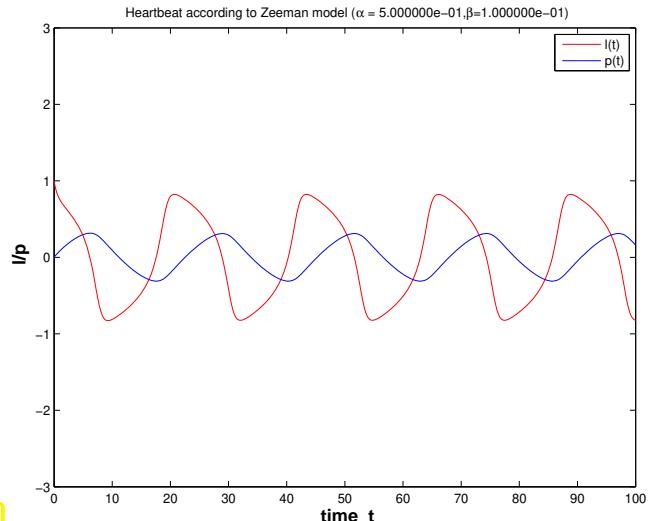


Fig. 370

Observation:  $\alpha \ll 1$  (bottom plots)  $\rightarrow$  ventricular fibrillation, a life-threatening condition.

**EXAMPLE 11.1.1.9 (Transient circuit simulation [Han02, Ch. 64])** In Chapter 1 and Chapter 8 we discussed circuit analysis as a source of linear and non-linear systems of equations, see Ex. 2.1.0.3 and Ex. 8.0.0.1. In the former example we admitted time-dependent currents and potentials, but dependence on time was confined to be “sinusoidal”. This enabled us to switch to frequency domain, see (2.1.0.6), which gave us a complex linear system of equations for the complex nodal potentials. Yet, this trick is only possible for *linear* circuits. In the general case, circuits have to be modelled by ODEs connecting time-dependent potentials and currents. This will be briefly explained now.

The approach is transient nodal analysis, *cf.* Ex. 2.1.0.3, based on the Kirchhoff current law, which reads for the node  $\bullet$  of the simple circuit drawn in Fig. 371

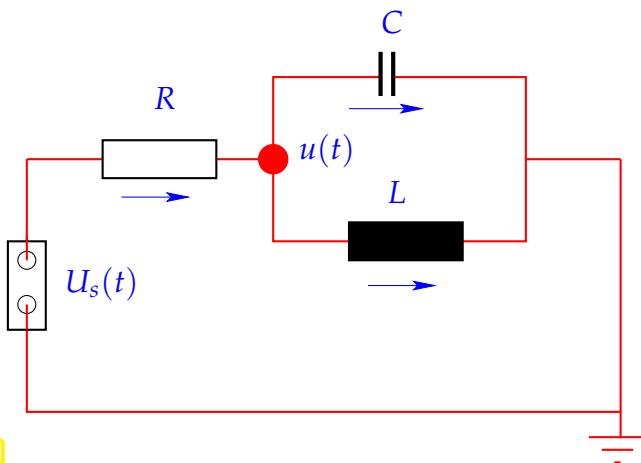
$$i_R(t) - i_L(t) - i_C(t) = 0 . \quad (11.1.1.10)$$

In addition we rely on known transient constitutive relations for basic linear circuit elements:

$$\text{resistor: } i_R(t) = R^{-1}u_R(t), \quad (11.1.1.11)$$

$$\text{capacitor: } i_C(t) = C \frac{du_C}{dt}(t), \quad (11.1.1.12)$$

$$\text{coil: } u_L(t) = L \frac{di_L}{dt}(t). \quad (11.1.1.13)$$



We assume that the source voltage  $U_s(t)$  is given. To apply nodal analysis to the circuit of Fig. 371 we differentiate (11.1.1.10) w.r.t.  $t$

$$\frac{di_R}{dt}(t) - \frac{di_L}{dt}(t) - \frac{di_C}{dt}(t) = 0,$$

and plug in the above constitutive relations for circuit elements:

$$\Rightarrow R^{-1} \frac{du_R}{dt}(t) - L^{-1}u_L(t) - C \frac{d^2u_C}{dt^2}(t) = 0.$$

We continue following the policy of nodal analysis and express all voltages by potential differences between nodes of the circuit.

$$u_R(t) = U_s(t) - u(t), \quad u_C(t) = u(t) - 0, \quad u_L(t) = u(t) - 0.$$

For this simple circuit there is only one node with unknown potential, see Fig. 371. Its time-dependent potential will be denoted by  $u(t)$  and this is the unknown of the model, a function of time satisfying the ordinary differential equation

$$R^{-1}(U_s(t) - u(t)) - L^{-1}u(t) - C \frac{d^2u}{dt^2}(t) = 0.$$

► This is an autonomous **2nd-order** ordinary differential equation:

$$Cu'' + R^{-1}u' + L^{-1}u = R^{-1}U_s. \quad (11.1.1.14)$$

The attribute “2nd-order” refers to the occurrence of a second derivative with respect to time. □

## 11.1.2 Theory of Initial-Value-Problems (IVPs)



*Supplementary literature.* [QSS00, Sect. 11.1], [DR08, Sect. 11.3]

**§11.1.2.1 (Initial value problems)** We start with an abstract mathematical description that also introduces key terminology:

A generic **Initial value problem** (IVP) for a **first-order ordinary differential equation** (ODE) ( $\rightarrow$  [Str09, Sect. 5.6], [DR08, Sect. 11.1]) can be stated as: find a function  $\mathbf{y} : I \rightarrow D$  that satisfies, cf. Def. 11.1.0.3,

$$\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}) \quad , \quad \mathbf{y}(t_0) = \mathbf{y}_0 . \quad (11.1.2.2)$$

- $\mathbf{f} : I \times D \mapsto \mathbb{R}^d \hat{=} \text{right hand side}$  (r.h.s.) ( $d \in \mathbb{N}$ ),
- $I \subset \mathbb{R} \hat{=} \text{(time)interval} \leftrightarrow \text{"time variable" } t$ ,
- $D \subset \mathbb{R}^d \hat{=} \text{state space/phase space} \leftrightarrow \text{"state variable" } \mathbf{y}$ ,
- $\Omega := I \times D \hat{=} \text{extended state space}$  (of tuples  $(t, \mathbf{y})$ ),
- $t_0 \in I \hat{=} \text{initial time}, \quad \mathbf{y}_0 \in D \hat{=} \text{initial state} \succ \text{initial conditions.}$

The time interval  $I$  may be finite or infinite. Frequently, the extended state space is not specified, but assumed to coincide with the maximal domain of definition of  $\mathbf{f}$ . Sometimes, the model suggests constraints on  $D$ , for instance, positivity of certain components that represent a density.  $\square$

**§11.1.2.3 (IVPs for autonomous ODEs)** Recall Def. 11.1.1.4: For an autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ , that is the right hand side  $\mathbf{f}$  does not depend on time  $t$ .

Hence, for autonomous ODEs we have  $I = \mathbb{R}$  and the right hand side function  $\mathbf{y} \mapsto \mathbf{f}(\mathbf{y})$  can be regarded as a stationary vector field (velocity field), see Fig. 364 or Fig. 367.

An important observation: If  $t \mapsto \mathbf{y}(t)$  is a solution of an autonomous ODE, then, for any  $\tau \in \mathbb{R}$ , also the shifted function  $t \mapsto \mathbf{y}(t - \tau)$  is a solution.

- $\succ$  For initial value problems for autonomous ODEs the initial time is irrelevant and therefore we can always make the “canonical choice  $t_0 = 0$ .

Autonomous ODEs naturally arise when modeling **time-invariant** systems or phenomena. All examples for Section 11.1.1 belong to this class.  $\square$

**§11.1.2.4 (Autonomization: Conversion into autonomous ODE)** In fact, autonomous ODEs already represent the general case, because every ODE can be converted into an autonomous one:

The idea is to include time as an extra  $d + 1$ -st component of an extended state vector  $\mathbf{z}(t)$ . This solution component has to grow linearly  $\leftrightarrow$  its temporal derivative must be  $= 1$

$$\mathbf{z}(t) := \begin{bmatrix} \mathbf{y}(t) \\ t \end{bmatrix} = \begin{bmatrix} \mathbf{z}' \\ z_{d+1} \end{bmatrix} : \quad \dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}) \leftrightarrow \dot{\mathbf{z}} = \mathbf{g}(\mathbf{z}), \quad \mathbf{g}(\mathbf{z}) := \begin{bmatrix} \mathbf{f}(z_{d+1}, \mathbf{z}') \\ 1 \end{bmatrix} .$$

This means  $\dot{z}_{d+1} = 1$  and implies  $z_{d+1}(t) = t + t_0$ , if  $t_0$  stands for the initial time in the original non-autonomous IVP.

- $\succ$  We restrict ourselves to autonomous ODEs in the remainder of this chapter.  $\square$

**Remark 11.1.2.5 (From higher order ODEs to first order systems)** [DR08, Sect. 11.2]

An ordinary differential equation of **order**  $n \in \mathbb{N}$  has the form

$$\boxed{\mathbf{y}^{(n)} = \mathbf{f}(t, \mathbf{y}, \dot{\mathbf{y}}, \dots, \mathbf{y}^{(n-1)})} . \quad (11.1.2.6)$$

where, with notations from § 11.1.2.1,  $\mathbf{f} : I \times D \times \cdots \times D \rightarrow \mathbb{R}^d$  is a function of times and  $n$  state arguments.

Notation: superscript  $(n) \doteq n$ -th temporal derivative  $t : \frac{d^n}{dt^n}$

Thus, no special treatment of higher order ODEs is necessary, because (11.1.2.6) can be turned into a 1st-order ODE (a system of size  $nd$ ) by adding all derivatives up to order  $n - 1$  as additional components to the state vector. This extended state vector  $\mathbf{z}(t) \in \mathbb{R}^{nd}$  is defined as

$$\mathbf{z}(t) := \begin{bmatrix} \mathbf{y}(t) \\ \mathbf{y}^{(1)}(t) \\ \vdots \\ \mathbf{y}^{(n-1)}(t) \end{bmatrix} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \vdots \\ \mathbf{z}_n \end{bmatrix} \in \mathbb{R}^{dn}: \quad (11.1.2.6) \Leftrightarrow \dot{\mathbf{z}} = \mathbf{g}(\mathbf{z}), \quad \mathbf{g}(\mathbf{z}) := \begin{bmatrix} \mathbf{z}_2 \\ \mathbf{z}_3 \\ \vdots \\ \mathbf{z}_n \\ \mathbf{f}(t, \mathbf{z}_1, \dots, \mathbf{z}_n) \end{bmatrix}. \quad (11.1.2.7)$$

Note that the extended system requires initial values  $\mathbf{y}(t_0), \dot{\mathbf{y}}(t_0), \dots, \mathbf{y}^{(n-1)}(t_0)$ :

For ODEs of order  $n \in \mathbb{N}$  well-posed initial value problems need to specify initial values for the first  $n - 1$  derivatives.

Now we review results about existence and uniqueness of solutions of initial value problems for first-order ODEs. These are surprisingly general and do not impose severe constraints on right hand side functions.

#### Definition 11.1.2.9. Lipschitz continuous function ( $\rightarrow$ [Str09, Def. 4.1.4])

Let  $\Theta := I \times D$ ,  $I \subset \mathbb{R}$  an interval,  $D \subset \mathbb{R}^d$  an open domain. A function  $\mathbf{f} : \Theta \mapsto \mathbb{R}^d$  is **Lipschitz continuous** (in the second argument) on  $\Theta$ , if

$$\exists L > 0: \|\mathbf{f}(t, \mathbf{w}) - \mathbf{f}(t, \mathbf{z})\| \leq L \|\mathbf{w} - \mathbf{z}\| \quad \forall (t, \mathbf{w}), (t, \mathbf{z}) \in \Theta. \quad (11.1.2.10)$$

#### Definition 11.1.2.11. Local Lipschitz continuity ( $\rightarrow$ [Str09, Def. 4.1.5])

Let  $\Omega := I \times D$ ,  $I \subset \mathbb{R}$  an interval,  $D \subset \mathbb{R}^d$  an open domain. A function  $\mathbf{f} : \Omega \mapsto \mathbb{R}^d$  is **locally Lipschitz continuous**, if for every  $(t, \mathbf{y}) \in \Omega$  there is a closed box  $B$  with  $(t, \mathbf{y}) \in B$  such that  $\mathbf{f}$  is Lipschitz continuous on  $B$ :

$$\begin{aligned} \forall (t, \mathbf{y}) \in \Omega: \quad & \exists \delta > 0, L > 0: \\ & \|\mathbf{f}(\tau, \mathbf{z}) - \mathbf{f}(\tau, \mathbf{w})\| \leq L \|\mathbf{z} - \mathbf{w}\| \\ & \forall \mathbf{z}, \mathbf{w} \in D: \|\mathbf{z} - \mathbf{y}\| \leq \delta, \|\mathbf{w} - \mathbf{y}\| \leq \delta, \forall \tau \in I: |\tau - t| \leq \delta. \end{aligned} \quad (11.1.2.12)$$

The property of local Lipschitz continuity means that the function  $(t, \mathbf{y}) \mapsto \mathbf{f}(t, \mathbf{y})$  has “locally finite slope” in  $\mathbf{y}$ .

**EXAMPLE 11.1.2.13 (A function that is not locally Lipschitz continuous)** The meaning of local Lipschitz continuity is best explained by giving an example of a function that fails to possess this property.

Consider the square root function  $t \mapsto \sqrt{t}$  on the *closed* interval  $[0, 1]$ . Its slope in  $t = 0$  is infinite and so

it is not locally Lipschitz continuous on  $[0, 1]$ .

However, if we consider the square root on the *open* interval  $]0, 1[$ , then it is locally Lipschitz continuous there.  $\dashv$

 Notation:  $D_y \mathbf{f} \doteq$  derivative of  $\mathbf{f}$  w.r.t. state variable, a Jacobian  $\in \mathbb{R}^{d,d}$  as defined in (8.2.2.8).

The next lemma gives a simple criterion for local Lipschitz continuity, which can be proved by the mean value theorem, *cf.* the proof of Lemma 8.2.2.9.

#### Lemma 11.1.2.14. Criterion for local Lipschitz continuity

If  $\mathbf{f}$  and  $D_y \mathbf{f}$  are continuous on the extended state space  $\Omega$ , then  $\mathbf{f}$  is locally Lipschitz continuous ( $\rightarrow$  Def. 11.1.2.11).

**Theorem 11.1.2.15. Theorem of Peano & Picard-Lindelöf** [Ama83, Satz II(7.6)], [Str09, Satz 6.5.1], [DR08, Thm. 11.10], [Han02, Thm. 73.1]

If the right hand side function  $\mathbf{f} : \hat{\Omega} \mapsto \mathbb{R}^d$  is locally Lipschitz continuous ( $\rightarrow$  Def. 11.1.2.11) then for all initial conditions  $(t_0, \mathbf{y}_0) \in \hat{\Omega}$  the IVP (11.1.2.2) has a solution  $\mathbf{y} \in C^1(J(t_0, \mathbf{y}_0), \mathbb{R}^d)$  with **maximal** (temporal) domain of definition  $J(t_0, \mathbf{y}_0) \subset \mathbb{R}$ .

In light of § 11.1.2.4 and Thm. 11.1.2.15 henceforth we mainly consider

$$\text{autonomous IVPs: } \dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) \quad , \quad \mathbf{y}(0) = \mathbf{y}_0 \quad , \quad (11.1.2.16)$$

with locally Lipschitz continuous ( $\rightarrow$  Def. 11.1.2.11) right hand side  $\mathbf{f}$ .

**§11.1.2.17 (Domain of definition of solutions of IVPs)** We emphasize a subtle message of Thm. 11.1.2.15.

Solutions of an IVP have an *intrinsic* maximal domain of definition

! Also note that the domain of definition/domain of existence  $J(t_0, \mathbf{y}_0)$  of the solution usually depends on the initial values  $(t_0, \mathbf{y}_0)$  !

Terminology: if  $J(t_0, \mathbf{y}_0) = I$  we say that the solution  $\mathbf{y} : I \mapsto \mathbb{R}^d$  is **global**.

Notation: for autonomous ODE we always have  $t_0 = 0$ , and therefore we write  $J(\mathbf{y}_0) := J(0, \mathbf{y}_0)$ .  $\dashv$

**EXAMPLE 11.1.2.18 (Finite-time blow-up)** Let us explain the still mysterious “maximal domain of definition” in statement of Thm. 11.1.2.15. It is related to the fact that every solution of an initial value problem (11.1.2.16) has its own largest possible time interval  $J(\mathbf{y}_0) \subset \mathbb{R}$  on which it is defined naturally.

As an example we consider the autonomous scalar ( $d = 1$ ) initial value problem, modeling “explosive growth” with a growth rate increasing linearly with the density:

$$\dot{y} = y^2 \quad , \quad y(0) = y_0 \in \mathbb{R} . \quad (11.1.2.19)$$

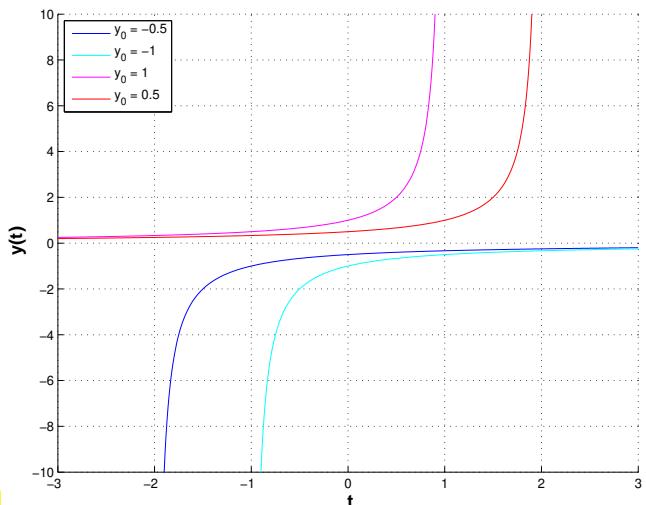
We choose  $I = D = \mathbb{R}$ . Clearly,  $y \mapsto y^2$  is locally Lipschitz-continuous, but only locally! Why not globally?

We find the solutions

$$y(t) = \begin{cases} \frac{1}{y_0^{-1}-t} & , \text{ if } y_0 \neq 0 , \\ 0 & , \text{ if } y_0 = 0 , \end{cases} \quad (11.1.2.20)$$

with domains of definition

$$J(y_0) = \begin{cases} ]-\infty, y_0^{-1}[ & , \text{ if } y_0 > 0 , \\ \mathbb{R} & , \text{ if } y_0 = 0 , \\ ]y_0^{-1}, \infty[ & , \text{ if } y_0 < 0 . \end{cases}$$



In this example, for  $y_0 > 0$  the solution experiences a **blow-up** in finite time and ceases to exist afterwards.

### 11.1.3 Evolution Operators

Now we examine a difficult but fundamental concept for time-dependent models stated by means of ODEs. For the sake of simplicity we restrict the discussion to autonomous IVPs

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) , \quad \mathbf{y}(0) = \mathbf{y}_0 , \quad (11.1.2.16)$$

with locally Lipschitz continuous right hand side and make the following assumption. A more general treatment is given in [DB02].

#### Assumption 11.1.3.1. Global solutions

Solutions of (11.1.2.16) are global:  $J(\mathbf{y}_0) = \mathbb{R}$  for all  $\mathbf{y}_0 \in D$ .

Now we return to the study of a generic ODE (11.1.0.2) instead of an IVP (11.1.2.2). We do this by temporarily changing the perspective: we fix a “time of interest”  $t \in \mathbb{R} \setminus \{0\}$  and follow all trajectories for the duration  $t$ . This induces a mapping of points in state space:

$$\geq \text{mapping } \Phi^t : \begin{cases} D & \mapsto D \\ \mathbf{y}_0 & \mapsto \mathbf{y}(t) \end{cases} , \quad t \mapsto \mathbf{y}(t) \text{ solution of IVP (11.1.2.16)} ,$$

This is a well-defined mapping of the state space into itself, by Thm. 11.1.2.15 and Ass. 11.1.3.1.

Now, we may also let  $t$  vary, which spawns a *family* of mappings  $\{\Phi^t\}$  of the state space into itself. However, it can also be viewed as a mapping with two arguments, a duration  $t$  and an initial state value  $\mathbf{y}_0$ !

### Definition 11.1.3.2. Evolution operator/mapping

Under Ass. 11.1.3.1 the mapping

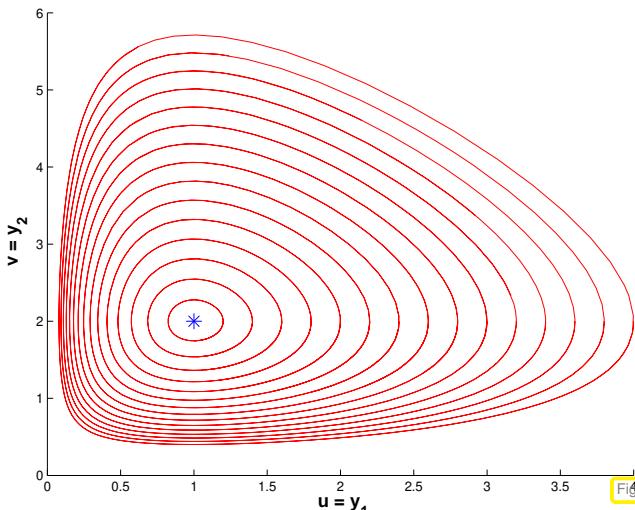
$$\Phi : \begin{cases} \mathbb{R} \times D & \mapsto D \\ (t, \mathbf{y}_0) & \mapsto \Phi^t \mathbf{y}_0 := \mathbf{y}(t) \end{cases},$$

where  $t \mapsto \mathbf{y}(t) \in C^1(\mathbb{R}, \mathbb{R}^d)$  is the unique (global) solution of the IVP  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ ,  $\mathbf{y}(0) = \mathbf{y}_0$ , is the **evolution operator**/mapping for the autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ .

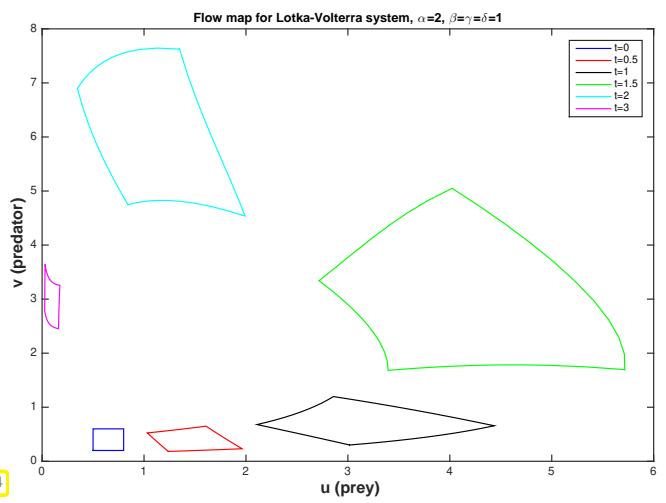
Note that  $t \mapsto \Phi^t \mathbf{y}_0$  describes the solution of  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  for  $\mathbf{y}(0) = \mathbf{y}_0$  (a trajectory). Therefore, by virtue of definition, we have

$$\frac{\partial \Phi}{\partial t}(t, \mathbf{y}) = \mathbf{f}(\Phi^t \mathbf{y}). \quad (11.1.3.3)$$

**EXAMPLE 11.1.3.4 (Evolution operator for Lotka-Volterra ODE (11.1.1.6))** For  $d = 2$  the action of an evolution operator can be visualized by tracking the movement of point sets in state space. Here this is done for the Lotka-Volterra ODE (11.1.1.6):



trajectories  $t \mapsto \Phi^t \mathbf{y}_0$



Think of  $\mathbf{y} \in \mathbb{R}^2 \mapsto \mathbf{f}(\mathbf{y}) \in \mathbb{R}^2$  as the velocity of the surface of a fluid. Specks of floating dust will be carried along by the fluid, patches of dust covering parts of the surface will move and deform over time. This can serve as a “mental image” of  $\Phi$ . ■

Given an evolution operator, we can recover the right-hand side function  $\mathbf{f}$  of the underlying autonomous ODE as  $\mathbf{f}(\mathbf{y}) = \frac{\partial \Phi}{\partial t}(0, \mathbf{y})$ : There is a one-to-one relationship between ODEs and their evolution operators, and those are the key objects behind an ODE.

An ODE “encodes” an evolution operator.

Understanding the concept of evolution operators is indispensable for numerical integration, that is the construction of numerical methods for the solution of IVPs for ODEs:

Numerical integration is concerned with the approximation of evolution operators.

**Remark 11.1.3.5 (Group property of autonomous evolutions)** Under Ass. 11.1.3.1 the evolution operator gives rise to a group of mappings  $D \mapsto D$ :

$$\Phi^s \circ \Phi^t = \Phi^{s+t} , \quad \Phi^{-t} \circ \Phi^t = Id \quad \forall t \in \mathbb{R} . \quad (11.1.3.6)$$

This is a consequence of the uniqueness theorem Thm. 11.1.2.15. It is also intuitive: following an evolution up to time  $t$  and then for some more time  $s$  leads us to the same final state as observing it for the whole time  $s+t$ .  $\square$

## 11.2 Introduction: Polygonal Approximation Methods

In this section we will see the first simple methods for the numerical integration (= solution) of initial-value problems (IVPs). We target an initial value problem (11.1.2.2) for a first-order ordinary differential equation

$$\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}) , \quad \mathbf{y}(t_0) = \mathbf{y}_0 . \quad (11.1.2.2)$$

As usual, the right hand side function  $\mathbf{f}$  may be given only in *procedural form*, for instance, in a C++ code as an functor object providing an evaluation operator

```
Eigen::VectorXd operator (double t, const Eigen::VectorXd &y) const;
```

see Rem. 5.1.0.8. An evaluation of  $\mathbf{f}$  may involve costly computations.

**§11.2.0.1 (Objectives of numerical integration)** Two basic tasks can be identified in the field of numerical integration = approximate solution of initial value problems for ODEs (Please distinguish from “numerical quadrature”, see Chapter 7.):

- (I) Given initial time  $t_0$ , final time  $T$ , and initial state  $\mathbf{y}_0$  compute an approximation of  $\mathbf{y}(T)$ , where  $t \mapsto \mathbf{y}(t)$  is the solution of (11.1.2.2). A corresponding function in C++ could look like

```
State solveivp(double t0, double T, State y0);
```

Here **State** is a type providing a fixed size or variable size vector  $\in \mathbb{R}^d$ , e.g.,

```
using State = Eigen::Matrix<double, statedim, 1>;
```

Here `statedim` is the dimension  $d$  of the state space that has to be known at compile time.

- (II) Output an approximate solution  $t \rightarrow \mathbf{y}_h(t)$  of (11.1.2.2) on  $[t_0, T]$  up to final time  $T \neq t_0$  for “all times”  $t \in [t_0, T]$  (in practice, of course, only for finitely many times  $t_0 = \tau_0 < \tau_1 < \tau_2 < \dots < \tau_{m-1} < \tau_m = T$ , consecutively)

```
std::vector<State>
solveivp(State y0, const std::vector<double> &tauvec);
```

This is the “plot solution” task.  $\square$

This section presents three methods that provide a piecewise linear, that is, “polygonal” approximation of solution trajectories  $t \mapsto \mathbf{y}(t)$ , cf. Ex. 5.1.0.12 for  $d = 1$ .

**§11.2.0.2 (Temporal mesh)** As in Section 6.5.1 the polygonal approximation in this section will be based on a (temporal) mesh ( $\rightarrow$  § 6.5.0.1)

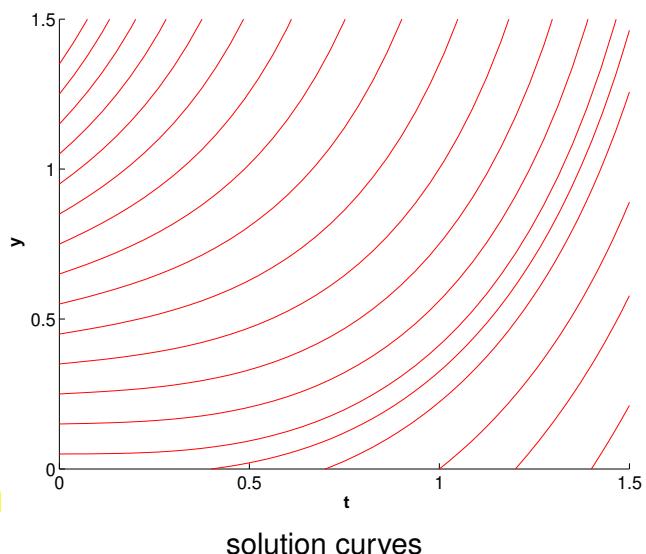
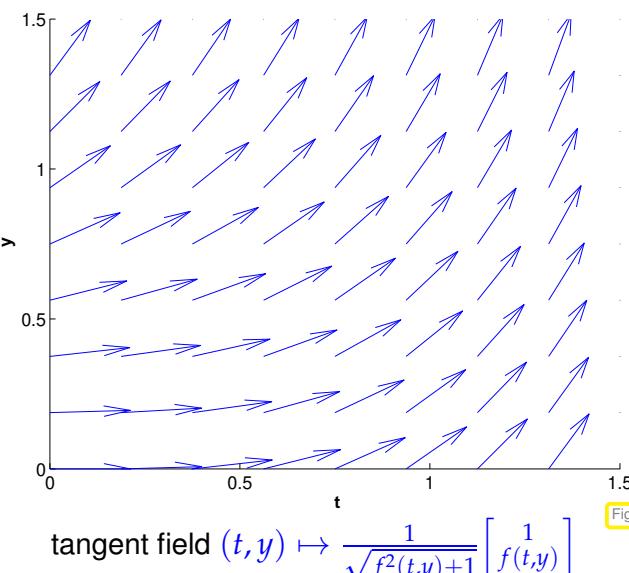
$$\mathcal{M} := \{t_0 < t_1 < t_2 < \dots < t_{N-1} < t_N := T\} \subset [t_0, T], \quad (11.2.0.3)$$

covering the time interval of interest between initial time  $t_0$  and final time  $T > t_0$ . We assume that the interval of interest is contained in the domain of definition of the solution of the IVP:  $[t_0, T] \subset J(t_0, \mathbf{y}_0)$ .  $\square$

## 11.2.1 Explicit Euler method

**EXAMPLE 11.2.1.1 (Tangent field and solution curves)** For  $d = 1$  polygonal methods can be constructed by geometric considerations in the  $t - y$  plane, a model for the extended state space. We explain this for the Riccati differential equation, a scalar ODE:

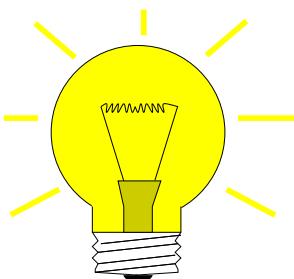
$$\dot{y} = f(t, y) := y^2 + t^2 \quad \triangleright \quad d = 1, \quad I, D = \mathbb{R}^+. \quad (11.2.1.2)$$



$$\text{tangent field } (t, y) \mapsto \frac{1}{\sqrt{f^2(t, y) + 1}} \begin{bmatrix} 1 \\ f(t, y) \end{bmatrix}$$

solution curves

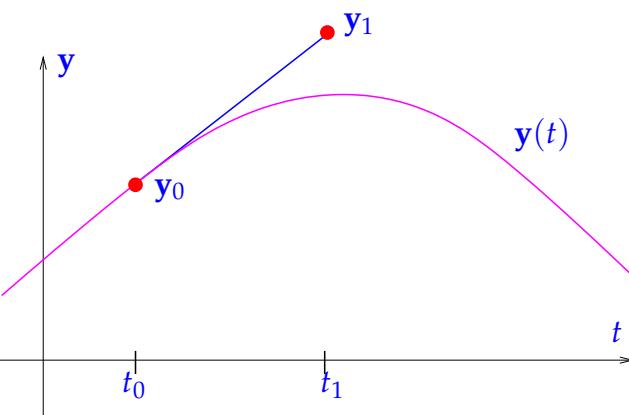
The solution curves run tangentially to the tangent field in each point of the extended state space.  $\square$



Idea:

“follow the tangents over short periods of time”

- ① **timestepping**: successive approximation of evolution on *mesh intervals*  $[t_{k-1}, t_k]$ ,  $k = 1, \dots, N$ ,  $t_N := T$ ,
- ② approximation of solution on  $[t_{k-1}, t_k]$  by **tangent line to solution trajectory through  $(t_{k-1}, \mathbf{y}_{k-1})$** .



**explicit Euler method** (Euler 1768)

$\triangleleft$  First step of explicit Euler method ( $d = 1$ ):

Slope of tangent  $= f(t_0, \mathbf{y}_0)$

$\mathbf{y}_1$  serves as initial value for next step!

See also [Han02, Ch. 74], [DR08, Alg. 11.4]

### EXAMPLE 11.2.1.3 (Visualization of explicit Euler method)

Temporal mesh

$$\mathcal{M} := \{t_j := j/5 : j = 0, \dots, 5\}.$$

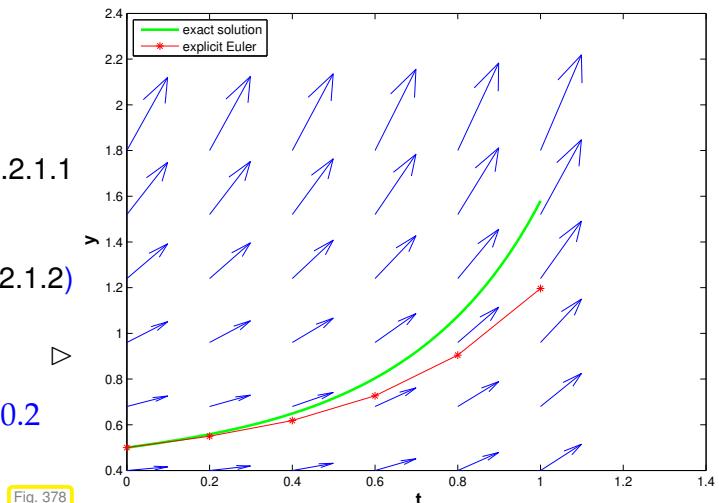
IVP for Riccati differential equation, see Ex. 11.2.1.1

$$\dot{y} = y^2 + t^2. \quad (11.2.1.2)$$

Here:  $y_0 = \frac{1}{2}$ ,  $t_0 = 0$ ,  $T = 1$ ,

— ≈ “Euler polygon” for uniform timestep  $h = 0.2$

↔ ≈ tangent field of Riccati ODE



↓

Formula: When applied to a general IVP of the form (11.1.2.2) the explicit Euler method generates a sequence  $(\mathbf{y}_k)_{k=0}^N$  by the recursion

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h_k \mathbf{f}(t_k, \mathbf{y}_k), \quad k = 0, \dots, N-1, \quad (11.2.1.4)$$

with local (size of) timestep (stepsize)  $h_k := t_{k+1} - t_k$ .

### Remark 11.2.1.5 (Explicit Euler method as difference scheme)

One can obtain (11.2.1.4) by approximating the derivative  $\frac{d}{dt}$  by a forward difference quotient on the (temporal) mesh  $\mathcal{M} := \{t_0, t_1, \dots, t_N\}$ :

$$\dot{\mathbf{y}}(t_k) \approx \frac{\mathbf{y}(t_k + h_k) - \mathbf{y}(t_k)}{h_k}$$

▶  $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}) \iff \frac{\mathbf{y}_{k+1} - \mathbf{y}_k}{h_k} = \mathbf{f}(t_k, \mathbf{y}_h(t_k)), \quad k = 0, \dots, N-1. \quad (11.2.1.6)$

**Difference schemes** follow a simple policy for the *discretization* of differential equations: replace all derivatives by difference quotients connecting solution values on a set of discrete points (the mesh). ↓

**Remark 11.2.1.7 (Output of explicit Euler method)** To begin with, the explicit Euler recursion (11.2.1.4) produces a sequence  $\mathbf{y}_0, \dots, \mathbf{y}_N$  of states. How does it deliver on the task (I) and (II) stated in § 11.2.0.1? By “geometric insight” we expect

$$\mathbf{y}_k \approx \mathbf{y}(t_k).$$

(As usual, we use the notation  $t \mapsto \mathbf{y}(t)$  for the exact solution of an IVP.)

Task (I): Easy, because  $\mathbf{y}_N$  already provides an approximation of  $\mathbf{y}(T)$ .

Task (II): The trajectory  $t \mapsto \mathbf{y}(t)$  is approximated by the piecewise linear function (“Euler polygon”)

$$\mathbf{y}_h : [t_0, t_N] \rightarrow \mathbb{R}^d, \quad \mathbf{y}_h(t) := \mathbf{y}_k \frac{t_{k+1} - t}{t_{k+1} - t_k} + \mathbf{y}_{k+1} \frac{t - t_k}{t_{k+1} - t_k} \quad \text{for } t \in [t_k, t_{k+1}], \quad (11.2.1.8)$$

see Fig. 378. This function can easily be sampled on any grid of  $[t_0, t_N]$ . In fact, it is the  $\mathcal{M}$ -piecewise linear interpolant of the data points  $(t_k, \mathbf{y}_k)$ ,  $k = 0, \dots, N$ , see Section 5.3.2).

The same considerations apply to the methods discussed in the next two sections and will not be repeated there.  $\square$

### 11.2.2 Implicit Euler method

Why forward difference quotient and not backward difference quotient? Let's try!

On (temporal) mesh  $\mathcal{M} := \{t_0, t_1, \dots, t_N\}$  we obtain

$$\dot{\mathbf{y}} = f(t, \mathbf{y}) \iff \frac{\mathbf{y}_{k+1} - \mathbf{y}_k}{h_k} = f(t_{k+1}, \mathbf{y}_h(t_{k+1})) , \quad k = 0, \dots, N-1 .$$

backward difference quotient

(11.2.2.1)

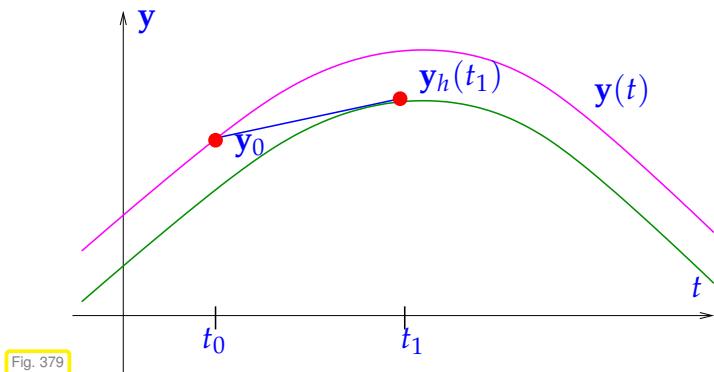
This leads to another simple timestepping scheme analogous to (11.2.1.4):

$$\boxed{\mathbf{y}_{k+1} = \mathbf{y}_k + h_k \mathbf{f}(t_{k+1}, \mathbf{y}_{k+1}) , \quad k = 0, \dots, N-1} , \quad (11.2.2.2)$$

with local timestep (stepsize)  $h_k := t_{k+1} - t_k$ .

(11.2.2.2) = implicit Euler method

Note: (11.2.2.2) requires solving a (possibly non-linear) system of equations to obtain  $\mathbf{y}_{k+1}$ !  
 (► Terminology “implicit”)



Geometry of implicit Euler method:

Approximate solution through  $(t_0, y_0)$  on  $[t_0, t_1]$  by

- straight line through  $(t_0, y_0)$

- with slope  $f(t_1, y_1)$

- ▷ —  $\hat{=}$  trajectory through  $(t_0, y_0)$ ,  
 —  $\hat{=}$  trajectory through  $(t_1, y_1)$ ,  
 —  $\hat{=}$  tangent at — in  $(t_1, y_1)$ .

#### Remark 11.2.2.3 (Feasibility of implicit Euler timestepping)

Issue: Is (11.2.2.2) well defined, that is, can we solve it for  $\mathbf{y}_{k+1}$  and is this solution unique?

Intuition: for small timesteps  $h > 0$  the right hand side of (11.2.2.2) is a “small perturbation of the identity”.

Formal: Consider an autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ , assume a continuously differentiable right hand side function  $\mathbf{f}, \mathbf{f} \in C^1(D, \mathbb{R}^d)$ , and regard (11.2.2.2) as an  $h$ -dependent non-linear system of equations:

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h_k \mathbf{f}(t_{k+1}, \mathbf{y}_{k+1}) \Leftrightarrow G(h, \mathbf{y}_{k+1}) = 0 \quad \text{with} \quad G(h, \mathbf{z}) := \mathbf{z} - h \mathbf{f}(t_{k+1}, \mathbf{z}) - \mathbf{y}_k .$$

To investigate the solvability of this non-linear equation we start with an observation about a partial derivative of  $G$ :

$$\frac{dG}{d\mathbf{z}}(h, \mathbf{z}) = \mathbf{I} - h D_{\mathbf{y}} \mathbf{f}(t_{k+1}, \mathbf{z}) \Rightarrow \frac{dG}{d\mathbf{z}}(0, \mathbf{z}) = \mathbf{I} .$$

In addition,  $G(0, \mathbf{y}_k) = 0$ . Next, recall the **implicit function theorem** [Str09, Thm. 7.8.1]:

#### Theorem 11.2.2.4. Implicit function theorem

Let  $G = G(\mathbf{x}, \mathbf{y})$  a continuously differentiable function of  $\mathbf{x} \in \mathbb{R}^k$  and  $\mathbf{y} \in \mathbb{R}^\ell$ , defined on the open set  $\Omega \subset \mathbb{R}^k \times \mathbb{R}^\ell$  with values in  $\mathbb{R}^\ell$ :  $G : \Omega \subset \mathbb{R}^k \times \mathbb{R}^\ell \rightarrow \mathbb{R}^\ell$ .

Assume that  $G$  has a zero in  $\mathbf{z}_0 := \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{y}_0 \end{bmatrix} \in \Omega$ ,  $\mathbf{x}_0 \in \mathbb{R}^k$ ,  $\mathbf{y}_0 \in \mathbb{R}^\ell$ :  $G(\mathbf{z}_0) = 0$ .

If the Jacobian  $\frac{\partial G}{\partial \mathbf{y}}(\mathbf{p}_0) \in \mathbb{R}^{\ell, \ell}$  is **invertible**, then there is an open neighborhood  $U$  of  $\mathbf{x}_0 \in \mathbb{R}^k$  and a continuously differentiable function  $\mathbf{g} : U \rightarrow \mathbb{R}^\ell$  such that

$$\mathbf{g}(\mathbf{x}_0) = \mathbf{y}_0 \quad \text{and} \quad G(\mathbf{x}, \mathbf{g}(\mathbf{x})) = 0 \quad \forall \mathbf{x} \in U .$$

For **sufficiently small**  $|h|$  it permits us to conclude that the equation  $G(h, \mathbf{z}) = 0$  defines a continuous function  $\mathbf{g} = \mathbf{g}(h)$  with  $\mathbf{g}(0) = \mathbf{y}_k$ .

➤ for **sufficiently small**  $h > 0$  the equation (11.2.2.2) has a unique solution  $\mathbf{y}_{k+1}$ . □

### 11.2.3 Implicit midpoint method

Beside using forward or backward difference quotients, the derivative  $\dot{\mathbf{y}}$  can also be approximated by the **symmetric difference quotient**, see also (5.2.3.18),

$$\dot{\mathbf{y}}(t) \approx \frac{\mathbf{y}(t+h) - \mathbf{y}(t-h)}{2h} , \quad h > 0 . \quad (11.2.3.1)$$

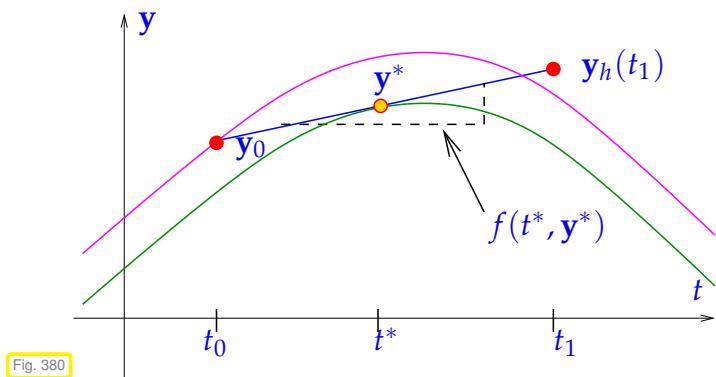
The idea is to apply this formula in  $t = \frac{1}{2}(t_k + t_{k+1})$  with  $h = h_k/2$ , which transforms the ODE into

$$\dot{\mathbf{y}} = f(t, \mathbf{y}) \iff \frac{\mathbf{y}_{k+1} - \mathbf{y}_k}{h_k} = f\left(\frac{1}{2}(t_k + t_{k+1}), \mathbf{y}_h\left(\frac{1}{2}(t_{k+1} + t_{k+1})\right)\right) , \quad k = 0, \dots, N-1 . \quad (11.2.3.2)$$

The trouble is that the value  $\mathbf{y}_h\left(\frac{1}{2}(t_{k+1} + t_{k+1})\right)$  does not seem to be available, unless we recall that the approximate trajectory  $t \mapsto \mathbf{y}_h(t)$  is supposed to be piecewise linear, which implies  $\mathbf{y}_h\left(\frac{1}{2}(t_{k+1} + t_{k+1})\right) = \frac{1}{2}(\mathbf{y}_h(t_k) + \mathbf{y}_h(t_{k+1}))$ . This gives the recursion formula for the **implicit midpoint method** in analogy to (11.2.1.4) and (11.2.2.2):

$$\boxed{\mathbf{y}_{k+1} = \mathbf{y}_k + h_k \mathbf{f}\left(\frac{1}{2}(t_k + t_{k+1}), \frac{1}{2}(\mathbf{y}_k + \mathbf{y}_{k+1})\right) , \quad k = 0, \dots, N-1} , \quad (11.2.3.3)$$

with local **timestep** (stepsize)  $h_k := t_{k+1} - t_k$ .



Implicit midpoint method: a geometric view:

Approximate trajectory through  $(t_0, \mathbf{y}_0)$  on  $[t_0, t_1]$  by

- straight line through  $(t_0, \mathbf{y}_0)$
  - with slope  $f(t^*, \mathbf{y}^*)$ , where  
 $t^* := \frac{1}{2}(t_0 + t_1)$ ,  $\mathbf{y}^* = \frac{1}{2}(\mathbf{y}_0 + \mathbf{y}_{h,t_1})$
- $\triangleq$
- $\hat{=}$  trajectory through  $(t_0, \mathbf{y}_0)$ ,
  - $\hat{=}$  trajectory through  $(t^*, \mathbf{y}^*)$ ,
  - $\hat{=}$  tangent at — in  $(t^*, \mathbf{y}^*)$ .

As in the case of (11.2.2.2), also (11.2.3.3) entails solving a (non-linear) system of equations in order to obtain  $\mathbf{y}_{k+1}$ . Rem. 11.2.2.3 also holds true in this case: for sufficiently small  $h$  (11.2.3.3) will have a unique solution  $\mathbf{y}_{k+1}$ , which renders the recursion well defined.

## 11.3 General single step methods

Now we fit the numerical schemes introduced in the previous section into a more general class of methods for the solution of (autonomous) initial value problems (11.1.2.16) for ODEs. Throughout we assume that all times considered belong to the domain of definition of the unique solution  $t \rightarrow \mathbf{y}(t)$  of (11.1.2.16), that is, for  $T > 0$  we take for granted  $[0, T] \subset J(\mathbf{y}_0)$  (temporal domain of definition of the solution of an IVP is explained in § 11.1.2.17).

### 11.3.1 Definition

**§11.3.1.1 (Discrete evolution operators)** Recall the Euler methods for autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ :

$$\begin{aligned} \text{explicit Euler: } \mathbf{y}_{k+1} &= \mathbf{y}_k + h_k \mathbf{f}(\mathbf{y}_k), \\ \text{implicit Euler: } \mathbf{y}_{k+1} &:= \mathbf{y}_k + h_k \mathbf{f}(\mathbf{y}_{k+1}), \quad h_k := t_{k+1} - t_k. \end{aligned}$$

Both formulas, for sufficiently small  $h_k$  ( $\rightarrow$  Rem. 11.2.2.3), provide a mapping

$$(\mathbf{y}_k, h_k) \mapsto \Psi(h, \mathbf{y}_k) := \mathbf{y}_{k+1}. \quad (11.3.1.2)$$

If  $\mathbf{y}_0$  is the initial value, then  $\mathbf{y}_1 := \Psi(h, \mathbf{y}_0)$  can be regarded as an approximation of  $\mathbf{y}(h)$ , the value returned by the evolution operator ( $\rightarrow$  Def. 11.1.3.2) for  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  applied to  $\mathbf{y}_0$  over the period  $h$ .  $\mathbf{y}(t_k)$ :

$$\mathbf{y}_1 = \Psi(h, \mathbf{y}_0) \iff \mathbf{y}(h) = \Phi^h \mathbf{y}_0 \Rightarrow \boxed{\Psi(h, \mathbf{y}) \approx \Phi^h \mathbf{y}}, \quad (11.3.1.3)$$

In a sense the polygonal approximation methods as based on approximations for the evolution operator associated with the ODE.

This is what every single step method does: it tries to approximate the evolution operator  $\Phi$  for an ODE by a mapping  $\Psi$  of the kind as described in (11.3.1.2).

→ A mapping  $\Psi$  as in (11.3.1.2) is called (a) **discrete evolution** operator.

☞ Notation: for discrete evolutions we often write  $\Psi^h \mathbf{y} := \Psi(h, \mathbf{y})$

**Remark 11.3.1.4 (Discretization)** The adjective “**discrete**” used above designates (components of) methods that attempt to approximate the solution of an IVP by a sequence of finitely many states. “**Discretization**” is the process of converting an ODE into a discrete model. This parlance is adopted for all procedures that reduce a “continuous model” involving ordinary or partial differential equations to a form with a finite number of unknowns.  $\square$

Above we identified the discrete evolutions underlying the polygonal approximation methods. Vice versa, a mapping  $\Psi$  as given in (11.3.1.2) defines a single step method.

**Definition 11.3.1.5. Single step method (for autonomous ODE)**  $\rightarrow$  [QSS00, Def. 11.2]

Given a **discrete evolution**  $\Psi : \Omega \subset \mathbb{R} \times D \mapsto \mathbb{R}^d$ , an initial state  $\mathbf{y}_0$ , and a temporal mesh  $\mathcal{M} := \{0 =: t_0 < t_1 < \dots < t_N := T\}$  the recursion

$$\mathbf{y}_{k+1} := \Psi(t_{k+1} - t_k, \mathbf{y}_k), \quad k = 0, \dots, N-1, \quad (11.3.1.6)$$

defines a **single-step method** (SSM) for the autonomous IVP  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ ,  $\mathbf{y}(0) = \mathbf{y}_0$  on the interval  $[0, T]$ .

- ☞ In a sense, a single step method defined through its associated discrete evolution does not approximate a concrete initial value problem, but tries to approximate an ODE in the form of its evolution operator.

In C++ a discrete evolution operator can be incarnated by a functor type offering an evaluation operator

```
State operator(double h, const State &y) const;
```

see § 11.2.0.1 for the **State** data type.

**Remark 11.3.1.7 (Discrete evolutions for non-autonomous ODEs)** The concept of single step method according to Def. 11.3.1.5 can be generalized to non-autonomous ODEs, which leads to recursions of the form:

$$\mathbf{y}_{k+1} := \Psi(t_k, t_{k+1}, \mathbf{y}_k), \quad k = 0, \dots, N-1,$$

for a discrete evolution operator  $\Psi$  defined on  $I \times I \times D$ .  $\square$

**§11.3.1.8 (Consistent single step methods)** All meaningful single step methods turn out to be modifications of the explicit Euler method (11.2.1.4).

**Consistent discrete evolution**

The discrete evolution  $\Psi$  defining a single step method according to Def. 11.3.1.5 and (11.3.1.6) for the autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  invariably is of the form

$$\Psi^h \mathbf{y} = \mathbf{y} + h \psi(h, \mathbf{y}) \quad \text{with} \quad \begin{aligned} \psi : I \times D &\rightarrow \mathbb{R}^d \text{ continuous,} \\ \psi(0, \mathbf{y}) &= \mathbf{f}(\mathbf{y}). \end{aligned} \quad (11.3.1.10)$$

**Definition 11.3.1.11. Consistent single step methods**

A single step method according to Def. 11.3.1.5 based on a discrete evolution of the form (11.3.1.10) is called **consistent** with the ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ .

**EXAMPLE 11.3.1.12 (Consistency of implicit midpoint method)** The discrete evolution  $\Psi$  and, hence, the function  $\psi = \psi(h, \mathbf{y})$  for the implicit midpoint method are defined only implicitly, of course. Thus, consistency cannot immediately be seen from a formula for  $\psi$ .

We examine consistency of the implicit midpoint method defined by

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h\mathbf{f}\left(\frac{1}{2}(t_k + t_{k+1}), \frac{1}{2}(\mathbf{y}_k + \mathbf{y}_{k+1})\right), \quad k = 0, \dots, N-1. \quad (11.2.3.3)$$

Assume that

- the right hand side function  $\mathbf{f}$  is smooth, at least  $\mathbf{f} \in C^1(D)$ ,
- and  $|h|$  is sufficiently small to guarantee the existence of a solution  $\mathbf{y}_{k+1}$  of (11.2.3.3), see Rem. 11.2.2.3.

The idea is to verify (11.3.1.10) by formal Taylor expansion of  $\mathbf{y}_{k+1}$  in  $h$ . To that end we plug (11.2.3.3) into itself and rely on Taylor expansion of  $\mathbf{f}$ :

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h\mathbf{f}\left(\frac{1}{2}(\mathbf{y}_k + \mathbf{y}_{k+1})\right) \stackrel{(11.2.3.3)}{=} \mathbf{y}_k + h\underbrace{\mathbf{f}\left(\mathbf{y}_k + \frac{1}{2}h\mathbf{f}\left(\frac{1}{2}(\mathbf{y}_k + \mathbf{y}_{k+1})\right)\right)}_{=\psi(h, \mathbf{y}_k)}.$$

Since, by the implicit function theorem Thm. 11.2.2.4,  $\mathbf{y}_{k+1}$  continuously depends on  $h$  and  $\mathbf{y}_k$ ,  $\psi(h, \mathbf{y}_k)$  has the desired properties, in particular  $\psi(0, \mathbf{y}) = \mathbf{f}(\mathbf{y})$  is clear.  $\square$

**Remark 11.3.1.13 (Notation for single step methods)** Many authors specify a single step method by writing down the first step for a general stepsize  $h$

$$\mathbf{y}_1 = (\text{implicit expression in } \mathbf{y}_0, h \text{ and } \mathbf{f}),$$

for instance, for the implicit midpoint rule

$$\mathbf{y}_1 = \mathbf{y}_0 + h\mathbf{f}\left(\frac{1}{2}(\mathbf{y}_0 + \mathbf{y}_1)\right).$$

Actually, this fixes the underlying discrete evolution. Also this course will sometimes adopt this practice.  $\square$

**§11.3.1.14 (Output of single step methods)** Here we resume and continue the discussion of Rem. 11.2.1.7 for general single step methods according to Def. 11.3.1.5. Assuming unique solvability of the systems of equations faced in each step of an **implicit** method, every single step method based on a mesh  $\mathcal{M} = \{0 = t_0 < t_1 < \dots < t_N := T\}$  produces a finite sequence  $(\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_N)$  of states, where the first agrees with the initial state  $\mathbf{y}_0$ .

We expect that the states provide a pointwise approximation of the solution trajectory  $t \rightarrow \mathbf{y}(t)$ :

$$\mathbf{y}_k \approx \mathbf{y}(t_k), \quad k = 1, \dots, N.$$

Thus task (I) from § 11.2.0.1, computing an approximation for  $\mathbf{y}(T)$ , is again easy: output  $\mathbf{y}_N$  as an approximation of  $\mathbf{y}(T)$ .

Task (II) from § 11.2.0.1, computing the solution trajectory, requires **interpolation** of the data points  $(t_k, \mathbf{y}_k)$  using some of the techniques presented in Chapter 5. The natural option is  $\mathcal{M}$ -piecewise polynomial interpolation, generalizing the polygonal approximation (11.2.1.8) used in Section 11.2.

Note that from the ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  the derivatives  $\dot{\mathbf{y}}_h(t_k) = \mathbf{f}(\mathbf{y}_k)$  are available without any further approximation. This facilitates **cubic Hermite interpolation** ( $\rightarrow$  Def. 5.4.1.1), which yields

$$\mathbf{y}_h \in C^1([0, T]): \quad \mathbf{y}_h|_{[x_{k-1}, x_k]} \in \mathcal{P}_3, \quad \mathbf{y}_h(t_k) = \mathbf{y}_k, \quad \frac{d\mathbf{y}_h}{dt}(t_k) = \mathbf{f}(\mathbf{y}_k).$$

Summing up, an approximate trajectory  $t \mapsto \mathbf{y}_h(t)$  is built in two stages:

- (i) Compute sequence  $(\mathbf{y}_k)_k$  by running the single step method.
- (ii) **Post-process** the obtained sequence, usually by applying interpolation, to get  $\mathbf{y}_h$ .

□

### 11.3.2 (Asymptotic) Convergence of Single-Step Methods

Of course, the accuracy of the solution sequence  $(\mathbf{y}_k)_k$  obtained by a particular single-step method ( $\rightarrow$  Def. 11.3.1.5) is a central concern. This motivates studying the dependence of suitable norms of the so-called discretization error on the choice of temporal mesh  $\mathcal{M}$ .



*Supplementary literature.* See [DR08, Sect. 11.5] and [QSS00, Sect. 11.3] for related presentations.

**§11.3.2.1 (Discretization error of single step methods)** Errors in numerical integration are called **discretization errors**, cf. Rem. 11.3.1.4.

Depending on the objective of numerical integration as stated in § 11.2.0.1 different (norms of) discretization errors are of interest:

- (I) If only the solution at final time  $T$  is sought, the relevant norm of the discretization error is

$$\epsilon_N := \|\mathbf{y}(T) - \mathbf{y}_N\|,$$

where  $\|\cdot\|$  is some vector norm on  $\mathbb{R}^d$ .

- (II) If we want to approximate the solution trajectory for (11.1.2.16) the discretization error is the function

$$t \mapsto \mathbf{e}(t) , \quad \mathbf{e}(t) := \mathbf{y}(t) - \mathbf{y}_h(t) ,$$

where  $t \mapsto \mathbf{y}_h(t)$  is the approximate trajectory obtained by post-processing, see § 11.3.1.14. In this case accuracy of the method is gauged by looking at norms of the function  $\mathbf{e}$ , see § 5.2.4.3 for examples.

- (III) Between (I) and (II) is the pointwise discretization error, which is the sequence (**grid function**)

$$\mathbf{e} : \mathcal{M} \rightarrow D , \quad \mathbf{e}_k := \mathbf{y}(t_k) - \mathbf{y}_k , \quad k = 0, \dots, N . \quad (11.3.2.2)$$

In this case one usually examines the maximum error in the mesh points

$$\|(\mathbf{e})\|_\infty := \max_{k \in \{1, \dots, N\}} \|\mathbf{e}_k\| ,$$

where  $\|\cdot\|$  is a suitable vector norm on  $\mathbb{R}^d$ , customarily the Euclidean vector norm.

□

**§11.3.2.3 (Asymptotic convergence of single step methods)** Once the discrete evolution  $\Psi$  associated with the ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  is specified, the single step method according to Def. 11.3.1.5 is fixed. The only way to control the accuracy of the solution  $\mathbf{y}_N$  or  $t \mapsto \mathbf{y}_h(t)$  is through the selection of the mesh  $\mathcal{M} = \{0 = t_0 < t_1 < \dots < t_N = T\}$ .

Hence we study convergence of single step methods for **families of meshes**  $\{\mathcal{M}_\ell\}$  and track the decay of (a norm) of the discretization error ( $\rightarrow$  § 11.3.2.1) as a function of the number  $N := |\mathcal{M}|$  of mesh points.

In other words, we examine **h-convergence**. We already did this in the case of piecewise polynomial interpolation in Section 6.5.1 and composite numerical quadrature in Section 7.4.

When investigating asymptotic convergence of single step methods we often resort to families of **equidistant** meshes of  $[0, T]$ :

$$\mathcal{M}_N := \{t_k := \frac{k}{N}T : k = 0, \dots, N\}. \quad (11.3.2.4)$$

We also call this the use of **uniform** timesteps of size  $h := \frac{T}{N}$ . □

### EXPERIMENT 11.3.2.5 (Speed of convergence of Euler methods)

The setting for this experiment is as follows:

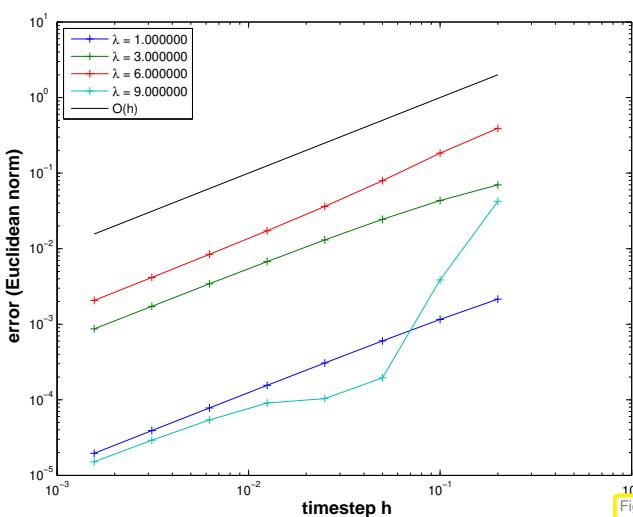
- ◆ We consider the following IVP for the logistic ODE, see Ex. 11.1.1.1

$$\dot{y} = \lambda y(1 - y), \quad y(0) = 0.01.$$

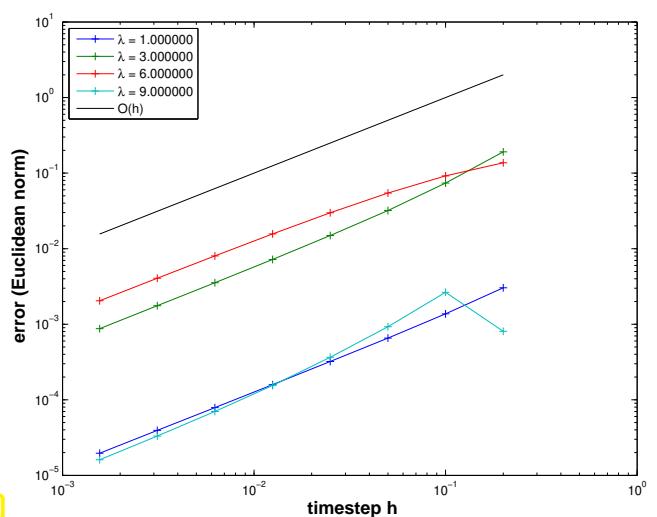
- ◆ We apply explicit and implicit Euler methods (11.2.1.4)/(11.2.2.2) with uniform timestep  $h = 1/N$ ,  $N \in \{5, 10, 20, 40, 80, 160, 320, 640\}$ .
- ◆ Monitored: Error at final time  $E(h) := |y(1) - y_N|$

We are mainly interested in the qualitative nature of the asymptotic convergence as  $h \rightarrow 0$  in the sense of the types of convergence introduced in Def. 6.1.2.7 with  $n$  replaced with  $h^{-1}$ . Abbreviating some error norm with  $T = T(h)$ , convergence can be classified as follows:

$$\begin{aligned} \exists p > 0: \quad T(h) \leq h^p &: \text{algebraic convergence, with order } p > 0, \quad \forall h > 0; \\ \exists 0 < q < 1: \quad T(h) \leq q^{1/h} &: \text{exponential convergence,} \end{aligned}$$



explicit Euler method



implicit Euler method

$\blacktriangleright O(N^{-1}) = O(h)$  algebraic convergence with order 1 in both cases for  $h \rightarrow 0$

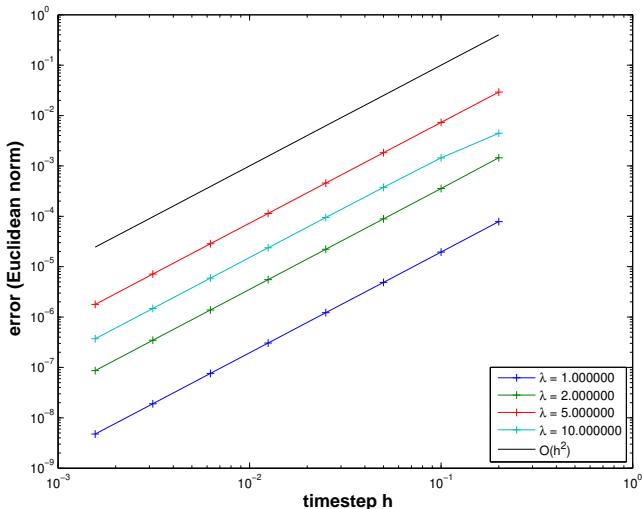


Fig. 383

### implicit midpoint method

Parlance: based on the observed rate of algebraic convergence, the two Euler methods are said to “converge with first order”, whereas the implicit midpoint method is called “second-order convergent”.

The observations made for polygonal timestepping methods reflect a general pattern:

#### Algebraic convergence of single step methods

Consider numerical integration of an initial value problem

$$\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}) , \quad \mathbf{y}(t_0) = \mathbf{y}_0 , \quad (11.1.2.2)$$

with sufficiently smooth right hand side function  $\mathbf{f} : I \times D \rightarrow \mathbb{R}^d$ .

Then conventional single step methods ( $\rightarrow$  Def. 11.3.1.5) will enjoy *algebraic convergence in the meshwidth*, more precisely, see [DR08, Thm. 11.25],

there is a  $p \in \mathbb{N}$  such that the sequence  $(\mathbf{y}_k)_k$  generated by the single step method for  $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$  on a mesh  $\mathcal{M} := \{t_0 < t_1 < \dots < t_N = T\}$  satisfies

$$\max_k \|\mathbf{y}_k - \mathbf{y}(t_k)\| \leq Ch^p \quad \text{for } h := \max_{k=1,\dots,N} |t_k - t_{k-1}| \rightarrow 0 , \quad (11.3.2.7)$$

with  $C > 0$  independent of  $\mathcal{M}$

#### Definition 11.3.2.8. Order of a single step method

The minimal integer  $p \in \mathbb{N}$  for which (11.3.2.7) holds for a single step method when applied to an ODE with (sufficiently) smooth right hand side, is called the **order** of the method.

As in the case of quadrature rules ( $\rightarrow$  Def. 7.3.0.1) their order is the principal intrinsic indicator for the “quality” of a single step method.

**§11.3.2.9 (Convergence analysis for the explicit Euler method [Han02, Ch. 74])** We consider the

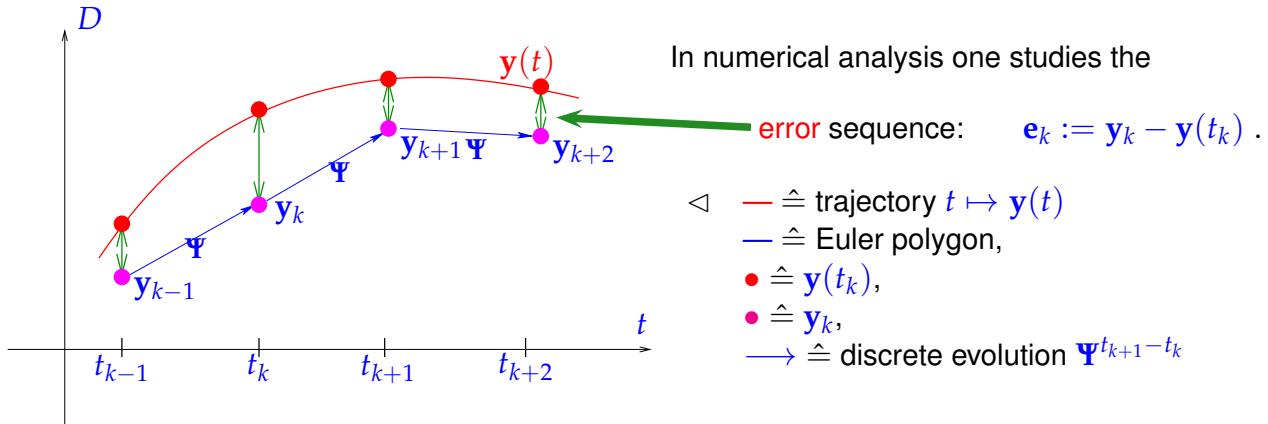
simplest single-step method, namely the explicit Euler method (11.2.1.4) on a mesh  $\mathcal{M} := \{0 = t_0 < t_1 < \dots < t_N = T\}$  for a generic autonomous IVP (11.1.2.2) with sufficiently smooth and (*globally*) Lipschitz continuous  $\mathbf{f}$ , that is,

$$\exists L > 0: \|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{z})\| \leq L \|\mathbf{y} - \mathbf{z}\| \quad \forall \mathbf{y}, \mathbf{z} \in D, \quad (11.3.2.10)$$

cf. Def. 11.1.2.11, and  $C^1$  exact solution  $t \mapsto \mathbf{y}(t)$ . Throughout we assume that solutions of  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  are defined on  $[0, T]$  for all initial states  $\mathbf{y}_0 \in D$ .

Recall the recursion defining the explicit Euler method

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h_k \mathbf{f}(\mathbf{y}_k), \quad k = 1, \dots, N-1. \quad (11.2.1.4)$$

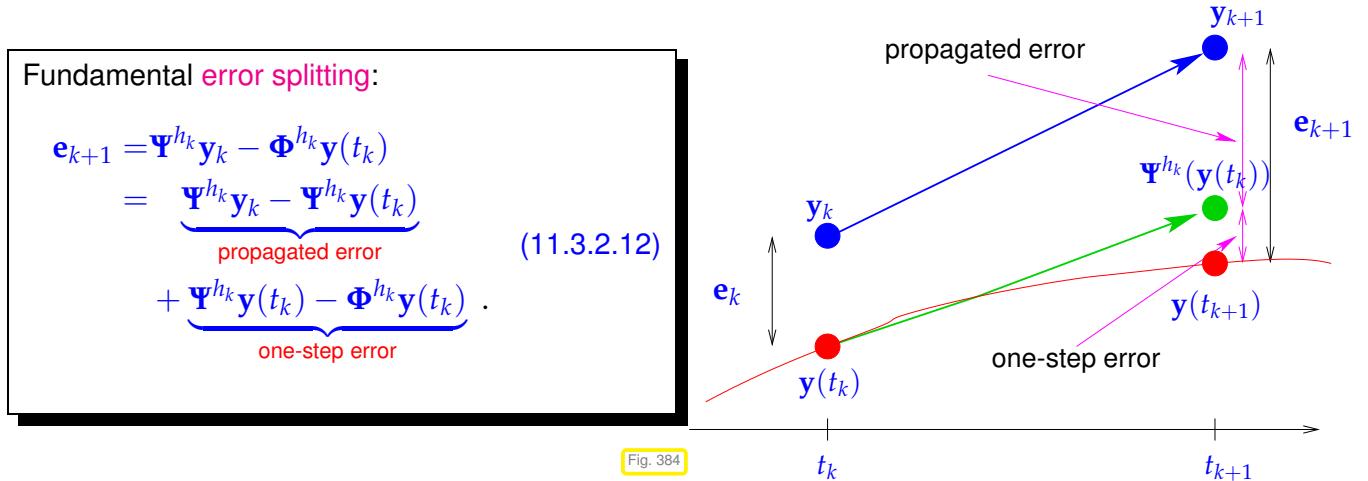


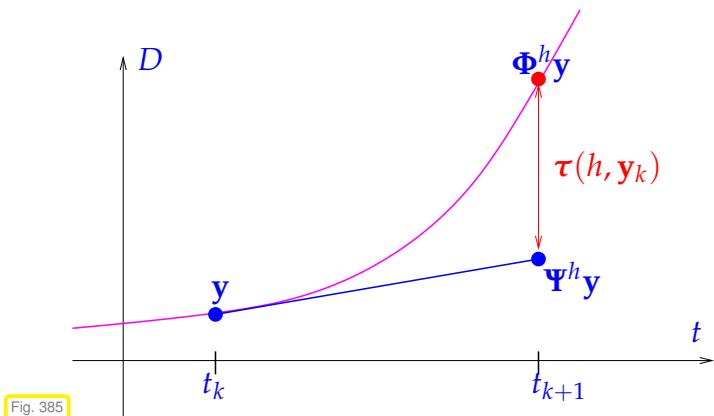
The approach to estimate  $\|\mathbf{e}_k\|$  follows a fundamental policy that comprises three key steps. To explain them we rely on the abstract concepts of the evolution operator  $\Phi$  associated with the ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  (→ Def. 11.1.3.2) and discrete evolution operator  $\Psi$  defining the explicit Euler single step method, see Def. 11.3.1.5:

$$(11.2.1.4) \Rightarrow \Psi^h \mathbf{y} = \mathbf{y} + h \mathbf{f}(\mathbf{y}). \quad (11.3.2.11)$$

We argue that in this context abstraction pays off, because it helps elucidate a general technique for the convergence analysis of single step methods.

### ① Abstract splitting of error:





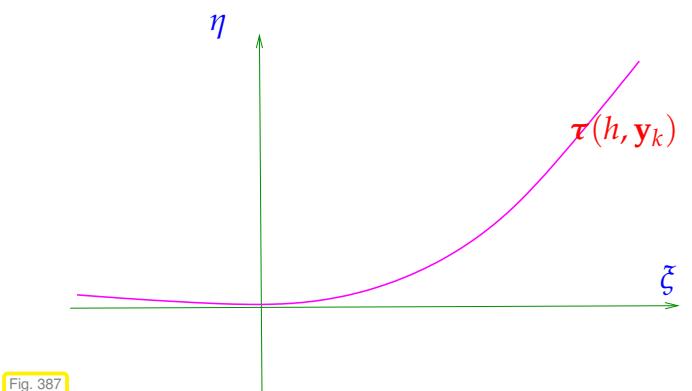
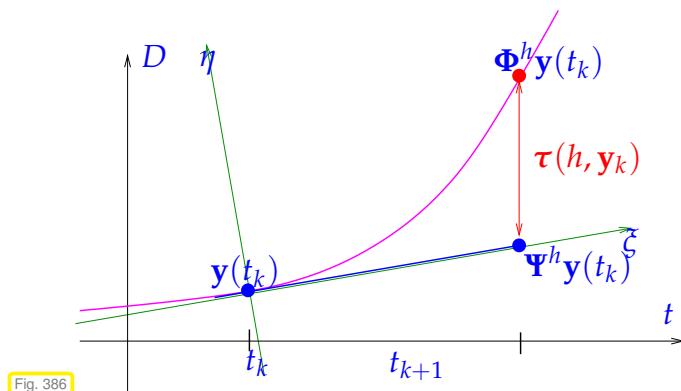
A generic **one-step error** expressed through continuous and discrete evolutions reads:

$$\tau(h, \mathbf{y}) := \Psi^h \mathbf{y} - \Phi^h \mathbf{y} . \quad (11.3.2.13)$$

▷ geometric visualisation of one-step error for explicit Euler method (11.2.1.4), cf. Fig. 377,  
 $h := t_{k+1} - t_k$   
—: solution trajectory through  $(t_k, \mathbf{y})$

## ② Estimate for one-step error $\tau(h_k, \mathbf{y}(t_k))$ :

Geometric considerations: distance of a smooth curve and its tangent shrinks as the square of the distance to the intersection point (curve locally looks like a parabola in the  $\xi - \eta$  coordinate system, see Fig. 387).



The geometric considerations can be made rigorous by analysis: recall Taylor's formula for the function  $\mathbf{y} \in C^{K+1}$  [Str09, Satz 5.5.1]:

$$\mathbf{y}(t+h) - \mathbf{y}(t) = \sum_{j=0}^K \mathbf{y}^{(j)}(t) \frac{h^j}{j!} + \underbrace{\int_t^{t+h} \mathbf{y}^{(K+1)}(\tau) \frac{(t+h-\tau)^K}{K!} d\tau}_{= \frac{\mathbf{y}^{(K+1)}(\xi)}{K!} h^{K+1}} , \quad (11.3.2.14)$$

for some  $\xi \in [t, t+h]$ . We conclude that, if  $\mathbf{y} \in C^2([0, T])$ , which is ensured for smooth  $\mathbf{f}$ , see Lemma 11.1.0.4, then

$$\mathbf{y}(t_{k+1}) - \mathbf{y}(t_k) = \dot{\mathbf{y}}(t_k) h_k + \frac{1}{2} \ddot{\mathbf{y}}(\xi_k) h_k^2 = \mathbf{f}(\mathbf{y}(t_k)) h_k + \frac{1}{2} \ddot{\mathbf{y}}(\xi_k) h_k^2 ,$$

for some  $t_k \leq \xi_k \leq t_{k+1}$ . This leads to an expression for the one-step error from (11.3.2.13)

$$\begin{aligned} \tau(h_k, \mathbf{y}(t_k)) &= \Psi^{h_k} \mathbf{y}(t_k) - \mathbf{y}(t_{k+1}) \\ &\stackrel{(11.3.2.11)}{=} \mathbf{y}(t_k) + h_k \mathbf{f}(\mathbf{y}(t_k)) - \mathbf{y}(t_k) - \mathbf{f}(\mathbf{y}(t_k)) h_k + \frac{1}{2} \ddot{\mathbf{y}}(\xi_k) h_k^2 \\ &= \frac{1}{2} \ddot{\mathbf{y}}(\xi_k) h_k^2 . \end{aligned} \quad (11.3.2.15)$$

Sloppily speaking, we observe

$$\boxed{\tau(h_k, \mathbf{y}(t_k)) = O(h_k^2)} \quad \text{uniformly for } h_k \rightarrow 0.$$

③ **Estimate for the propagated error** from (11.3.2.12)

$$\begin{aligned} \|\Psi^{h_k} \mathbf{y}_k - \Psi^{h_k} \mathbf{y}(t_k)\| &= \|\mathbf{y}_k + h_k \mathbf{f}(\mathbf{y}_k) - \mathbf{y}(t_k) - h_k \mathbf{f}(\mathbf{y}(t_k))\| \\ &\stackrel{(11.3.2.10)}{\leq} (1 + Lh_k) \|\mathbf{y}_k - \mathbf{y}(t_k)\|. \end{aligned} \quad (11.3.2.16)$$


③ Obtain *recursion* for error norms  $\epsilon_k := \|\mathbf{e}_k\|$  by  $\triangle$ -inequality:

$$\epsilon_{k+1} \leq (1 + h_k L) \epsilon_k + \rho_k, \quad \rho_k := \frac{1}{2} h_k^2 \max_{t_k \leq \tau \leq t_{k+1}} \|\ddot{\mathbf{y}}(\tau)\|. \quad (11.3.2.17)$$

Taking into account  $\epsilon_0 = 0$ , this leads to

$$\epsilon_k \leq \sum_{l=1}^k \prod_{j=1}^{l-1} (1 + Lh_j) \rho_l, \quad k = 1, \dots, N. \quad (11.3.2.18)$$

Use the elementary estimate  $(1 + Lh_j) \leq \exp(Lh_j)$  (by convexity of exponential function):

$$(11.3.2.18) \Rightarrow \epsilon_k \leq \sum_{l=1}^k \prod_{j=1}^{l-1} \exp(Lh_j) \cdot \rho_l = \sum_{l=1}^k \exp(L \sum_{j=1}^{l-1} h_j) \rho_l.$$

Note:  $\sum_{j=1}^{l-1} h_j \leq T$  for final time  $T$  and conclude

$$\begin{aligned} \epsilon_k &\leq \exp(LT) \sum_{l=1}^k \rho_l \leq \exp(LT) \max_k \frac{\rho_k}{h_k} \sum_{l=1}^k h_l \leq T \exp(LT) \max_{l=1, \dots, k} h_l \cdot \max_{t_0 \leq \tau \leq t_k} \|\ddot{\mathbf{y}}(\tau)\|. \\ \Rightarrow \|\mathbf{y}_k - \mathbf{y}(t_k)\| &\leq T \exp(LT) \max_{l=1, \dots, k} h_l \cdot \max_{t_0 \leq \tau \leq t_k} \|\ddot{\mathbf{y}}(\tau)\|. \end{aligned} \quad (11.3.2.19)$$

We can summarize the insight gleaned through this theoretical analysis as follows:

Total error arises from accumulation of propagated one-step errors!

First conclusions from (11.3.2.19):

- ◆ error bound  $= O(h)$ ,  $h := \max_l h_l$  ( $\Rightarrow$  1st-order algebraic convergence)
- ◆ Error bound grows exponentially with the length  $T$  of the integration interval.  $\perp$

**§11.3.2.20 (One-step error and order of a single step method)** In the analysis of the global discretization error of the explicit Euler method in § 11.3.2.9 a one-step error of size  $O(h_k^2)$  led to a total error of  $O(h)$  through the effect of error accumulation over  $N \approx h^{-1}$  steps. This relationship remains valid for

almost all single step methods:

Consider an IVP (11.1.2.2) with solution  $t \mapsto \mathbf{y}(t)$  and a single step method defined by the discrete evolution  $\Psi$  ( $\rightarrow$  Def. 11.3.1.5). If the *one-step error along the solution trajectory* satisfies ( $\Phi$  is the evolution map associated with the ODE, see Def. 11.1.3.2)

$$\|\Psi^h \mathbf{y}(t) - \Phi^h \mathbf{y}(t)\| \leq Ch^{p+1} \quad \forall h \text{ sufficiently small}, t \in [0, T],$$

for some  $p \in \mathbb{N}$  and  $C > 0$ , then, usually,

$$\max_k \|\mathbf{y}_k - \mathbf{y}(t_k)\| \leq \bar{C} h_M^p,$$

with  $\bar{C} > 0$  independent of the temporal mesh  $M$ .

A rigorous statement as a theorem would involve some particular assumptions on  $\Psi$ , which we do not want to give here. These assumptions are satisfied, for instance, for all the methods presented in the sequel.  $\square$

## 11.4 Explicit Runge-Kutta Methods



*Supplementary literature.* [DR08, Sect. 11.6], [Han02, Ch. 76], [QSS00, Sect. 11.8]

So far we only know first and second order methods from 11.2: the explicit and implicit Euler method (11.2.1.4) and (11.2.2.2), respectively, are of first order, the implicit midpoint rule of second order. We observed this in Exp. 11.3.2.5 and it can be proved rigorously for all three methods adapting the arguments of § 11.3.2.9.

Thus, barring the impact of roundoff, the low-order polygonal approximation methods are guaranteed to achieve any prescribed accuracy provided that the mesh is fine enough. Why should we need any other timestepping schemes?

**Remark 11.4.0.1 (Rationale for high-order single step methods *cf.* [DR08, Sect. 11.5.3])** We argue that the use of higher-order timestepping methods is highly *advisable for the sake of efficiency*. The reasoning is very similar to that of Rem. 7.3.0.49, when we considered numerical quadrature. The reader is advised to study that remark again.

As we saw in § 11.3.2.3 error bounds for single step methods for the solution of IVPs will inevitably feature unknown constants “ $C > 0$ ”. Thus they do **not** give useful information about the discretization error for a concrete IVP and mesh. Hence, it is too ambitious to ask how many timesteps are needed so that  $\|\mathbf{y}(T) - \mathbf{y}_N\|$  stays below a prescribed bound, *cf.* the discussion in the context of numerical quadrature.

However, an easier question can be answered by *asymptotic estimates* like (11.3.2.7):

What extra computational effort buys a prescribed *reduction of the error*?

(also recall the considerations in Section 8.3.3!)

The usual concept of “computational effort” for single step methods ( $\rightarrow$  Def. 11.3.1.5) is as follows

- Computational effort**
- ~ total number of  $\mathbf{f}$ -evaluations for approximately solving the IVP,
  - ~ number of timesteps, if evaluation of discrete evolution  $\Psi^h$  ( $\rightarrow$  Def. 11.3.1.5) requires fixed number of  $\mathbf{f}$ -evaluations,
  - ~  $h^{-1}$ , in the case of uniform timestep size  $h > 0$  (equidistant mesh (11.3.2.4)).

Now, let us consider a single step method of order  $p \in \mathbb{N}$ , employed with a uniform timestep  $h_{\text{old}}$ . We focus on the maximal discretization error in the mesh points, see § 11.3.2.1. As in (7.3.0.50) we assume that the asymptotic error bounds are *sharp*:

$$\text{err}(h) \approx Ch^p \quad \text{for small meshwidth } h > 0 ,$$

with a “generic constant”  $C > 0$  independent of the mesh.

Goal:  $\frac{\text{err}(h_{\text{new}})}{\text{err}(h_{\text{old}})} \stackrel{!}{=} \frac{1}{\rho}$  for reduction factor  $\rho > 1$ .

$$(11.3.2.7) \Rightarrow \frac{h_{\text{new}}^p}{h_{\text{old}}^p} \stackrel{!}{=} \frac{1}{\rho} \Leftrightarrow h_{\text{new}} = \rho^{-1/p} h_{\text{old}} .$$

For single step method of order  $p \in \mathbb{N}$

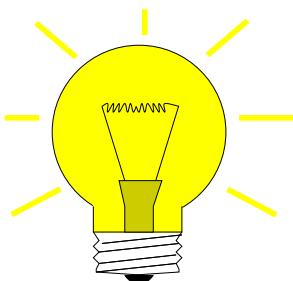
increase effort by factor  $\rho^{1/p}$  reduce error by factor  $\rho > 1$

☞ the larger the order  $p$ , the less effort for a prescribed reduction of the error!

We remark that another (minor) rationale for using higher-order methods is to curb impact of roundoff errors ( $\rightarrow$  Section 1.5.3) accumulating during timestepping [DR08, Sect. 11.5.3].

**§11.4.0.2 (Bootstrap construction of explicit single step methods)** Now we will build a class of methods that are explicit and achieve orders  $p > 2$ . The starting point is a simple *integral equation* satisfied by any solution  $t \mapsto \mathbf{y}(t)$  of an initial value problems for the ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ :

$$\text{IVP: } \begin{aligned} \dot{\mathbf{y}}(t) &= \mathbf{f}(t, \mathbf{y}(t)) , \\ \mathbf{y}(t_0) &= \mathbf{y}_0 \end{aligned} \Rightarrow \mathbf{y}(t_1) = \mathbf{y}_0 + \int_{t_0}^{t_1} \mathbf{f}(\tau, \mathbf{y}(\tau)) d\tau$$



Idea: approximate the integral by means of  $s$ -point quadrature formula ( $\rightarrow$  Section 7.1, defined on the reference interval  $[0, 1]$ ) with nodes  $c_1, \dots, c_s$ , weights  $b_1, \dots, b_s$ .

$$\mathbf{y}(t_1) \approx \mathbf{y}_1 = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{f}(t_0 + c_i h, \boxed{\mathbf{y}(t_0 + c_i h)}) , \quad h := t_1 - t_0 . \quad (11.4.0.3)$$

Obtain these values by **bootstrapping**

“Bootstrapping” = use the same idea in a simpler version to get  $\mathbf{y}(t_0 + c_i h)$ , noting that these values can be replaced by other approximations obtained by methods already constructed (this approach will be elucidated in the next example).

What error can we afford in the approximation of  $\mathbf{y}(t_0 + c_i h)$  (under the assumption that  $\mathbf{f}$  is Lipschitz continuous)? We take the cue from the considerations in § 11.3.2.9.

Goal: aim for one-step error bound  $\mathbf{y}(t_1) - \mathbf{y}_1 = O(h^{p+1})$

Note that there is a factor  $h$  in front of the quadrature sum in (11.4.0.3). Thus, our goal can already be achieved, if only

$\mathbf{y}(t_0 + c_i h)$  is approximated up to an error  $O(h^p)$ ,

again, because in (11.4.0.3) a factor of size  $h$  multiplies  $\mathbf{f}(t_0 + c_i, \mathbf{y}(t_0 + c_i h))$ .

This is accomplished by a less accurate discrete evolution than the one we are about to build. Thus, we can construct discrete evolutions of higher and higher order, in turns, starting with the explicit Euler method. All these methods will be **explicit**, that is,  $\mathbf{y}_1$  can be computed directly from point values of  $\mathbf{f}$ .  $\square$

**EXAMPLE 11.4.0.4 (Simple Runge-Kutta methods by quadrature & bootstrapping)** Now we apply the bootstrapping idea outlined above. We write  $\mathbf{k}_\ell \in \mathbb{R}^d$  for the approximations of  $\mathbf{y}(t_0 + c_i h)$ .

- Quadrature formula = trapezoidal rule (7.2.0.5):

$$Q(f) = \frac{1}{2}(f(0) + f(1)) \leftrightarrow s=2: c_1 = 0, c_2 = 1, b_1 = b_2 = \frac{1}{2}, \quad (11.4.0.5)$$

and  $\mathbf{y}(t_1)$  approximated by explicit Euler step (11.2.1.4)

$$\mathbf{k}_1 = \mathbf{f}(t_0, \mathbf{y}_0), \quad \mathbf{k}_2 = \mathbf{f}(t_0 + h, \mathbf{y}_0 + h\mathbf{k}_1), \quad \mathbf{y}_1 = \mathbf{y}_0 + \frac{h}{2}(\mathbf{k}_1 + \mathbf{k}_2). \quad (11.4.0.6)$$

(11.4.0.6) = **explicit trapezoidal method** (for numerical integration of ODEs).

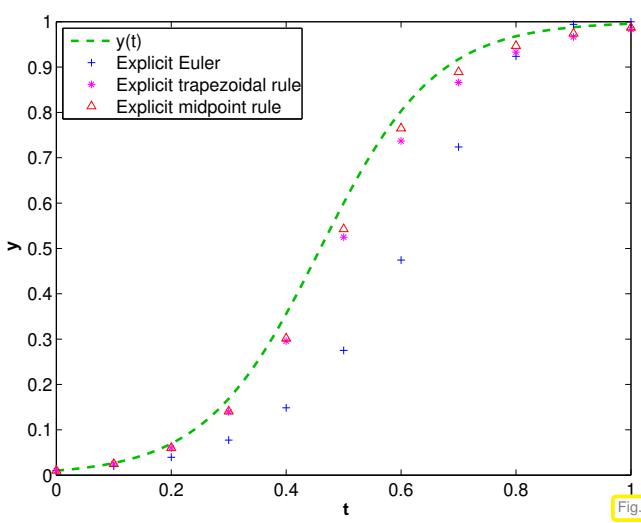
- Quadrature formula  $\rightarrow$  simplest Gauss quadrature formula = midpoint rule ( $\rightarrow$  Ex. 7.2.0.3) &  $\mathbf{y}(\frac{1}{2}(t_1 + t_0))$  approximated by explicit Euler step (11.2.1.4)

$$\mathbf{k}_1 = \mathbf{f}(t_0, \mathbf{y}_0), \quad \mathbf{k}_2 = \mathbf{f}(t_0 + \frac{h}{2}, \mathbf{y}_0 + \frac{h}{2}\mathbf{k}_1), \quad \mathbf{y}_1 = \mathbf{y}_0 + h\mathbf{k}_2. \quad (11.4.0.7)$$

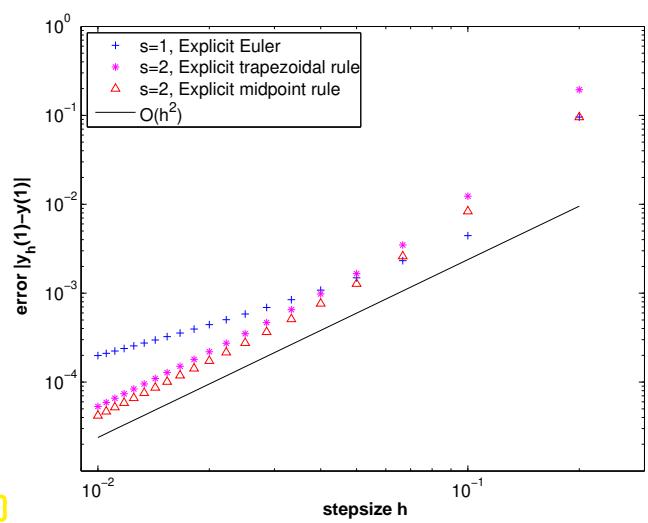
(11.4.0.7) = **explicit midpoint method** (for numerical integration of ODEs) [DR08, Alg. 11.18].  $\square$

**EXAMPLE 11.4.0.8 (Convergence of simple Runge-Kutta methods)** We perform an empiric study of the order of the explicit single step methods constructed in Ex. 11.4.0.4.

- ◆ IVP:  $\dot{y} = 10y(1 - y)$  (logistic ODE (11.1.1.2)),  $y(0) = 0.01$ ,  $T = 1$ ,
- ◆ Explicit single step methods, uniform timestep  $h$ .



$y_h(j/10)$ ,  $j = 1, \dots, 10$  for explicit RK-methods



Errors at final time  $y_h(1) - y(1)$

Observation: obvious *algebraic convergence* in meshwidth  $h$  with integer rates/orders:

- explicit trapezoidal rule (11.4.0.6) → order 2
- explicit midpoint rule (11.4.0.7) → order 2

This is what one expects from the considerations in Ex. 11.4.0.4. □

The formulas that we have obtained follow a general pattern:

#### Definition 11.4.0.9. Explicit Runge-Kutta method

For  $b_i, a_{ij} \in \mathbb{R}$ ,  $c_i := \sum_{j=1}^{i-1} a_{ij}$ ,  $i, j = 1, \dots, s$ ,  $s \in \mathbb{N}$ , an **s-stage explicit Runge-Kutta single step method** (RK-SSM) for the ODE  $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$ ,  $\mathbf{f} : \Omega \rightarrow \mathbb{R}^d$ , is defined by ( $\mathbf{y}_0 \in D$ )

$$\mathbf{k}_i := \mathbf{f}\left(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^{i-1} a_{ij} \mathbf{k}_j\right), \quad i = 1, \dots, s, \quad \mathbf{y}_1 := \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i.$$

The vectors  $\mathbf{k}_i \in \mathbb{R}^d$ ,  $i = 1, \dots, s$ , are called **increments**,  $h > 0$  is the size of the timestep.

Recall Rem. 11.3.1.13 to understand how the discrete evolution for an explicit Runge-Kutta method is specified in this definition by giving the formulas for the first step. This is a convention widely adopted in the literature about numerical methods for ODEs. Of course, the increments  $\mathbf{k}_i$  have to be computed anew in each timestep.

The implementation of an **s**-stage explicit Runge-Kutta single step method according to Def. 11.4.0.9 is straightforward: The increments  $\mathbf{k}_i \in \mathbb{R}^d$  are computed successively, starting from  $\mathbf{k}_1 = \mathbf{f}(t_0 + c_1 h, \mathbf{y}_0)$ .

► Only **s**  $\mathbf{f}$ -evaluations and AXPY operations (→ Section 1.3.2) are required.

In books and research articles a particular way to write down the coefficients characterizing RK-SSMs is widely used:

#### Butcher scheme notation for explicit RK-SSM

Shorthand notation for (explicit) Runge-Kutta methods [DR08, (11.75)]  
Butcher scheme  $\triangleright$   $\begin{array}{c|cc} \mathbf{c} & \mathfrak{A} \\ \hline & \mathbf{b}^T \end{array} :=$   
 (Note:  $\mathfrak{A}$  is strictly lower triangular  $s \times s$ -matrix)

$$\begin{array}{c|cccc} c_1 & 0 & \cdots & & 0 \\ c_2 & a_{21} & \ddots & & \vdots \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & & \ddots & \vdots \\ c_s & a_{s1} & \cdots & a_{s,s-1} & 0 \\ \hline b_1 & & \cdots & b_{s-1} & b_s \end{array} \quad (11.4.0.11)$$

Now we restrict ourselves to the case of an autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ . Matching Def. 11.4.0.9 and Def. 11.3.1.5, we see that the discrete evolution induced by an explicit Runge-Kutta single-step method is

$$\Psi^h \mathbf{y} = \mathbf{y} + h \sum_{i=1}^s b_i \mathbf{k}_i, \quad h \in \mathbb{R}, \quad \mathbf{y} \in D, \quad (11.4.0.12)$$

where the increments  $\mathbf{k}_i$  are defined by the increment equations

$$\mathbf{k}_i := \mathbf{f}(\mathbf{y} + h \sum_{j=1}^{i-1} a_{ij} \mathbf{k}_j) .$$

In line with (11.3.1.10), this discrete evolution can be written as

$$\Psi^h = \mathbf{y} + h \psi(h, \mathbf{y}) , \quad \psi(h, \mathbf{y}) = \sum_{i=1}^s b_i \mathbf{k}_i .$$

Is this discrete evolution consistent in the sense of § 11.3.1.8, that is, does  $\psi(0, \mathbf{y}) = \mathbf{f}(\mathbf{y})$  hold? If  $h = 0$ , the increment equations yield

$$h = 0 \Rightarrow \mathbf{k}_1 = \dots = \mathbf{k}_s = \mathbf{f}(\mathbf{y}) . \quad \blacktriangleright \quad \psi(0, \mathbf{y}) = \left( \sum_{i=1}^s b_i \right) \mathbf{f}(\mathbf{y}) .$$

#### Corollary 11.4.0.13. Consistent Runge-Kutta single step methods

A Runge-Kutta single step method according to Def. 11.4.0.9 is *consistent* ( $\rightarrow$  Def. 11.3.1.11) with the ODE  $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$ , if and only if

$$\sum_{i=1}^s b_i = 1 .$$

**Remark 11.4.0.14 (RK-SSM and quadrature rules)** Note that in Def. 11.4.0.9 the coefficients  $C_i$  and  $b_i$ ,  $i \in \{1, \dots, s\}$ , can be regarded as nodes and weights of a quadrature formula ( $\rightarrow$  Def. 7.1.0.1) on  $[0, 1]$ : apply the explicit Runge-Kutta single step method to “ODE”  $\dot{\mathbf{y}} = f(t)$ ,  $f \in C^0([0, 1])$ , on  $[[0, 1]]$  with timestep  $h = 1$  and initial value  $\mathbf{y}(0)$ , with exact solution

$$\dot{\mathbf{y}}(t) = f(t) , \quad \mathbf{y}(0) = 0 \Rightarrow \mathbf{y}(t) = \int_0^t f(\tau) d\tau .$$

Then the formulas of Def. 11.4.0.9 reduce to

$$\mathbf{y}_1 = 0 + \sum_{i=1}^s b_i f(c_i) \approx \int_0^1 f(\tau) d\tau .$$

Recall that the quadrature rule with these weights and nodes  $c_j$  will have order  $\geq 1$  ( $\rightarrow$  Def. 7.3.0.1), if the weights add up to 1!

**EXAMPLE 11.4.0.15 (Butcher schemes for some explicit RK-SSM [DR08, Sect. 11.6.1])** The following explicit Runge-Kutta single step methods are often mentioned in literature.

- Explicit Euler method (11.2.1.4):
- |   |   |
|---|---|
| 0 | 0 |
| 1 |   |
- $\geq$  order = 1

- explicit trapezoidal rule (11.4.0.6):
- |               |               |               |
|---------------|---------------|---------------|
| 0             | 0             | 0             |
| 1             | 1             | 0             |
| $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
- $\geq$  order = 2

- explicit midpoint rule (11.4.0.7):

$$\begin{array}{c|cc} 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \hline 0 & 1 \end{array} \quad \Rightarrow \quad \text{order} = 2$$

- Classical 4th-order RK-SSM:

$$\begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ \hline 1 & 0 & 0 & 1 & 0 \\ \hline & \frac{1}{6} & \frac{2}{6} & \frac{2}{6} & \frac{1}{6} \end{array} \quad \Rightarrow \quad \text{order} = 4$$

- Kutta's 3/8-rule:

$$\begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ \frac{2}{3} & -\frac{1}{3} & 1 & 0 & 0 \\ \hline 1 & 1 & -1 & 1 & 0 \\ \hline & \frac{1}{8} & \frac{3}{8} & \frac{3}{8} & \frac{1}{8} \end{array} \quad \Rightarrow \quad \text{order} = 4$$

Hosts of (explicit) Runge-Kutta methods can be found in the literature, see for example the [Wikipedia page](#). They are stated in the form of Butcher schemes (11.4.0.11) most of the time.

**Remark 11.4.0.16 (Construction of higher order Runge-Kutta single step methods)** Runge-Kutta single step methods of order  $p > 2$  are not found by bootstrapping as in Ex. 11.4.0.4, because the resulting methods would have quite a lot of stages compared to their order.

Rather one derives **order conditions** yielding large non-linear systems of equations for the coefficients  $a_{ij}$  and  $b_i$  in Def. 11.4.0.9, see [DB02, Sect .4.2.3] and [HLW06, Ch. III]. This approach is similar to the construction of a Gauss quadrature rule in Ex. 7.3.0.14. Unfortunately, the systems of equations are very difficult to solve and no universal recipe is available. Nevertheless, through massive use of symbolic computation, Runge-Kutta methods of order up to 19 have been constructed in this way.

**Remark 11.4.0.17 (“Butcher barriers” for explicit RK-SSM)** The following table gives lower bounds for the number of stages needed to achieve order  $p$  for an explicit Runge-Kutta method.

order $p$	1	2	3	4	5	6	7	8	$\geq 9$
minimal no. $s$ of stages	1	2	3	4	6	7	9	11	$\geq p+3$

No general formula is has been discovered. What is known is that for explicit Runge-Kutta single step methods according to Def. 11.4.0.9

$$\text{order } p \leq \text{number } s \text{ of stages of RK-SSM}$$

**§11.4.0.18 (EIGEN-compatible adaptive explicit embedded Runge-Kutta integrator)** An implementation of an explicit embedded Runge-Kutta single-step method with adaptive stepsize control for solving an autonomous IVP is provided by the utility class **ode45**. The terms “embedded” and “adaptive” will be explained in Section 11.5.

The class is templated with two type parameters:

- (i) **StateType**: type for vectors in state space  $V$ , e.g. a fixed size vector type of EIGEN:  
`Eigen::Matrix<double, N, 1>`, where  $N$  is an integer constant § 11.2.0.1.
- (ii) **RhsType**: a functor type, see Section 0.3.3, for the right hand side function  $f$ ; must match **StateType**, default type provided.

### C++ code 11.4.0.19: Runge-Kutta-Fehlberg 4(5) numerical integrator class

```

2  template <class StateType ,
3      class RhsType = std::function<StateType (const StateType &) >>
4  class ode45 {
5  public:
6      // Idle constructor
7      ode45(const RhsType & rhs) : f(rhs) {}
8      // Main timestepping routine
9      template<class NormFunc = decltype(_norm<StateType>) >
10     std::vector< std::pair<StateType , double>>
11     solve(const StateType & y0,double T,const NormFunc & norm = _norm<StateType>);
12     // Print statistics and options of this class instance.
13     void print();
14     struct Options {
15         bool save_init = true;    // Set true if you want to save the initial
16         data
17         bool fixed_stepsize = false; // TODO: Set true if you want a fixed step
18         size
19         unsigned int max_iterations = 5000;
20         double min_dt = -1.; // Set the minimum step size (-1 for none)
21         double max_dt = -1.; // Set the maximum step size (-1 for none)
22         double initial_dt = -1.; // Set an initial step size
23         double start_time = 0; // Set a starting time
24         double rtol = 1e-6; // Relative tolerance for the error.
25         double atol = 1e-8; // Absolute tolerance for the error.
26         bool do_statistics = false; // Set to true before solving to save
27             statistics
28         bool verbose = false; // Print more output.
29     } options;
30     // Data structure for usage statistics.
31     // Available after a call of solve(), if do_statistics is set to true.
32     struct Statistics {
33         unsigned int cycles = 0; // Number of loops (sum of all accepted and
34             rejected steps)
35         unsigned int steps = 0; // Number of actual time steps performed
36             (accepted step)
37         unsigned int rejected_steps = 0; // Number of rejected steps per step
38         unsigned int funcalls = 0; // Function calls
39     } statistics;
40 };
41 
```

The functor for the right hand side  $f : D \subset V \rightarrow V$  of the ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  is specified as an argument of the constructor.

The single-step numerical integrator is invoked by the templated method

```

template<class NormFunc = decltype(_norm<StateType>) >
std::vector< std::pair<StateType , double>>
solve(const StateType & y0,double T,const NormFunc & norm =
      _norm<StateType>);
```

The following arguments have to be supplied:

1.  $y_0$ : the initial value  $\mathbf{y}_0$

2.  $T$ : the final time  $T$ , initial time  $t_0 = 0$  is assumed, because the class can deal with autonomous ODEs only, recall § 11.1.2.3.
3. `norm`: a functor returning a suitable norm for a state vector. Defaults to EIGEN's maximum vector norm.

The method returns a vector of 2-tuples  $(\mathbf{y}_k, t_k)$  (note the order!),  $k = 0, \dots, N$ , of temporal mesh points  $t_k$ ,  $t_0 = 0$ ,  $t_N = T$ , see § 11.2.0.2, and approximate states  $\mathbf{y}_k \approx \mathbf{y}(t_k)$ , where  $t \mapsto \mathbf{y}(t)$  stands for the exact solution of the initial value problem.

The next self-explanatory code snippet uses the numerical integrator class `ode45` for solving a scalar autonomous ODE.

#### C++ code 11.4.0.20: Invocation of adaptive embedded Runge-Kutta-Fehlberg integrator

```

2 int main(void)
3 {
4     // Types to be used for a scalar ODE with state space ℝ
5     using StateType = double;
6     using RhsType = std::function<StateType(StateType)>;
7     // Logistic differential equation (11.1.1.2)
8     RhsType f = [] (StateType y) { return 5*y*(1-y); };
9     StateType y0 = 0.2;    // Initial value
10    // Exact solution of IVP, see (11.1.1.3)
11    auto y = [y0](double t) { return y0/(y0+(1-y0)*exp(-5*t)); };
12    // State space ℝ, simple modulus supplies norm
13    auto normFunc = [] (StateType x){ return fabs(x); };
14
15    // Invoke explicit Runge-Kutta method with stepsize control
16    ode45<StateType,RhsType> integrator(f);
17    // Do the timestepping with adaptive strategy of Section 11.5
18    std::vector<std::pair<StateType, double>> states =
19        integrator.solve(y0, 1.0, normFunc);
20    // Output information accumulation during numerical integration
21    integrator.options.do_statistics = true; integrator.print();
22
23    for (auto state : states)
24        std::cout << "t = " << state.second << ", y = " << state.first
25            << ", |err| = " << fabs(state.first - y(state.second)) << std::endl;
}
```

## 11.5 Adaptive Stepsize Control

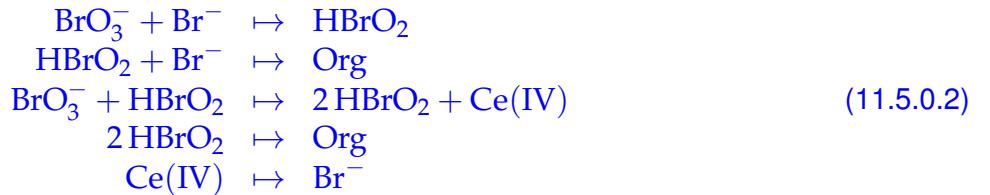
In Section 7.5, in the context of numerical quadrature, we learned about an **a-posteriori** way to adjust the mesh underlying a composite quadrature rule to the integrand: During the computation we estimate the local quadrature error by comparing the approximations obtained by using quadrature formulas of different order. The same policy for adapting the integration mesh is very popular in the context of numerical integration, too. Since the size  $h_k := t_{k+1} - t_k$  of the cells of the temporal mesh is also called the **timestep size**, this kind of a-posteriori mesh adaptation is also known as stepsize control.



*Supplementary literature.* [DR08, Sect. 11.7], [QSS00, Sect. 11.8.2]

**EXAMPLE 11.5.0.1 (Oregonator reaction)** Chemical reaction kinetics is a field where ODE based models are very common. This example presents a famous reaction with extremely abrupt dynamics. Refer to [Han02, Ch. 62] for more information about the ODE-based modelling of kinetics of chemical reactions.

This is a special case of an “oscillating” Zhabotinski-Belousov reaction [Gra02]:



$$\begin{aligned}
 y_1 := c(\text{BrO}_3^-): \quad \dot{y}_1 &= -k_1 y_1 y_2 - k_3 y_1 y_3, \\
 y_2 := c(\text{Br}^-): \quad \dot{y}_2 &= -k_1 y_1 y_2 - k_2 y_2 y_3 + k_5 y_5, \\
 y_3 := c(\text{HBrO}_2): \quad \dot{y}_3 &= k_1 y_1 y_2 - k_2 y_2 y_3 + k_3 y_1 y_3 - 2k_4 y_3^2, \\
 y_4 := c(\text{Org}): \quad \dot{y}_4 &= k_2 y_2 y_3 + k_4 y_3^2, \\
 y_5 := c(\text{Ce(IV)}): \quad \dot{y}_5 &= k_3 y_1 y_3 - k_5 y_5, \\
 \end{aligned} \tag{11.5.0.3}$$

with (non-dimensionalized) reaction constants:

$$k_1 = 1.34, \quad k_2 = 1.6 \cdot 10^9, \quad k_3 = 8.0 \cdot 10^3, \quad k_4 = 4.0 \cdot 10^7, \quad k_5 = 1.0.$$



periodic chemical reaction [Video 1](#), [Video 2](#)

MATLAB simulation with initial state  $y_1(0) = 0.06$ ,  $y_2(0) = 0.33 \cdot 10^{-6}$ ,  $y_3(0) = 0.501 \cdot 10^{-10}$ ,  $y_4(0) = 0.03$ ,  $y_5(0) = 0.24 \cdot 10^{-7}$ :

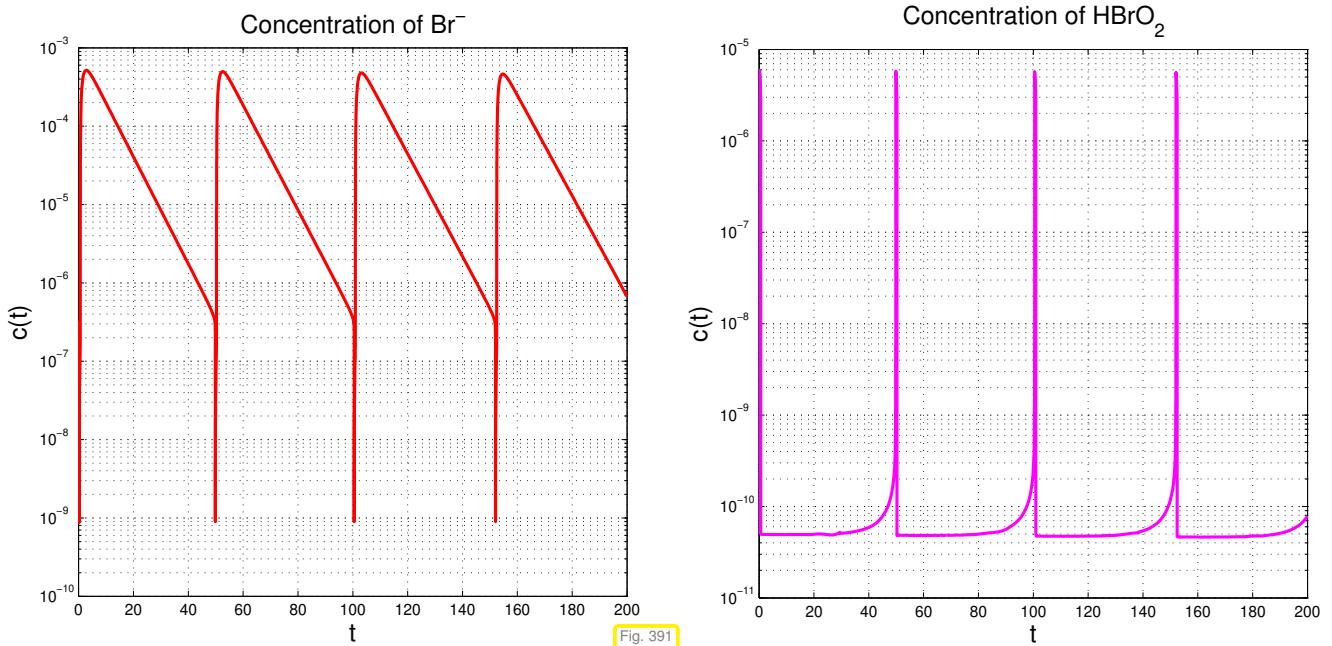


Fig. 390

Fig. 391

We observe a strongly non-uniform behavior of the solution in time.

This is very common with evolutions arising from practical models (circuit models, chemical reaction models, mechanical systems)

### EXAMPLE 11.5.0.4 (Blow-up)

We return to the “explosion ODE” of Ex. 11.1.2.18 and consider the scalar autonomous IVP:

$$\dot{y} = y^2, \quad y(0) = y_0 > 0.$$

►  $y(t) = \frac{y_0}{1 - y_0 t}, \quad t < 1/y_0.$

As we have seen a solution exists only for finite time and then suffers a **Blow-up**, that is,  $\lim_{t \rightarrow 1/y_0} y(t) = \infty$

:  $J(y_0) = ]-\infty, 1/y_0]$ !

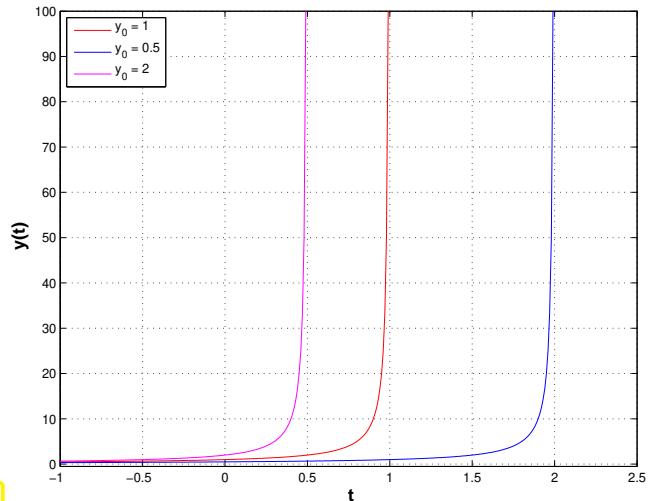
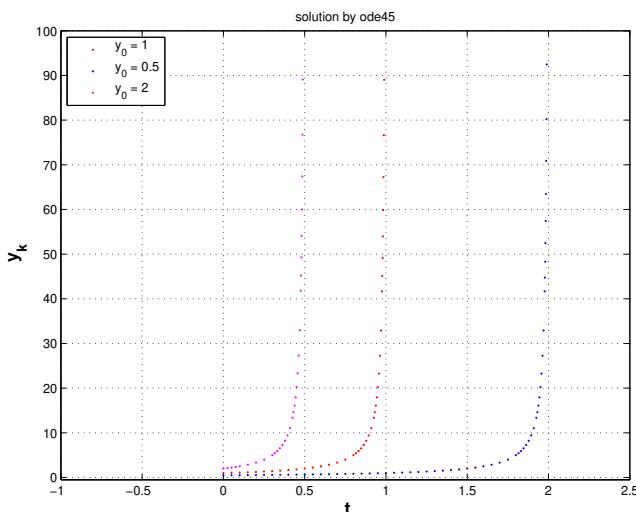


Fig. 392

How to choose temporal mesh  $\{t_0 < t_1 < \dots < t_{N-1} < t_N\}$  for single step method in case  $J(y_0)$  is not known, even worse, if it is not clear a priori that a blow up will happen?

Just imagine: what will result from equidistant explicit Euler integration (11.2.1.4) applied to the above IVP?



Simulation with the numerical integrator **ode45** from Code 12.0.0.3.

Fig. 393

We observe: **ode45** manages to reduce stepsize more and more as it approaches the singularity of the solution! How can it accomplish this feat! □

Key challenge (discussed for autonomous ODEs below):

How to choose a *good temporal mesh*  $\{0 = t_0 < t_1 < \dots < t_{N-1} < t_N\}$   
for a given single step method applied to a concrete IVP?

What does “good” mean ?

Be efficient!

Be accurate!

Stepsize adaptation for single step methods	
Objective:	$N$ as small as possible & $\max_{k=1,\dots,N} \ \mathbf{y}(t_k) - \mathbf{y}_k\  < \text{TOL}$ or $\ \mathbf{y}(T) - \mathbf{y}_N\  < \text{TOL}$ , $\text{TOL} = \text{tolerance}$
Policy:	Try to curb/balance one-step error by <ul style="list-style-type: none"> <li>◆ adjusting current stepsize <math>h_k</math>,</li> <li>◆ predicting suitable next timestep <math>h_{k+1}</math></li> </ul>
Tool:	Local-in-time one-step error estimator ( <i>a posteriori</i> , based on $\mathbf{y}_k, h_{k-1}$ )

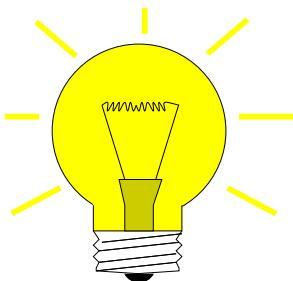
Why local-in-time timestep control (based on estimating only the one-step error)?

Consideration: If a small time-local error in a single timestep leads to large error  $\|\mathbf{y}_k - \mathbf{y}(t_k)\|$  at later times, then local-in-time timestep control is powerless about it and will not even notice!!

Nevertheless, local-in-time timestep control is used almost exclusively,

- ☞ because we do not want to discard past timesteps, which could amount to tremendous waste of computational resources,
- ☞ because it is inexpensive and it works for many practical problems,
- ☞ because there is no reliable method that can deliver guaranteed accuracy for general IVP.

**§11.5.0.6 (Local-in-time error estimation)** We “recycle” heuristics already employed for adaptive quadrature, see Section 7.5, § 7.5.0.10. There we tried to get an idea of the local quadrature error by comparing two approximations of different order. Now we pursue a similar idea over a single timestep.



Idea:

*Estimation of one-step error*

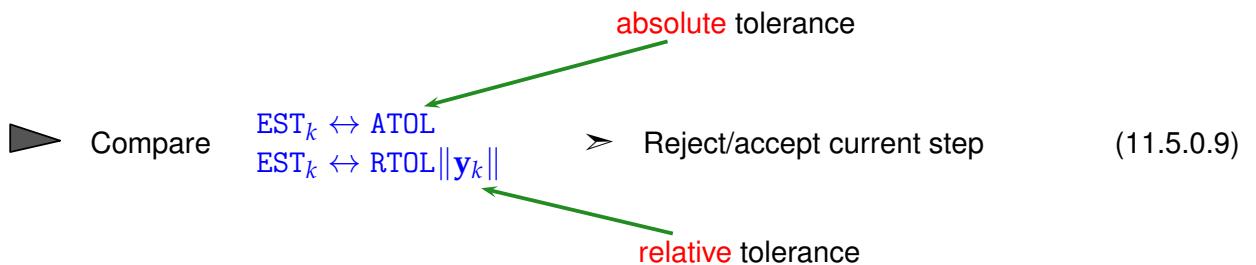
Compare results for two discrete evolutions  $\Psi^h, \tilde{\Psi}^h$  of different order over current timestep  $h$ :

If  $\text{Order}(\tilde{\Psi}) > \text{Order}(\Psi)$ , then we expect

$$\underbrace{\Phi^h \mathbf{y}(t_k) - \Psi^h \mathbf{y}(t_k)}_{\text{one-step error}} \approx \text{EST}_k := \tilde{\Psi}^h \mathbf{y}(t_k) - \Psi^h \mathbf{y}(t_k). \quad (11.5.0.7)$$

Heuristics for concrete  $h$

**§11.5.0.8 (Temporal mesh refinement)** We take for granted a local error estimate  $\text{EST}_k$ .



For a similar use of absolute and relative tolerances see Section 8.1.2: termination criteria for iterations, in particular (8.1.2.3).

☞ Simple algorithm:

- $EST_k < \max\{ATOL, \|y_k\|RTOL\}$ : Carry out next timestep (stepsize  $h$ )  
Use larger stepsize (e.g.,  $\alpha h$  with some  $\alpha > 1$ ) for following step (\*)
- $EST_k > \max\{ATOL, \|y_k\|RTOL\}$ : Repeat current step with smaller stepsize  $< h$ , e.g.,  $\frac{1}{2}h$

Rationale for (\*): if the current stepsize guarantees sufficiently small one-step error, then it might be possible to obtain a still acceptable one-step error with a larger timestep, which would enhance efficiency (fewer timesteps for total numerical integration). This should be tried, since timestep control will usually provide a safeguard against undue loss of accuracy.

The following C++ code implements a wrapper function `odeintadapt()` for a general adaptive single-step method according to the policy outlined above. The arguments are meant to pass the following information:

- `Psilow`, `Psihigh`: functors passing discrete evolution operators for autonomous ODE of different order, type `@(y, h)`, expecting a state (usually a column vector) as first argument, and a stepsize as second,
- $T$ : final time  $T > 0$ ,
- $y_0$ : initial state  $y_0$ ,
- $h_0$ : stepsize  $h_0$  for the first timestep
- `reltol`, `abstol`: relative and absolute tolerances, see (11.5.0.9),
- $h_{\min}$ : minimal stepsize, timestepping terminates when stepsize control  $h_k < h_{\min}$ , which is relevant for detecting blow-ups or collapse of the solution.

#### C++ code 11.5.0.10: Simple local stepsize control for single step methods

```

2 // Auxiliary function: default norm for an EIGEN vector type
3 template <class State>
4 double _norm(const State &y) { return y.norm(); }
5
6 template <class DiscEvolOp, class State, class NormFunc = decltype(_norm<State>) >
7 std::vector<std::pair<double, State>>
8 odeintadapt(DiscEvolOp &&Psi_low, DiscEvolOp &&Psi_high,
9             const State &y0, double T, double h0,
10            double reltol, double abstol, double hmin,
11            NormFunc &norm = _norm<State>){
12     double t = 0; // initial time  $t_0 = 0$ 
13     State y = y0; // current state
14     double h = h0; // timestep to start with
15     std::vector<std::pair<double, State>> states; // vector of times/computed

```

```

    states:  $(t_k, \mathbf{y}_k)_k$ 
16 states.push_back({t, y}); // initial time and state
17
18 while ((states.back().first < T) && (h >= hmin)) { //
19     State yh = Psihigh(h, y); // high order discrete evolution  $\tilde{\Psi}^h$ 
20     State yH = Psilow(h, y); // low order discrete evolution  $\Psi^h$ 
21     double est = norm(yH-yh); // local error estimate  $EST_k$ 
22
23     if (est < max(reltol*norm(y), abstol)) { // step accepted
24         y = yh; // use high order approximation
25         t = t+min(T-t,h); // next time  $t_k$ 
26         states.push_back({t,y}); //
27         h = 1.1*h; // try with increased stepsize
28     }
29     else // step rejected
30         h = h/2; // try with half the stepsize
31 }
32 // Numerical integration has ground to a halt !
33 if (h < hmin) {
34     cerr << "Warning: Failure at t=" << states.back().first
35     << ". Unable to meet integration tolerances without reducing the step "
36     << "size below the smallest value allowed (" << hmin << ") at time t." << endl;
37 }
38 return states;
39 }
```

Comments on Code 11.5.0.10:

- line 18: check whether final time is reached or timestepping has ground to a halt ( $h_k < h_{\min}$ ).
- line 19, 20: advance state by low and high order integrator.
- line 21: compute norm of estimated error, see (11.5.0.7).
- line 23: make comparison (11.5.0.9) to decide whether to accept or reject local step.
- line 26, 27: step accepted, update state and current time and suggest 1.1 times the current stepsize for next step.
- line 30 step rejected, try again with half the stepsize.
- Return value is a vector of pairs consisting of
  - times  $t \leftrightarrow$  temporal mesh  $t_0 < t_1 < t_2 < \dots < t_N < T$ , where  $t_N < T$  indicated premature termination (collapse, blow-up),
  - states  $y \leftrightarrow$  sequence  $(\mathbf{y}_k)_{k=0}^N$ .

**Remark 11.5.0.11 (Estimation of “wrong” error?)** We face the same conundrum as in the case of adaptive numerical quadrature, see Rem. 7.5.0.17:

! By the heuristic considerations, see (11.5.0.7) it seems that  $EST_k$  measures the one-step error for the low-order method  $\Psi$  and that we should use  $\mathbf{y}_{k+1} = \Psi^{h_k} \mathbf{y}_k$ , if the timestep is accepted.

However, it would be foolish not to use the better value  $\mathbf{y}_{k+1} = \tilde{\Psi}^{h_k} \mathbf{y}_k$ , since it is available for free. This is what is done in every implementation of adaptive methods, also in Code 11.5.0.10, and this choice can

be justified by control theoretic arguments [DB02, Sect. 5.2].

**EXPERIMENT 11.5.0.12 (Simple adaptive stepsize control)** We test adaptive timestepping routine from Code 11.5.0.10 for a scalar IVP and compare the estimated local error and true local error.

- ◆ IVP for ODE  $\dot{y} = \cos(\alpha y)^2$ ,  $\alpha > 0$ , solution  $y(t) = \arctan(\alpha(t - c))/\alpha$  for  $y(0) \in [-\pi/2, \pi/2]$
- ◆ Simple adaptive timestepping based on explicit Euler (11.2.1.4) and explicit trapezoidal rule (11.4.0.6)

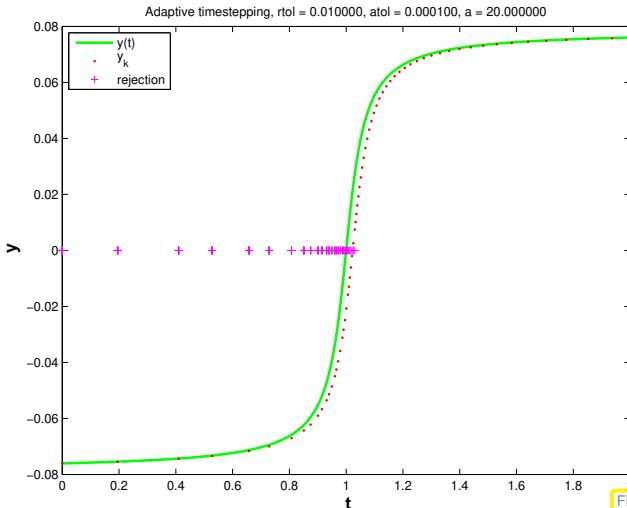


Fig. 394

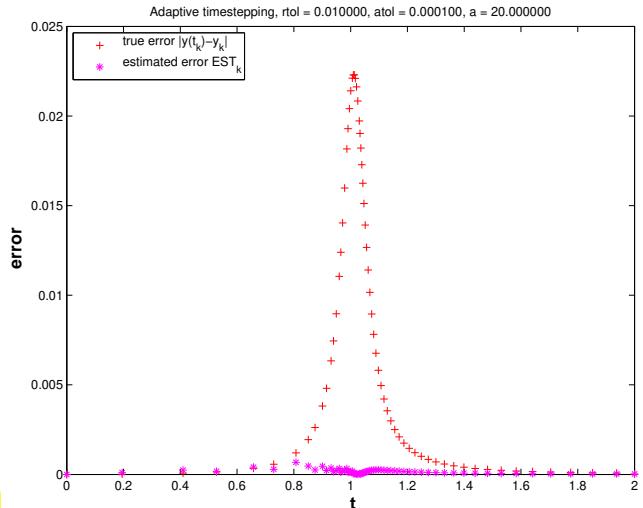


Fig. 395

Statistics: 66 timesteps, 131 rejected timesteps

Observations:

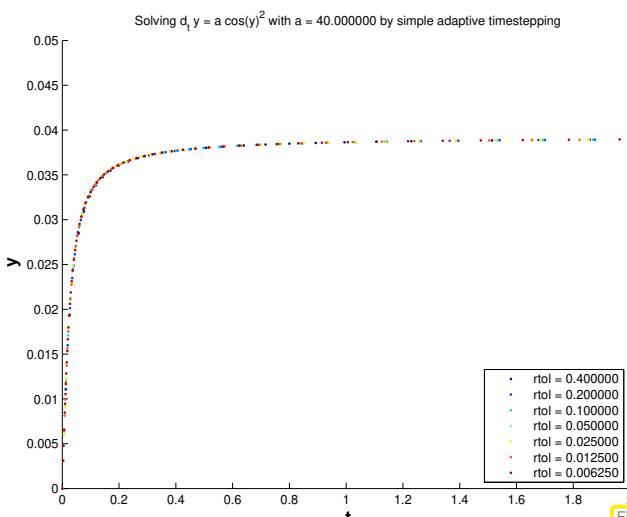
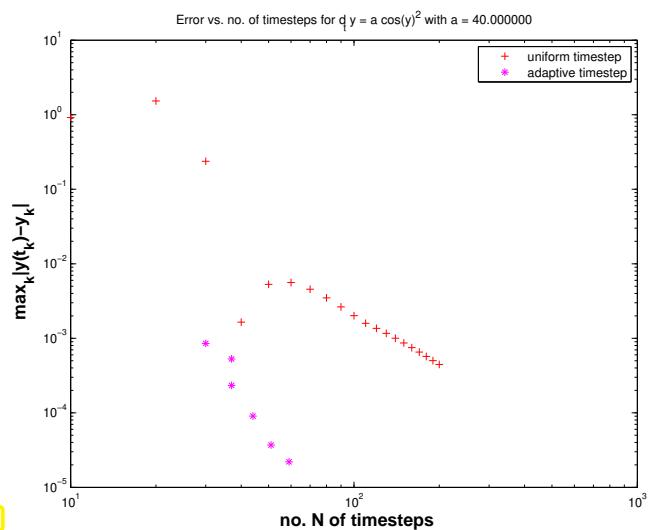
- ☞ Adaptive timestepping well resolves local features of solution  $y(t)$  at  $t = 1$
- ☞ Estimated error (an estimate for the one-step error) and true error are **not** related! To understand this recall Rem. 11.5.0.11.

**EXPERIMENT 11.5.0.13 (Gain through adaptivity → Exp. 11.5.0.12)** In this experiment we want to explore whether adaptive timestepping is worth while, as regards reduction of computational effort without sacrificing accuracy.

We retain the simple adaptive timestepping from previous experiment Exp. 11.5.0.12 and also study the same IVP.

New: initial state  $y(0) = 0$ !

Now we examine the dependence of the maximal discretization error in mesh points on the computational effort. The latter is proportional to the number of timesteps.

Solutions  $(y_k)_k$  for different values of  $rtol$ 

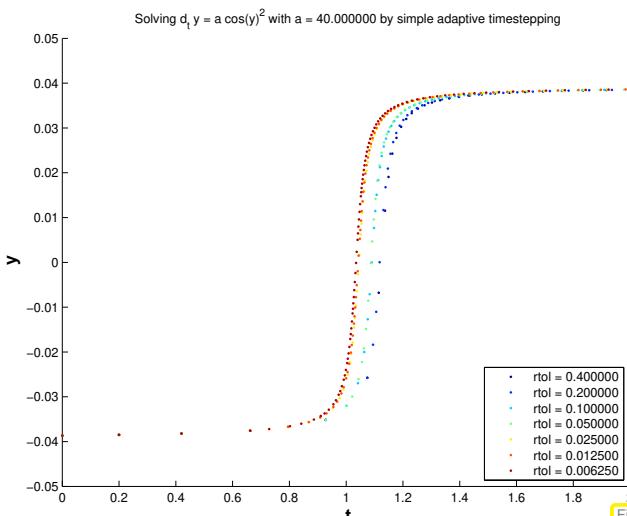
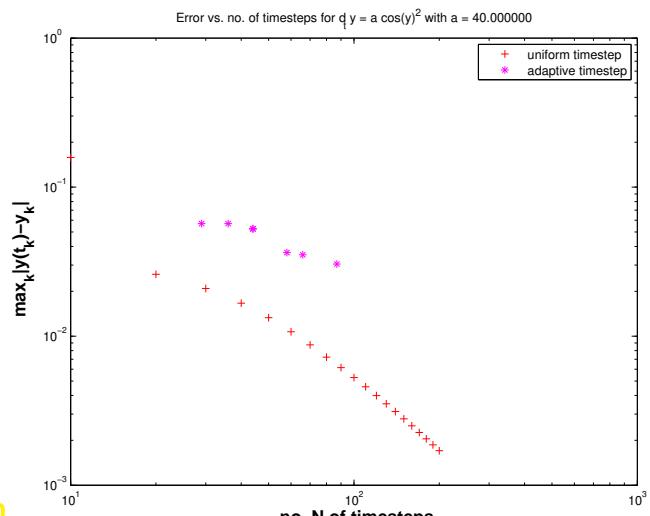
Observations:

- Adaptive timestepping achieves much better accuracy for a fixed computational effort.

### EXPERIMENT 11.5.0.14 (“Failure” of adaptive timestepping → Exp. 11.5.0.13)

Same ODE and simple adaptive timestepping as in previous experiment Exp. 11.5.0.13.

$$\dot{y} = \cos^2(\alpha y) \Rightarrow y(t) = \arctan(\alpha(t - c)) / \alpha, y(0) \in [-\frac{\pi}{2\alpha}, -\frac{\pi}{2\alpha}],$$

for  $\alpha = 40$ .Now: initial state  $y(0) = -0.0386 \approx \frac{\pi}{2\alpha}$  as in Exp. 11.5.0.12Solutions  $(y_k)_k$  for different values of  $rtol$ 

Observations:

- Adaptive timestepping leads to larger errors at the same computational cost as uniform timestepping !

Explanation: the position of the steep step of the solution has a sensitive dependence on an initial value, if  $y(0) \approx \frac{\pi}{2\alpha}$ :

$$y(t) = \frac{1}{\alpha} \arctan(\alpha(t + \tan(y_0/\alpha))), \text{ step at } \approx -\tan(y_0/\alpha).$$

Hence, small local errors in the initial timesteps will lead to large errors at around time  $t \approx 1$ . The stepsize control is mistaken in condoning these small one-step errors in the first few steps and, therefore, incurs huge errors later.

However, the perspective of **backward error analysis** ( $\rightarrow$  § 1.5.5.18) rehabilitates adaptive stepsize control in this case: it gives us a numerical solution that is very close to the exact solution of the ODE with slightly perturbed initial state  $y_0$ .  $\square$

**Remark 11.5.0.15 (Refined local stepsize control  $\rightarrow$  [DR08, Sect. 11.7])**

The above algorithm (Code 11.5.0.10) is simple, but the rule for increasing/shrinking of timestep “squanders” the information contained in  $\text{EST}_k : \text{TOL}$ :

More ambitious goal !	When $\text{EST}_k > \text{TOL}$ : stepsize <b>adjustment</b> better $h_k = ?$
	When $\text{EST}_k < \text{TOL}$ : stepsize <b>prediction</b> good $h_{k+1} = ?$

Assumption: At our disposal are two discrete evolutions:

- ◆  $\Psi$  with  $\text{order}(\Psi) = p$  ( $\rightarrow$  “low order” single step method)
- ◆  $\tilde{\Psi}$  with  $\text{order}(\tilde{\Psi}) > p$  ( $\rightarrow$  “higher order” single step method)

These are the same building blocks as for the simple adaptive strategy employed in Code 11.5.0.10 (passed as arguments  $\text{Psilow}, \text{Psihigh}$  there).

Asymptotic expressions for one-step error for  $h \rightarrow 0$ :

$$\begin{aligned}\Psi^{h_k} \mathbf{y}(t_k) - \Phi^{h_k} \mathbf{y}(t_k) &= ch^{p+1} + O(h_k^{p+2}), \\ \tilde{\Psi}^{h_k} \mathbf{y}(t_k) - \Phi^{h_k} \mathbf{y}(t_k) &= O(h_k^{p+2}),\end{aligned}\tag{11.5.0.16}$$

with some (unknown)  $c > 0$ .

Why  $h^{p+1}$ ? Remember estimate (11.3.2.15) from the error analysis of the explicit Euler method: we also found  $O(h^2)$  there for the one-step error of a single step method of order 1.

Heuristics: the timestep  $h_k$  is small  $\rightarrow$  “higher order terms”  $O(h_k^{p+2})$  can be ignored.

$$\begin{aligned}\Psi^{h_k} \mathbf{y}(t_k) - \Phi^{h_k} \mathbf{y}(t_k) &\doteq ch_k^{p+1} + O(h_k^{p+2}), \\ \tilde{\Psi}^{h_k} \mathbf{y}(t_k) - \Phi^{h_k} \mathbf{y}(t_k) &\doteq O(h_k^{p+2}).\end{aligned}\Rightarrow \text{EST}_k \doteq ch_k^{p+1}.\tag{11.5.0.17}$$

☞ notation:  $\doteq$  equality up to higher order terms in  $h_k$

$$\text{EST}_k \doteq ch_k^{p+1} \Rightarrow c \doteq \frac{\text{EST}_k}{h_k^{p+1}}.\tag{11.5.0.18}$$

Available in algorithm, see (11.5.0.7)

For the sake of *accuracy* (stipulates “ $\text{EST}_k < \text{TOL}$ ”) & *efficiency* (favors “ $>$ ”) we aim for

$$\text{EST}_k \stackrel{!}{=} \text{TOL} := \max\{\text{ATOL}, \|\mathbf{y}_k\| \text{RTOL}\}.\tag{11.5.0.19}$$

What timestep  $h_*$  can actually achieve (11.5.0.19), if we “believe” in (11.5.0.17) (and, therefore, in (11.5.0.18))?

$$(11.5.0.18) \& (11.5.0.19) \Rightarrow \text{TOL} = \frac{\text{EST}_k}{h_k^{p+1}} h_*^{p+1}.$$

► “Optimal timestep”: (stepsize prediction)

$$h_* = h_k^{p+1} \sqrt{\frac{\text{TOL}}{\text{EST}_k}}. \quad (11.5.0.20)$$

adjusted stepsize (**R**)  
& suggested stepsize  
(**A**)

(Reject): In case  $\text{EST}_k > \text{TOL}$  ➤ repeat step with stepsize  $h_*$ .

(Accept): If  $\text{EST}_k \leq \text{TOL}$  ➤ use  $h_*$  as stepsize for next step.

### C++ code 11.5.0.21: Refined local stepsize control for single step methods

```

2 // Auxiliary function: default norm for an EIGEN vector type
3 template <class State>
4 double _norm(const State &y) { return y.norm(); }
5
6 template <class DiscEvolOp, class State, class NormFunc = decltype(_norm<State>) >
7 std::vector<std::pair<double, State>>
8 odeintssctrl(DiscEvolOp &&Psilow, unsigned int p, DiscEvolOp &&Psihigh,
9           const State &y0, double T, double h0,
10          double reltol, double abstol, double hmin,
11          NormFunc &norm = _norm<State> ) {
12     double t = 0; // initial time  $t_0 = 0$ 
13     State y = y0; // current state, initialized here
14     double h = h0; // timestep to start with
15     std::vector<std::pair<double, State>> states; // vector  $(t_k, y_k)_k$ 
16     states.push_back({t, y});
17
18     // Main timestepping loop
19     while ((states.back().first < T) && (h >= hmin)) { //
20         State yh = Psihigh(h, y); // high order discrete evolution  $\tilde{\Psi}^h$ 
21         State yH = Psilow(h, y); // low order discrete evolution  $\Psi^h$ 
22         double est = norm(yH - yh); //  $\leftrightarrow \text{EST}_k$ 
23         double tol = max(reltol * norm(y), abstol); // effective tolerance
24
25         // Optimal stepsize according to (11.5.0.20)
26         h = h * max(0.5, min(2., pow(tol / est, 1. / (p + 1)))); //
27         if (est < tol) // step accepted
28             states.push_back({t = t + min(T - t, h), y = yh}); // store next approximate
29             state
30     }
31     if (h < hmin) {
32         cerr << "Warning: Failure at t="
33         << states.back().first
34         << ". Unable to meet integration tolerances without reducing the step"
35         << " size below the smallest value allowed (" << hmin << ") at time t." << endl;
36     }
37     return states;
}

```

Comments on Code 11.5.0.21 (see comments on Code 11.5.0.10 for more explanations):

- Input arguments as for Code 11.5.0.10, except for  $p \hat{=} \text{order of lower order discrete evolution}$ .
- line 26: compute presumably better local stepsize according to (11.5.0.20),
- line 27: decide whether to repeat the step or advance,
- line 27: extend output arrays if current step has not been rejected.

**Remark 11.5.0.22 (Stepsize control in `ode45`)** In Code 12.0.0.3 you have seen the interface for the numerical integrator class `ode45`. Its name already indicates the orders of the pair of single step methods used for adaptive stepsize control:

$$\Psi \hat{=} \text{RK-method of order } 4 \quad \tilde{\Psi} \hat{=} \text{RK-method of order } 5$$

`ode45`

Specifying tolerances for `ode45` is done by setting the values of the fields `rtol` and `atol` of the options data member of `ode45`.

The possibility to pass tolerances to numerical integrators based on adaptive timestepping may tempt one into believing that they allow to control the accuracy of the solutions. However, as is clear from Rem. 11.5.0.15, these tolerances are solely applied to local error estimates and, inherently, have nothing to do with global discretization errors, see Exp. 11.5.0.12.

### No global error control through local-in-time adaptive timestepping

The absolute/relative tolerances imposed for local-in-time adaptive timestepping do *not* allow to predict accuracy of solution!

**Remark 11.5.0.24 (Embedded Runge-Kutta methods)** For higher order RK-SSM with a considerable number of stages computing different sets of increments ( $\rightarrow$  Def. 11.4.0.9) for two methods of different order just for the sake of local-in-time stepsize control would mean incommensurate effort.

Embedding idea: Use two RK-SSMs based on the **same increments**, that is, built with the same coefficients  $a_{ij}$ , but different weights  $b_i$ , see Def. 11.4.0.9 for the formulas, and **different orders**  $p$  and  $p + 1$ .

Butcher scheme for **embedded explicit Runge-Kutta methods**

(Lower order scheme has weights  $\hat{b}_i$ .)

$$\begin{array}{c|ccccc} c & \alpha \\ \hline b^T & & & & \\ \hline \hat{b}^T & & & & \end{array} := \begin{array}{c|ccccc} c_1 & a_{11} & \cdots & a_{1s} \\ \vdots & \vdots & & \vdots \\ c_s & a_{s1} & \cdots & a_{ss} \\ \hline b_1 & \cdots & b_s \\ \hline \hat{b}_1 & \cdots & \hat{b}_s \end{array} .$$

**EXAMPLE 11.5.0.25 (Commonly used embedded explicit Runge-Kutta methods)** The following two embedded RK-SSM, presented in the form of their extended Butcher schemes, provided single step methods of orders 4 & 5.

$0$					
$\frac{1}{3}$	$\frac{1}{3}$				
$\frac{1}{3}$	$\frac{1}{6}$	$\frac{1}{6}$			
$\frac{1}{2}$	$\frac{1}{8}$	$0$		$\frac{3}{8}$	
$1$	$\frac{1}{2}$	$0$	$-\frac{3}{2}$	$2$	
$y_1$	$\frac{1}{6}$	$0$	$0$	$\frac{2}{3}$	$\frac{1}{6}$
$\widehat{y}_1$	$\frac{1}{10}$	$0$	$\frac{3}{10}$	$\frac{2}{5}$	$\frac{1}{5}$

Merson's embedded RK-SSM

$0$	$\frac{1}{2}$
$\frac{1}{2}$	$0$
$\frac{1}{2}$	$\frac{1}{2}$
$1$	$0$
$1$	$0$
$\frac{3}{4}$	$\frac{5}{32}$
$\frac{3}{4}$	$\frac{7}{32}$
$\frac{3}{4}$	$\frac{13}{32}$
$\frac{3}{4}$	$-\frac{1}{32}$
$y_1$	$\frac{1}{6}$
$y_1$	$\frac{1}{3}$
$y_1$	$\frac{1}{3}$
$\widehat{y}_1$	$-\frac{1}{2}$
$\widehat{y}_1$	$\frac{7}{3}$
$\widehat{y}_1$	$\frac{7}{3}$
$\widehat{y}_1$	$\frac{13}{6}$
$\widehat{y}_1$	$-\frac{16}{3}$

Fehlberg's embedded RK-SSM

**EXAMPLE 11.5.0.26 (Adaptive timestepping for mechanical problem)** We test the effect of adaptive stepsize control in MATLAB for the equations of motion describing the planar movement of a point mass in a conservative force field  $\mathbf{x} \in \mathbb{R}^2 \mapsto F(\mathbf{x}) \in \mathbb{R}^2$ : Let  $t \mapsto \mathbf{y}(t) \in \mathbb{R}^2$  be the trajectory of point mass (in the plane).

From Newton's law:  $\ddot{\mathbf{y}} = F(\mathbf{y}) := -\frac{2\mathbf{y}}{\|\mathbf{y}\|_2^2}$ . (11.5.0.27)

As in Rem. 11.1.2.5 we can convert the second-order ODE (11.5.0.27) into an equivalent 1st-order ODE by introducing the **velocity**  $\mathbf{v} := \dot{\mathbf{y}}$  as an extra solution component:

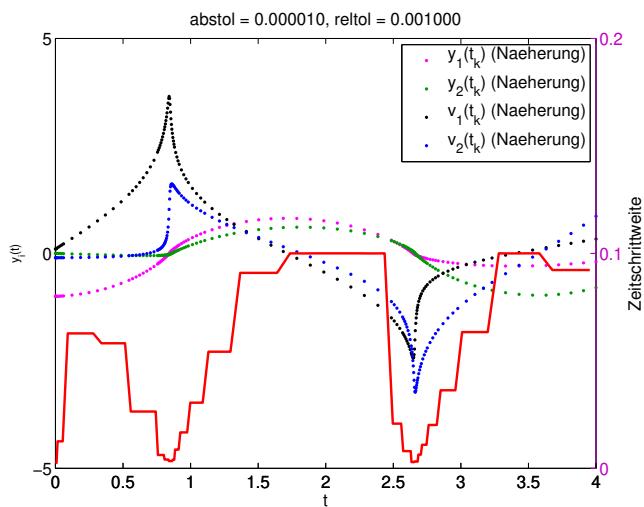
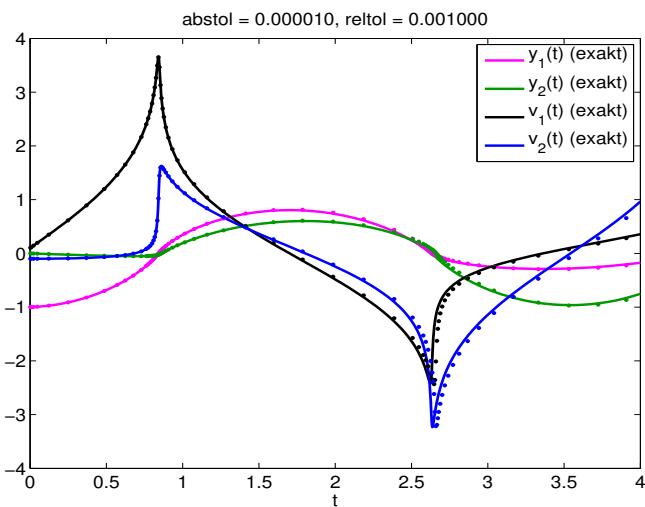
$$(11.5.0.27) \quad \Rightarrow \quad \begin{bmatrix} \dot{\mathbf{y}} \\ \dot{\mathbf{v}} \end{bmatrix} = \begin{bmatrix} \mathbf{v} \\ -\frac{2\mathbf{y}}{\|\mathbf{y}\|_2^2} \end{bmatrix}. \quad (11.5.0.28)$$

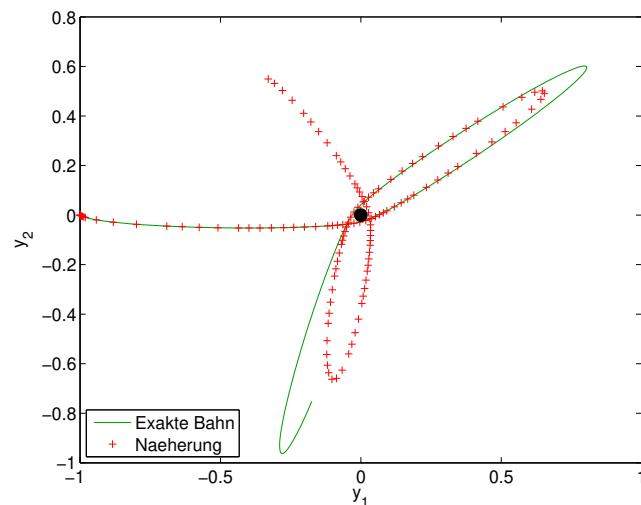
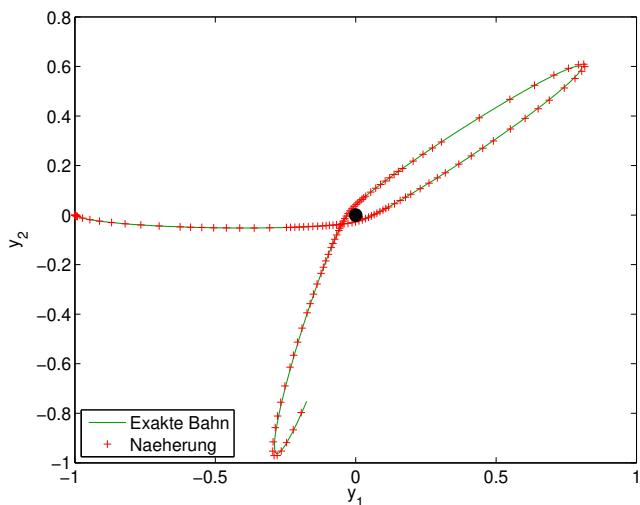
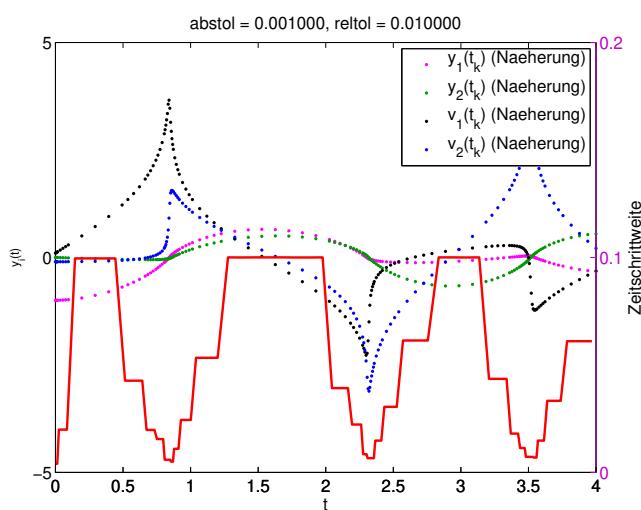
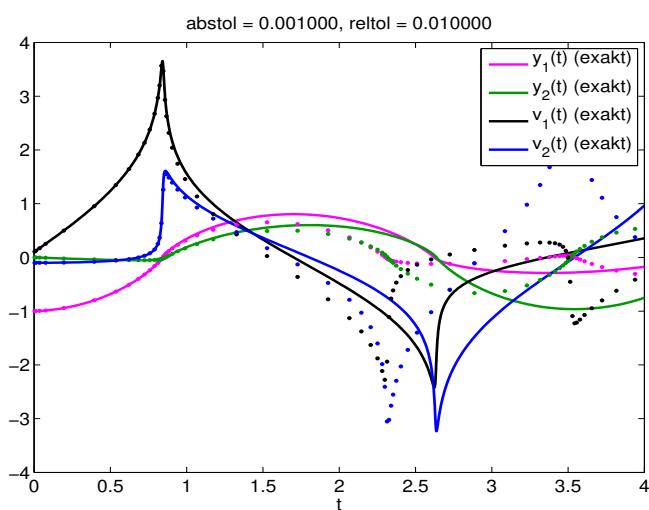
The following initial values used in the experiment:

$$\mathbf{y}(0) := \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \quad \mathbf{v}(0) := \begin{bmatrix} 0.1 \\ -0.1 \end{bmatrix}$$

Adaptive numerical integration with **ode45** from Code 12.0.0.3 with

- ```
① options.rtol = 0.001,options.atol = 1.0E-5,  
② options.rtol = 0.01,options.atol = 1.0E-3,
```





Observations:

- Fast changes in solution components captured by adaptive approach through very small timesteps.
- Completely wrong solution, if tolerance reduced slightly.

In this example we face a rather **sensitive dependence** of the trajectories on initial states or intermediate states. Small perturbations at one instance in time can have a massive impact on the solution at later times. Local stepsize control is powerless about preventing this. □

## Summary and Learning Outcomes

After having studied this chapter you should

- know the concept of evolution operator for an ODE and its relationship with solutions of associated initial value problems.
- be able to convert higher-order and non-autonomous ODEs into the form  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ .
- know about discrete evolutions and how they induce single step methods (SSMs).

- remember that single step method converge asymptotically algebraically for stepsize  $h \rightarrow 0$  and that the rate of convergence is called the order of the SSM.
- know the general form of explicit Runge-Kutta methods and Butcher schemes.
- understand when adaptive timestep control is essential for meaningful numerical integration of initial value problems.
- be able to describe the policy to time-local adaptive timestep control for embedded Runge-Kutta methods.

# Bibliography

- [Ama83] H. Amann. *Gewöhnliche Differentialgleichungen*. 1st. Berlin: Walter de Gruyter, 1983 (cit. on pp. 647, 653).
- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 647, 650, 651, 653, 657, 664, 666, 670–674, 677, 685).
- [Dea80] M.A.B. Deakin. “Applied catastrophe theory in the social and biological sciences”. In: *Bulletin of Mathematical Biology* 42.5 (1980), pp. 647–679 (cit. on p. 648).
- [DB02] P. Deuflhard and F. Bornemann. *Scientific Computing with Ordinary Differential Equations*. 2nd ed. Vol. 42. Texts in Applied Mathematics. New York: Springer, 2002 (cit. on pp. 646, 654, 675, 683).
- [Gra02] C.R. Gray. “An analysis of the Belousov-Zhabotinski reaction”. In: *Rose-Hulman Undergraduate Math Journal* 3.1 (2002) (cit. on p. 678).
- [HLW06] E. Hairer, C. Lubich, and G. Wanner. *Geometric numerical integration*. 2nd ed. Vol. 31. Springer Series in Computational Mathematics. Heidelberg: Springer, 2006 (cit. on pp. 646, 647, 675).
- [HNW93] E. Hairer, S.P. Norsett, and G. Wanner. *Solving Ordinary Differential Equations I. Nonstiff Problems*. 2nd ed. Berlin, Heidelberg, New York: Springer-Verlag, 1993 (cit. on p. 646).
- [HW11] E. Hairer and G. Wanner. *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic Problems*. Vol. 14. Springer Series in Computational Mathematics. Berlin: Springer-Verlag, 2011 (cit. on p. 646).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 647, 649, 653, 657, 666, 670, 678).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 650, 662, 664, 670, 677).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on pp. 646, 647, 651–653, 660, 668).

# Chapter 12

## Single Step Methods for Stiff Initial Value Problems



*Supplementary literature.* [DR08, Sect. 11.9]

Explicit Runge-Kutta methods with stepsize control ( $\rightarrow$  Section 11.5) seem to be able to provide approximate solutions for any IVP with good accuracy provided that tolerances are set appropriately.

Does this mean that everything settled about numerical integration?

**EXAMPLE 12.0.0.1 (Explicit adaptive RK-SSM for stiff IVP)** In this example we will witness the near failure of a high-order adaptive explicit Runge-Kutta method for a simple scalar autonomous ODE.

$$\text{IVP considered: } \dot{y} = \lambda y^2(1 - y), \quad \lambda := 500, \quad y(0) = \frac{1}{100}. \quad (12.0.0.2)$$

This is a logistic ODE as introduced in Ex. 11.1.1.1. We try to solve it by means of an explicit adaptive embedded Runge-Kutta-Fehlberg method ( $\rightarrow$  Rem. 11.5.0.24) using the class **ode45** from § 11.4.0.18 (Preprocessor switch MATLABCOEFF activated).

### C++ code 12.0.0.3: Solving (12.0.0.2) with class **ode45**

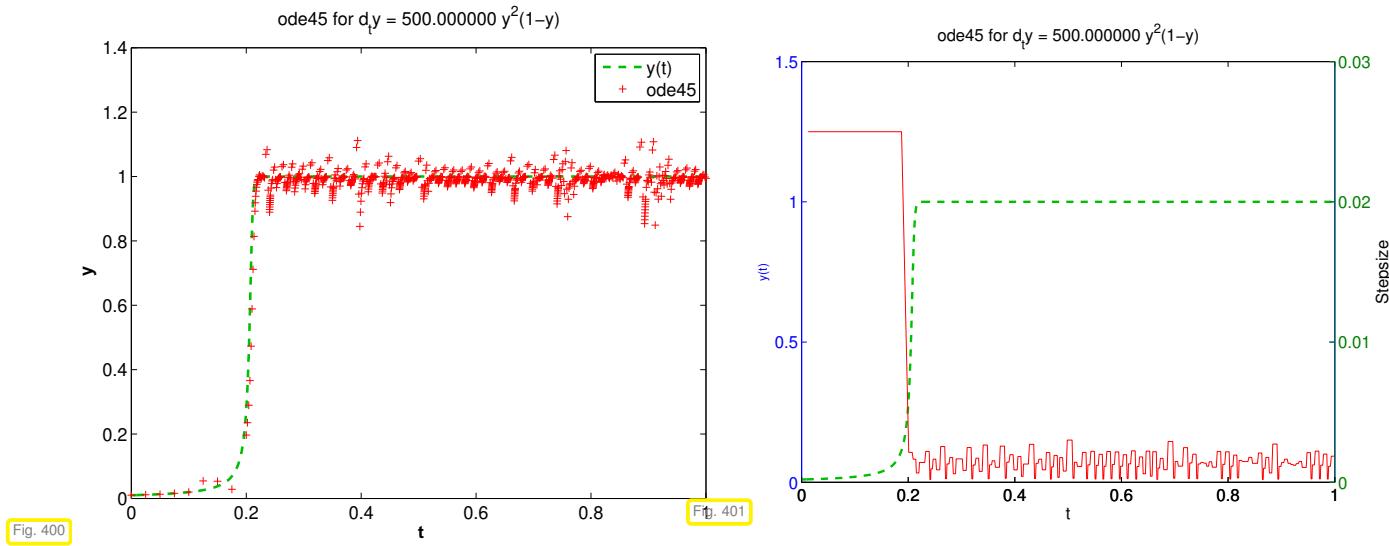
```
2 // Types to be used for a scalar ODE with state space ℝ
3 using StateType = double;
4 using RhsType = std::function<StateType(StateType)>;
5 // Logistic differential equation (11.1.1.2)
6 double lambda = 500.0;
7 RhsType f = [lambda](StateType y) { return lambda*y*y*(1-y); };
8 StateType y0 = 0.01; // Initial value, will create a STIFF IVP
9 // State space ℝ, simple modulus supplies norm
10 auto normFunc = [] (StateType x){ return fabs(x); };

11
12 // Invoke explicit Runge-Kutta method with stepsize control
13 ode45<StateType, RhsType> integrator(f);
14 // Set rather loose tolerances
15 integrator.options.rtol = 0.1;
16 integrator.options.atol = 0.001;
17 integrator.options.min_dt = 1E-18;
18 std::vector<std::pair<StateType, double>> states =
    integrator.solve(y0, 1.0, normFunc);
19 // Output information accumulation during numerical integration
20 integrator.options.do_statistics = true; integrator.print();
```

Statistics of the integrator run



|                                             |      |
|---------------------------------------------|------|
| <code>1</code> — number of steps :          | 183  |
| <code>2</code> — number of rejected steps : | 185  |
| <code>3</code> — function calls :           | 1302 |

Stepsize control of **ode45** running amok!

❓ The solution is virtually constant from  $t > 0.2$  and, nevertheless, the integrator uses tiny timesteps until the end of the integration interval. Why this crazy behavior?

## Contents

|                                                                                       |     |
|---------------------------------------------------------------------------------------|-----|
| <b>12.1 Model Problem Analysis</b> . . . . .                                          | 693 |
| <b>12.2 Stiff Initial Value Problems</b> . . . . .                                    | 706 |
| <b>12.3 Implicit Runge-Kutta Single Step Methods</b> . . . . .                        | 710 |
| 12.3.1 The implicit Euler method for stiff IVPs . . . . .                             | 711 |
| 12.3.2 Collocation single step methods . . . . .                                      | 712 |
| 12.3.3 General implicit RK-SSMs . . . . .                                             | 715 |
| 12.3.4 Model Problem Analysis for Implicit Runge-Kutta Singlge-Step Method (IRK-SSMs) | 717 |
| <b>12.4 Semi-implicit Runge-Kutta Methods</b> . . . . .                               | 723 |
| <b>12.5 Splitting methods</b> . . . . .                                               | 725 |

## 12.1 Model Problem Analysis

Fortunately, full insight into the observations made in Ex. 12.0.0.1 can already be gleaned from a scalar linear model problem that is extremely easy to analyze.



*Supplementary literature.* See also [Han02, Ch. 77], [QSS00, Sect. 11.3.3].

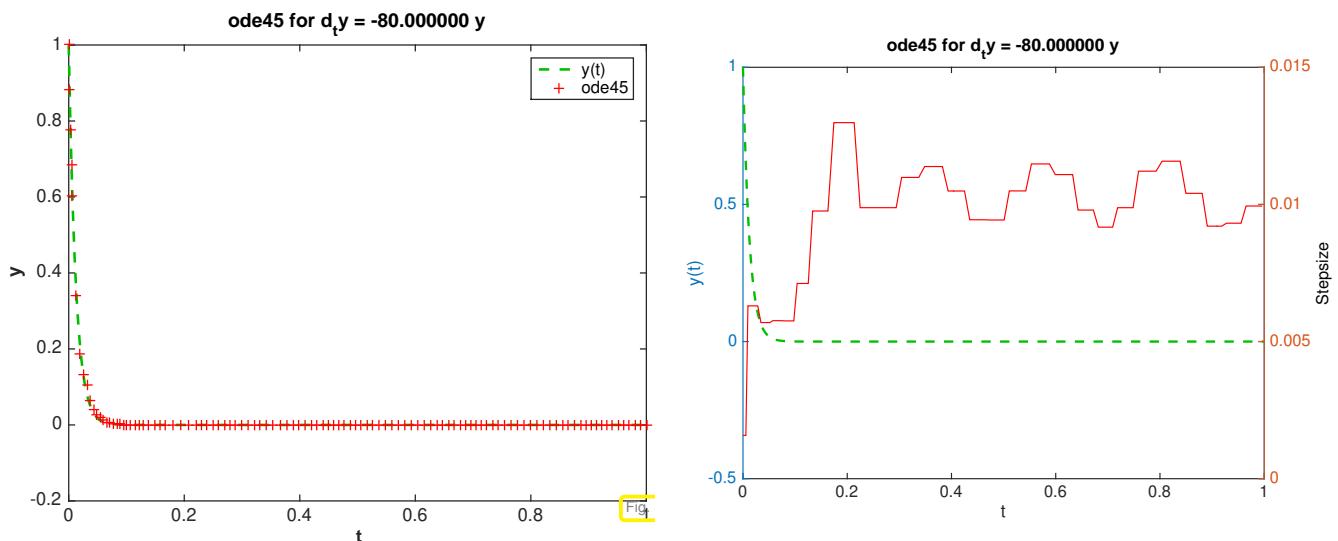
**EXPERIMENT 12.1.0.1 (Adaptive explicit RK-SSM for scalar linear decay ODE)** To rule out that what we observed in Ex. 12.0.0.1 might have been a quirk of the IVP (12.0.0.2) we conduct the same investigations for the simple **linear, scalar, autonomous** IVP

$$\dot{y} = \lambda y \quad , \quad \lambda := -80 \quad , \quad y(0) = 1 . \quad (12.1.0.2)$$

We use the **ode45** class to solve (12.1.0.2) with the same parameters as in Code 12.0.0.3.



|                                             |     |
|---------------------------------------------|-----|
| <code>1</code> — number of steps :          | 33  |
| <code>2</code> — number of rejected steps : | 32  |
| <code>3</code> — function calls :           | 231 |



Observation: Though  $y(t) \approx 0$  for  $t > 0.1$ , the integrator keeps on using “unreasonably small” timesteps even then.

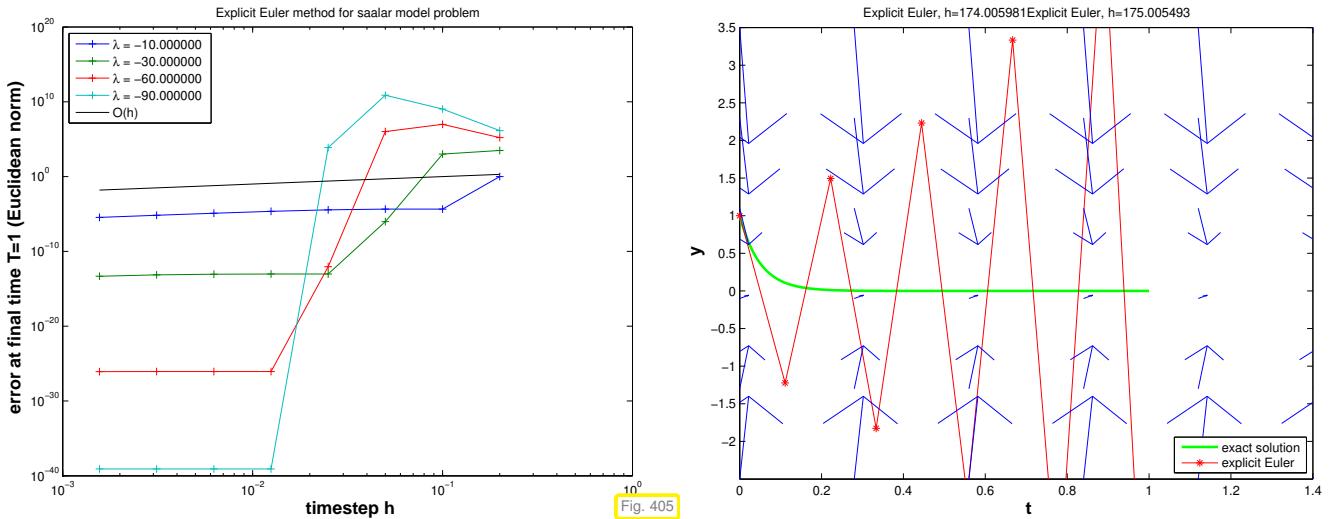
In this section we will discover a simple explanation for the startling behavior of **ode45** in Ex. 12.0.0.1.

**EXAMPLE 12.1.0.3 (Blow-up of explicit Euler method)** The simplest explicit RK-SSM is the explicit Euler method, see Section 11.2.1. We know that it should converge like  $O(h)$  for meshwidth  $h \rightarrow 0$ . In this example we will see that this may be true only for sufficiently small  $h$ , which may be extremely small.

- ◆ We consider the IVP for the scalar linear decay ODE:

$$\dot{y} = f(y) := \lambda y, \quad \lambda \ll 0, \quad y(0) = 1.$$

- ◆ We apply the explicit Euler method (11.2.1.4) with uniform timestep  $h = 1/N$ ,  $N \in \{5, 10, 20, 40, 80, 160, 320, 640\}$ .



$\lambda \ll 0$ : blow-up of  $y_k$  for large timestep  $h$

$\lambda = 20$ :  $\text{---} \hat{=} y(t)$ ,  $\text{---} \hat{=} \text{Euler polygon}$

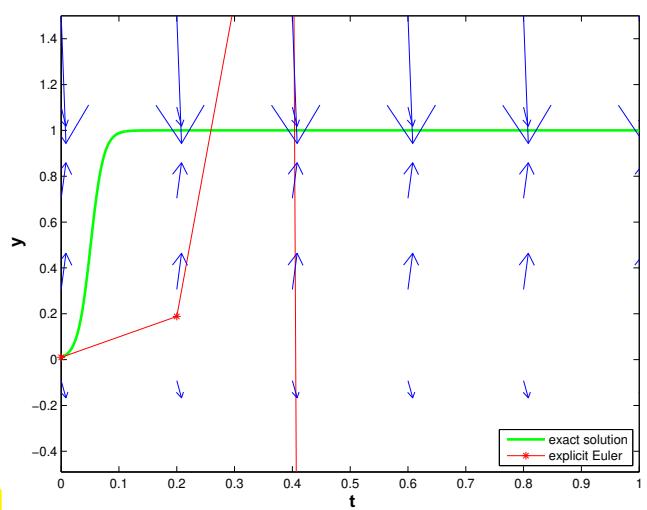
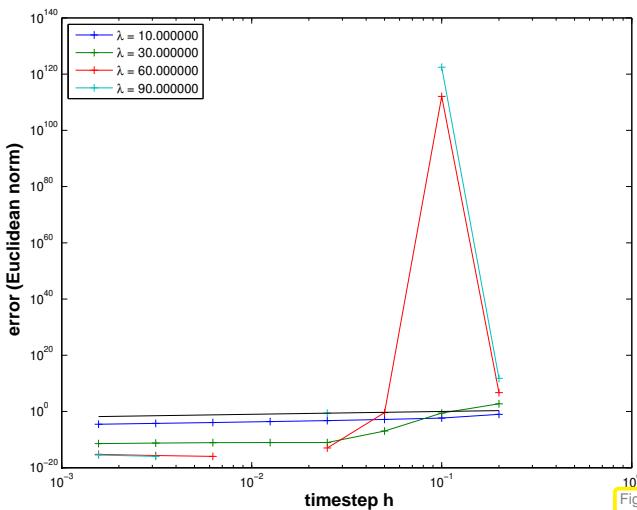
Explanation: From Fig. 405 we draw the geometric conclusion that, if  $h$  is “large in comparison with  $\lambda^{-1}$ ”, then the approximations  $y_k$  may miss the stationary point  $y = 0$  due to overshooting.

This leads to a sequence  $(y_k)_k$  with exponentially increasing oscillations.

- ◆ Now we look at an IVP for the logistic ODE, see Ex. 11.1.1.1:

$$\dot{y} = f(y) := \lambda y(1 - y), \quad y(0) = 0.01.$$

- ◆ As before, we apply the explicit Euler method (11.2.1.4) with uniform timestep  $h = 1/N$ ,  $N \in \{5, 10, 20, 40, 80, 160, 320, 640\}$ .



$\lambda$  large: blow-up of  $y_k$  for large timestep  $h$

For large timesteps  $h$  we also observe oscillatory blow-up of the sequence  $(y_k)_k$ .

#### Deeper analysis:

For  $y \approx 1$ :  $f(y) \approx \lambda(1-y) \Rightarrow$  If  $y(t_0) \approx 1$ , then the solution of the IVP will behave like the solution of  $\dot{y} = \lambda(1-y)$ , which is a linear ODE. Similarly,  $z(t) := 1 - y(t)$  will behave like the solution of the “decay equation”  $\dot{z} = -\lambda z$ . Thus, around the stationary point  $y = 1$  the explicit Euler method behaves like it did for  $\dot{y} = \lambda y$  in the vicinity of the stationary point  $y = 0$ ; it grossly overshoots. □

**§12.1.0.4 (Linear model problem analysis: explicit Euler method)** The phenomenon observed in the two previous examples is accessible to a remarkably simple rigorous analysis: Motivated by the considerations in Ex. 12.1.0.3 we study the explicit Euler method (11.2.1.4) for the

$$\text{linear model problem: } \dot{y} = \lambda y, \quad y(0) = y_0, \quad \text{with } \lambda \ll 0, \quad (12.1.0.5)$$

which has *exponentially decaying* exact solution

$$y(t) = y_0 \exp(\lambda t) \rightarrow 0 \quad \text{for } t \rightarrow \infty.$$

Recall the recursion for the explicit Euler with uniform timestep  $h > 0$  method for (12.1.0.5):

$$(11.2.1.4) \text{ for } f(y) = \lambda y: \quad y_{k+1} = y_k(1 + \lambda h). \quad (12.1.0.6)$$

We easily get a closed form expression for the approximations  $y_k$ :

►  $y_k = y_0(1 + \lambda h)^k \Rightarrow |y_k| \rightarrow \begin{cases} 0 & , \text{ if } \lambda h > -2 \quad (\text{qualitatively correct}), \\ \infty & , \text{ if } \lambda h < -2 \quad (\text{qualitatively wrong}). \end{cases}$

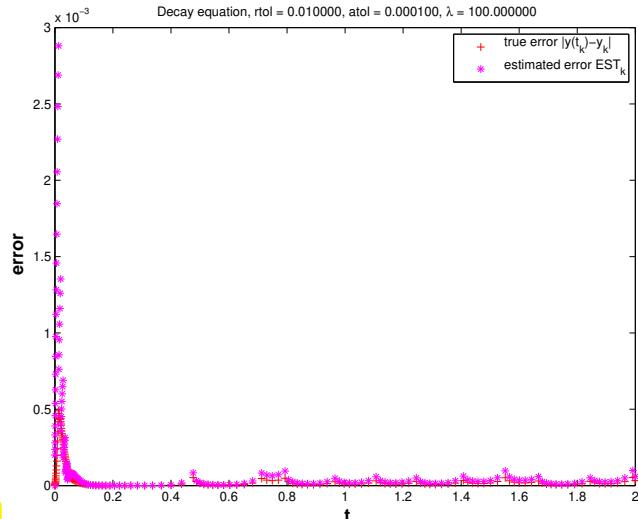
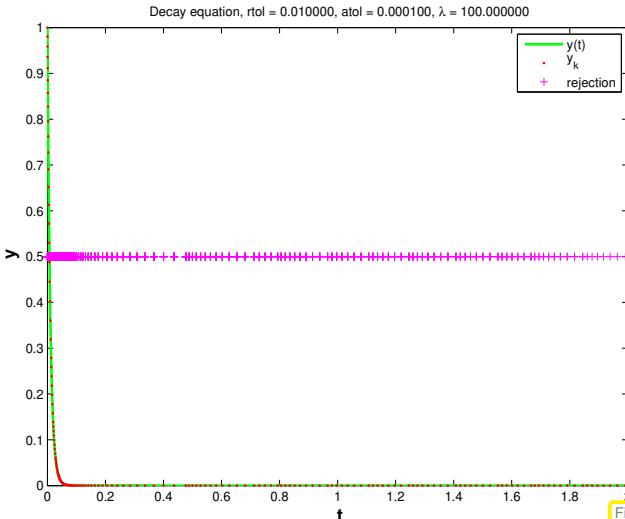
#### Observed: timestep constraint

Only if  $|\lambda| h < 2$  we obtain a decaying solution by the explicit Euler method!

Could it be that the timestep control is desperately trying to enforce the qualitatively correct behavior of the numerical solution in Ex. 12.1.0.3? Let us examine how the simple stepsize control of Code 11.5.0.10 fares for model problem (12.1.0.5):

**EXAMPLE 12.1.0.8 (Simple adaptive timestepping for fast decay)** In this example we let a transparent adaptive timestep struggle with “overshooting”:

- ◆ “Linear model problem IVP”:  $\dot{y} = \lambda y$ ,  $y(0) = 1$ ,  $\lambda = -100$
- ◆ Simple adaptive timestepping method as in Exp. 11.5.0.12, see Code 11.5.0.10. Timestep control based on the pair of 1st-order explicit Euler method and 2nd-order explicit trapezoidal method.



Observation: in fact, stepsize control enforces small timesteps even if  $y(t) \approx 0$  and persistently triggers rejections of timesteps. This is necessary to prevent overshooting in the Euler method, which contributes to the estimate of the one-step error.

We see the purpose of *stepsize control thwarted*, because after only a very short time the solution is almost zero and then, in fact, large timesteps should be chosen. □

Are these observations a particular “flaw” of the explicit Euler method? Let us study the behavior of another simple explicit Runge-Kutta method applied to the linear model problem.

**EXAMPLE 12.1.0.9 (Explicit trapzoidal method for decay equation → [DR08, Ex. 11.29])**

Recall recursion for the explicit trapezoidal method derived in Ex. 11.4.0.4:

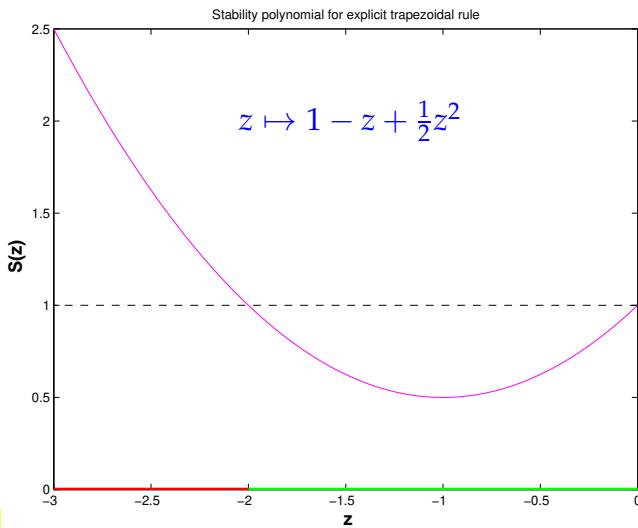
$$\mathbf{k}_1 = \mathbf{f}(t_0, \mathbf{y}_0), \quad \mathbf{k}_2 = \mathbf{f}(t_0 + h, \mathbf{y}_0 + h\mathbf{k}_1), \quad \mathbf{y}_1 = \mathbf{y}_0 + \frac{h}{2}(\mathbf{k}_1 + \mathbf{k}_2). \quad (11.4.0.6)$$

Apply it to the model problem (12.1.0.5), that is, the scalar autonomous ODE with right hand side function  $\mathbf{f}(y) = f(y) = \lambda y$ ,  $\lambda < 0$ :

$$\blacktriangleright \quad k_1 = \lambda y_0, \quad k_2 = \lambda(y_0 + hk_1) \quad \Rightarrow \quad y_1 = \underbrace{(1 + \lambda h + \frac{1}{2}(\lambda h)^2)}_{=:S(h\lambda)} y_0. \quad (12.1.0.10)$$

- The sequence of approximations generated by the explicit trapezoidal rule can be expressed in closed form as

$$y_k = S(h\lambda)^k y_0, \quad k = 0, \dots, N. \quad (12.1.0.11)$$



**§12.1.0.13 (Model problem analysis for general explicit Runge-Kutta single step methods)** We generalize the approach taken in Ex. 12.1.0.9 and apply an explicit Runge-Kutta method ( $\rightarrow$  Def. 11.4.0.9) encoded by the Butcher scheme  $\begin{array}{c|cc} c & \mathfrak{A} \\ \hline & b^T \end{array}$  to the autonomous scalar linear ODE (12.1.0.5) ( $\dot{y} = \lambda y$ ). We write down the equations for the increments and  $y_1$  from Def. 11.4.0.9 for  $f(y) := \lambda y$  and then convert the resulting system of equations into matrix form:

$$\begin{aligned} \blacktriangleright \quad k_i &= \lambda(y_0 + h \sum_{j=1}^{i-1} a_{ij} k_j), & \Rightarrow \quad \begin{bmatrix} \mathbf{I} - z\mathfrak{A} & \mathbf{0} \\ -z\mathbf{b}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{k} \\ y_1 \end{bmatrix} &= y_0 \begin{bmatrix} \mathbf{1} \\ 1 \end{bmatrix}, \\ y_1 &= y_0 + h \sum_{i=1}^s b_i k_i \end{aligned} \quad (12.1.0.14)$$

where  $\mathbf{k} \in \mathbb{R}^s \doteq$  denotes the vector  $[k_1, \dots, k_s]^\top / \lambda$  of increments, and  $z := \lambda h$ . Next we apply block Gaussian elimination ( $\rightarrow$  Rem. 2.3.1.11) to solve for  $y_1$  and obtain

$$y_1 = S(z)y_0 \quad \text{with} \quad S(z) := 1 + z\mathbf{b}^T(\mathbf{I} - z\mathfrak{A})^{-1}\mathbf{1}. \quad (12.1.0.15)$$

Alternatively we can express  $y_1$  through determinants appealing to Cramer's rule,

$$y_1 = y_0 \frac{\det \begin{bmatrix} \mathbf{I} - z\mathfrak{A} & \mathbf{1} \\ -z\mathbf{b}^T & 1 \end{bmatrix}}{\det \begin{bmatrix} \mathbf{I} - z\mathfrak{A} & \mathbf{0} \\ -z\mathbf{b}^T & 1 \end{bmatrix}} \Rightarrow S(z) = \det(\mathbf{I} - z\mathfrak{A} + z\mathbf{1}\mathbf{b}^T), \quad (12.1.0.16)$$

and note that  $\mathfrak{A}$  is a strictly lower triangular matrix, which means that  $\det(\mathbf{I} - z\mathfrak{A}) = 1$ . Thus we have proved the following theorem.

**Theorem 12.1.0.17. Stability function of explicit Runge-Kutta methods** → [Han02, Thm. 77.2], [QSS00, Sect. 11.8.4]

The discrete evolution  $\Psi_\lambda^h$  of an explicit  $s$ -stage Runge-Kutta single step method (→ Def. 11.4.0.9) with Butcher scheme  $\begin{array}{c|cc} \mathbf{c} & \mathfrak{A} \\ \hline & \mathbf{b}^T \end{array}$  (see (11.4.0.11)) for the ODE  $\dot{\mathbf{y}} = \lambda \mathbf{y}$  amounts to a multiplication with the number

$$\Psi_\lambda^h = S(\lambda h) \Leftrightarrow y_1 = S(\lambda h)y_0 ,$$

where  $S$  is the **stability function** (SF)

$$S(z) := 1 + z\mathbf{b}^T(\mathbf{I} - z\mathfrak{A})^{-1}\mathbf{1} = \det(\mathbf{I} - z\mathfrak{A} + z\mathbf{1}\mathbf{b}^T) , \quad \mathbf{1} := [1, \dots, 1]^\top \in \mathbb{R}^s . \quad (12.1.0.18)$$

□

**EXAMPLE 12.1.0.19 (Stability functions of explicit Runge-Kutta single step methods)** From Thm. 12.1.0.17 and their Butcher schemes we can instantly compute the stability functions of explicit RK-SSM. We do this for a few methods whose Butcher schemes were listed in Ex. 11.4.0.15

- Explicit Euler method (11.2.1.4):  $\begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array} \Rightarrow S(z) = 1 + z .$

- Expl. trapezoidal method (11.4.0.6):  $\begin{array}{c|cc} 0 & 0 & 0 \\ \hline 1 & 1 & 0 \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array} \Rightarrow S(z) = 1 + z + \frac{1}{2}z^2 .$

- Classical RK4 method:  $\begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ \hline \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ \hline & \frac{1}{6} & \frac{2}{6} & \frac{2}{6} & \frac{1}{6} \end{array} \Rightarrow S(z) = 1 + z + \frac{1}{2}z^2 + \frac{1}{6}z^3 + \frac{1}{24}z^4 .$

These examples confirm an immediate consequence of the determinant formula for the stability function  $S(z)$ .

#### Corollary 12.1.0.20. Polynomial stability function of explicit RK-SSM

For a consistent (→ Def. 11.3.1.11)  $s$ -stage explicit Runge-Kutta single step method according to Def. 11.4.0.9 the stability function  $S$  defined by (12.1.0.18) is a non-constant **polynomial** of degree  $\leq s$ , that is,  $S \in \mathcal{P}_s$ .

□

**Remark 12.1.0.21 (Stability function and exponential function)** Let us compare the two evolution operators:

- $\Phi \hat{=} \text{evolution operator}$  (→ Def. 11.1.3.2) for  $\dot{\mathbf{y}} = \lambda \mathbf{y}$ ,

- $\Psi \hat{=} \text{discrete evolution operator} (\rightarrow \S 11.3.1.1)$  for an  $s$ -stage Runge-Kutta single step method.

$$\Phi^h y = e^{\lambda h} y \longleftrightarrow \Psi^h y = S(\lambda h) y .$$

In light of that  $\Psi$  is supposed to be an approximation for  $\Phi$ ,  $\Psi \approx \Phi$ , see (11.3.1.3), we expect that

$$S(z) \approx \exp(z) \quad \text{for small } |z| . \quad (12.1.0.22)$$

A more precise statement is made by the following lemma:

**Lemma 12.1.0.23. Stability function as approximation of exp for small arguments**

Let  $S$  denote the stability function of an  $s$ -stage explicit Runge-Kutta single step method of order  $q \in \mathbb{N}$ . Then

$$|S(z) - \exp(z)| = O(|z|^{q+1}) \quad \text{for } |z| \rightarrow 0 . \quad (12.1.0.24)$$

This means that the lowest  $q + 1$  coefficients of  $S(z)$  must be equal to the first coefficients of the exponential series:

$$S(z) = \sum_{j=0}^q \frac{1}{j!} z^j + z^{q+1} p(z) \quad \text{with some } p \in \mathcal{P}_{s-q-1} .$$

In order to match the first  $q$  terms of the exponential series, we need at least  $S \in \mathcal{P}_q$ , which entails a minimum of  $q$  stages.

**Corollary 12.1.0.25. Stages limit order of explicit RK-SSM**

An explicit  $s$ -stage RK-SSM has maximal order  $q \leq s$ .

§12.1.0.26 (Stability induced timestep constraint) In § 12.1.0.13 we established that for the sequence  $(y_k)_{k=0}^\infty$  produced by an explicit Runge-Kutta single step method applied to the linear scalar model ODE  $\dot{y} = \lambda y$ ,  $\lambda \in \mathbb{R}$ , with uniform timestep  $h > 0$  holds

$$y_{k+1} = S(\lambda h) y_k \Rightarrow y_k = S(\lambda h^k) y_0 .$$



$$\begin{aligned} (y_k)_{k=0}^\infty \text{ non-increasing} &\Leftrightarrow |S(\lambda h)| \leq 1 , \\ (y_k)_{k=0}^\infty \text{ exponentially increasing} &\Leftrightarrow |S(\lambda h)| > 1 . \end{aligned} \quad (12.1.0.27)$$

where  $S = S(z)$  is the stability function of the RK-SSM as defined in (12.1.0.18).

Invariably polynomials tend to  $\pm\infty$  for large (in modulus) arguments:

$$\forall S \in \mathcal{P}_s, S \neq \text{const} : \lim_{|z| \rightarrow \infty} S(z) = \infty \quad \text{uniformly} . \quad (12.1.0.28)$$

So, for any  $\lambda \neq 0$  there will be a threshold  $h_{\max} > 0$  so that  $|y_k| \rightarrow \infty$  as  $|h| > h_{\max}$ .

Reversing the argument we arrive at a **timestep constraint**, as already observed for the explicit Euler methods in § 12.1.0.4.

Only if one ensures that  $|\lambda h|$  is sufficiently small, one can avoid exponentially increasing approximations  $y_k$  (qualitatively wrong for  $\lambda < 0$ ) when applying an explicit RK-SSM to the model problem (12.1.0.5) with uniform timestep  $h > 0$ .

For  $\lambda \ll 0$  this stability induced timestep constraint may force  $h$  to be much *smaller than required by demands on accuracy*: in this case timestepping becomes **inefficient**.  $\square$

**Remark 12.1.0.29 (Stepsize control detects instability)** Ex. 12.0.0.1, Ex. 12.1.0.8 send the message that local-in-time stepsize control as discussed in Section 11.5 selects timesteps that avoid blow-up, with a hefty price tag however in terms of computational cost and poor accuracy.  $\square$

Objection: simple linear scalar IVP (12.1.0.5) may be an oddity rather than a model problem: the weakness of explicit Runge-Kutta methods discussed above may be just a peculiar response to an unusual situation. Let us extend our investigations to **systems of linear ODEs**,  $d > 1$ .

**§12.1.0.30 (Systems of linear ordinary differential equations)** A generic linear ordinary differential equation on state space  $\mathbb{R}^d$  has the form

$$\dot{\mathbf{y}} = \mathbf{My} \quad \text{with a matrix } \mathbf{M} \in \mathbb{R}^{d,d}. \quad (12.1.0.31)$$

As explained in [NS02, Sect. 8.1], (12.1.0.31) can be solved by **diagonalization**: If we can find a *regular* matrix  $\mathbf{V} \in \mathbb{C}^{d,d}$  such that

$$\mathbf{MV} = \mathbf{VD} \quad \text{with diagonal matrix } \mathbf{D} = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_d \end{bmatrix} \in \mathbb{C}^{d,d}, \quad (12.1.0.32)$$

then the 1-parameter family of global solutions of (12.1.0.31) is given by

$$\mathbf{y}(t) = \mathbf{V} \begin{bmatrix} \exp(\lambda_1 t) & & 0 \\ & \ddots & \\ 0 & & \exp(\lambda_d t) \end{bmatrix} \mathbf{V}^{-1} \mathbf{y}_0, \quad \mathbf{y}_0 \in \mathbb{R}^d. \quad (12.1.0.33)$$

The columns of  $\mathbf{V}$  are a basis of **eigenvectors** of  $\mathbf{M}$ , the  $\lambda_j \in \mathbb{C}$ ,  $j = 1, \dots, d$  are the associated **eigenvalues** of  $\mathbf{M}$ , see Def. 9.1.0.1.

The idea behind diagonalization is the transformation of (12.1.0.31) into  $d$  *decoupled* scalar linear ODEs:

$$\dot{\mathbf{y}} = \mathbf{My} \quad \xrightarrow{\mathbf{z}(t) := \mathbf{V}^{-1}\mathbf{y}(t)} \quad \dot{\mathbf{z}} = \mathbf{Dz} \quad \leftrightarrow \quad \begin{array}{l} \dot{z}_1 = \lambda_1 z_1 \\ \vdots \\ \dot{z}_d = \lambda_d z_d \end{array}, \quad \text{since } \mathbf{M} = \mathbf{VDV}^{-1}.$$

The formula (12.1.0.33) can be generalized to

$$\mathbf{y}(t) = \exp(\mathbf{Mt})\mathbf{y}_0 \quad \text{with matrix exponential} \quad \exp(\mathbf{B}) := \sum_{k=0}^{\infty} \frac{1}{k!} \mathbf{B}^k, \quad \mathbf{B} \in \mathbb{C}^{d,d}. \quad (12.1.0.34)$$

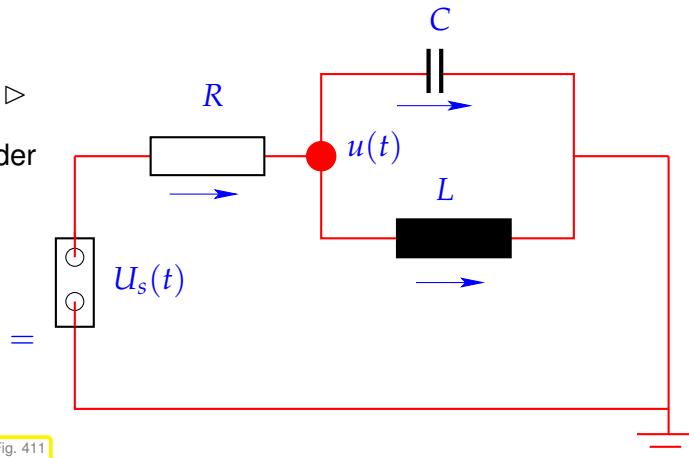
**EXAMPLE 12.1.0.35 (Transient simulation of RLC-circuit)** This example demonstrates the diagonalization of a linear system of ODEs.

Consider circuit from Ex. 11.1.1.9

Transient nodal analysis leads to the second-order linear ODE

$$\ddot{u} + \alpha \dot{u} + \beta u = g(t),$$

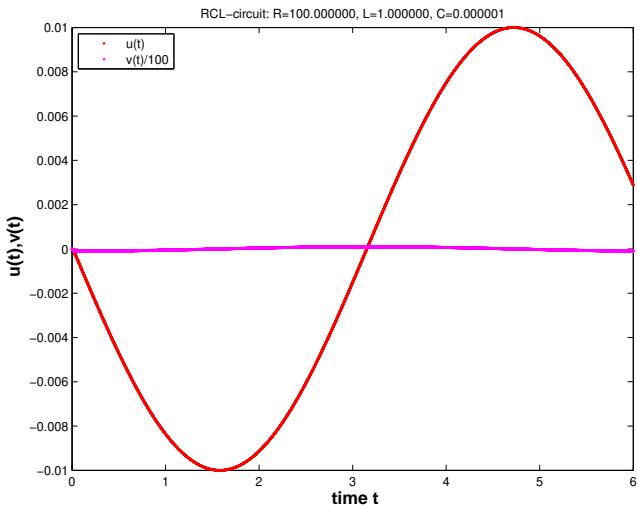
with coefficients  $\alpha := (RC)^{-1}$ ,  $\beta = (LC)^{-1}$ ,  $g(t) = \alpha U_s(t)$ .



We transform it to a linear 1st-order ODE as in Rem. 11.1.2.5 by introducing  $v := \dot{u}$  as additional solution component:

$$\underbrace{\begin{bmatrix} \dot{u} \\ v \end{bmatrix}}_{=: \dot{y}} = \underbrace{\begin{bmatrix} 0 & 1 \\ -\beta & -\alpha \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}}_{=: f(t, y)} - \underbrace{\begin{bmatrix} 0 \\ g(t) \end{bmatrix}}_{=: f(t, y)}, \quad \text{with } \beta \gg \alpha \gg 1 \text{ in usual settings.}$$

We integrate IVPs for this ODE by means of the adaptive integrator **ode45** from § 11.4.0.18.



$R = 100\Omega$ ,  $L = 1\text{H}$ ,  $C = 1\mu\text{F}$ ,  $U_s(t) = 1\text{V} \sin(t)$ ,  $u(0) = v(0) = 0$  ("switch on")

**ode45** statistics:

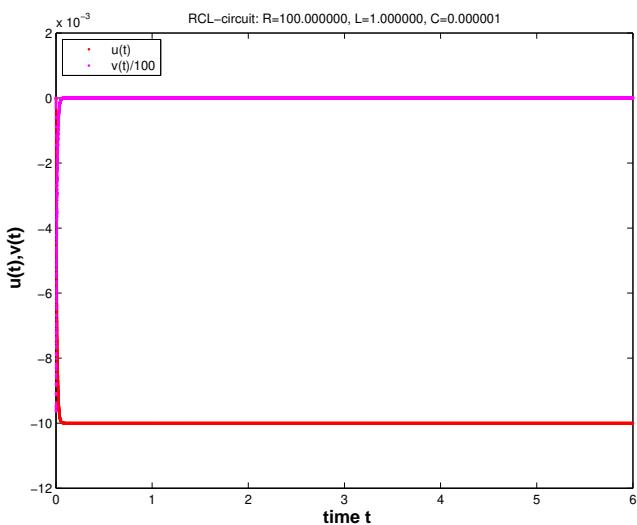
17897 successful steps

1090 failed attempts

113923 function evaluations

Inefficient: way more timesteps than required for resolving smooth solution, cf. remark in the end of § 12.1.0.26.

Maybe the time-dependent right hand side due to the time-harmonic excitation severely affects **ode45**? Let us try a constant exciting voltage:



$R = 100\Omega$ ,  $L = 1\text{H}$ ,  $C = 1\mu\text{F}$ ,  $U_s(t) = 1\text{V}$ ,  $u(0) = v(0) = 0$  ("switch on")

**ode45** statistics:

17901 successful steps

1210 failed attempts

114667 function evaluations

Tiny timesteps despite virtually constant solution!

We make the same observation as in Ex. 12.0.0.1, Ex. 12.1.0.8: the local-in-time stepsize control of **ode45**

(→ Section 11.5) enforces extremely small timesteps though the solution almost constant except at  $t = 0$ .

To understand the structure of the solutions for this transient circuit example, let us apply the diagonalization technique from § 12.1.0.30 to the linear ODE

$$\dot{\mathbf{y}} = \underbrace{\begin{bmatrix} 0 & 1 \\ -\beta & -\alpha \end{bmatrix}}_{=: \mathbf{M}} \mathbf{y} , \quad \mathbf{y}(0) = \mathbf{y}_0 \in \mathbb{R}^2 . \quad (12.1.0.36)$$

Above we face the situation  $\beta \gg \frac{1}{4}\alpha^2 \gg 1$ .

We can obtain the general solution of  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$ ,  $\mathbf{M} \in \mathbb{R}^{2,2}$ , by diagonalization of  $\mathbf{M}$  (if possible):

$$\mathbf{MV} = \mathbf{M}(\mathbf{v}_1, \mathbf{v}_2) = (\mathbf{v}_1, \mathbf{v}_2) \begin{bmatrix} \lambda_1 & \\ & \lambda_2 \end{bmatrix} . \quad (12.1.0.37)$$

where  $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^2 \setminus \{0\}$  are the eigenvectors of  $\mathbf{M}$ ,  $\lambda_1, \lambda_2$  are the eigenvalues of  $\mathbf{M}$ , see Def. 9.1.0.1. The latter are the roots of the characteristic polynomial  $t \mapsto \chi(t) := t^2 + \alpha t + \beta$  in  $\mathbb{C}$ , and we find

$$\lambda_{1/2} = \frac{1}{2}(-\alpha \pm D) , \quad D := \begin{cases} \sqrt{\alpha^2 - 4\beta} & , \text{ if } \alpha^2 \geq 4\beta , \\ i\sqrt{4\beta - \alpha^2} & , \text{ if } \alpha^2 < 4\beta . \end{cases}$$

Note that the eigenvalues have a large (in modulus) negative real part and a non-vanishing imaginary part in the setting of the experiment.

Then we transform  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$  into *decoupled* scalar linear ODEs:

$$\dot{\mathbf{y}} = \mathbf{M}\mathbf{y} \Leftrightarrow \mathbf{V}^{-1}\dot{\mathbf{y}} = \mathbf{V}^{-1}\mathbf{M}\mathbf{V}(\mathbf{V}^{-1}\mathbf{y}) \stackrel{\mathbf{z}(t) := \mathbf{V}^{-1}\mathbf{y}(t)}{\Leftrightarrow} \dot{\mathbf{z}} = \begin{bmatrix} \lambda_1 & \\ & \lambda_2 \end{bmatrix} \mathbf{z} . \quad (12.1.0.38)$$

This yields the general solution of the ODE  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$ , see also [Str09, Sect. 5.6]:

$$\mathbf{y}(t) = A\mathbf{v}_1 \exp(\lambda_1 t) + B\mathbf{v}_2 \exp(\lambda_2 t) , \quad A, B \in \mathbb{R} . \quad (12.1.0.39)$$

Note:  $t \mapsto \exp(\lambda_i t)$  is general solution of the ODE  $\dot{z}_i = \lambda_i z_i$ . □

**§12.1.0.40 (“Diagonalization” of explicit Euler method)** Recall discrete evolution of explicit Euler method (11.2.1.4) for ODE  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$ ,  $\mathbf{M} \in \mathbb{R}^{d,d}$ :

$$\Psi^h \mathbf{y} = \mathbf{y} + h\mathbf{M}\mathbf{y} \Leftrightarrow \mathbf{y}_{k+1} = \mathbf{y}_k + h\mathbf{M}\mathbf{y}_k .$$

As in § 12.1.0.30 we assume that  $\mathbf{M}$  can be diagonalized, that is (12.1.0.32) holds:  $\mathbf{V}^{-1}\mathbf{M}\mathbf{V} = \mathbf{D}$  with a diagonal matrix  $\mathbf{D} \in \mathbb{C}^{d,d}$  containing the eigenvalues of  $\mathbf{M}$  on its diagonal. Next, apply the **decoupling by diagonalization** idea to the recursion of the explicit Euler method.

$$\mathbf{V}^{-1}\mathbf{y}_{k+1} = \mathbf{V}^{-1}\mathbf{y}_k + h\mathbf{V}^{-1}\mathbf{M}\mathbf{V}(\mathbf{V}^{-1}\mathbf{y}_k) \stackrel{\mathbf{z}_k := \mathbf{V}^{-1}\mathbf{y}_k}{\Leftrightarrow} \underbrace{(\mathbf{z}_{k+1})_i = (\mathbf{z}_k)_i + h\lambda_i(\mathbf{z}_k)_i}_{\triangleq \text{explicit Euler step for } \dot{z}_i = \lambda_i z_i} . \quad (12.1.0.41)$$

Crucial insight:

The explicit Euler method generates uniformly bounded solution sequences  $(\mathbf{y}_k)_{k=0}^\infty$  for  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$  with diagonalizable matrix  $\mathbf{M} \in \mathbb{R}^{d,d}$  with eigenvalues  $\lambda_1, \dots, \lambda_d$ , if and only if it generates uniformly bounded sequences for all the scalar ODEs  $\dot{z} = \lambda_i z$ ,  $i = 1, \dots, d$ .

So far we conducted the model problem analysis under the premises  $\lambda < 0$ .

However, in Ex. 12.1.0.35 we face  $\lambda_{1/2} = -\frac{1}{2}(\alpha \pm i\sqrt{4\beta - \alpha^2})$  (complex eigenvalues!). Let us now examine how the explicit Euler method and even general explicit RK-methods respond to them.

**EXAMPLE 12.1.0.42 (Explicit Euler method for damped oscillations)** Consider linear model IVP (12.1.0.5) for  $\lambda \in \mathbb{C}$ :

$$\operatorname{Re} \lambda < 0 \Rightarrow \text{exponentially decaying solution } y(t) = y_0 \exp(\lambda t),$$

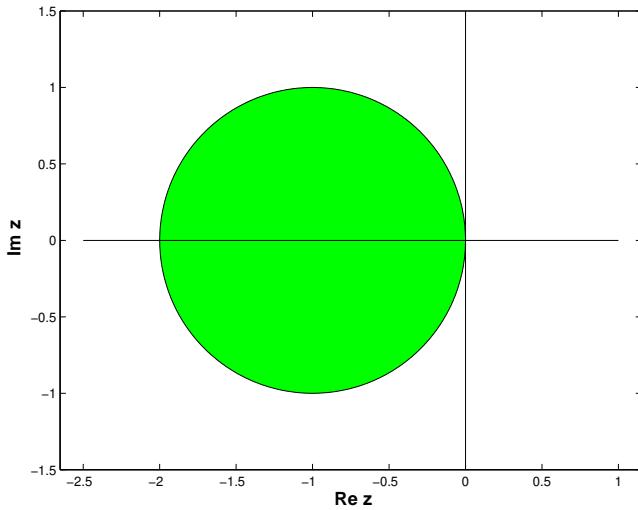
because  $|\exp(\lambda t)| = \exp(\operatorname{Re} \lambda \cdot t)$ .

The **model problem analysis** from Ex. 12.1.0.3, Ex. 12.1.0.9 can be extended verbatim to the case of  $\lambda \in \mathbb{C}$ . It yields the following insight for the explicit Euler method and  $\lambda \in \mathbb{C}$ :

The sequence generated by the explicit Euler method (11.2.1.4) for the model problem (12.1.0.5) satisfies

$$y_{k+1} = y_k(1 + h\lambda) \quad \Rightarrow \quad \lim_{k \rightarrow \infty} y_k = 0 \Leftrightarrow |1 + h\lambda| < 1. \quad (12.1.0.6)$$

↑  
timestep constraint to get decaying (discrete) solution !



$$\triangleleft \{z \in \mathbb{C}: |1 + z| < 1\}$$

The green region of the complex plane marks values for  $\lambda h$ , for which the explicit Euler method will produce exponentially decaying solutions.

Fig. 414

Now we can conjecture what happens in Ex. 12.1.0.35: the eigenvalues  $\lambda_{1/2} = -\frac{1}{2}\alpha \pm i\sqrt{\beta - \frac{1}{4}\alpha^2}$  of  $\mathbf{M}$  have a very large (in modulus) negative real part. Since **ode45** can be expected to behave as if it integrates  $\dot{z} = \lambda_2 z$ , it faces a severe timestep constraint, if exponential blow-up is to be avoided, see Ex. 12.1.0.3. Thus stepsize control must resort to tiny timesteps. □

**§12.1.0.43 (Extended model problem analysis for explicit Runge-Kutta single step methods)** We apply an explicit  $s$ -stage RK-SSM ( $\rightarrow$  Def. 11.4.0.9) described by the Butcher scheme  $\begin{array}{c|cc} c & \alpha \\ \hline b^T & \end{array}$  to the autonomous linear ODE  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$ ,  $\mathbf{M} \in \mathbb{C}^{d,d}$ , and obtain (for the first step with timestep size  $h > 0$ )

$$\mathbf{k}_\ell = \mathbf{M}(\mathbf{y}_0 + h \sum_{j=1}^{s-1} a_{\ell j} \mathbf{k}_j), \quad \ell = 1, \dots, s, \quad \mathbf{y}_1 = \mathbf{y}_0 + h \sum_{\ell=1}^s b_\ell \mathbf{k}_\ell. \quad (12.1.0.44)$$

Now assume that  $\mathbf{M}$  can be diagonalized, that is (12.1.0.32) holds:  $\mathbf{V}^{-1}\mathbf{M}\mathbf{V} = \mathbf{D}$  with a diagonal matrix  $\mathbf{D} \in \mathbb{C}^{d,d}$  containing the eigenvalues  $\lambda_i \in \mathbb{C}$  of  $\mathbf{M}$  on its diagonal. Then apply the substitutions

$$\hat{\mathbf{k}}_\ell := \mathbf{V}^{-1}\mathbf{k}_\ell, \quad \ell = 1, \dots, s, \quad \hat{\mathbf{y}}_k := \mathbf{V}^{-1}\mathbf{y}_k, \quad k = 0, 1,$$

to (12.1.0.44), which yield

$$\hat{\mathbf{k}}_\ell = \mathbf{D}(\hat{\mathbf{y}}_0 + h \sum_{j=1}^{s-1} a_{\ell j} \hat{\mathbf{k}}_j), \quad \ell = 1, \dots, s, \quad \hat{\mathbf{y}}_1 = \hat{\mathbf{y}}_0 + h \sum_{\ell=1}^s b_\ell \hat{\mathbf{k}}_\ell. \quad (12.1.0.45)$$

↑

$$(\hat{\mathbf{k}}_\ell)_i = \lambda_i ((\mathbf{y}_0)_i + h \sum_{j=1}^{s-1} a_{\ell j} (\hat{\mathbf{k}}_j)_i), \quad (\hat{\mathbf{y}}_1)_i = (\hat{\mathbf{y}}_0)_i + h \sum_{\ell=1}^s b_\ell (\hat{\mathbf{k}}_\ell)_i, \quad i = 1, \dots, d. \quad (12.1.0.46)$$

We infer that, if  $(\mathbf{y}_k)_k$  is the sequence produced by an explicit RK-SSM applied to  $\dot{\mathbf{y}} = \mathbf{My}$ , then

$$\mathbf{y}_k = \mathbf{V} \begin{bmatrix} y_k^{[1]} & & 0 \\ & \ddots & \\ 0 & & y_k^{[d]} \end{bmatrix} \mathbf{V}^{-1},$$

where  $(y_k^{[i]})_k$  is the sequence generated by the same RK-SSM with the same sequence of timesteps for the IVP  $\dot{y} = \lambda_i y, y(0) = (\mathbf{V}^{-1} \mathbf{y}_0)_i$ .

The RK-SSM generates uniformly bounded solution sequences  $(\mathbf{y}_k)_{k=0}^\infty$  for  $\dot{\mathbf{y}} = \mathbf{My}$  with diagonalizable matrix  $\mathbf{M} \in \mathbb{R}^{d,d}$  with eigenvalues  $\lambda_1, \dots, \lambda_d$ , if and only if it generates uniformly bounded sequences for all the scalar ODEs  $\dot{z} = \lambda_i z, i = 1, \dots, d$ .

### Stability analysis: reduction to scalar case

Understanding the behavior of RK-SSM for autonomous scalar linear ODEs  $\dot{y} = \lambda y$  with  $\lambda \in \mathbb{C}$  is enough to predict their behavior for general autonomous linear systems of ODEs.

From the considerations of § 12.1.0.26 we deduce the following fundamental result.

#### Theorem 12.1.0.48. (Absolute) stability of explicit RK-SSM for linear systems of ODEs

The sequence  $(\mathbf{y}_k)_k$  of approximations generated by an explicit RK-SSM ( $\rightarrow$  Def. 11.4.0.9) with stability function  $S$  (defined in (12.1.0.18)) applied to the linear autonomous ODE  $\dot{\mathbf{y}} = \mathbf{My}$ ,  $\mathbf{M} \in \mathbb{C}^{d,d}$ , with uniform timestep  $h > 0$  decays exponentially for every initial state  $\mathbf{y}_0 \in \mathbb{C}^d$ , if and only if  $|S(\lambda_i h)| < 1$  for all eigenvalues  $\lambda_i$  of  $\mathbf{M}$ .

Please note that

$$\operatorname{Re} \lambda_i < 0 \quad \forall i \in \{1, \dots, d\} \implies \|\mathbf{y}(t)\| \rightarrow 0 \quad \text{for } t \rightarrow \infty,$$

for any solution of  $\dot{\mathbf{y}} = \mathbf{My}$ . This is obvious from the representation formula (12.1.0.33).  $\square$

**§12.1.0.49 (Region of (absolute) stability of explicit RK-SSM)** We consider an explicit Runge-Kutta single step method with stability function  $S$  for the model linear scalar IVP  $\dot{y} = \lambda y, y(0) = y_0, \lambda \in \mathbb{C}$ . From Thm. 12.1.0.17 we learn that for uniform stepsize  $h > 0$  we have  $y_k = S(\lambda h)^k y_0$  and conclude that

$$y_k \rightarrow 0 \quad \text{for } k \rightarrow \infty \iff |S(\lambda h)| < 1. \quad (12.1.0.50)$$

Hence, the modulus  $|S(\lambda h)|$  tells us for which combinations of  $\lambda$  and stepsize  $h$  we achieve exponential decay  $y_k \rightarrow \infty$  for  $k \rightarrow \infty$ , which is the desirable behavior of the approximations for  $\operatorname{Re} \lambda < 0$ .

#### Definition 12.1.0.51. Region of (absolute) stability

Let the discrete evolution  $\Psi$  for a single step method applied to the scalar linear ODE  $\dot{y} = \lambda y$ ,  $\lambda \in \mathbb{C}$ , be of the form

$$\Psi^h y = S(z)y, \quad y \in \mathbb{C}, h > 0 \quad \text{with} \quad z := h\lambda \quad (12.1.0.52)$$

and a function  $S : \mathbb{C} \rightarrow \mathbb{C}$ . Then the **region of (absolute) stability** of the single step method is given by

$$\mathcal{S}_\Psi := \{z \in \mathbb{C} : |S(z)| < 1\} \subset \mathbb{C}.$$

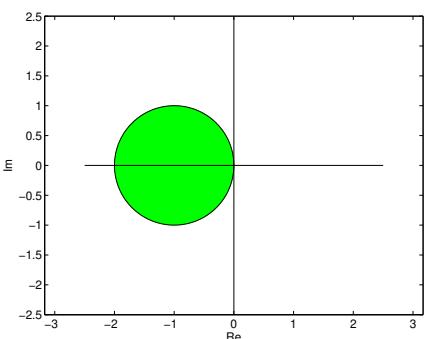
Of course, by Thm. 12.1.0.17, in the case of explicit RK-SSM the function  $S$  will coincide with their **stability function** from (12.1.0.18).

We can easily combine the statement of Thm. 12.1.0.48 with the concept of a region of stability and conclude that an explicit RK-SSM will generate exponentially decaying solutions for the linear ODE  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$ ,  $\mathbf{M} \in \mathbb{C}^{d,d}$ , for every initial state  $\mathbf{y}_0 \in \mathbb{C}^d$ , if and only if  $\lambda_i h \in \mathcal{S}_\Psi$  for all eigenvalues  $\lambda_i$  of  $\mathbf{M}$ .

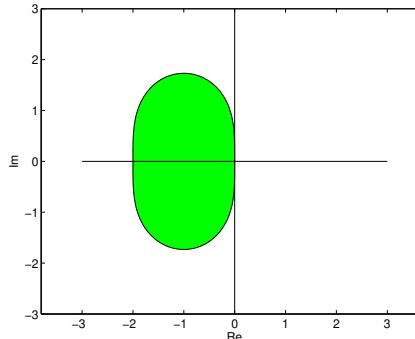
Adopting the arguments of § 12.1.0.26 we conclude from Cor. 12.1.0.20 that

- ◆ the regions of (absolute) stability of explicit RK-SSM are **bounded**,
- ◆ a **timestep constraint** depending on the eigenvalues of  $\mathbf{M}$  is necessary to have a guaranteed exponential decay RK-solutions for  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$ .

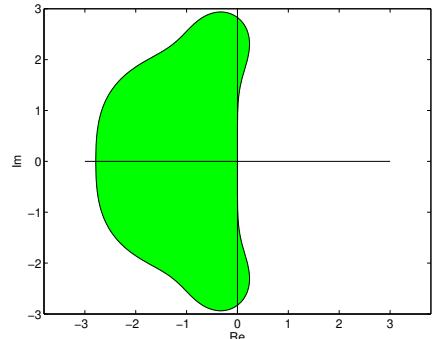
**EXAMPLE 12.1.0.53 (Regions of stability of some explicit RK-SSM)** The green domains  $\subset \mathbb{C}$  depict the bounded regions of stability for some RK-SSM from Ex. 11.4.0.15.



$\mathcal{S}_\Psi$ : explicit Euler (11.2.1.4)



$\mathcal{S}_\Psi$ : explicit trapezoidal method



$\mathcal{S}_\Psi$ : classical RK4 method

In general we have for a consistent RK-SSM ( $\rightarrow$  Def. 11.3.1.11) that their stability functions satisfy  $S(z) = 1 + z + O(z^2)$  for  $z \rightarrow 0$ . Therefore,  $\mathcal{S}_\Psi \neq \emptyset$  and the imaginary axis will be tangent to  $\mathcal{S}_\Psi$  in  $z = 0$ .  $\square$

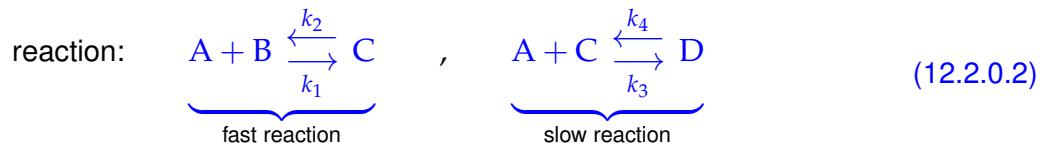
## 12.2 Stiff Initial Value Problems



*Supplementary literature.* [QSS00, Sect. 11.10]

This section will reveal that the behavior observed in Ex. 12.0.0.1 and Ex. 12.1.0.3 is typical for a large class of problems and that the model problem (12.1.0.5) really represents a “generic case”. This justifies the attention paid to linear model problem analysis in Section 12.1.

**EXAMPLE 12.2.0.1 (Kinetics of chemical reactions → [Han02, Ch. 62])** In Ex. 11.5.0.1 we already saw an ODE model for the dynamics of a chemical reaction. Now we study an abstract reaction.



Vastly different reaction constants:

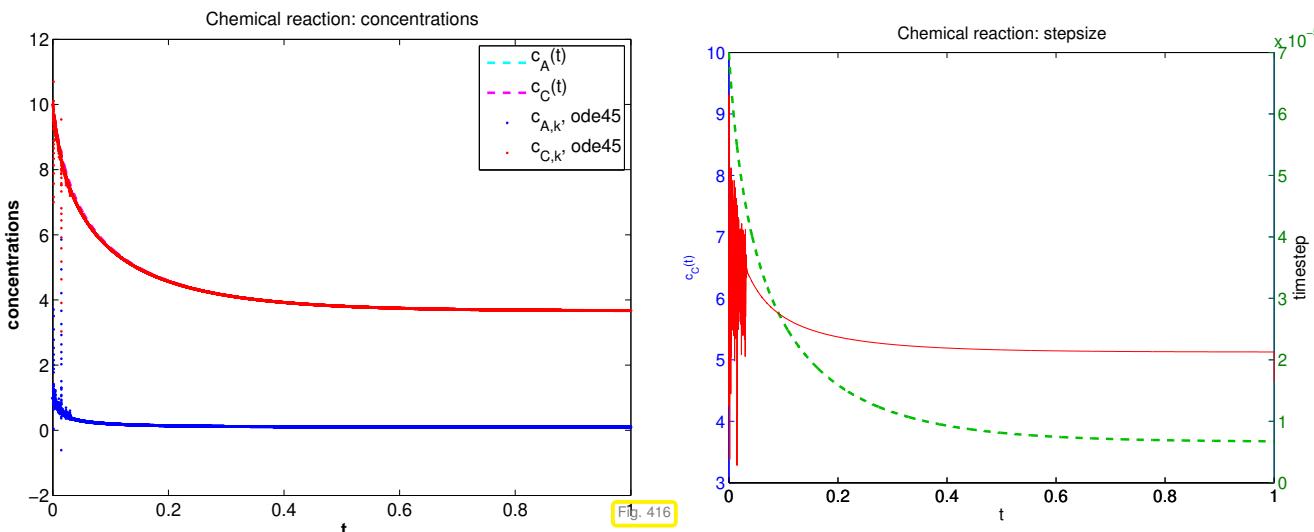
$$k_1, k_2 \gg k_3, k_4$$

► If  $c_A(0) > c_B(0)$  ➤ 2nd reaction determines overall long-term reaction dynamics

Mathematical model: non-linear ODE involving concentrations  $\mathbf{y}(t) = [c_A(t), c_B(t), c_C(t), c_D(t)]^\top$

$$\dot{\mathbf{y}} := \frac{d}{dt} \begin{bmatrix} c_A \\ c_B \\ c_C \\ c_D \end{bmatrix} = \mathbf{f}(\mathbf{y}) := \begin{bmatrix} -k_1 c_A c_B + k_2 c_C - k_3 c_A c_C + k_4 c_D \\ -k_1 c_A c_B + k_2 c_C \\ k_1 c_A c_B - k_2 c_C - k_3 c_A c_C + k_4 c_D \\ k_3 c_A c_C - k_4 c_D \end{bmatrix}. \quad (12.2.0.3)$$

Concrete choice of parameters:  $t_0 = 0$ ,  $T = 1$ ,  $k_1 = 10^4$ ,  $k_2 = 10^3$ ,  $k_3 = 10$ ,  $k_4 = 1$ , initial value  $\mathbf{y}_0 = [1, 1, 10, 0]^\top$ .



Observations: After a **fast initial transient** phase, the solution shows only slow dynamics. Nevertheless, the explicit adaptive integrator used for this simulation insists on using a tiny timestep. It behaves very much like **ode45** in Ex. 12.0.0.1.

**EXAMPLE 12.2.0.4 (Strongly attractive limit cycle)** We consider the non-linear Autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  with

$$\mathbf{f}(\mathbf{y}) := \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \mathbf{y} + \lambda(1 - \|\mathbf{y}\|^2) \mathbf{y}, \quad (12.2.0.5)$$

on the state space  $D = \mathbb{R}^2 \setminus \{0\}$ .

For  $\lambda = 0$ , the initial value problem  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ ,  $\mathbf{y}(0) = \begin{bmatrix} \cos \varphi \\ \sin \varphi \end{bmatrix}$ ,  $\varphi \in \mathbb{R}$  has the solution

$$\mathbf{y}(t) = \begin{bmatrix} \cos(t - \varphi) \\ \sin(t - \varphi) \end{bmatrix}, \quad t \in \mathbb{R}. \quad (12.2.0.6)$$

For this solution we have  $\|\mathbf{y}(t)\|_2 = 1$  for all times.

- (12.2.0.6) provides a solution even for  $\lambda \neq 0$ , if  $\|\mathbf{y}(0)\|_2 = 1$ , because in this case the term  $\lambda(1 - \|\mathbf{y}\|^2) \mathbf{y}$  will never become non-zero on the solution trajectory.

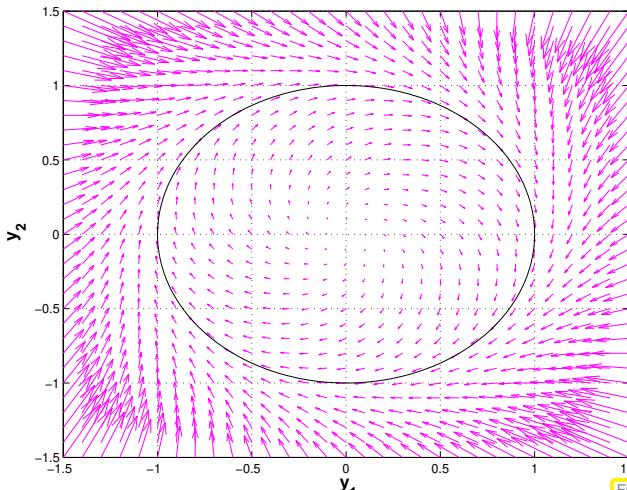
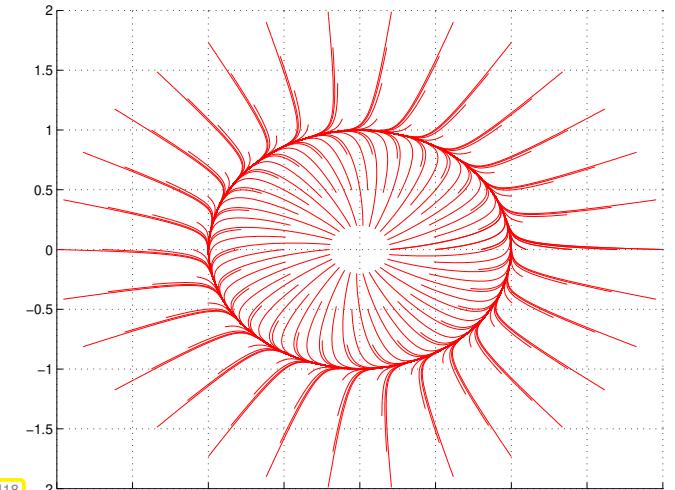


Fig. 417

solution trajectories ( $\lambda = 10$ )

We study the response of **ode45** introduced in § 11.4.0.18 to different choice of  $\lambda$  with initial state  $\mathbf{y}_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ . According to the above considerations this initial state should completely “hide the impact of  $\lambda$  from our view”.

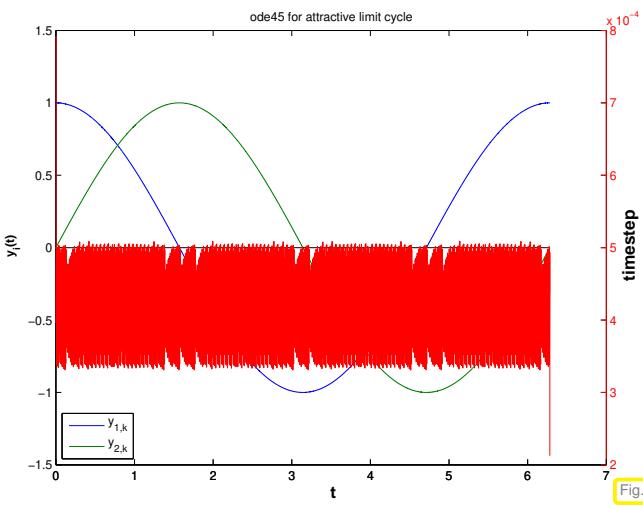
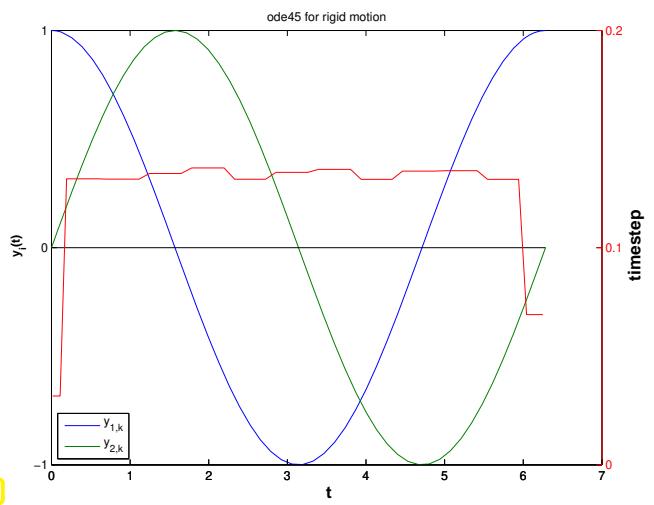


Fig. 419

many (3794) steps ( $\lambda = 1000$ )accurate solution with few steps ( $\lambda = 0$ )

Confusing observation: we have  $\|\mathbf{y}_0\| = 1$ , which implies  $\|\mathbf{y}(t)\| = 1 \quad \forall t!$

Thus, the term of the right hand side, which is multiplied by  $\lambda$  will always vanish on the exact solution trajectory, which stays on the unit circle.

Nevertheless, **ode45** is forced to use tiny timesteps by the *mere presence* of this term! □

We want to find criteria that allow to predict the massive problems haunting explicit single step methods in the case of the *non-linear* IVP of Ex. 12.0.0.1, Ex. 12.2.0.1, and Ex. 12.2.0.4. Recall that for *linear* IVPs of the form  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$ ,  $\mathbf{y}(0) = \mathbf{y}_0$ , the model problem analysis of Section 12.1 tells us that, given knowledge of the region of stability of the timestepping scheme, the eigenvalues of the matrix  $\mathbf{M} \in \mathbb{C}^{d,d}$  provide full information about timestep constraint we are going to face. Refer to Thm. 12.1.0.48 and § 12.1.0.49.

The ODEs we saw in Ex. 12.2.0.1 and Ex. 12.2.0.4 are *non-linear*. Yet, the entire stability analysis of Section 12.1 was based on linear ODEs. Thus, we need to extend the stability analysis to non-linear ODEs.

We start with a “phenomenological notion”, just a keyword to refer to the kind of difficulties presented by the IVPs of Ex. 12.0.0.1, Ex. 12.2.0.1, Ex. 12.1.0.8, and Ex. 12.2.0.4.

#### Notion 12.2.0.7. Stiff IVP

An initial value problem is called **stiff**, if stability imposes much tighter timestep constraints on *explicit single step methods* than the accuracy requirements.

**§12.2.0.8 (Linearization of ODEs)** Linear ODEs, though very special, are highly relevant as “local model” for general ODEs:

We consider a general autonomous ODE

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) \quad , \quad \mathbf{f} : D \subset \mathbb{R}^d \rightarrow \mathbb{R}^d .$$

As usual, we assume  $\mathbf{f}$  to be  $C^2$ -smooth and that it enjoys local Lipschitz continuity ( $\rightarrow$  Def. 11.1.2.11) on  $D$  so that unique solvability of IVPs is guaranteed by Thm. 11.1.2.15.

We fix a state  $\mathbf{y}^* \in D$ ,  $D$  the state space, write  $t \mapsto \mathbf{y}(t)$  for the solution with  $\mathbf{y}(0) = \mathbf{y}^*$ . We set  $\mathbf{z}(t) = \mathbf{y}(t) - \mathbf{y}^*$ , which satisfies

$$\mathbf{z}(0) = 0 \quad , \quad \dot{\mathbf{z}} = \mathbf{f}(\mathbf{y}^* + \mathbf{z}) = \mathbf{f}(\mathbf{y}^*) + D\mathbf{f}(\mathbf{y}^*)\mathbf{z} + R(\mathbf{y}^*, \mathbf{z}) \quad , \quad \text{with} \quad \|R(\mathbf{y}^*, \mathbf{z})\| = O(\|\mathbf{z}\|^2) .$$

This is obtained by Taylor expansion of  $\mathbf{f}$  at  $\mathbf{y}^*$ , see [Str09, Satz 7.5.2]. Hence, in a neighborhood of a state  $\mathbf{y}^*$  on a solution trajectory  $t \mapsto \mathbf{y}(t)$ , the deviation  $\mathbf{z}(t) = \mathbf{y}(t) - \mathbf{y}^*$  satisfies

$$\dot{\mathbf{z}} \approx \mathbf{f}(\mathbf{y}^*) + D\mathbf{f}(\mathbf{y}^*)\mathbf{z} . \tag{12.2.0.9}$$

► The short-time evolution of  $\mathbf{y}$  with  $\mathbf{y}(0) = \mathbf{y}^*$  is approximately governed by the **affine-linear ODE**

$$\dot{\mathbf{y}} = \mathbf{M}(\mathbf{y} - \mathbf{y}^*) + \mathbf{b} \quad , \quad \mathbf{M} := D\mathbf{f}(\mathbf{y}^*) \in \mathbb{R}^{d,d} \quad , \quad \mathbf{b} := \mathbf{f}(\mathbf{y}^*) \in \mathbb{R}^d . \tag{12.2.0.10}$$

In the scalar case we have come across this linearization already in Ex. 12.1.0.3. □

**§12.2.0.11 (Linearization of explicit Runge-Kutta single step methods)** We consider one step a general  $s$ -stage RK-SSM according to Def. 11.4.0.9 for the autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ , with smooth right hand side function  $\mathbf{f} : D \subset \mathbb{R}^d \rightarrow \mathbb{R}^d$ :

$$\mathbf{k}_i = \mathbf{f}(\mathbf{y}_0 + h \sum_{j=1}^{i-1} a_{ij} \mathbf{k}_j) \quad , \quad i = 1, \dots, s \quad , \quad \mathbf{y}_1 = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i .$$

We perform linearization at  $\mathbf{y}^* := \mathbf{y}_0$  and ignore all terms at least quadratic in the timestep size  $h$  (this is indicated by the  $\approx$  symbol):

$$\mathbf{k}_i \approx \mathbf{f}(\mathbf{y}^*) + D\mathbf{f}(\mathbf{y}^*)h \sum_{j=1}^{i-1} a_{ij}\mathbf{k}_j, \quad i = 1, \dots, s, \quad \mathbf{y}_1 = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i.$$

The defining equations for the same RK-SSM applied to

$$\dot{\mathbf{z}} = \mathbf{M}(\mathbf{z}) + \mathbf{b}, \quad \mathbf{M} := D\mathbf{f}(\mathbf{y}^*) \in \mathbb{R}^{d,d}, \quad \mathbf{b} := \mathbf{f}(\mathbf{y}^*),$$

which agrees with (12.2.0.10) after substitution  $\mathbf{z}(t) - \mathbf{y}(t) - \mathbf{y}^*$ , are

$$\mathbf{k}_i \approx \mathbf{b} + \mathbf{M}h \sum_{j=1}^{i-1} a_{ij}\mathbf{k}_j, \quad i = 1, \dots, s, \quad \mathbf{y}_1 = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i.$$

We find that for small timesteps

the discrete evolution of the RK-SSM for  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  in the state  $\mathbf{y}^*$  is close to the discrete evolution of the same RK-SSM applied to the linearization (12.2.0.10) of the ODE in  $\mathbf{y}^*$ .

By straightforward manipulations of the defining equations of an explicit RK-SSM we find that, if

- $(\mathbf{y}_k)_k$  is the sequence of states generated by the RK-SSM applied to the affine-linear ODE  $\dot{\mathbf{y}} = \mathbf{M}(\mathbf{y} - \mathbf{y}_0) + \mathbf{b}$ ,  $\mathbf{M} \in \mathbb{C}^{d,d}$  regular,
- $(\mathbf{w}_k)_k$  is the sequence of states generated by the same RK-SSM applied to the linear ODE  $\dot{\mathbf{w}} = \mathbf{M}\mathbf{w}$  and  $\mathbf{w}_0 := \mathbf{M}^{-1}\mathbf{b}$ , then

$$\mathbf{w}_k = \mathbf{y}_k - \mathbf{y}_0 + \mathbf{M}^{-1}\mathbf{b}.$$

- The analysis of the behavior of an RK-SSM for an affine-linear ODE can be reduced to understanding its behavior for a linear ODE with the same matrix.

Combined with the insights from § 12.1.0.43 this means that

for small timestep the behavior of an explicit RK-SSM applied to  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  close to the state  $\mathbf{y}^*$  is determined by the eigenvalues of the Jacobian  $D\mathbf{f}(\mathbf{y}^*)$ .

In particular, if  $D\mathbf{f}(\mathbf{y}^*)$  has at least one eigenvalue whose modulus is large, then an exponential drift-off of the approximate states  $\mathbf{y}_k$  away from  $\mathbf{y}^*$  can only be avoided for sufficiently small timestep, again a **timestep constraint**.

### How to distinguish stiff initial value problems

An initial value problem for an autonomous ODE  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  will probably be stiff, if, for substantial periods of time,

$$\min\{\operatorname{Re} \lambda : \lambda \in \sigma(D\mathbf{f}(\mathbf{y}(t)))\} \ll 0, \tag{12.2.0.13}$$

$$\text{and} \quad \max\{\operatorname{Re} \lambda : \lambda \in \sigma(D\mathbf{f}(\mathbf{y}(t)))\} \lesssim 0, \tag{12.2.0.14}$$

where  $t \mapsto \mathbf{y}(t)$  is the solution trajectory and  $\sigma(\mathbf{M})$  is the spectrum of the matrix  $\mathbf{M}$ , see

Def. 9.1.0.1.

The condition (12.2.0.14) has to be read as “the real parts of all eigenvalues are below a bound with small modulus”. If this is not the case, then the exact solution will experience blow-up. It will change drastically over very short periods of time and small timesteps will be required anyway in order to resolve this.  $\square$

**EXAMPLE 12.2.0.15 (Predicting stiffness of non-linear IVPs)**

- ① We consider the IVP from Ex. 12.0.0.1:

$$\text{IVP considered: } \dot{y} = f(y) := \lambda y^2(1 - y), \quad \lambda := 500, \quad y(0) = \frac{1}{100}.$$

We find

$$f'(y) = \lambda(2y - 3y^2) \Rightarrow f'(1) = -\lambda.$$

Hence, in case  $\lambda \gg 1$  as in Fig. 401, we face a stiff problem close to the stationary state  $y = 1$ . The observations made in Fig. 401 exactly match this prediction.

- ② The solution of the IVP from Ex. 12.2.0.4

$$\dot{\mathbf{y}} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \mathbf{y} + \lambda(1 - \|\mathbf{y}\|^2) \mathbf{y}, \quad \|\mathbf{y}_0\|_2 = 1. \quad (12.2.0.5)$$

satisfies  $\|\mathbf{y}(t)\|_2 = 1$  for all times. Using the product rule (8.4.1.10) of multi-dimensional differential calculus, we find

$$\begin{aligned} D\mathbf{f}(\mathbf{y}) &= \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} + \lambda \left( -2\mathbf{y}\mathbf{y}^\top + (1 - \|\mathbf{y}\|_2^2) \mathbf{I} \right). \\ \rightarrow \sigma(D\mathbf{f}(\mathbf{y})) &= \left\{ -\lambda - \sqrt{\lambda^2 - 1}, -\lambda + \sqrt{\lambda^2 - 1} \right\}, \text{ if } \|\mathbf{y}\|_2 = 1. \end{aligned}$$

Thus, for  $\lambda \gg 1$ ,  $D\mathbf{f}(\mathbf{y}(t))$  will always have an eigenvalue with large negative real part, whereas the other eigenvalue is close to zero: the IVP is stiff.  $\square$

**Remark 12.2.0.16 (Characteristics of stiff IVPs)** Often one can already tell from the expected behavior of the solution of an IVP, which is often clear from the modeling context, that one has to brace for stiffness.

Typical features of stiff IVPs:

- ◆ Presence of **fast transients** in the solution, see Ex. 12.1.0.3, Ex. 12.1.0.35,
- ◆ Occurrence of **strongly attractive** fixed points/limit cycles, see Ex. 12.2.0.4

## 12.3 Implicit Runge-Kutta Single Step Methods

**Explicit** Runge-Kutta single step method cannot escape tight timestep constraints for stiff IVPs that may render them inefficient, see § 12.1.0.49. In this section we are going to augment the class of Runge-Kutta methods by timestepping schemes that can cope well with stiff IVPs.



*Supplementary literature.* [DR08, Sect. 11.6.2], [QSS00, Sect. 11.8.3]

### 12.3.1 The implicit Euler method for stiff IVPs

**EXAMPLE 12.3.1.1 (Euler methods for stiff decay IVP)** We revisit the setting of Ex. 12.1.0.3 and again consider Euler methods for the decay IVP

$$\dot{y} = \lambda y , \quad y(0) = 1 , \quad \lambda < 0 .$$

We apply both the explicit Euler method (11.2.1.4) and the implicit Euler method (11.2.2.2) with uniform timesteps  $h = 1/N$ ,  $N \in \{5, 10, 20, 40, 80, 160, 320, 640\}$  and monitor the error at final time  $T = 1$  for different values of  $\lambda$ .

Explicit Euler method (11.2.1.4)

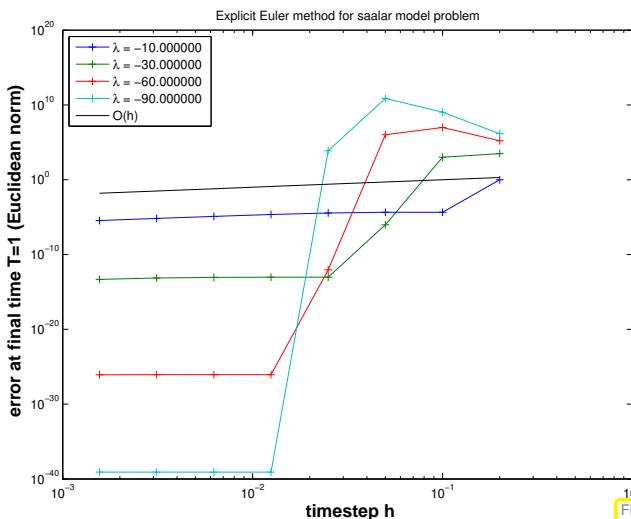
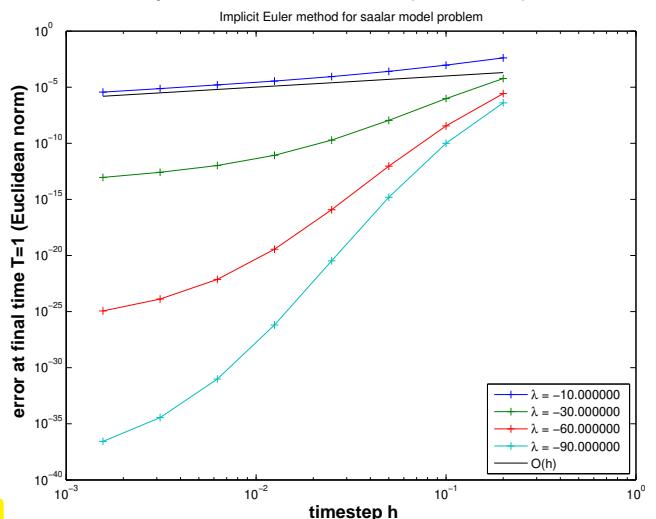


Fig. 421

Implicit Euler method (11.2.2.2)



$\lambda$  large: blow-up of  $y_k$  for large timestep  $h$

$\lambda$  large: stable for all timesteps  $h > 0$  !

We observe onset of convergence of the implicit Euler method already for large timesteps  $h$ . □

#### §12.3.1.2 (Linear model problem analysis: implicit Euler method)

We follow the considerations of

§ 12.1.0.4 and consider the *implicit* Euler method (11.2.2.2) for the

$$\text{linear model problem: } \dot{y} = \lambda y , \quad y(0) = y_0 , \quad \text{with } \operatorname{Re} \lambda \ll 0 , \quad (12.1.0.5)$$

with *exponentially decaying* (maybe oscillatory for  $\operatorname{Im} \lambda \neq 0$ ) exact solution

$$y(t) = y_0 \exp(\lambda t) \rightarrow 0 \quad \text{for } t \rightarrow \infty .$$

The recursion of the implicit Euler method for (12.1.0.5) is defined by

$$(11.2.2.2) \text{ for } f(y) = \lambda y \Rightarrow y_{k+1} = y_k + \lambda h y_{k+1} \quad k \in \mathbb{N}_0 . \quad (12.3.1.3)$$

► generated sequence  $y_k := \left( \frac{1}{1 - \lambda h} \right)^k y_0 . \quad (12.3.1.4)$

$\Rightarrow \boxed{\operatorname{Re} \lambda < 0 \Rightarrow \lim_{k \rightarrow \infty} y_k = 0 \quad \forall h > 0 !} \quad (12.3.1.5)$

**No** timestep constraint: qualitatively correct behavior of  $(y_k)_k$  for  $\operatorname{Re} \lambda < 0$  and **any**  $h > 0$ !

As in § 12.1.0.40 this analysis can be extended to linear systems of ODEs  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$ ,  $\mathbf{M} \in \mathbb{C}^{d,d}$ , by means of **diagonalization**.

As in § 12.1.0.30 and § 12.1.0.40 we assume that  $\mathbf{M}$  can be diagonalized, that is (12.1.0.32) holds:  $\mathbf{V}^{-1}\mathbf{M}\mathbf{V} = \mathbf{D}$  with a diagonal matrix  $\mathbf{D} \in \mathbb{C}^{d,d}$  containing the eigenvalues of  $\mathbf{M}$  on its diagonal. Next, apply the **decoupling by diagonalization** idea to the recursion of the implicit Euler method.

$$\mathbf{V}^{-1}\mathbf{y}_{k+1} = \mathbf{V}^{-1}\mathbf{y}_k + h \underbrace{\mathbf{V}^{-1}\mathbf{M}\mathbf{V}}_{=\mathbf{D}} (\mathbf{V}^{-1}\mathbf{y}_{k+1}) \stackrel{\mathbf{z}_k := \mathbf{V}^{-1}\mathbf{y}_k}{\Leftrightarrow} \underbrace{(\mathbf{z}_{k+1})_i = \frac{1}{1 - \lambda_i h} (\mathbf{z}_k)_i}_{\triangleq \text{implicit Euler step for } \dot{z}_i = \lambda_i z_i} . \quad (12.3.1.6)$$

Crucial insight:

For any timestep, the implicit Euler method generates exponentially decaying solution sequences  $(\mathbf{y}_k)_{k=0}^{\infty}$  for  $\dot{\mathbf{y}} = \mathbf{M}\mathbf{y}$  with diagonalizable matrix  $\mathbf{M} \in \mathbb{R}^{d,d}$  with eigenvalues  $\lambda_1, \dots, \lambda_d$ , if  $\operatorname{Re} \lambda_i < 0$  for all  $i = 1, \dots, d$ .

Thus we expect that the implicit Euler method will not face stability induced timestep constraints for stiff problems ( $\rightarrow$  Notion 12.2.0.7).  $\square$

## 12.3.2 Collocation single step methods

Unfortunately the implicit Euler method is of first order only, see Exp. 11.3.2.5. This section presents an algorithm for designing higher order single step methods generalizing the implicit Euler method.

Setting: We consider the general ordinary differential equation  $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$ ,  $\mathbf{f} : I \times D \rightarrow \mathbb{R}^d$  locally Lipschitz continuous, which guarantees the local existence of unique solutions of initial value problems, see Thm. 11.1.2.15.

We define the single step method through specifying the first step  $\mathbf{y}_0 = \mathbf{y}(t_0) \rightarrow \mathbf{y}_1 \approx \mathbf{y}(t_1)$ , where  $\mathbf{y}_0 \in D$  is the initial step at initial time  $t_0 \in I$ . We assume that the exact solution trajectory  $t \mapsto \mathbf{y}(t)$  exists on  $[t_0, t_1]$ . Use as a timestepping scheme on a temporal mesh ( $\rightarrow$  § 11.2.0.2) in the sense of Def. 11.3.1.5 is straightforward.

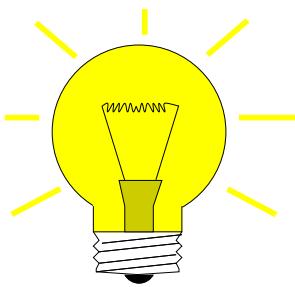
### §12.3.2.1 (Collocation principle)

#### Abstract collocation idea

Collocation is a paradigm for the **discretization** ( $\rightarrow$  Rem. 11.3.1.4) of **differential equations**:

- (I) Write the discrete solution  $u_h$ , a function, as linear combination of  $N \in \mathbb{N}$  sufficiently smooth (basis) functions  $\geq N$  unknown coefficients.
- (II) Demand that  $u_h$  satisfies the differential equation at  $N$  points/times  $\geq N$  equations.

We apply this policy to the differential equation  $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$  on  $[t_0, t_1]$ :



- Idea:
- ➊ Approximate  $t \mapsto \mathbf{y}(t)$ ,  $t \in [t_0, t_1]$ , by a function  $t \mapsto \mathbf{y}_h(t) \in V$ ,  $V$  an  $d \cdot (s+1)$ -dimensional trial space  $V$  of functions  $[t_0, t_1] \mapsto \mathbb{R}^d \rightarrow$  Item (I).
  - ➋ Fix  $\mathbf{y}_h \in V$  by imposing collocation conditions  

$$\mathbf{y}_h(t_0) = \mathbf{y}_0 \quad , \quad \dot{\mathbf{y}}_h(\tau_j) = \mathbf{f}(\tau_j, \mathbf{y}_h(\tau_j)) \quad , \quad j = 1, \dots, s \quad , \quad (12.3.2.3)$$
for collocation points  $t_0 \leq \tau_1 < \dots < \tau_s \leq t_1 \rightarrow$  Item (II).
  - ➌ Choose  $\mathbf{y}_1 := \mathbf{y}_h(t_1)$ .

Our choice (the “standard option”):

(Componentwise) polynomial trial space  $V = (\mathcal{P}_s)^d$

Recalling  $\dim \mathcal{P}_s = s+1$  from Thm. 5.2.1.2 we see that our choice makes the number  $N := d(s+1)$  of collocation conditions match the dimension of the trial space  $V$ .

Now we want to derive a concrete representation for the polynomial  $\mathbf{y}_h$ . We draw on concepts introduced in Section 5.2.2. We define the collocation points as

$$\tau_j := t_0 + c_j h \quad , \quad j = 1, \dots, s \quad , \quad \text{for } 0 \leq c_1 < c_2 < \dots < c_s \leq 1 \quad , \quad h := t_1 - t_0 \quad .$$

Let  $\{L_j\}_{j=1}^s \subset \mathcal{P}_{s-1}$  denote the set of Lagrange polynomials of degree  $s-1$  associated with the node set  $\{c_j\}_{j=1}^s$ , see (5.2.2.4). They satisfy  $L_j(c_i) = \delta_{ij}$ ,  $i, j = 1, \dots, s$  and form a basis of  $\mathcal{P}_{s-1}$ .

In each of its  $d$  components, the derivative  $\dot{\mathbf{y}}_h$  is a polynomial of degree  $s-1$ :  $\dot{\mathbf{y}} \in (\mathcal{P}_{s-1})^d$ . Hence, it has the following representation, compare (5.2.2.6).

$$\dot{\mathbf{y}}_h(t_0 + \tau h) = \sum_{j=1}^s \dot{\mathbf{y}}_h(t_0 + c_j h) L_j(\tau) \quad . \quad (12.3.2.4)$$

As  $\tau_j = t_0 + c_j h$ , the comcollocation conditions make it possible to replace  $\dot{\mathbf{y}}_h(c_j h)$  with an expression in the right hand side function  $\mathbf{f}$ :

$$(12.3.2.3) \quad \dot{\mathbf{y}}_h(t_0 + \tau h) = \sum_{j=1}^s \mathbf{k}_j L_j(\tau) \quad \text{with “coefficients”} \quad \mathbf{k}_j := \mathbf{f}(t_0 + c_j h, \mathbf{y}_h(t_0 + c_j h)) \quad .$$

Next we integrate and use  $\mathbf{y}_h(t_0) = \mathbf{y}_0$

$$\dot{\mathbf{y}}_h(t_0 + \tau h) = \mathbf{y}_0 + h \sum_{j=1}^s \mathbf{k}_j \int_0^\tau L_j(\zeta) d\zeta \quad .$$

This yields the following formulas for the computation of  $\mathbf{y}_1$ , which characterize the  $s$ -stage collocation single step method induced by the (normalized) collocation points  $c_j \in [0, 1]$ ,  $j = 1, \dots, s$ .

$$\mathbf{k}_i = \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^s a_{ij} \mathbf{k}_j) \quad , \quad \text{where} \quad a_{ij} := \int_0^{c_i} L_j(\tau) d\tau \quad , \quad (12.3.2.5)$$

$$\mathbf{y}_1 := \mathbf{y}_h(t_1) = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i \quad , \quad b_i := \int_0^1 L_i(\tau) d\tau \quad .$$

Note that, since arbitrary  $\mathbf{y}_0 \in D$ ,  $t_0, t_1 \in I$  were admitted, this defines a discrete evolution  $\Psi : I \times I \times D \rightarrow \mathbb{R}^d$  by  $\Psi^{t_0, t_1} \mathbf{y}_0 := \mathbf{y}_h(t_1)$ .  $\square$

**Remark 12.3.2.6 (Implicit nature of collocation single step methods)** Note that (12.3.2.5) represents a generically non-linear system of  $s \cdot d$  equations for the  $s \cdot d$  components of the vectors  $\mathbf{k}_i$ ,  $i = 1, \dots, s$ . Usually, it will not be possible to obtain  $\mathbf{k}_i$  by a fixed number of evaluations of  $\mathbf{f}$ . For this reason the single step methods defined by (12.3.2.5) are called **implicit**.

With similar arguments as in Rem. 11.2.2.3 one can prove that for sufficiently small  $|t_1 - t_0|$  a unique solution for  $\mathbf{k}_1, \dots, \mathbf{k}_s$  can be found.  $\square$

**§12.3.2.7 (Collocation single step methods and quadrature)** Clearly, in the case  $d = 1$ ,  $f(t, \mathbf{y}) = f(t)$ ,  $\mathbf{y}_0 = 0$  the computation of  $\mathbf{y}_1$  boils down to the evaluation of a quadrature formula on  $[t_0, t_1]$ , because from (12.3.2.5) we get

$$\mathbf{y}_1 = h \sum_{i=1}^s b_i f(t_0 + c_i h), \quad b_i := \int_0^1 L_i(\tau) d\tau, \quad (12.3.2.8)$$

which is a polynomial quadrature formula (7.2.0.2) on  $[0, 1]$  with nodes  $c_j$  transformed to  $[t_0, t_1]$  according to (7.1.0.5).  $\square$

**EXPERIMENT 12.3.2.9 (Empiric Convergence of collocation single step methods)** We consider the scalar logistic ODE (11.1.1.2) with parameter  $\lambda = 10$  ( $\rightarrow$  only mildly stiff), initial state  $y_0 = 0.01$ ,  $T = 1$ .

Numerical integration by timestepping with uniform timestep  $h$  based on collocation single step method (12.3.2.5).

①

Equidistant collocation points,  $c_j = \frac{j}{s+1}$ ,  $j = 1, \dots, s$ .

We observe **algebraic convergence** with the empiric rates

- $s = 1 : p = 1.96$
- $s = 2 : p = 2.03$
- $s = 3 : p = 4.00$
- $s = 4 : p = 4.04$

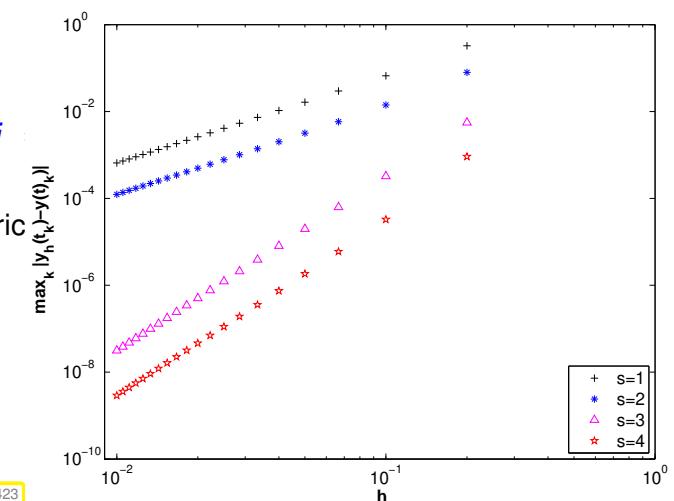


Fig. 423

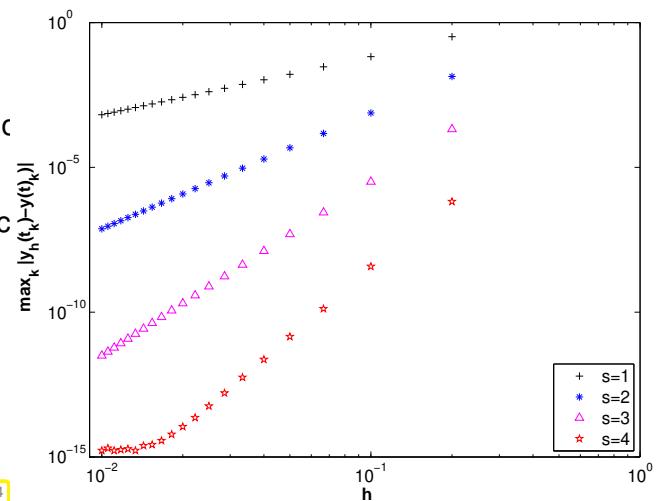
In this case we conclude the following (empiric) order ( $\rightarrow$  Def. 11.3.2.8) of the collocation single step method:

$$(\text{empiric}) \text{ order} = \begin{cases} s & \text{for even } s, \\ s+1 & \text{for odd } s. \end{cases}$$

- ① Gauss points in  $[0, 1]$  as normalized collocative points  $c_j, j = 1, \dots, s$ .

We observe **algebraic convergence** with the empirical rates

$$\begin{aligned}s = 1 &: p = 1.96 \\ s = 2 &: p = 4.01 \\ s = 3 &: p = 6.00 \\ s = 4 &: p = 8.02\end{aligned}$$



Obviously, for the (empirical) order ( $\rightarrow$  Def. 11.3.2.8) of the **Gauss collocation single step method** holds

$$(\text{empirical order}) = 2s.$$

Note that the 1-stage Gauss collocation single step method is the implicit midpoint method from Section 11.2.3.  $\square$

**§12.3.2.10 (Order of collocation single step method)** What we have observed in Exp. 12.3.2.9 reflects a fundamental result on collocation single step methods as defined in (12.3.2.5).

**Theorem 12.3.2.11. Order of collocation single step method [DB02, Satz .6.40]**

Provided that  $\mathbf{f} \in C^p(I \times D)$ , the order ( $\rightarrow$  Def. 11.3.2.8) of an  $s$ -stage collocation single step method according to (12.3.2.5) agrees with the order ( $\rightarrow$  Def. 7.3.0.1) of the quadrature formula on  $[0, 1]$  with nodes  $c_j$  and weights  $b_j, j = 1, \dots, s$ .

- By Thm. 7.3.0.23 the  $s$ -stage **Gauss collocation single step method** whose nodes  $c_j$  are chosen as the  $s$  Gauss points on  $[0, 1]$  is of order  $2s$ .  $\square$

### 12.3.3 General implicit RK-SSMs

The notations in (12.3.2.5) have deliberately been chosen to allude to Def. 11.4.0.9. In that definition it takes only letting the sum in the formula for the increments run up to  $s$  to capture (12.3.2.5).

**Definition 12.3.3.1. General Runge-Kutta single step method (cf. Def. 11.4.0.9)**

For  $b_i, a_{ij} \in \mathbb{R}, c_i := \sum_{j=1}^s a_{ij}, i, j = 1, \dots, s, s \in \mathbb{N}$ , an  $s$ -stage Runge-Kutta single step method (RK-SSM) for the IVP (11.1.2.2) is defined by

$$\mathbf{k}_i := \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^s a_{ij} \mathbf{k}_j), \quad i = 1, \dots, s, \quad \mathbf{y}_1 := \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i.$$

As before, the  $\mathbf{k}_i \in \mathbb{R}^d$  are called **increments**.

Note: computation of increments  $\mathbf{k}_i$  may now require the solution of (*non-linear*) systems of equations of size  $s \cdot d$  ( $\rightarrow$  “implicit” method, cf. Rem. 12.3.2.6)

### General Butcher scheme notation for RK-SSM

Shorthand notation for Runge-Kutta methods

**Butcher scheme**

Note: now  $\mathfrak{A}$  can be a general  $s \times s$ -matrix.

$$\frac{\mathbf{c} \mid \mathfrak{A}}{\mathbf{b}^T} := \begin{array}{c|ccc} c_1 & a_{11} & \cdots & a_{1s} \\ \vdots & \vdots & & \vdots \\ c_s & a_{s1} & \cdots & a_{ss} \\ \hline b_1 & b_1 & \cdots & b_s \end{array}. \quad (12.3.3.3)$$

Summary: terminology for Runge-Kutta single step methods:

- strict lower triangular matrix
- explicit Runge-Kutta method, Def. 11.4.0.9
- lower triangular matrix
- diagonally-implicit Runge-Kutta method (DIRK)

Many of the techniques and much of the theory discussed for explicit RK-SSMs carry over to general (implicit) Runge-Kutta single step methods:

- Sufficient condition for consistence from Cor. 11.4.0.13
- Algebraic convergence for meshwidth  $h \rightarrow 0$  and the related concept of order ( $\rightarrow$  Def. 11.3.2.8)
- Embedded methods and algorithms for adaptive stepsize control from Section 11.5

**Remark 12.3.3.4 (Stage form equations for increments)** In Def. 12.3.3.1 instead of the increments we can consider as unknowns the so-called **stages**

$$\mathbf{g}_i := h \sum_{j=1}^s a_{ij} \mathbf{k}_j, \quad i = 1, \dots, s, \quad \Leftrightarrow \quad \mathbf{k}_i = \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + \mathbf{g}_i). \quad (12.3.3.5)$$

This leads to the equivalent defining equations in “stage form” for an implicit RK-SSM

$$\begin{aligned} \mathbf{k}_i &:= \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^s a_{ij} \mathbf{k}_j), \quad i = 1, \dots, s, \\ &\Downarrow \\ \mathbf{g}_i &= h \sum_{j=1}^s a_{ij} \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + \mathbf{g}_j), \quad \mathbf{y}_1 = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + \mathbf{g}_i). \end{aligned} \quad (12.3.3.6)$$

In terms of implementation there is no difference: Also the **stage equations** (12.3.3.6) are usually solved by means of Newton’s method, see next remark.  $\dashv$

**Remark 12.3.3.7 (Solving the increment equations for implicit RK-SSMs)** We reformulate the increment equations in stage form (12.3.3.6) as a non-linear system of equations in standard form  $F(\mathbf{x}) = 0$ . Unknowns are the total  $s \cdot d$  components of the stage vectors  $\mathbf{g}_i, i = 1, \dots, s$  as defined in (12.3.3.5).

$$\begin{aligned} \mathbf{g} &= [\mathbf{g}_1, \dots, \mathbf{g}_s]^\top \in \mathbb{R}^{s \cdot d}, \\ \mathbf{g}_i &:= h \sum_{j=1}^s a_{ij} \mathbf{f}(t_0 + c_j h, \mathbf{y}_0 + \mathbf{g}_j) \end{aligned} \quad \Rightarrow \quad F(\mathbf{g}) = \mathbf{g} - h(\mathfrak{A} \otimes \mathbf{I}_d) \begin{bmatrix} \mathbf{f}(t_0 + c_1 h, \mathbf{y}_0 + \mathbf{g}_1) \\ \vdots \\ \mathbf{f}(t_0 + c_s h, \mathbf{y}_0 + \mathbf{g}_s) \end{bmatrix} \stackrel{!}{=} \mathbf{0},$$

where  $\mathbf{I}_d$  is the  $d \times d$  identity matrix and  $\otimes$  designates the Kronecker product introduced in Def. 1.4.3.7.

We compute an approximate solution of  $\mathbf{F}(\mathbf{g}) = \mathbf{0}$  iteratively by means of the simplified Newton method presented in Rem. 8.4.1.39. This is a Newton method with “frozen Jacobian”. As  $\mathbf{g} \rightarrow \mathbf{0}$  for  $\mathbf{h} \rightarrow \mathbf{0}$ , we choose zero as initial guess:

$$\mathbf{g}^{(k+1)} = \mathbf{g}^{(k)} - \mathbf{D}\mathbf{F}(\mathbf{0})^{-1}\mathbf{F}(\mathbf{g}^{(k)}) \quad k = 0, 1, 2, \dots, \quad \mathbf{g}^{(0)} = \mathbf{0}. \quad (12.3.3.8)$$

with the Jacobian

$$\mathbf{D}\mathbf{F}(\mathbf{0}) = \begin{bmatrix} \mathbf{I} - ha_{11} \frac{\partial \mathbf{f}}{\partial \mathbf{y}}(t_0, \mathbf{y}_0) & \cdots & -ha_{1s} \frac{\partial \mathbf{f}}{\partial \mathbf{y}}(t_0, \mathbf{y}_0) \\ \vdots & \ddots & \vdots \\ -ha_{s1} \frac{\partial \mathbf{f}}{\partial \mathbf{y}}(t_0, \mathbf{y}_0) & \cdots & \mathbf{I} - ha_{ss} \frac{\partial \mathbf{f}}{\partial \mathbf{y}}(t_0, \mathbf{y}_0) \end{bmatrix} \in \mathbb{R}^{sd, sd}. \quad (12.3.3.9)$$

Obviously,  $\mathbf{D}\mathbf{F}(\mathbf{0}) \rightarrow \mathbf{I}$  for  $\mathbf{h} \rightarrow \mathbf{0}$ . Thus,  $\mathbf{D}\mathbf{F}(\mathbf{0})$  will be regular for sufficiently small  $\mathbf{h}$ .

In each step of the simplified Newton method we have to solve a linear system of equations with coefficient matrix  $\mathbf{D}\mathbf{F}(\mathbf{0})$ . If  $s \cdot d$  is large, an efficient implementation has to reuse the LU-decomposition of  $\mathbf{D}\mathbf{F}(\mathbf{0})$ , see Code 8.4.1.40 and Rem. 2.5.0.10.  $\square$

## 12.3.4 Model Problem Analysis for Implicit Runge-Kutta Single-Step Method (IRK-SSMs)

**Model problem analysis** for general Runge-Kutta single step methods ( $\rightarrow$  Def. 12.3.3.1) runs parallel to that for explicit RK-methods as elaborated in Section 12.1, § 12.1.0.13. Familiarity with the techniques and results of this section is assumed. The reader is asked to recall the concept of **stability function** from Thm. 12.1.0.17, the **diagonalization technique** from § 12.1.0.43, and the definition of **region of (absolute) stability** from Def. 12.1.0.51.

We apply the implicit RK-SSM according to Def. 12.3.3.1 to the autonomous linear scalar ODE  $\dot{\mathbf{y}} = \lambda \mathbf{y}$ ,  $\lambda \in \mathbb{C}$ , and utterly parallel to the considerations in § 12.1.0.13 we obtain

$$\begin{array}{l} \blacktriangleright \quad k_i = \lambda(y_0 + h \sum_{j=1}^{i-1} a_{ij} k_j), \\ \quad \quad \quad \Rightarrow \quad \begin{bmatrix} \mathbf{I} - z\mathfrak{A} & \mathbf{0} \\ -z\mathbf{b}^T & 1 \end{bmatrix} \begin{bmatrix} \mathbf{k} \\ \mathbf{y}_1 \end{bmatrix} = y_0 \begin{bmatrix} \mathbf{1} \\ 1 \end{bmatrix}, \\ \quad \quad \quad \mathbf{y}_1 = y_0 + h \sum_{i=1}^s b_i k_i \end{array} \quad (12.1.0.14)$$

where  $\mathbf{k} \in \mathbb{R}^s \hat{=} \text{denotes the vector } [k_1, \dots, k_s]^T / \lambda$  of increments, and  $z := \lambda h$ . As in § 12.1.0.13 we can eliminate the increments and obtain an expression for  $\mathbf{y}_1$ :

$$\mathbf{y}_1 = S(z)\mathbf{y}_0 \quad \text{with} \quad S(z) := 1 + z\mathbf{b}^T(\mathbf{I} - z\mathfrak{A})^{-1}\mathbf{1}. \quad (12.1.0.15)$$

### Theorem 12.3.4.1. Stability function of Runge-Kutta methods, cf. Thm. 12.1.0.17

The discrete evolution  $\Psi_\lambda^h$  of an  $s$ -stage Runge-Kutta single step method ( $\rightarrow$  Def. 12.3.3.1) with Butcher scheme  $\begin{array}{c|cc} \mathbf{c} & \mathfrak{A} \\ \hline & \mathbf{b}^T \end{array}$  (see (12.3.3.3)) for the ODE  $\dot{\mathbf{y}} = \lambda \mathbf{y}$  is given by a multiplication with

$$S(z) := \underbrace{1 + z\mathbf{b}^T(\mathbf{I} - z\mathfrak{A})^{-1}\mathbf{1}}_{\text{stability function}} = \frac{\det(\mathbf{I} - z\mathfrak{A} + z\mathbf{1}\mathbf{b}^T)}{\det(\mathbf{I} - z\mathfrak{A})}, \quad z := \lambda h, \quad \mathbf{1} = [1, \dots, 1]^T \in \mathbb{R}^s.$$

**EXAMPLE 12.3.4.2 (Regions of stability for simple implicit RK-SSM)** We determine the Butcher schemes (12.3.3.3) for simple implicit RK-SSM and apply the formula from Thm. 12.3.4.1 to compute their stability functions.

- Implicit Euler method:

$$\begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array} \Rightarrow S(z) = \frac{1}{1-z}.$$

- Implicit midpoint method:

$$\begin{array}{c|cc} \frac{1}{2} & \frac{1}{2} \\ \hline & 1 \end{array} \Rightarrow S(z) = \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z}.$$

Their regions of stability  $\mathcal{S}_\Psi$  as defined in Def. 12.1.0.51,

$$\mathcal{S}_\Psi := \{z \in \mathbb{C}: |S(z)| < 1\} \subset \mathbb{C},$$

can easily found from the respective stability functions:

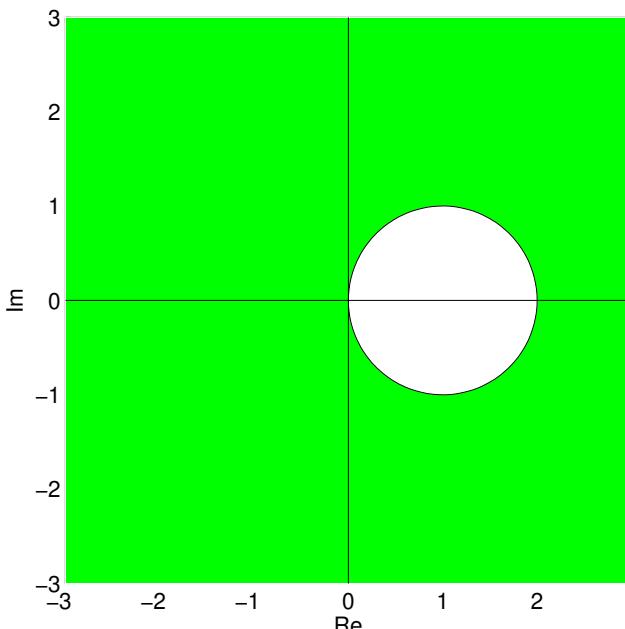
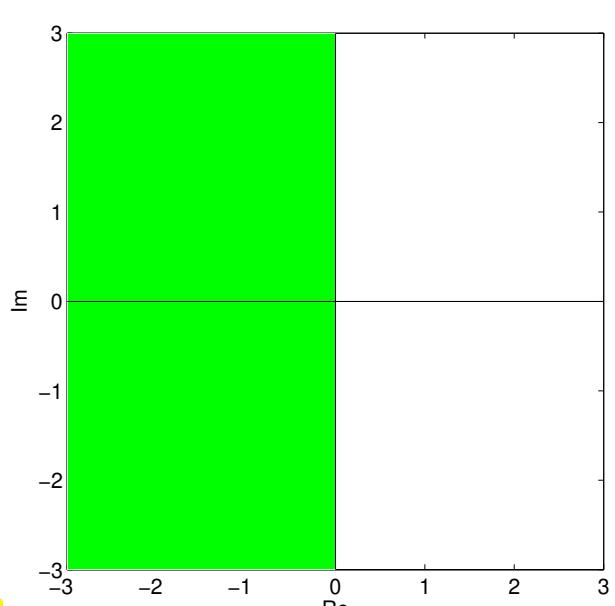


Fig. 425

$\mathcal{S}_\Psi$ : implicit Euler method (11.2.2.2)



$\mathcal{S}_\Psi$ : implicit midpoint method (11.2.3.3)

We see that in both cases  $|S(z)| < 1$ , if  $\text{Re } z < 0$ . □

From the determinant formula for the stability function  $S(z)$  we can conclude a generalization of Cor. 12.1.0.20.

#### Corollary 12.3.4.3. Rational stability function of explicit RK-SSM

For a consistent ( $\rightarrow$  Def. 11.3.1.11)  $s$ -stage general Runge-Kutta single step method according to Def. 12.3.3.1 the stability function  $S$  is a non-constant rational function of the form  $S(z) = \frac{P(z)}{Q(z)}$  with polynomials  $P \in \mathcal{P}_s$ ,  $Q \in \mathcal{P}_s$ .

Of course, a rational function  $z \mapsto S(z)$  can satisfy  $\lim_{|z| \rightarrow \infty} |S(z)| < 1$  as we have seen in Ex. 12.3.4.2. As a consequence, the region of stability for implicit RK-SSM need not be bounded.

**§12.3.4.4 (A-stability)** A general RK-SSM with stability function  $S$  applied to the scalar linear IVP  $\dot{y} = \lambda y$ ,  $y(0) = y_0 \in \mathbb{C}$ ,  $\lambda \in \mathbb{C}$ , with uniform timestep  $h > 0$  will yield the sequence  $(y_k)_{k=0}^\infty$  defined by

$$y_k = S(z)^k y_0 \quad , \quad z = \lambda h . \quad (12.3.4.5)$$

Hence, the next property of a RK-SSM guarantees that the sequence of approximations decays exponentially whenever the exact solution of the model problem IVP (12.1.0.5) does so.

#### Definition 12.3.4.6. A-stability of a Runge-Kutta single step method

A Runge-Kutta single step method with stability function  $S$  is **A-stable**, if

$$\mathbb{C}^- := \{z \in \mathbb{C} : \operatorname{Re} z < 0\} \subset \mathcal{S}_\Psi . \quad (\mathcal{S}_\Psi \hat{=} \text{region of stability Def. 12.1.0.51})$$

From Ex. 12.3.4.2 we conclude that both the implicit Euler method and the implicit midpoint method are A-stable.

A-stable Runge-Kutta single step methods will not be affected by stability induced timestep constraints when applied to **stiff** IVP ( $\rightarrow$  Notion 12.2.0.7).

**§12.3.4.7 (“Ideal” region of stability)** In order to reproduce the qualitative behavior of the exact solution, a single step method when applied to the scalar linear IVP  $\dot{y} = \lambda y$ ,  $y(0) = y_0 \in \mathbb{C}$ ,  $\lambda \in \mathbb{C}$ , with uniform timestep  $h > 0$ ,

- should yield an *exponentially decaying* sequence  $(y_k)_{k=0}^\infty$ , whenever  $\operatorname{Re} \lambda < 0$ ,
- should produce an *exponentially increasing* sequence  $(y_k)_{k=0}^\infty$ , whenever  $\operatorname{Re} \lambda > 0$ .

Thus, in light of (12.3.4.5), we agree that the stability if

$$\text{“ideal” region of stability is } \boxed{\mathcal{S}_\Psi = \mathbb{C}^-} . \quad (12.3.4.8)$$

Are there RK-SSMs that can boast of an ideal region of stability?

Regions of stability of Gauss collocation single step methods, see Exp. 12.3.2.9:

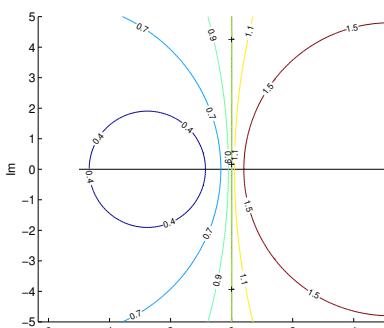


Fig. 427

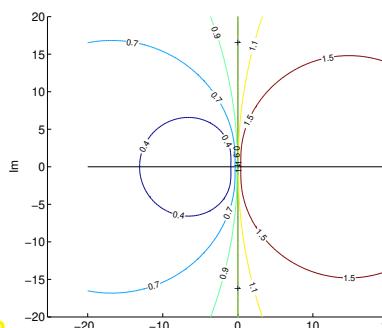


Fig. 428

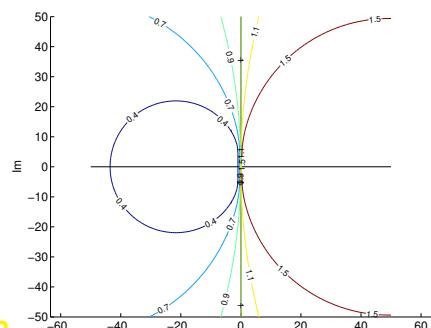


Fig. 429

Implicit midpoint method

 $s = 2$  (order 4) $s = 4$  (order 8)

Level lines for  $|S(z)|$  for Gauss collocation methods

**Theorem 12.3.4.9. Region of stability of Gauss collocation single step methods [DB02, Satz 6.44]**

$s$ -stage **Gauss collocation single step methods** defined by (12.3.2.5) with the nodes  $c_s$  given by the  $s$  Gauss points on  $[0, 1]$ , feature the “ideal” stability domain:

$$\mathcal{S}_\Psi = \mathbb{C}^- . \quad (12.3.4.8)$$

In particular, all Gauss collocation single step methods are A-stable.

**EXPERIMENT 12.3.4.10 (Implicit RK-SSMs for stiff IVP)** We consider the stiff IVP

$$\dot{y} = -\lambda y + \beta \sin(2\pi t), \quad \lambda = 10^6, \beta = 10^6, \quad y(0) = 1,$$

whose solution essentially is the smooth function  $t \mapsto \sin(2\pi t)$ . Applying the criteria (12.2.0.13) and (12.2.0.14) we immediately see that this IVP is extremely stiff.

We solve it with different implicit RK-SSM on  $[0, 1]$  with large uniform timestep  $h = \frac{1}{20}$ .

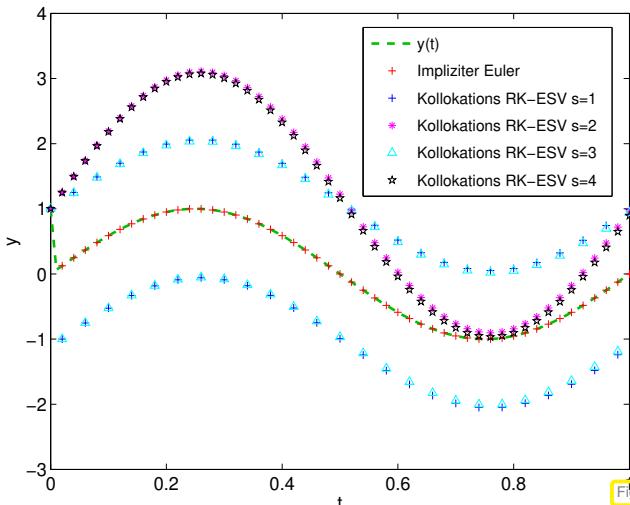


Fig. 430

Solutions by RK-SSMs

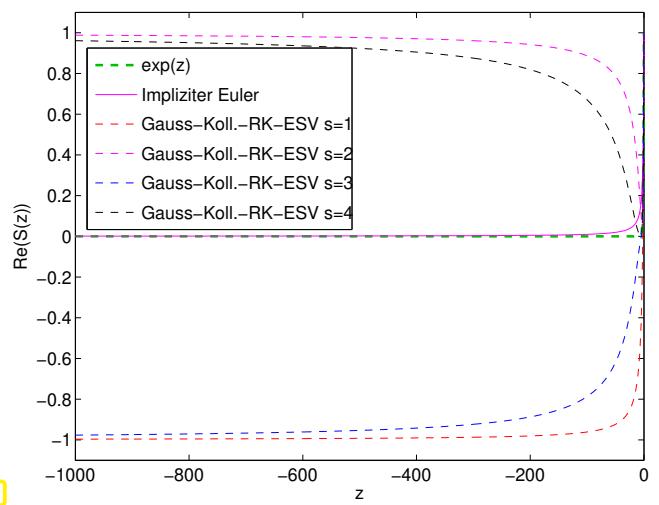


Fig. 431

Stability functions on  $\mathbb{C}^-$ 

We observe that Gauss collocation RK-SSMs incur a huge discretization error, whereas the simple implicit Euler method provides a perfect approximation!

Explanation: The stability functions for Gauss collocation RK-SSMs satisfy

$$\lim_{|z| \rightarrow \infty} |S(z)| = 1 .$$

Hence, when they are applied to  $\dot{y} = \lambda y$  with extremely large (in modulus)  $\lambda < 0$ , they will produce sequences that decay only very slowly or even oscillate, which misses the very rapid decay of the exact solution. The stability function for the implicit Euler method is  $S(z) = (1 - z)^{-1}$  and satisfies  $\lim_{|z| \rightarrow \infty} S(z) = 0$ , which will mean a fast exponential decay of the  $y_k$ . □

**§12.3.4.11 (L-stability)** In light of what we learned in the previous experiment we can now state what we expect from the stability function of a Runge-Kutta method that is suitable for stiff IVP ( $\rightarrow$  Notion 12.2.0.7>):

**Definition 12.3.4.12. L-stable Runge-Kutta method** → [Han02, Ch. 77]

A Runge-Kutta method (→ Def. 12.3.3.1) is **L-stable/asymptotically stable**, if its stability function (→ Thm. 12.3.4.1) satisfies

$$(i) \quad \operatorname{Re} z < 0 \Rightarrow |S(z)| < 1 , \quad (12.3.4.13)$$

$$(ii) \quad \lim_{\operatorname{Re} z \rightarrow -\infty} S(z) = 0 . \quad (12.3.4.14)$$

Remember: L-stable : $\Leftrightarrow$  A-stable & “ $S(-\infty) = 0$ ”

**Remark 12.3.4.15 (Necessary condition for L-stability of Runge-Kutta methods)**

Consider a Runge-Kutta single step method (→ Def. 12.3.3.1) described by the Butcher scheme  $\begin{array}{c|cc} c & \mathfrak{A} \\ \hline b^T & \end{array}$ .

Assume that  $\mathfrak{A} \in \mathbb{R}^{s,s}$  is regular, which can be fulfilled only for an implicit RK-SSM.

For a rational function  $S(z) = \frac{P(z)}{Q(z)}$  the limit for  $|z| \rightarrow \infty$  exists and can easily be expressed by the leading coefficients of the polynomials  $P$  and  $Q$ :

$$\text{Thm. 12.3.4.1} \Rightarrow S(-\infty) = 1 - \mathbf{b}^T \mathfrak{A}^{-1} \mathbf{1} . \quad (12.3.4.16)$$

► If  $\mathbf{b}^T = (\mathfrak{A})_{:,j}^T$  (row of  $\mathfrak{A}$ )  $\Rightarrow S(-\infty) = 0$ . (12.3.4.17)

Butcher scheme (12.3.3.3) for L-stable RK-methods, see Def. 12.3.4.12

$$\triangleright \begin{array}{c|cc} c & \mathfrak{A} \\ \hline b^T & \end{array} := \begin{array}{c|cccc} c_1 & a_{11} & \cdots & a_{1s} \\ \vdots & \vdots & & \vdots \\ c_{s-1} & a_{s-1,1} & \cdots & a_{s-1,s} \\ \hline 1 & b_1 & \cdots & b_s \\ \hline & b_1 & \cdots & b_s \end{array} .$$

A closer look at the coefficient formulas of (12.3.2.5) reveals that the algebraic condition (12.3.4.17) will automatically satisfied for a collocation single step method with  $c_s = 1$ !

**EXAMPLE 12.3.4.18 (L-stable implicit Runge-Kutta methods)** There is a family of  $s$ -point quadrature formulas on  $[0, 1]$  with a node located in 1 and (maximal) order  $2s - 1$ : **Gauss-Radau formulas**. They induce the **L-stable** Gauss-Radau collocation single step methods of order  $2s - 1$  according to Thm. 12.3.2.11.

$$\begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array}$$

Implicit Euler method

$$\begin{array}{c|ccc} & \frac{1}{3} & \frac{5}{12} & -\frac{1}{12} \\ \hline 1 & \frac{3}{4} & \frac{1}{4} & \\ \hline & \frac{3}{4} & \frac{1}{4} & \end{array}$$

Radau RK-SSM, order 3

$$\begin{array}{c|cccc} & \frac{4-\sqrt{6}}{10} & \frac{88-7\sqrt{6}}{360} & \frac{296-169\sqrt{6}}{1800} & \frac{-2+3\sqrt{6}}{225} \\ \hline & \frac{4+\sqrt{6}}{10} & \frac{296+169\sqrt{6}}{1800} & \frac{88+7\sqrt{6}}{360} & \frac{-2-3\sqrt{6}}{225} \\ & 1 & \frac{16-\sqrt{6}}{36} & \frac{16+\sqrt{6}}{36} & \frac{1}{9} \\ \hline & \frac{16-\sqrt{6}}{36} & \frac{16+\sqrt{6}}{36} & & \frac{1}{9} \end{array}$$

Radau RK-SSM, order 5

The stability functions of  $s$ -stage Gauss-Radau collocation SSMs are rational functions of the form

$$S(z) = \frac{P(z)}{Q(z)}, \quad P \in \mathcal{P}_{s-1}, Q \in \mathcal{P}_s.$$

Beware that also " $S(\infty) = 0$ ", which means that Gauss-Radau methods when applied to problems with fast exponential blow-up may produce a spurious decaying solution.

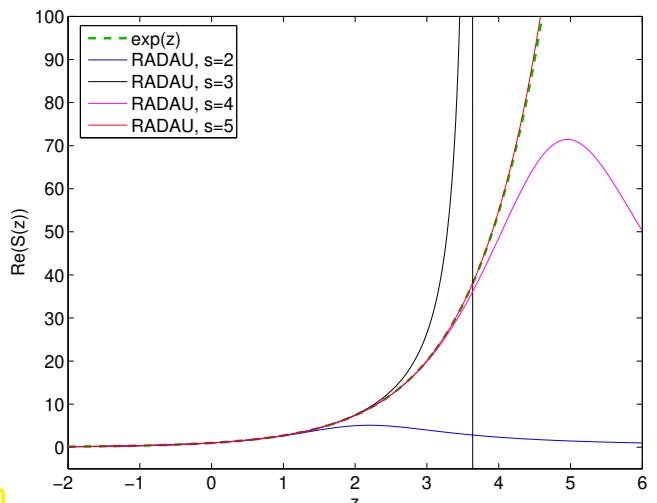
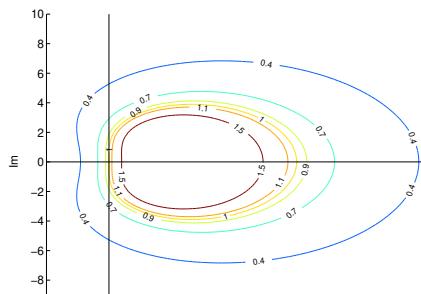


Fig. 432

Level lines of stability functions of  $s$ -stage Gauss-Radau collocation SSMs:



The fourth-order Gauss method is already so accurate that potential overshoots when approaching  $y = 1$  are damped fast enough.  $\square$

## 12.4 Semi-implicit Runge-Kutta Methods



*Supplementary literature.* [Han02, Ch. 80]



The equations fixing the increments  $\mathbf{k}_i \in \mathbb{R}^d$ ,  $i = 1, \dots, s$ , for an  $s$ -stage implicit RK-method constitute a (Non-)linear system of equations with  $s \cdot d$  unknowns.



Expensive iterations needed to find  $\mathbf{k}_i$ ?

Remember that we compute approximate solutions anyway, and the increments *are weighted with the stepsize  $h \ll 1$* , see Def. 12.3.3.1. So there is no point in determining them with high accuracy!



Idea: Use only a fixed *small* number of Newton steps to solve for the  $\mathbf{k}_i$ ,  $i = 1, \dots, s$ .

Extreme case: use only a single Newton step! Let's try.

### EXAMPLE 12.4.0.1 (Semi-implicit Euler single-step method)

- We consider an Initial value problem for logistic ODE, see Ex. 11.1.1.1

$$\dot{y} = \lambda y(1 - y) , \quad y(0) = 0.1 , \quad \lambda = 5 .$$

- We use the implicit Euler method (11.2.2.2) with uniform timestep  $h = 1/n$ ,  $n \in \{5, 8, 11, 17, 25, 38, 57, 85, 128, 192, 288, 432, 649, 973, 1460, 2189, 3284, 4926, 7389\}$ .

& approximate computation of  $y_{k+1}$  by  
1 Newton step with initial guess  $y_k$

= semi-implicit Euler method

- Measured error  $\text{err} = \max_{j=1, \dots, n} |y_j - y(t_j)|$

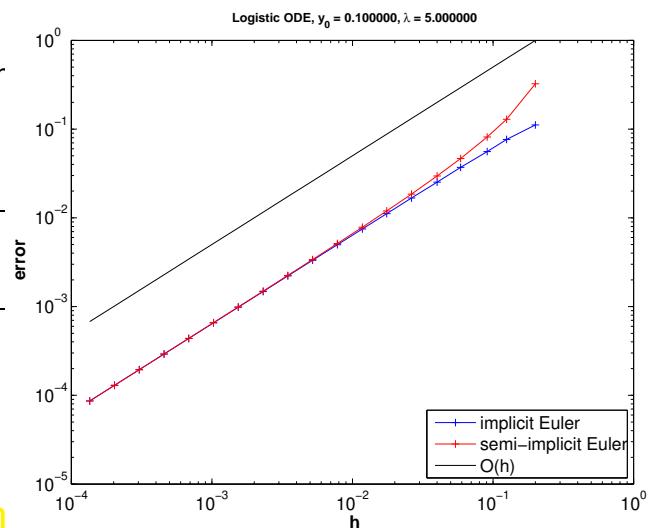


Fig. 436

From (11.2.2.2) with timestep  $h > 0$

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h\mathbf{f}(\mathbf{y}_{k+1}) \Leftrightarrow F(\mathbf{y}_{k+1}) := \mathbf{y}_{k+1} - h\mathbf{f}(\mathbf{y}_{k+1}) - \mathbf{y}_k = 0 .$$

One Newton step (8.4.1.1) applied to  $F(\mathbf{y}) = 0$  with initial guess  $\mathbf{y}_k$  yields

$$\mathbf{y}_{k+1} = \mathbf{y}_k - D\mathbf{f}(\mathbf{y}_k)^{-1}F(\mathbf{y}_k) = \mathbf{y}_k + (\mathbf{I} - hD\mathbf{f}(\mathbf{y}_k))^{-1}h\mathbf{f}(\mathbf{y}_k) .$$

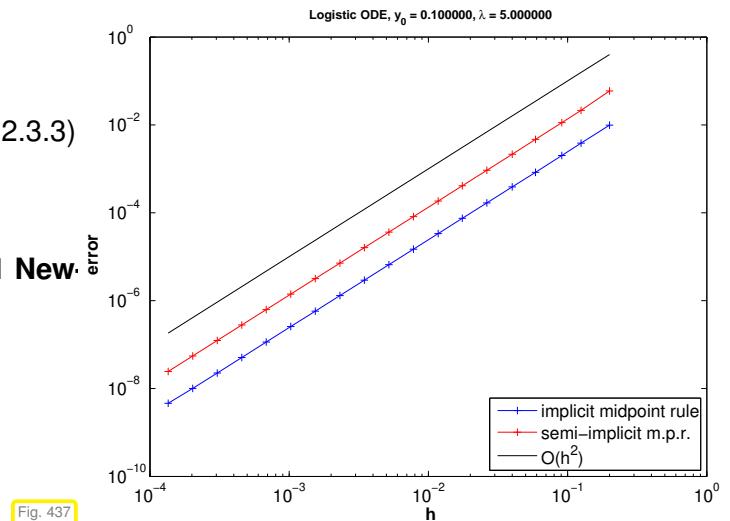
Note: for linear ODE with  $\mathbf{f}(\mathbf{y}) = \mathbf{A}\mathbf{y}$ ,  $\mathbf{A} \in \mathbb{R}^{d,d}$ , we recover the original implicit Euler method!

Observation: Approximate evaluation of defining equation for  $y_{k+1}$  preserves 1st order convergence.

- Now, implicit midpoint method (11.2.3.3) uniform timestep  $h = 1/n$  as above

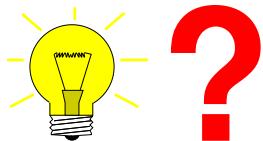
& approximate computation of  $y_{k+1}$  by **1 New-ton step**, initial guess  $y_k$

- Fehlermass  $\text{err} = \max_{j=1,\dots,n} |y_j - y(t_j)|$



We still observe second-order convergence! □

Try: Use **linearized increment equations** for implicit RK-SSM



$$\mathbf{k}_i := \mathbf{f}(\mathbf{y}_0 + h \sum_{j=1}^s a_{ij} \mathbf{k}_j), \quad i = 1, \dots, s$$

$$\mathbf{k}_i = \mathbf{f}(\mathbf{y}_0) + h D\mathbf{f}(\mathbf{y}_0) \left( \sum_{j=1}^s a_{ij} \mathbf{k}_j \right), \quad i = 1, \dots, s. \quad (12.4.0.2)$$

The good news is that all results about stability derived from model problem analysis (→ Section 12.1) remain valid despite linearization of the increment equations:

Linearization does nothing for linear ODEs ▷ stability function (→ Thm. 12.3.4.1) not affected!

The bad news is that the preservation of the order observed in Ex. 12.4.0.1 will no longer hold in the general case.

#### EXPERIMENT 12.4.0.3 (Convergence of naive semi-implicit Radau method)

- We consider an IVP for the logistic ODE from Ex. 11.1.1.1:

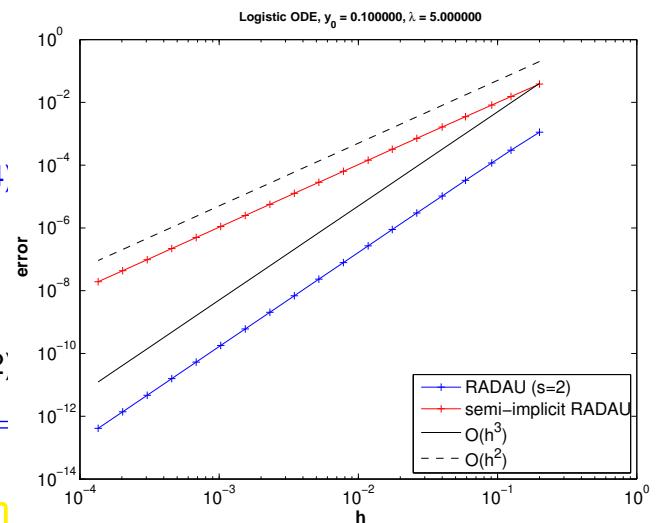
$$\dot{y} = \lambda y(1 - y), \quad y(0) = 0.1, \quad \lambda = 5.$$

- ◆ 2-stage Radau RK-SSM, Butcher scheme

$$\begin{array}{c|cc} \frac{1}{3} & \frac{5}{12} & -\frac{1}{12} \\ 1 & \frac{3}{4} & \frac{1}{4} \\ \hline & \frac{3}{4} & \frac{1}{4} \end{array}, \quad (12.4.0.4)$$

order = 3, see Ex. 12.3.4.18.

- ◆ Increments from linearized equations (12.4.0.2)
- ◆ We monitor the error through  $\text{err} = \max_{j=1,\dots,n} |y_j - y(t_j)|$



Loss of order due to linearization !

**§12.4.0.5 (Rosenbrock-Wanner methods)** We have just seen that the simple linearization according to (12.4.0.2) will degrade the order of implicit RK-SSMs and leads to a substantial loss of accuracy. This is not an option.

Yet, the idea behind (12.4.0.2) has been refined. One does not start from a known RK-SSM, but introduces general coefficients for structurally linear increment equations.

► Class of  $s$ -stage semi-implicit (linearly implicit) Runge-Kutta methods (Rosenbrock-Wanner (ROW) methods):

$$\begin{aligned} (\mathbf{I} - ha_{ii}\mathbf{J})\mathbf{k}_i &= \mathbf{f}(\mathbf{y}_0 + h \sum_{j=1}^{i-1} (a_{ij} + d_{ij})\mathbf{k}_j) - h\mathbf{J} \sum_{j=1}^{i-1} d_{ij}\mathbf{k}_j, \quad \mathbf{J} = \mathbf{D}\mathbf{f}(\mathbf{y}_0), \\ \mathbf{y}_1 &:= \mathbf{y}_0 + \sum_{j=1}^s b_j \mathbf{k}_j. \end{aligned} \quad (12.4.0.6)$$

Then the coefficients  $a_{ij}$ ,  $d_{ij}$ , and  $b_i$  are determined from order conditions by solving large non-linear systems of equations.

In each step  $s$  linear systems with coefficient matrices  $\mathbf{I} - ha_{ii}\mathbf{J}$  have to be solved. For methods used in practice one often demands that  $a_{ii} = \gamma$  for all  $i = 1, \dots, s$ . As a consequence, we have to solve  $s$  linear systems with the same coefficient matrix  $\mathbf{I} - h\gamma\mathbf{J} \in \mathbb{R}^{d,d}$ , which permits us to reuse LU-factorizations, see Rem. 2.5.0.10. □

## 12.5 Splitting methods

**§12.5.0.1 (Splitting idea: composition of partial evolutions)** Many relevant ordinary differential equations feature a right hand side function that is the sum to two (or more) terms. Consider an autonomous IVP with a right hand side function that can be split in an additive fashion:

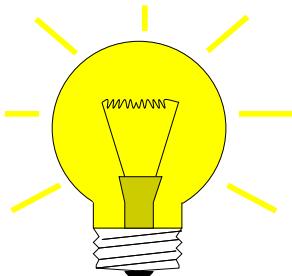
$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) + \mathbf{g}(\mathbf{y}), \quad \mathbf{y}(0) = \mathbf{y}_0, \quad (12.5.0.2)$$

with  $\mathbf{f} : D \subset \mathbb{R}^d \mapsto \mathbb{R}^d$ ,  $\mathbf{g} : D \subset \mathbb{R}^d \mapsto \mathbb{R}^d$  “sufficiently smooth”, locally Lipschitz continuous ( $\rightarrow$  Def. 11.1.2.11).

Let us introduce the evolution operators ( $\rightarrow$  Def. 11.1.3.2) for both summands:

$$\begin{array}{ll} \text{(Continuous) evolution maps:} & \Phi_f^t \leftrightarrow \text{ODE } \dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}), \\ & \Phi_g^t \leftrightarrow \text{ODE } \dot{\mathbf{y}} = \mathbf{g}(\mathbf{y}). \end{array}$$

Temporarily we assume that both  $\Phi_f^t$ ,  $\Phi_g^t$  are available in the form of analytic formulas or highly accurate approximations.

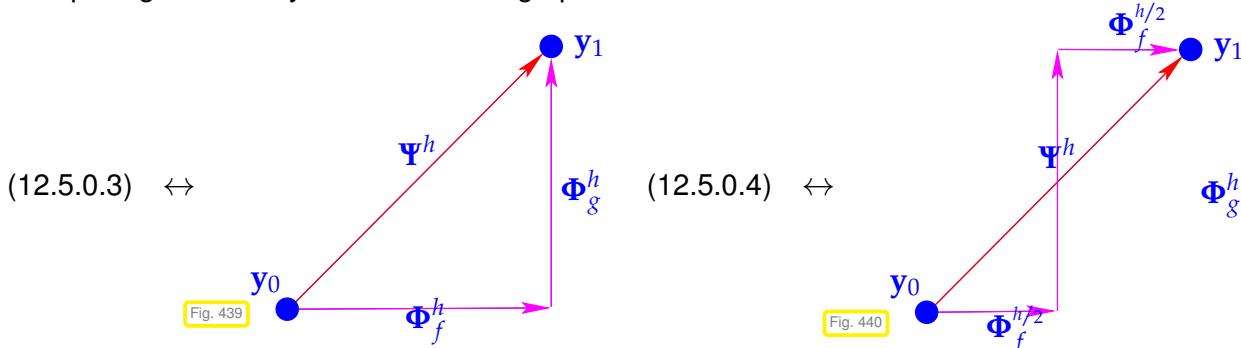


Idea: Build single step methods ( $\rightarrow$  Def. 11.3.1.5) based on the following discrete evolutions

$$\text{Lie-Trotter splitting: } \Psi^h = \Phi_g^h \circ \Phi_f^h, \quad (12.5.0.3)$$

$$\text{Strang splitting: } \Psi^h = \Phi_f^{h/2} \circ \Phi_g^h \circ \Phi_f^{h/2}. \quad (12.5.0.4)$$

These splittings are easily remembered in graphical form:



Note that over many timesteps the Strang splitting approach is not more expensive than Lie-Trotter splitting, because the actual implementation of (12.5.0.4) should be done as follows:

$$\begin{aligned} \mathbf{y}_{1/2} &:= \Phi_f^{h/2}, & \mathbf{y}_1 &:= \Phi_g^h \mathbf{y}_{1/2}, \\ \mathbf{y}_{3/2} &:= \Phi_f^h \mathbf{y}_1, & \mathbf{y}_2 &:= \Phi_g^h \mathbf{y}_{3/2}, \\ \mathbf{y}_{5/2} &:= \Phi_f^h \mathbf{y}_2, & \mathbf{y}_3 &:= \Phi_g^h \mathbf{y}_{5/2}, \\ &\vdots & &\vdots, \end{aligned}$$

because  $\Phi_f^{h/2} \circ \Phi_f^{h/2} = \Phi_f^h$ . This means that a Strang splitting SSM differs from a Lie-Trotter splitting SSM in the first and the last step only.  $\dashv$

**EXAMPLE 12.5.0.5 (Convergence of simple splitting methods)** We consider the following IVP whose right hand side function is the sum of two functions for which the ODEs can be solved analytically:

$$\dot{y} = \underbrace{\lambda y(1-y)}_{=:f(y)} + \underbrace{\sqrt{1-y^2}}_{=:g(y)}, \quad y(0) = 0.$$

►  $\Phi_f^t y = \frac{1}{1 + (y^{-1} - 1)e^{-\lambda t}}, \quad t > 0, y \in [0, 1] \quad (\text{logistic ODE (11.1.1.2)})$

►  $\Phi_g^t y = \begin{cases} \sin(t + \arcsin(y)) & , \text{if } t + \arcsin(y) < \frac{\pi}{2}, \\ 1 & , \text{else,} \end{cases} \quad t > 0, y \in [0, 1].$

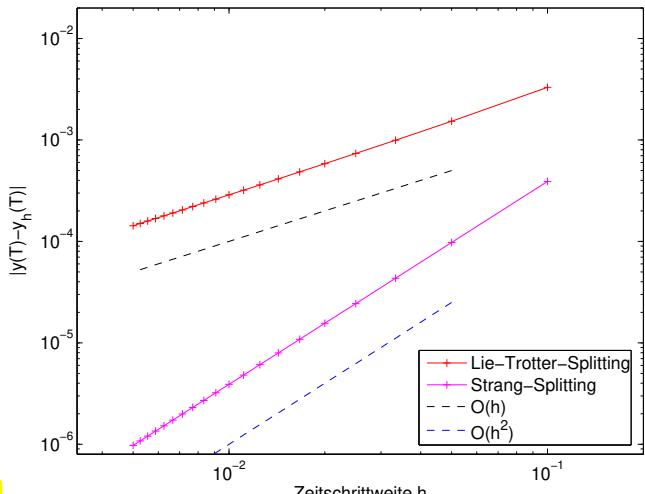


Fig. 441

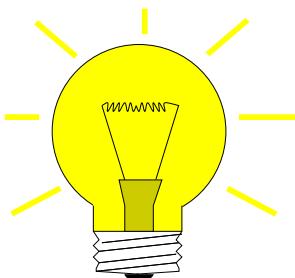
We observe algebraic convergence of the two splitting methods, order 1 for (12.5.0.3), oder 2 for (12.5.0.4). □

The observation made in Ex. 12.5.0.5 reflects a general truth:

#### Theorem 12.5.0.6. Order of simple splitting methods

Die single step methods defined by (12.5.0.3) or (12.5.0.4) are of order ( $\rightarrow$  Def. 11.3.2.8) 1 and 2, respectively.

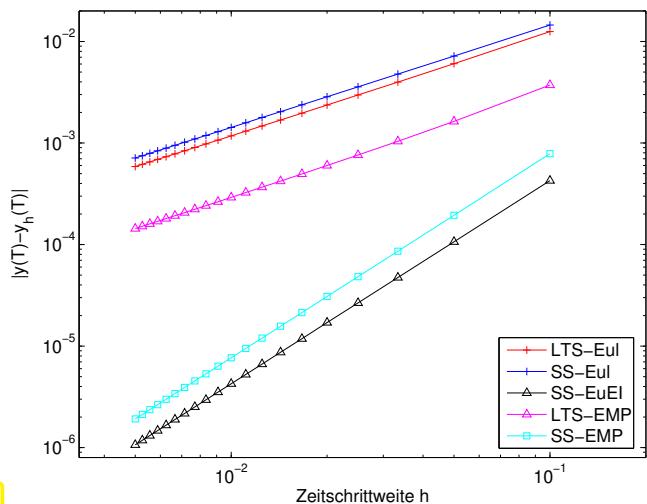
**§12.5.0.7 (Inexact splitting methods)** Of course, the assumption that  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$  and  $\dot{\mathbf{y}} = \mathbf{g}(\mathbf{y})$  can be solved exactly will hardly ever be met. However, it should be clear that a “sufficiently accurate” approximation of the evolution maps  $\Phi_g^h$  and  $\Phi_f^h$  is all we need



Idea: In (12.5.0.3)/(12.5.0.4) replace

exact evolutions  $\longrightarrow$  discrete evolutions  
 $\Phi_g^h, \Phi_f^h \longrightarrow \Psi_g^h, \Psi_f^h$

**EXAMPLE 12.5.0.8 (Convergence of inexact simple splitting methods)** Again we consider the IVP of Ex. 12.5.0.5 and inexact splitting methods based on different single step methods for the two ODE corresponding to the summands.



|         |                                                                                                                                 |
|---------|---------------------------------------------------------------------------------------------------------------------------------|
| LTS-Eul | explicit Euler method (11.2.1.4) → $\Psi_{h,g}^h$ ,<br>$\Psi_{h,f}^h$ + Lie-Trotter splitting (12.5.0.3)                        |
| SS-Eul  | explicit Euler method (11.2.1.4) → $\Psi_{h,g}^h$ ,<br>$\Psi_{h,f}^h$ + Strang splitting (12.5.0.4)                             |
| SS-EuEl | Strang splitting (12.5.0.4): Explizites Euler method (11.2.1.4) ○ exact evolution $\Phi_g^h$ ○ implicit Euler method (11.2.2.2) |
| LTS-EMP | explicit midpoint method (11.2.3.3) → $\Psi_{h,g}^h$ , $\Psi_{h,f}^h$ + Lie-Trotter splitting (12.5.0.3)                        |
| SS-EMP  | explicit midpoint method (11.4.0.7) → $\Psi_{h,g}^h$ , $\Psi_{h,f}^h$ + Strang splitting (12.5.0.4)                             |

- The order of splitting methods may be (but need not be) limited by the order of the SSMs used for  $\Phi_f^h$ ,  $\Phi_g^h$ .

**§12.5.0.9 (Application of splitting methods)** In the following situation the use splitting methods seems advisable:

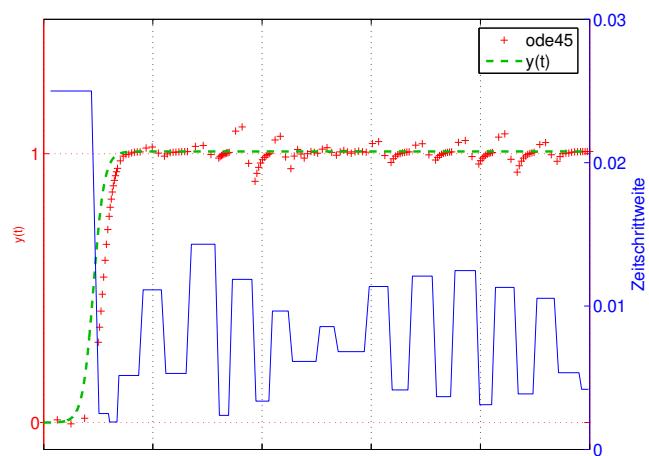
### “Splittable” ODEs

$\dot{\mathbf{y}} = f(\mathbf{y}) + g(\mathbf{y})$  “difficult” :  $\dot{\mathbf{y}} = f(\mathbf{y}) \rightarrow$  stiff, but with an analytic solution  
(e.g., stiff → Section 12.2)       $\dot{\mathbf{y}} = g(\mathbf{y})$  “easy”, amenable to explicit integration.

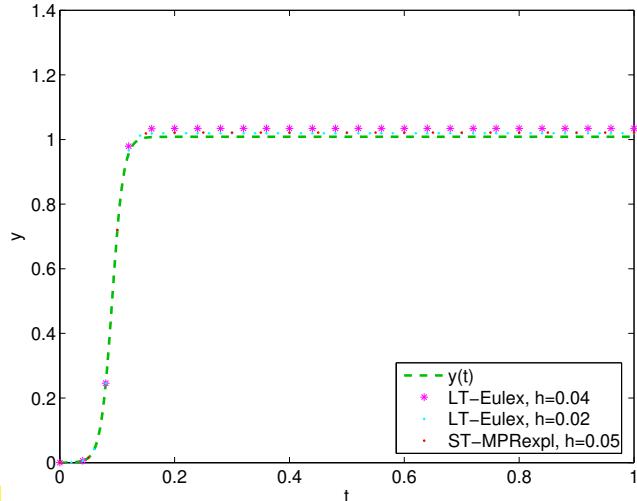
**EXPERIMENT 12.5.0.11 (Splitting off stiff components)** Recall Ex. 12.0.0.1 and the IVP studied there:

$$\text{AWP } \dot{y} = \lambda y(1-y) + \alpha \sin(y), \quad \lambda = 100, \quad \alpha = 1, \quad y(0) = 10^{-4}.$$

small perturbation



Solution by **ode45**, see Ex. 12.0.0.1



Total number of timesteps

|                          |     |
|--------------------------|-----|
| ode45:                   | 152 |
| LT-Eulex, $h = 0.04$ :   | 25  |
| LT-Eulex, $h = 0.02$ :   | 50  |
| ST-MPRexpl, $h = 0.05$ : | 20  |

Details of the methods:

- LT-Eulex:  $\dot{y} = \lambda y(1 - y) \rightarrow$  exact evolution,  $\dot{y} = \alpha \sin y \rightarrow$  expl. Euler (11.2.1.4) & Lie-Trotter splitting (12.5.0.3)
- ST-MPRExpl:  $\dot{y} = \lambda y(1 - y) \rightarrow$  exacte evolution,  $\dot{y} = \alpha \sin y \rightarrow$  expl. midpoint rule (11.4.0.7) & Strang splitting (12.5.0.4)

We observe that this splitting scheme can cope well with the stiffness of the problem, because the stiff term on the right hand side is integrated exactly.  $\square$

**EXAMPLE 12.5.0.12 (Splitting linear and local terms)** In the numerical treatment of partial differential equation one commonly encounters ODEs of the form

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) := -\mathbf{A}\mathbf{y} + \begin{bmatrix} g(y_1) \\ \vdots \\ g(y_d) \end{bmatrix}, \quad \mathbf{A} = \mathbf{A}^\top \in \mathbb{R}^{d,d} \text{ positive definite } (\rightarrow \text{Def. 1.1.2.6}), \quad (12.5.0.13)$$

with state space  $D = \mathbb{R}^d$ , where  $\lambda_{\min}(\mathbf{A}) \approx 1$ ,  $\lambda_{\max}(\mathbf{A}) \approx d^2$ , and the derivative of  $g : \mathbb{R} \rightarrow \mathbb{R}$  is bounded. Then IVPs for (12.5.0.13) will be stiff, since the Jacobian

$$\mathbf{D}\mathbf{f}(\mathbf{y}) = -\mathbf{A} + \begin{bmatrix} g'(y_1) & & \\ & \ddots & \\ & & g'(y_d) \end{bmatrix} \in \mathbb{R}^{d,d}$$

will have eigenvalues “close to zero” and others that are large (in modulus) and negative. Hence,  $\mathbf{D}\mathbf{f}(\mathbf{y})$  will satisfy the criteria (12.2.0.13) and (12.2.0.14) for any state  $\mathbf{y} \in \mathbb{R}^d$ .

The natural splitting is

$$\mathbf{f}(\mathbf{y}) = \mathbf{g}(\mathbf{y}) + \mathbf{q}(\mathbf{y}) \quad \text{with} \quad \mathbf{g}(\mathbf{y}) := -\mathbf{A}\mathbf{y}, \quad \mathbf{q}(\mathbf{y}) := \begin{bmatrix} g(y_1) \\ \vdots \\ g(y_d) \end{bmatrix}.$$

- For the **linear ODE**  $\dot{\mathbf{y}} = \mathbf{g}(\mathbf{y})$  we have to use and L-stable ( $\rightarrow$  Def. 12.3.4.12) single step method, for instance a second-order *implicit* Runge-Kutta method. Its increments can be obtained by solving a *linear system of equations*, whose coefficient matrix will be the same for every step, if uniform timesteps are used.
- The ODE  $\dot{\mathbf{y}} = \mathbf{q}(\mathbf{y})$  boils down to **decoupled** scalar ODEs  $\dot{y}_j = g(y_j)$ ,  $j = 1, \dots, d$ . For them we can use an inexpensive *explicit* RK-SSM like the explicit trapezoidal method (11.4.0.6). According to our assumptions on  $g$  these ODEs are not haunted by stiffness.

$\square$

## Summary and Learning Outcomes

In this chapter you have learned

- the notion of a stiff initial-value problem,
- how to detect stiffness in systems of ODEs,
- that explicit (Runge-Kutta) single-step methods struggle with severe stability-induced timestep constraints for stiff IVPs,

- the concept of a stability function and how it can be used to predict stability-induced timestep constraints,
- what is an implicit Runge-Kutta single-step method,
- the collocation construction of implicit Runge-Kutta single-step method based on a quadrature rule,
- the idea underlying semi-implicit Runge-Kutta single-step methods,
- the splitting approach to ODEs whose right-hand-side function is the sum of two different terms.

# Bibliography

- [DR08] W. Dahmen and A. Reusken. *Numerik für Ingenieure und Naturwissenschaftler*. Heidelberg: Springer, 2008 (cit. on pp. 692, 696, 710).
- [DB02] P. Deuflhard and F. Bornemann. *Scientific Computing with Ordinary Differential Equations*. 2nd ed. Vol. 42. Texts in Applied Mathematics. New York: Springer, 2002 (cit. on pp. 715, 720).
- [Han02] M. Hanke-Bourgeois. *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*. Mathematische Leitfäden. Stuttgart: B.G. Teubner, 2002 (cit. on pp. 693, 698, 706, 721–723).
- [NS02] K. Nipp and D. Stoffer. *Lineare Algebra*. 5th ed. Zürich: vdf Hochschulverlag, 2002 (cit. on p. 700).
- [QSS00] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Vol. 37. Texts in Applied Mathematics. New York: Springer, 2000 (cit. on pp. 693, 698, 706, 710).
- [Ran15] Joachim Rang. “Improved traditional Rosenbrock-Wanner methods for stiff ODEs and DAEs”. In: *J. Comput. Appl. Math.* 286 (2015), pp. 128–144. DOI: [10.1016/j.cam.2015.03.010](https://doi.org/10.1016/j.cam.2015.03.010).
- [Str09] M. Struwe. *Analysis für Informatiker*. Lecture notes, ETH Zürich. 2009 (cit. on pp. 702, 708).

# Index

- LU*-decomposition
  - existence, 137
- $L^2$ -inner product, 488
- $h$ -convergence, 469
- $p$ -convergence, 471
- (Asymptotic) complexity, 68
- (Size of) best approximaton error, 413
- Fill-in, 172
- Preconditioner, 674
- BLAS
  - axpy, 64
- C++
  - move semantics, 27
- EIGEN: triangularView, 72
- MATLAB: cumsum, 73
- PYTHON: reshape, 52
- 3-term recursion
  - for Chebychev polynomials, 428
  - for Legendre polynomials, 493
- 3-term recusion
  - orthogonal polynomials, 448
- 5-points-star-operator, 325
- a posteriori
  - adaptive quadrature, 507
- a posteriori adaptive, 427
- a posteriori error bound, 527
- a posteriori termination, 525
- a priori
  - adaptive quadrature, 507
- a priori termination, 525
- A-inner product, 655
- A-orthogonal, 665
- A-stability of a Runge-Kutta single step method, 762
- A-stable single step method, 762
- Absolute and relative error, 83
- absolute error, 83
- absolute tolerance, 525, 554, 721
- adaptive
  - a posteriori, 427
- adaptive multigrid quadrature, 509
- adaptive quadrature, 507
- a posteriori, 507
- a priori, 507
- Adding EPS to 1, 85
- AGM, 524
- Aitken-Neville scheme, 354
- algebraic convergence, 418, 706
- algebraic dependence, 69
- algebraically equivalent, 91
- aliasing, 460
- alternation theorem, 452
- Analyticity of a complex valued function, 424
- approximation
  - uniform, 407
- approximation error, 407
- arrow matrix, 173
- Ass: “Axiom” of roundoff analysis, 84
- Ass: Analyticity of interpoland, 425
- Ass: Global solutions, 694
- Ass: Sampling in a period, 392
- Ass: Self-adjointness of multiplication operator, 446
- asymptotic complexity, 68
  - sharp bounds, 68
- asymptotic rate of linear convergence, 534
- augmented normal equations, 266
- autonomization, 690
- Autonomous ODE, 686
- AXPY operation, 667
- axpy operation, 64
- back substitution, 118
- backward error analysis, 104
- backward substitution, 131
- Bandbreite
  - Zeilen-, 177
- banded matrix, 177
- bandwidth, 177
  - lower, 177
  - minimizing, 181
  - upper, 177
- barycentric interpolation formula, 352
- basis
  - cosine, 329

- orthonormal, 593
- sine, 323
- trigonometric, 285
- Belousov-Zhabotinsky reaction, 719
- bending energy, 384
- Bernstein approximant, 411
- Bernstein polynomials, 411
- Besetzungsmuster, 183
- best approximation
  - uniform, 452
- best approximation error, 413, 444
- bicg, 681
- BiCGStab, 681
- bisection, 537
- BLAS, 61
- block LU-decomposition, 132
- block matrix multiplication, 60
- blow-up, 720
- blurring operator, 302
- Boundary edge, 163
- Broyden
  - quasi-Newton method, 573
- Broyden-Verfahren
  - convergence monitor, 574
- Butcher scheme, 714, 759
- cache miss, 53
- cache thrashing, 53
- cancellation, 86, 88
- capacitance, 110
- capacitor, 110
- cardinal
  - spline, 387
- cardinal basis, 343, 347
- cardinal basis function, 387
- cardinal interpolant, 387, 480
- Cauchy product
  - of power series, 277
- Cauchy-Schwarz inequality, 422
- Causal channel/filter, 272
- causal filter, 271
- CCS format, 159
- cell
  - of a mesh, 467
- CG
  - convergence, 672
  - preconditioned, 675
  - termination criterion, 668
- CG = conjugate gradient method, 662
- CG algorithm, 667
- chain rule, 557
- channel, 271
- Characteristic parameters of IEEE floating point numbers, 82
- characteristic polynomial, 591
- Chebychev expansion, 437
- Chebychev nodes, 430, 431
- Chebychev polynomials, 428, 670
  - 3-term recursion, 428
- Chebychev-interpolation, 427
- chemical reaction kinetics, 748
- Cholesky decomposition
  - costs, 189
- circuit simulation
  - transient, 688
- circulant matrix, 280
- Classical Runge-Kutta method
  - Butcher scheme, 716
- Clenshaw algorithm, 438
- cluster analysis, 259
- coefficient matrix, 109
- coil, 110
- collocation, 755
- collocation conditions, 755
- collocation points, 755
- collocation single step methods, 755
- column
  - of a matrix, 39
- column major matrix format, 49
- column sum norm, 103
- column transformation, 60
- combinatorial graph Laplacian, 165
- complexity
  - asymptotic, 68
  - linear, 71
  - of SVD, 240
- composite quadrature formulas, 499
- Compressed Column Storage (CCS), 159
- compressed row storage, 159
- Compressed Row Storage (CRS), 159
- computational cost
  - Gaussian elimination, 120
- computational costs
  - LU-decomposition, 130
  - QR-decomposition, 227
- Computational effort, 67
- computational effort, 67, 551
  - eigenvalue computation, 595
- concave
  - data, 370
  - function, 371
- Condition (number) of a matrix, 114
- condition number

- of a matrix, 114
- spectral, 662
- conjugate gradient method, 662
- consistency
  - of iterative methods, 519
  - fixed point iteration, 529
- Consistency of fixed point iterations, 529
- Consistency of iterative methods, 519
- Consistent single step methods, 703
- constant
  - Lebesgue, 432
- constitutive relations, 110, 340
- constrained least squares, 265
- Contractive mapping, 532
- Convergence, 519
- convergence
  - algebraic, 418, 706
  - asymptotic, 548
  - exponential, 418, 422, 433, 706
  - global, 519
  - iterative method, 519
  - linear, 520
  - linear in Gauss-Newton method, 583
  - local, 519
  - numerical quadrature, 480
  - quadratic, 524
  - rate, 520
- convergence monitor, 574
  - of Broyden method, 574
- convex
  - data, 370
  - function, 371
- Convex/concave data, 370
- convex/concave function, 371
- convolution
  - discrete, 271, 276
  - discrete periodic, 278
  - of sequences, 277
- Corollary: “Optimality” of CG iterates, 670
- Corollary: Best approximant by orthogonal projection, 444
- Corollary: Composition of orthogonal transformations, 215
- Corollary: Consistent Runge-Kutta single step methods, 715
- Corollary: Continuous local Lagrange interpolants, 468
- Corollary: Dimension of  $\mathcal{P}_{2n}^T$ , 394
- Corollary: Euclidean matrix norm and eigenvalues, 103
- Corollary: Invariance of order under affine trans-
- formation, 484
- Corollary: Lagrange interpolation as linear mapping, 348
- Corollary: ONB representation of best approximant, 445
- Corollary: Periodicity of Fourier transforms, 308
- Corollary: Piecewise polynomials Lagrange interpolation operator, 468
- Corollary: Polynomial stability function of explicit RK-SSM, 740
- Corollary: Principal axis transformation, 592
- Corollary: Rational stability function of explicit RK-SSM, 762
- Corollary: Smoothness of cubic Hermite polynomial interpolant, 374
- Corollary: Stages limit order of explicit RK-SSM, 741
- Corollary: Uniqueness of least squares solutions, 201
- Corollary: Uniqueness of QR-factorization, 215
- Correct rounding, 83
- cosine
  - basis, 329
  - transform, 329
- cosine matrix, 329
- cosine transform, 329
- costs
  - Cholesky decomposition, 189
- Crout's algorithm, 129
- CRS, 159
- CRS format
  - diagonal, 159
- cubic complexity, 69
- cubic Hermite interpolation, 350, 374, 376
- Cubic Hermite polynomial interpolant, 374
- cubic spline interpolation
  - error estimates, 473
- cyclic permutation, 175
- damped Newton method, 567
- damping factor, 569
- data fitting, 400, 578
  - linear, 401
  - polynomial, 403
- data interpolation, 338
- deblurring, 301
- definite, 102
- dense matrix, 156
- derivative
  - in vector spaces, 556
- Derivative of functions between vector spaces, 556

descent methods, 655  
 destructor, 28  
 DFT, 282, 286  
     two-dimensional, 298  
 Diagonal dominance, 185  
 diagonal matrix, 40  
 diagonalization  
     for solving linear ODEs, 742  
     of a matrix, 592  
 diagonalization of local translation invariant linear operators, 325  
 Diagonally dominant matrix, 185  
 diagonally implicit Runge-Kutta method, 759  
 difference quotient, 88  
     backward, 699  
     forward, 698  
     symmetric, 700  
 difference scheme, 698  
 differential, 556  
 dilation, 393  
 direct power method, 604  
 DIRK-SSM, 759  
 discrete  $L^2$ -inner product, 446, 449  
 Discrete convolution, 276  
 discrete convolution, 271, 276  
 discrete evolution, 702  
 discrete Fourier transform, 282, 286  
 Discrete periodic convolution, 278  
 discrete periodic convolution, 278  
 discretization  
     of a differential equation, 702  
 discretization error, 704  
 discriminant formula, 86, 92  
 divided differences, 362  
 domain of definition, 517  
 domain specific language (DSL), 42  
 dot product, 54  
 double nodes, 350  
 double precision, 81  
 DSL: domain specific language, 42  
 economical singular value decomposition, 237  
 efficiency, 551  
 Eigen, 42  
     arrays, 45  
     data types, 42  
     initialisation, 44  
     sparse matrices, 159  
 eigen  
     accessing matrix entries, 44  
 eigenspace, 591  
 eigenvalue, 591  
     generalized, 593  
 eigenvalue problem  
     generalized, 593  
 eigenvalues and eigenvectors, 591  
 eigenvector, 591  
     generalized, 593  
 electric circuit, 109, 516  
     resonant frequencies, 587  
 element  
     of a matrix, 39  
 elementary arithmetic operations, 79, 84  
 elimination matrix, 126  
 embedded Runge-Kutta methods, 728  
 Energy norm, 655  
 entry  
     of a matrix, 39  
 envelope  
     matrix, 177  
 Equation  
     non-linear, 517  
 equidistant mesh, 467  
 equidistribution principle  
     for quadrature error, 508  
 equivalence  
     of norms, 102  
 Equivalence of norms, 521  
 ergodicity, 603  
 error  
     absolute, 83  
     relative, 83  
 error estimator  
     a posteriori, 527  
 error indicator  
     for extrapolation, 360  
 Euler method  
     explicit, 697  
     implicit, 699  
     implicit, stability function, 761  
     semi implicit, 767  
 Euler polygon, 698  
 Euler's formula, 393  
 Euler's iteration, 546  
 evolution operator, 694  
 Evolution operator/mapping, 694  
 expansion  
     asymptotic, 357  
 explicit Euler method, 697  
     Butcher scheme, 715  
 explicit midpoint rule  
     Butcher scheme, 715  
     for ODEs, 713

- Explicit Runge-Kutta method, 714
- explicit Runge-Kutta method, 714
- explicit trapezoidal rule
  - Butcher scheme, 715
  - exponential convergence, 433
  - extended normal equations, 209
  - extended state space
    - of an ODE, 690
  - extrapolation, 357
- Families of sparse matrices, 157
- fast Fourier transform, 316
- FFT, 316
- fill-in, 172
- filter
  - high pass, 293
  - low pass, 293
- Finding out EPS in C++, 84
- Finite channel/filter, 272
- finite filter, 271
- Fitted polynomial, 450
- fixed point, 529
- fixed point form, 529
- fixed point iteration, 528
- fixed point iteration
  - consistency, 529
  - Newton's method, 564
- floating point number, 80
- floating point numbers, 79, 80
- forward elimination, 118
- forward substitution, 131
- Fourier
  - matrix, 285
  - Fourier coefficient, 311
  - Fourier modes, 460
  - Fourier series, 308, 459
  - Fourier transform, 308
    - discrete, 282, 286
  - fractional order of convergence, 548
  - frequency domain, 110
  - frequency filtering, 289
  - Frobenius norm, 248
  - full-rank condition, 201
  - function
    - concave, 371
    - convex, 371
  - function object, 30
  - function representation, 341
- Gauss collocation single step method, 758
- Gauss Quadrature, 484
- Gauss-Legendre quadrature formulas, 492
- Gauss-Newton method, 581
- Gauss-Radau quadrature formulas, 765
- Gauss-Seidel preconditioner, 676
- Gaussian elimination, 117
  - block version, 124
  - by rank-1 modifications, 123
  - for non-square matrices, 122
- general least squares problem, 205
- generalization error, 404
- Generalized condition number of a matrix, 206
- Generalized Lagrange polynomials, 350
- Generalized solution of a linear system of equations, 203
- Givens rotation, 218
- Givens-Rotation, 221, 233
- global solution
  - of an IVP, 693
- GMRES, 680
- Golub-Welsch algorithm, 494
- gradient, 557, 657
- Gradient and Hessian, 557
- Gram-Schmidt
  - Orthonormalisierung, 645
- Gram-Schmidt orthogonalisation, 76, 211
- Gram-Schmidt orthogonalization, 445, 645, 666
- Gram-Schmidt orthonormalization, 632
- graph partitioning, 617
- grid, 467
- grid cell, 467
- grid function, 326
- grid interval, 467
- Halley's iteration, 546
- harmonic mean, 378
- hat function, 342
- heartbeat model, 687
- Hermite interpolation
  - cubic, 350
- Hermitian matrix, 41
- Hermitian/symmetric matrices, 41
- Hessenberg matrix, 230
- Hessian, 41
- Hessian matrix, 557
- high pass filter, 293
- Hilbert matrix, 90
- homogeneous, 102
- Hooke's law, 625
- Horner scheme, 346
- Householder reflection, 216
- I/O-complexity, 68
- identity matrix, 40

- IEEE standard 754, 81
- ill conditioned, 115
- ill-conditioned problem, 105
- image segmentation, 607
- image space, 196
- Image space and kernel of a matrix, 111
- implicit differentiation, 561
- implicit Euler method, 699
- implicit function theorem, 700
- implicit midpoint method, 700, 701
- Impulse response, 271
- impulse response, 271
  - of a filter, 271
- in place, 129, 130
- in situ, 124, 130
- increment equations
  - linearized, 768
- increments
  - Runge-Kutta, 714, 759
- inductance, 110
- inductor, 110
- inexact splitting methods, 771
- inf, 81
- infinity, 81
- initial guess, 519, 528
- initial value problem
  - stiff, 750
- initial value problem (IVP), 690
- Inner product, 442
- inner product
  - A-, 655
- intermediate value theorem, 537
- interpolant
  - piecewise linear, 372
- interpolation
  - barycentric formula, 352
  - Chebychev, 427
  - complete cubic spline, 383
  - cubic Hermite, 376
  - Hermite, 350
  - Lagrange, 347
  - natural cubic spline, 383
  - periodic cubic spline, 384
  - spline cubic, 381
  - spline cubic, locality, 387
  - spline shape preserving, 388
    - trigonometric, 392
  - interpolation operator, 344
  - interpolation problem, 340
  - interpolation scheme, 340
  - inverse interpolation, 549
- inverse iteration, 618
  - preconditioned, 620
- inverse matrix, 111
- Invertible matrix, 111
- invertible matrix, 111, 112
- iteration, 518
  - Halley's, 546
  - Euler's, 546
  - quadratical inverse interpolation, 546
- iteration function, 519, 528
- iterative method, 518
  - convergence, 519
- IVP, 690
- Jacobi preconditioner, 676
- Jacobian, 533, 553, 556
- kernel, 196
- kinetics
  - of chemical reaction, 748
- Kirchhoff (current) law, 110
- knots
  - spline, 380
- Konvergenz
  - Algebraische, Quadratur, 497
- Kronecker product, 74
- Kronecker symbol, 39
- Krylov space, 664
  - for Ritz projection, 640
- L-stable, 764
- L-stable Runge-Kutta method, 764
- Lagrange function, 265
- Lagrange interpolation approximation scheme, 416
- Lagrange multiplier, 265
- Lagrangian (interpolation polynomial) approximation scheme, 417
- Lagrangian multiplier, 265
- lambda function, 21, 30
- Landau symbol, 68, 418
- Landau-O, 68
- Lapack, 122
- leading coefficient
  - of polynomial, 345
- Least squares
  - with linear constraint, 265
- least squares
  - total, 263
- least squares problem, 203
- Least squares solution, 196
- least-squares problem

- non-linear, 579
- least-squares solution
  - non-linear, 579
- Lebesgue
  - constant, 432
- Lebesgue constant, 368, 423
- Legendre polynomials, 491
- Lemma:  $\mathbf{r}_k \perp U_k$ , 663
- Lemma: Absolute conditioning of polynomial interpolation, 367
- Lemma: Affine pullbacks preserve polynomials, 415
- Lemma: Bases for Krylov spaces in CG, 666
- Lemma: Cholesky decomposition, 188
- Lemma: Criterion for local Lipschitz continuity, 692
- Lemma: Cubic convergence of modified Newton methods, 545
- Lemma: Decay of Fourier coefficients, 462
- Lemma: Diagonal dominance and definiteness, 187
- Lemma: Diagonalization of circulant matrices, 286
- Lemma: Equivalence of Gaussian elimination and LU-factorization, 140
- Lemma: Error representation for polynomial Lagrange interpolation, 421
- Lemma: Exact quadrature by equidistant trapezoidal rule, 505
- Lemma: Existence of LU-decomposition, 127
- Lemma: Existence of LU-factorization with pivoting, 137
- Lemma: Formula for Euclidean norm of a Hermitian matrix, 103
- Lemma: Fourier coefficients of derivatives, 461
- Lemma: Gershgorin circle theorem, 592
- Lemma: Group of regular diagonal/triangular matrices, 58
- Lemma: Higher order local convergence of fixed point iterations, 535
- Lemma: Interpolation error estimates for exponentially decaying Fourier coefficients, 465
- Lemma: Kernel and range of (Hermitian) transposed matrices, 201
- Lemma: LU-factorization of diagonally dominant matrices, 185
- Lemma: Ncut and Rayleigh quotient ( $\rightarrow$  [SHM00]), 611
- Lemma: Necessary conditions for s.p.d., 41
- Lemma: Positivity of Gauss-Legendre quadrature weights, 492
- Lemma: Properties of cosine matrix, 329
- Lemma: Properties of Fourier matrix, 286
- Lemma: Properties of the sine matrix, 323
- Lemma: Quadrature error estimates for  $C^r$ -integrands, 496
- Lemma: Quadrature formulas from linear interpolation schemes, 480
- Lemma: Residual formula for quotients, 425
- Lemma: S.p.d. LSE and quadratic minimization problem, 655
- Lemma: Sherman-Morrison-Woodbury formula, 155
- Lemma: Similarity and spectrum  $\rightarrow$  [Gut09, Thm. 9.7], [DR08, Lemma 7.6], [NS02, Thm. 7.2], 592
- Lemma: Smoothness of solutions of ODEs, 685
- Lemma: Stability function as approximation of  $\exp$  for small arguments, 740
- Lemma: Sufficient condition for linear convergence of fixed point iteration, 534
- Lemma: Sufficient condition for local linear convergence of fixed point iteration, 533
- Lemma: SVD and Euclidean matrix norm, 245
- Lemma: SVD and rank of a matrix  $\rightarrow$  [NS02, Cor. 9.7], 238
- Lemma: Taylor expansion of inverse distance function, 626
- Lemma: Theory of Arnoldi process, 646
- Lemma: Transformation of norms under affine pullbacks, 415
- Lemma: Tridiagonal Ritz projection from CG residuals, 642
- Lemma: Unique solvability of linear least squares fitting problem, 403
- Lemma: Uniqueness of orthonormal polynomials, 447
- Lemma: Zeros of Legendre polynomials, 492
- Levinson algorithm, 334
- Lie-Trotter splitting, 770
- limit cycle, 749
- limiter, 377
- line search, 656
- Linear channel/filter, 272
- linear complexity, 69, 71
- Linear convergence, 520
- linear correlation, 256
- linear data fitting, 401
- linear electric circuit, 109
- linear filter, 271
- Linear interpolation operator, 344

linear operator, 344  
     diagonalization, 325  
 linear ordinary differential equation, 590  
 linear regression, 69  
 linear system of equations, 109  
     multiple right hand sides, 122  
 Lipschitz continuous function, 691, 692  
 Lloyd-Max algorithm, 259  
 Local and global convergence, 519  
 local Lagrange interpolation, 467  
 local linearization, 553  
 locality  
     of interpolation, 386  
 logistic differential equation, 686  
 Lotka-Volterra ODE, 686  
 low pass filter, 293  
 lower triangular matrix, 40  
 LU-decomposition  
     blocked, 132  
     computational costs, 130  
     envelope aware, 179  
     existence, 127  
     in place, 130  
 LU-decomposition/LU-factorization, 127  
 LU-factorization  
     envelope aware, 179  
     of sparse matrices, 171  
     with pivoting, 135  
 machine number, 80  
     exponent, 80  
 machine numbers, 79, 81  
     distribution, 80  
     extremal, 80  
 Machine numbers/floating point numbers, 80  
 machine precision, 84  
 mantissa, 80  
 Markov chain, 332, 598  
     stationary distribution, 599  
 mass matrix, 626  
 Matrix  
     adjoint, 40  
     Hermitian, 592  
     Hermitian transposed, 40  
     normal, 592  
     skew-Hermitian, 592  
     transposed, 40  
     unitary, 592  
 matrix  
     banded, 177  
     condition number, 114  
     dense, 156  
     diagonal, 40  
     envelope, 177  
     Fourier, 285  
     Hermitian, 41  
     Hessian, 557  
     lower triangular, 40  
     normalized, 40  
     orthogonal, 211  
     positive definite, 41  
     positive semi-definite, 41  
     rank, 112  
     sine, 323  
     sparse, 156  
     storage formats, 49  
     structurally symmetric, 181  
     symmetric, 41  
     tridiagonal, 177  
     unitary, 211  
     upper triangular, 40  
 matrix block, 39  
 matrix compression, 247  
 Matrix envelope, 177  
 matrix exponential, 742  
 matrix factorization, 125, 126  
 Matrix norm, 102  
     matrix norm, 102  
         column sums, 103  
         row sums, 103  
 matrix storage  
     envelope oriented, 181  
 maximum likelihood, 197  
 member function, 18  
 mesh, 467  
     equidistant, 467  
     in time, 702  
     temporal, 696  
 mesh adaptation, 509  
 mesh refinement, 509  
 mesh width, 467  
 Method  
     Quasi-Newton, 571  
 method, 18  
 midpoint method  
     implicit, stability function, 761  
 midpoint rule, 482, 713  
 Milne rule, 483  
 min-max theorem, 613  
 minimal residual methods, 680  
 model function, 539  
 model reduction, 408  
 Modellfunktionsverfahren, 539

modification techniques, 229  
 modified Newton method, 545  
 monomial representation  
     of a polynomial, 345  
 monomials, 345  
 monotonic data, 370  
 Moore-Penrose pseudoinverse, 205  
 move semantics, 27  
 multi-point methods, 539, 546  
 multiplicity  
     geometric, 591  
     of an interpolation node, 349  
  
 NaN, 81  
 Ncut, 609  
 nested  
     subspaces, 663  
 nested spaces, 410  
 Newton  
     basis, 361  
     damping, 569  
     damping factor, 569  
     monotonicity test, 569  
     simplified method, 563  
 Newton correction, 553  
     simplified, 567  
 Newton iteration, 553  
     numerical Differentiation, 564  
     termination criterion, 565  
 Newton method  
     1D, 539  
     damped, 567  
     local quadratic convergence, 564  
     modified, 545  
 Newton's law of motion, 626  
 nodal analysis, 109, 516  
     transient, 688  
 nodal polynomial, 427  
 nodal potentials, 110  
 node  
     double, 350  
     for interpolation, 347  
     in electric circuit, 109  
     multiple, 349  
     multiplicity, 349  
     of a mesh, 467  
     quadrature, 478  
 nodes, 347  
     Chebychev, 431  
     Chebychev nodes, 432  
         for interpolation, 347  
 non-linear least-squares problem, 579  
  
 Non-linear least-squares solution, 579  
 non-linear least-squares solution, 579  
 non-normalized numbers, 81  
 Norm, 102  
 norm, 102  
      $L^1$ , 367  
      $L^2$ , 367  
      $\infty$ -, 102  
     1-, 102  
     energy-, 655  
     Euclidean, 102  
     Frobenius norm, 248  
     of matrix, 102  
     Sobolev semi-, 423  
     supremum, 367  
 normal equations, 198, 443  
     augmented, 266  
     extended, 209  
     with constraint, 266  
 normalization, 604  
 Normalized cut, 609  
 normalized lower triangular matrix, 127  
 normalized triangular matrix, 40  
 not a number, 81  
 nullspace, 196  
 Nullstellenbestimmung  
     Modellfunktionsverfahren, 539  
 Numerical differentiation  
     roundoff, 89  
 numerical Differentiation  
     Newton iteration, 564  
     numerical differentiation, 88  
     numerical quadrature, 476  
     numerical rank, 240, 243  
  
 ODE, 690  
     scalar, 697  
 Ohmic resistor, 110  
 one-point methods, 539  
 one-step error, 708  
 order  
     of an ODE, 691  
     of quadrature formula, 484  
 Order of a quadrature rule, 484  
 Order of a single step method, 707  
 order of convergence, 523  
     fractional, 548  
 ordinary differential equation  
     linear, 590  
 ordinary differential equation (ODE), 690  
 oregonator, 719  
 orthogonal matrix, 211

- orthogonal polynomials, 491
- orthogonal projection, 444
- Orthogonality, 442
- Orthonormal basis, 445
- orthonormal basis, 445, 593
- Orthonormal polynomials, 447
- overfitting, 404
- overflow, 81, 85
- overloading
  - of functions, 18
  - of operators, 19
- page rank, 597
  - stochastic simulation, 598
- parameter estimation, 193
- PARDISO, 169
- partial pivoting, 135
- pattern
  - of a matrix, 57
- PCA, 251
- PCG, 675
- Peano
  - Theorem of, 692
- penalization, 614
- penalty parameter, 614
- periodic
  - function, 392
- periodic sequence, 277
- permutation, 136
- Permutation matrix, 136
- permutation matrix, 136, 231
- perturbation lemma, 114
- Petrov-Galerkin condition, 681
- phase space
  - of an ODE, 690
- phenomenological model, 687
- Picard-Lindelöf
  - Theorem of, 692
- Piecewise cubic Hermite interpolant (with exact slopes), 471
- PINVIT, 620
- Pivot
  - choice of, 134
  - pivot, 118, 119
  - pivot row, 118, 119
  - pivoting, 133
- Planar triangulation, 161
- point spread function, 301
- polynomial
  - characteristic, 591
  - generalized Lagrange, 350
  - Lagrange, 347
- polynomial fitting, 403
- polynomial interpolation
  - existence and uniqueness, 348
  - generalized, 349
- polynomial space, 345
- positive definite
  - criteria, 41
  - matrix, 41
- potentials
  - nodal, 110
- power series, 424
- power spectrum
  - of a signal, 294
- preconditioned CG method, 675
- preconditioned inverse iteration, 620
- preconditioner, 674
- preconditioning, 673
- predator-prey model, 686
- principal axis, 259
- principal axis transformation, 659
- principal component, 256
- principal component analysis (PCA), 251
- principal minor, 131
- principal orthogonal decomposition (POD), 257
- problem
  - ill conditioned, 115
  - ill-conditioned, 105
  - sensitivity, 112
  - well conditioned, 115
- procedural form, 476
- product rule, 557
- propagated error, 708
- pullback, 414, 478
- Punkt
  - stationär, 687
- pwer method
  - direct, 604
- Python, 48
- QR algorithm, 593
- QR-algorithm with shift, 594
- QR-decomposition, 78, 213
  - computational costs, 227
- QR-factorization, QR-decomposition, 217
- quadratic complexity, 69
- quadratic convergence, 536
- quadratic eigenvalue problem, 587
- quadratic functional, 655
- quadratic inverse interpolation, 550
- quadratical inverse interpolation, 546
- quadrature
  - adaptive, 507

polynomial formulas, 481  
 quadrature formula  
     order, 484  
 Quadrature formula/quadrature rule, 478  
 quadrature node, 478  
 quadrature numerical, 476  
 quadrature point, 478  
 quadrature weight, 478  
 quasi-linear system, 559  
 Quasi-Newton method, 571, 573  
  
 Radau RK-method  
     order 3, 765  
     order 5, 765  
 radiative heat transfer, 278  
 range, 196  
 rank  
     column rank, 112  
     computation, 240  
     numerical, 240, 243  
     of a matrix, 112  
     row rank, 112  
 Rank of a matrix, 112  
 rank-1 modification, 123, 229  
 rank-1-matrix, 71  
 rank-1-modification, 154, 573  
 rate  
     of algebraic convergence, 418, 706  
     of convergence, 520  
 Rayleigh quotient, 605, 613  
 Rayleigh quotient iteration, 619  
 Region of (absolute) stability, 747  
 regular matrix, 111  
 Regular refinement of a planar triangulation, 165  
 relative error, 83  
 relative tolerance, 525, 554, 721  
 rem:Fspec, 286  
 Residual, 139  
 residual quantity, 621  
 Riccati differential equation, 696, 698  
 Riemann sum, 311  
 right hand side  
     of an ODE, 690  
 right hand side vector, 109  
 rigid body mode, 628  
 Ritz projection, 635, 639  
 Ritz value, 635  
 Ritz vector, 635  
 root of unity, 284  
 roots of unity, 504  
 rounding, 83  
  
 rounding up, 83  
 roundoff  
     for numerical differentiation, 89  
 row  
     of a matrix, 39  
 row major matrix format, 49  
 ROW methods, 769  
 row sum norm, 103  
 row transformation, 60, 117, 125, 126  
 Runge's example, 366  
 Runge-Kutta  
     increments, 714, 759  
 Runge-Kutta method, 714, 759  
     L-stable, 764  
 Runge-Kutta methods  
     embedded, 728  
     semi-implicit, 766  
     stability function, 739, 761  
  
 saddle point problem, 265  
     matrix form, 266  
 scalar ODE, 697  
 scaling  
     of a matrix, 58  
 scheme  
     Horner, 346  
 Schur  
     Komplement, 132  
 Schur complement, 132, 151, 154  
 scientific notation, 80  
 secant condition, 572  
 secant method, 547, 550, 572  
 segmentation  
     of an image, 607  
 semi-implicit Euler method, 767  
 seminorm, 423  
 sensitive dependence, 105  
 sensitivity  
     of a problem, 112  
 shape  
     preservation, 373  
     preserving spline interpolation, 388  
 Sherman-Morrison-Woodbury formula, 155  
 shifted inverse iteration, 618  
 signal  
     time-discrete, 270  
 similarity  
     of matrices, 592  
 similarity function  
     for image segmentation, 608  
 similarity transformations, 592  
 similary transformation

- unitary, 594
- Simpson rule, 482
- sine
  - basis, 323
  - matrix, 323
  - transform, 323
- Sine transform, 322
- single precision, 81
- single step method
  - A-stability, 762
- Single-step method, 702
- single-step method, 702
- singular value decomposition, 235, 237
- Singular value decomposition (SVD), 237
- slopes
  - for cubic Hermite interpolation, 374
- Smoothed triangulation, 164
- Solution of an ordinary differential equation, 685
- Space of trigonometric polynomials, 393
- Sparse matrix, 157
- sparse matrix, 156
  - COO format, 158
  - LU-factorization, 171
  - triplet format, 158
- sparse matrix storage formats, 158
- spectral condition number, 662
- spectral partitioning, 617
- spectral radius, 591
- spectrum, 591
  - of a matrix, 659
- spline, 380
  - cardinal, 387
  - complete cubic, 383
  - cubic, 381
  - cubic, locality, 387
  - knots, 380
  - natural cubic, 383
  - periodic cubic, 384
  - physical, 385
  - shape preserving interpolation, 388
- Splines, 380
- splitting
  - Lie-Trotter, 770
  - Strang, 770
- splitting methods, 769
  - inexact, 771
- spy, 57, 58
- stability function
  - of explicit Runge-Kutta methods, 739
  - of Runge-Kutta methods, 761
- stable
  - algorithm, 104
  - numerically, 104
- Stable algorithm, 104
- stages, 759
- state space
  - of an ODE, 690
- stationary distribution, 599
- steepest descent, 656
- Stiff IVP, 750
- stiffness matrix, 626
- stochastic matrix, 599
- stochastic simulation of page rank, 598
- stopping rule, 524
- Strang splitting, 770
- Strassen's algorithm, 70
- Structurally symmetric matrix, 181
- structurally symmetric matrix, 181
- sub-matrix, 39
- sub-multiplicative, 102
- subspace correction, 663
- subspace iteration
  - for direct power method, 637
- subspaces
  - nested, 663
- SuperLU, 169
- surrogate function, 408
- SVD, 235, 237
- symmetric matrix, 41
- Symmetric positive definite (s.p.d.) matrices, 41
- symmetry
  - structural, 181
- system matrix, 109
- system of equations
  - linear, 109
- tangent field, 696
- Taylor expansion, 535
- Taylor polynomial, 410
- Taylor series, 424
- Taylor's formula, 410
- template, 19
- tensor product, 54
- tent function, 342
- Teoplitz matrices, 331
- termination criterion, 524
  - ideal, 525
  - Newton iteration, 565
  - residual based, 526
- Theorem: → [Han02, Thm. 25.4], 620
- Theorem:  $L^2$ -error estimate for trigonometric interpolation, 462

- Theorem:  $L^\infty$  polynomial best approximation estimate, 414
- Theorem: (Absolute) stability of explicit RK-SSM for linear systems of ODEs, 746
- Theorem: 3-term recursion for Chebychev polynomials, 428
- Theorem: 3-term recursion for orthogonal polynomials, 448
- Theorem: Uniform approximation by polynomials, 411
- Theorem: Courant-Fischer min-max theorem  
→ [GOL89], 613
- Theorem: Banach's fixed point theorem, 532
- Theorem: Best low rank approximation, 249
- Theorem: Bound for spectral radius, 591
- Theorem: Chebychev alternation theorem, 453
- Theorem: Commuting matrices have the same eigenvectors, 283
- Theorem: Composition of analytic functions, 426
- Theorem: Conditioning of LSEs, 113
- Theorem: Convergence of gradient method/steepest descent, 662
- Theorem: Convergence of approximation by cubic Hermite interpolation, 472
- Theorem: Convergence of CG method, 671
- Theorem: Convergence of direct power method  
→ [DR08, Thm. 25.1], 607
- Theorem: Convolution theorem, 288
- Theorem: Cost for solving triangular systems, 147
- Theorem: Cost of Gaussian elimination, 147
- Theorem: Criteria for invertibility of matrix, 112
- Theorem: Dimension of space of polynomials, 345
- Theorem: Dimension of spline space, 381
- Theorem: Divergent polynomial interpolants, 420
- Theorem: Envelope and fill-in, 178
- Theorem: Equivalence of all norms on finite dimensional vector spaces, 521
- Theorem: Existence & uniqueness of generalized Lagrange interpolation polynomials, 350
- Theorem: Existence & uniqueness of Lagrange interpolation polynomial, 348
- Theorem: Existence of  $n$ -point quadrature formulas of order  $2n$ , 490
- Theorem: Existence of least squares solutions, 197
- Theorem: Exponential convergence of trigonometric interpolation for analytic interpolants, 465
- Theorem: Exponential decay of Fourier coefficients of analytic functions, 463
- Theorem: Formula for generalized solution, 205
- Theorem: Gaussian elimination for s.p.d. matrices, 186
- Theorem: Gram-Schmidt orthonormalization, 445
- Theorem: Implicit function theorem, 700
- Theorem: Isometry property of Fourier transform, 313
- Theorem: Kernel and range of  $\mathbf{A}^\top \mathbf{A}$ , 201
- Theorem: Least squares solution of data fitting problem, 402
- Theorem: Local quadratic convergence of Newton's method, 565
- Theorem: Local shape preservation by piecewise linear interpolation, 373
- Theorem: Maximal order of  $n$ -point quadrature rule, 486
- Theorem: Mean square (semi-)norm/Inner product (semi-)norm, 442
- Theorem: Mean square norm best approximation through normal equations, 443
- Theorem: Minimax property of the Chebychev polynomials, 429
- Theorem: Monotonicity preservation of limited cubic Hermite interpolation, 379
- Theorem: Obtaining least squares solutions by solving normal equations, 198
- Theorem: Optimality of natural cubic spline interpolant, 385
- Theorem: Order of collocation single step method, 758
- Theorem: Order of simple splitting methods, 771
- Theorem: Positivity of Clenshaw-Curtis weights, 483
- Theorem: Preservation of Euclidean norm, 211
- Theorem: Property of linear, monotonicity preserving interpolation into  $C^1$ , 379
- Theorem: Pseudoinverse and SVD, 244
- Theorem: QR-decomposition, 214
- Theorem: QR-decomposition “preserves bandwidth”, 221
- Theorem: Quadrature error estimate for quadrature rules with positive weights, 496
- Theorem: Rayleigh quotient, 613
- Theorem: Region of stability of Gauss collocation single step methods, 763
- Theorem: Representation of interpolation error, 420
- Theorem: Residue theorem, 424

- Theorem: Schur's lemma, 592  
Theorem: Sensitivity of full-rank linear least squares problem, 206  
Theorem: singular value decomposition → [NS02, Thm. 9.6], [Gut09, Thm. 11.1], 235  
Theorem: Solution of POD problem, 258  
Theorem: Span property of G.S. vectors, 212  
Theorem: Stability function of explicit Runge-Kutta methods, 739  
Theorem: Stability function of general Runge-Kutta methods, 761  
Theorem: Stability of Gaussian elimination with partial pivoting, 140  
Theorem: Stability of Householder QR [Hig02, Thm. 19.4], 222  
Theorem: Sufficient order conditions for quadrature rules, 485  
Theorem: Taylor's formula, 535  
Theorem: Theorem of Peano & Picard-Lindelöf [Ama83, Satz II(7.6)], [Str09, Satz 6.5.1], [DR08, Thm. 11.10], [Han02, Thm. 73.1], 692  
Time-invariant channel/filter, 272  
time-invariant filter, 271  
timestep (size), 698  
timestep constraint, 741  
timestepping, 697  
Toeplitz matrix, 333  
Toeplitz solvers
  - fast algorithms, 336
tolerance, 526
  - absolute, 721
  - absoute, 525, 554
  - for adaptive timestepping for ODEs, 721
  - for termination, 525
  - realitive, 721
  - relative, 525, 554
total least squares, 263  
trajectory, 687  
transform
  - cosine, 329
  - fast Fourier, 316
  - sine, 323
transformation matrix, 60  
trapezoidal rule, 482, 504, 712
  - for ODEs, 713
trend, 251  
trial space
  - for collocation, 755
triangle inequality, 102  
triangular linear systems, 151  
triangulation, 161  
tridiagonal matrix, 177  
trigonometric basis, 285  
trigonometric interpolation, 392, 457  
Trigonometric polynomial, 312  
trigonometric polynomial, 312  
trigonometric polynomials, 393, 457  
trigonometric transformations, 322  
tripled format, 158  
truss structure
  - vibrations, 624
trust region method, 584  
Types of asymptotic convergence of approximation schemes, 418  
Types of matrices, 40  
UMFPACK, 169  
underflow, 81, 85  
uniform approximation, 407  
uniform best approximation, 452  
Uniform convergence
  - of Fourier series, 309
unit vector, 39  
Unitary and orthogonal matrices, 211  
unitary matrix, 211  
unitary similary transformation, 594  
upper Hessenberg matrix, 646  
upper triangular matrix, 40, 118, 126  
Vandermonde matrix, 349  
variational calculus, 385  
vector field, 690  
vectorization
  - of a matrix, 51
Vieta's formula, 92  
Weddle rule, 483  
weight
  - quadrature, 478
weight function, 446  
weighted  $L^2$ -inner product, 446  
well conditioned, 115  
Young's modulus, 625  
Zerlegung
  - LU, 131
zero padding, 281, 334, 398

# List of Symbols

|                                                                                                                                       |     |
|---------------------------------------------------------------------------------------------------------------------------------------|-----|
| $(\mathbf{A})_{i,j} \triangleq$ reference to entry $a_{ij}$ of matrix $\mathbf{A}$ , 39                                               | 196 |
| $(\mathbf{A})_{k:l,r:s} \triangleq$ reference to submatrix of $\mathbf{A}$ spanning rows $k, \dots, l$ and columns $r, \dots, s$ , 39 |     |
| $(\mathbf{x})_i \triangleq$ $i$ -th component of vector $\mathbf{x}$ , 38                                                             |     |
| $(x_k) *_n (y_k) \triangleq$ discrete periodic convolution, 278                                                                       |     |
| $C(I) \triangleq$ space of continuous functions $I \rightarrow \mathbb{R}$ , 366                                                      |     |
| $C^1([a, b]) \triangleq$ space of continuously differentiable functions $[a, b] \mapsto \mathbb{R}$ , 374                             |     |
| $J(t_0, \mathbf{y}_0) \triangleq$ maximal domain of definition of a solution of an IVP, 692                                           |     |
| $O \triangleq$ zero matrix, 40                                                                                                        |     |
| $O(\cdot) \triangleq$ Landau symbol, 68                                                                                               |     |
| $V^\perp \triangleq$ orthogonal complement of a subspace, 201                                                                         |     |
| $\mathcal{E} \triangleq$ expected value of a random variable, 332                                                                     |     |
| $\mathcal{P}_n^T \triangleq$ space of trigonometric polynomials of degree $n$ , 393                                                   |     |
| $\mathcal{R}_k(m, n) \triangleq$ set of rank- $k$ matrices, 248                                                                       |     |
| $D\Phi \triangleq$ Jacobian of $\Phi : D \mapsto \mathbb{R}^n$ at $\mathbf{x} \in D$ , 533                                            |     |
| $D_y \mathbf{f} \triangleq$ Derivative of $\mathbf{f}$ w.r.t.. $y$ (Jacobian), 692                                                    |     |
| $\text{EPS} \triangleq$ machine precision, 84                                                                                         |     |
| $\text{Eig}\mathbf{A}\lambda \triangleq$ eigenspace of $\mathbf{A}$ for eigenvalue $\lambda$ , 591                                    |     |
| $\mathcal{R}\mathbf{A} \triangleq$ range/column space of matrix $\mathbf{A}$ , 238                                                    |     |
| $\mathcal{N}(\mathbf{A}) \triangleq$ kernel/nullspace of a matrix, 111, 196                                                           |     |
| $\mathcal{N}\mathbf{A} \triangleq$ nullspace of matrix $\mathbf{A}$ , 238                                                             |     |
| $\mathcal{K}_l(\mathbf{A}, \mathbf{z}) \triangleq$ Krylov subspace, 664                                                               |     |
| $L_{\mathcal{T}} \triangleq$ Lagrangian (interpolation polynomial) approximation scheme on node set $\mathcal{T}$ , 417               |     |
| $\ \mathbf{Ax} - \mathbf{b}\ _2 \rightarrow \min \triangleq$ minimize $\ \mathbf{Ax} - \mathbf{b}\ _2$ , 203                          |     |
| $\ \mathbf{A}\ _F^2$ , 248                                                                                                            |     |
| $\ \mathbf{x}\ _A \triangleq$ energy norm induced by s.p.d. matrix $\mathbf{A}$ , 655                                                 |     |
| $\ \cdot\  \triangleq$ Euclidean norm of a vector $\in \mathbb{K}^n$ , 76                                                             |     |
| $\ \cdot\  \triangleq$ norm on vector space, 102                                                                                      |     |
| $\ f\ _{L^\infty(I)}$ , 367                                                                                                           |     |
| $\ f\ _{L^1(I)}$ , 367                                                                                                                |     |
| $\ f\ _{L^2(I)}^2$ , 367                                                                                                              |     |
| $\mathcal{P}_k$ , 345                                                                                                                 |     |
| $\Psi^h \mathbf{y} \triangleq$ discretei evolution for autonomous ODE, 702                                                            |     |
| $\mathcal{R}(\mathbf{A}) \triangleq$ image/range space of a matrix, 111,                                                              |     |
| $(\cdot, \cdot)_V \triangleq$ inner product on vector space $V$ , 442                                                                 |     |
| $\mathcal{S}_{d,\mathcal{M}}$ , 380                                                                                                   |     |
| $\mathbf{A}^\dagger \triangleq$ Moore-Penrose pseudoinverse of $\mathbf{A}$ , 205                                                     |     |
| $\mathbf{A}^\top \triangleq$ transposed matrix, 40                                                                                    |     |
| $\mathbf{I} \triangleq$ identity matrix, 40                                                                                           |     |
| $\mathbf{h} * \mathbf{x} \triangleq$ discrete convolution of two vectors, 276                                                         |     |
| $\mathbf{x} *_n \mathbf{y} \triangleq$ discrete periodic convolution of vectors, 278                                                  |     |
| $\bar{z} \triangleq$ complex conjugation, 40                                                                                          |     |
| $\mathbb{C}^- := \{z \in \mathbb{C}: \text{Re } z < 0\}$ , 762                                                                        |     |
| $\mathbb{K} \triangleq$ generic field of numbers, either $\mathbb{R}$ or $\mathbb{C}$ , 38                                            |     |
| $\mathbb{K}_*^{n,n} \triangleq$ set of invertible $n \times n$ matrices, 112                                                          |     |
| $\mathbb{M} \triangleq$ set of machine numbers, 79                                                                                    |     |
| $\delta_{ij} \triangleq$ Kronecker symbol, 39, 347                                                                                    |     |
| $\delta_{ij} \triangleq$ Kronecker symbol, 271                                                                                        |     |
| $\imath \triangleq$ imaginary unit, “ $\imath := \sqrt{-1}$ ”, 110                                                                    |     |
| $\kappa(\mathbf{A}) \triangleq$ spectral condition number, 662                                                                        |     |
| $\lambda_{\mathcal{T}} \triangleq$ Lebesgue constant for Lagrange interpolation on node set $\mathcal{T}$ , 368                       |     |
| $\lambda_{\max} \triangleq$ largest eigenvalue (in modulus), 662                                                                      |     |
| $\lambda_{\min} \triangleq$ smallest eigenvalue (in modulus), 662                                                                     |     |
| $\mathbf{1} = [1, \dots, ]^\top$ , 739                                                                                                |     |
| $\text{Ncut}(\mathcal{X}) \triangleq$ normalized cut of subset of weighted graph, 609                                                 |     |
| $\text{argmin} \triangleq$ (global) minimizer of a functional, 656                                                                    |     |
| $\text{cond}(\mathbf{A})$ , 114                                                                                                       |     |
| $\text{cut}(\mathcal{X}) \triangleq$ cut of subset of weighted graph, 609                                                             |     |
| $\text{diag}(d_1, \dots, d_n) \triangleq$ square diagonal matrix, 40                                                                  |     |
| $\text{dist}_{\ \cdot\ }(x, V) \triangleq$ distance of an element of a normed vector spcace from set $V$ , 413                        |     |
| $\text{env}(\mathbf{A})$ , 177                                                                                                        |     |
| $\text{lsq}(\mathbf{A}, \mathbf{b}) \triangleq$ set of least squares solutions of $\mathbf{ax} = \mathbf{b}$ , 196                    |     |
| $\text{nnz}$ , 157                                                                                                                    |     |
| $\text{rank}(\mathbf{A}) \triangleq$ rank of matrix $\mathbf{A}$ , 112                                                                |     |
| $\text{sgn} \triangleq$ sign function, 377                                                                                            |     |
| $\text{vec}(\mathbf{A}) \triangleq$ vectorization of a matrix, 51                                                                     |     |
| $\text{weight}(\mathcal{X}) \triangleq$ connectivity of subset of weighted graph, 609                                                 |     |
| $\overline{\phantom{x}}$ $\triangleq$ complex conjugation, 442                                                                        |     |
| $\overline{m}(\mathbf{A})$ , 177                                                                                                      |     |

- $\rho(\mathbf{A}) \doteq$  spectral radius of  $\mathbf{A} \in \mathbb{K}^{n,n}$ , 591  
 $\rho_{\mathbf{A}}(\mathbf{u}) \doteq$  Rayleigh quotient, 605  
 $\mathbf{f} \doteq$  right hand side of an ODE, 690  
 $\sigma(\mathbf{A}) \doteq$  spectrum of matrix  $\mathbf{A}$ , 591  
 $\sigma(\mathbf{M})$  *hat*= spectrum of matrix  $\mathbf{M}$ , 659  
 $\tilde{x}$ , 83  
 $\mathbb{S}^1 \doteq$  unit circle in the complex plane, 394  
 $m(\mathbf{A})$ , 177  
 $\widehat{\tilde{f}}_j \doteq$   $j$ -th Fourier coefficient of periodic function  $f$ ,  
                           311  
 $f^{(k)} \doteq$   $k$ -th derivative of function  $f : I \subset \mathbb{R} \rightarrow \mathbb{K}$ ,  
                           410  
 $f^{(k)} \doteq$   $k$  derivative of  $f$ , 99  
 $m(\mathbf{A})$ , 177  
 $y[t_i, \dots, t_{i+k}] \doteq$  divided difference, 363  
 $\|\mathbf{x}\|_1$ , 102  
 $\|\mathbf{x}\|_2$ , 102  
 $\|\mathbf{x}\|_\infty$ , 102  
 $\cdot \doteq$  Derivative w.r.t. time  $t$ , 685  
  
 TOL tolerance, 721

# Examples and Remarks

- LU*-decomposition of sparse matrices, 171  
 $L^2$ -error estimates for polynomial interpolation, 422  
 $h$ -adaptive numerical quadrature, 512  
 $p$ -convergence of piecewise polynomial interpolation, 471  
'**auto**' in EIGEN codes, 46  
(Nearly) singular LSE in shifted inverse iteration, 618  
(Relative) point locations from distances, 195  
 $L^2([-1,1])$ -orthogonal polynomials → [Han02, Bsp. 33.2], 448  
BLAS calling conventions, 65  
EIGEN-based code: debug mode and release mode, 47  
EIGEN in use, 47  
General non-linear systems of equations, 517  
Complex step differentiation [LM67], 100  
'Partial LU-decompositions' of principal minors, 131  
"Annihilating" orthogonal transformations in 2D, 215  
"Behind the scenes" of **MyVector** arithmetic, 33  
"Butcher barriers" for explicit RK-SSM, 716  
"Failure" of adaptive timestepping, 725  
"Fast" matrix multiplication, 70  
"Full-rank condition", 579  
"Low" and "high" frequencies, 292  
"Squeezed" DFT of a periodically truncated signal, 305  
"**auto**" considered harmful, 23  
"**Overfitting**", 404  
 $\mathbf{B} = \mathbf{B}^H$  s.p.d. mit Cholesky-Zerlegung, 593  
L-stable implicit Runge-Kutta methods, 765  
fft  
    Efficiency, 314  
2-norm from eigenvalues, 662  
3-Term recursion for Legendre polynomials, 493  
Analytic solution of homogeneous linear ordinary differential equations, 590  
Convergence of PINVIT, 622  
Convergence of subspace variant of direct power method, 638  
Direct power method, 605  
Eigenvalue computation with Arnoldi process, 650  
Impact of roundoff on Lanczos process, 644  
Lagrange polynomials for uniformly spaced nodes, 347  
Lanczos process for eigenvalue computation, 643  
Page rank algorithm, 597  
Power iteration with Ritz projection, 636  
qr based orthogonalization, 633  
Rayleigh quotient iteration, 619  
Resonances of linear electrical circuits, 587  
Ritz projections onto Krylov space, 640  
Runtimes of `eig`, 595  
Stability of Arnoldi process, 649  
Subspace power iteration with orthogonal projection, 630  
Vibrations of a truss structure, 624  
A data type designed for of interpolation problem, 344  
A function that is not locally Lipschitz continuous, 692  
A posteriori error bound for linearly convergent iteration, 527  
A posteriori termination criterion for linearly convergent iterations, 527  
A posteriori termination criterion for plain CG, 668  
A priori and a posteriori choice of optimal interpolation nodes, 427  
A special quasi-linear system of equations, 560  
Accessing matrix data as a vector, 49  
Adapted Newton method, 544  
Adaptive explicit RK-SSM for scalar linear decay ODE, 734  
Adaptive quadrature in PYTHON, 514  
Adaptive timestepping for mechanical problem, 729  
Adding `EPS` to 1, 85  
Affine invariance of Newton method, 555

- Algorithm for cluster analysis, 259  
 Angles in a triangulation, 194  
 Application of modified Newton methods, 546  
 Approximate computaton of Fourier coefficients, 506  
 Approximation by discrete polynomial fitting, 451  
 Arnoldi process Ritz projection, 646  
 Asymptotic behavior of Lagrange interpolation error, 417  
 Asymptotic complexity of Householder QR-factorization, 223  
 Auxiliary construction for shape preserving quadratic spline interpolation, 390  
 Bad behavior of global polynomial interpolants, 371  
 Banach's fixed point theorem, 532  
 Bernstein approximants, 412  
 Block LU-factorization, 132  
 Block Gaussian elimination, 124  
 Blow-up, 720  
 Blow-up of explicit Euler method, 735  
 Blow-up solutions of vibration equations, 623  
 Bound for *asymptotic* rate of linear convergence, 534  
 Breakdown of associativity, 83  
 Broyden method for a large non-linear system, 577  
 Broyden's quasi-Newton method: convergence, 573  
 Butcher scheme for some explicit RK-SSM, 715  
 Calling BLAS routines from C/C++, 65  
 Cancellation during the computation of relative errors, 90  
 Cancellation in decimal system, 88  
 Cancellation in Gram-Schmidt orthogonalisation, 90  
 Cancellation when evaluating difference quotients, 88  
 Cancellation: roundoff error analysis, 91  
 Cardinal shape preserving quadratic spline, 391  
 CG convergence and spectrum, 672  
 Characteristics of stiff IVPs, 753  
 Chebychev interpolation errors, 432  
 Chebychev interpolation of analytic function, 436  
 Chebychev interpolation of analytic functions, 434  
 Chebychev nodes, 432  
 Chebychev polynomials on arbitrary interval, 431  
 Chebychev representation of built-in functions, 441  
 Chebychev vs equidistant nodes, 431  
 Choice of quadrature weights, 485  
 Class PolyEval, 364  
 Classification from measured data, 251  
 Clenshaw-Curtis quadrature rules, 483  
 Commonly used embedded explicit Runge-Kutta methods, 729  
 Communicating special properties of system matrices in EIGEN, 147  
 Composite quadrature and piecewise polynomial interpolation, 501  
 Composite quadrature rules vs. global quadrature rules, 503  
 Computation of nullspace and image space of matrices, 241  
 Computational effort for eigenvalue computations, 595  
 Computing Gauss nodes and weights, 494  
 Computing the zeros of a quadratic polynomial, 86  
 Conditioning and relative error, 143  
 Conditioning of conventional row transformations, 222  
 Conditioning of normal equations, 207  
 Conditioning of the extended normal equations, 209  
 Connection with linear least squares problems Chapter 3, 444  
 Consistency of implicit midpoint method, 703  
 Consistent right hand side vectors are highly improbable, 196  
 Constitutive relations from measurements, 340  
 Construction of higher order Runge-Kutta single step methods, 716  
 Contiguous arrays in C++, 25  
 Convergence monitors, 574  
 Convergence of CG as iterative solver, 669  
 Convergence of Fourier sums, 310  
 Convergence of global quadrature rules, 497  
 Convergence of gradient method, 660  
 Convergence of Hermite interpolation, 472  
 Convergence of Hermite interpolation with exact slopes, 471  
 Convergence of inexact simple splitting methods, 771  
 Convergence of Krylov subspace methods for non-symmetric system matrix, 682  
 Convergence of naive semi-implicit Radau method, 768  
 Convergence of Newton's method in 2D, 564  
 Convergence of quadratic inverse interpolation,

- 550  
 Convergence of Remez algorithm, 456  
 Convergence of secant method, 548  
 Convergence of simple Runge-Kutta methods, 713  
 Convergence of simple splitting methods, 770  
 Convergence rates for CG method, 672  
 Convergence theory for PCG, 676  
 Convex least squares functional, 203  
 Convolution of sequences, 277  
 Cosine transforms for compression, 331  
 Damped Broyden method, 575  
 Damped Newton method, 570  
 Data points confined to a subspace, 257  
 Deblurring by DFT, 301  
 Decay conditions for bi-infinite signals, 309  
 Decimal floating point numbers, 80  
 Derivative of Euclidean norm, 558  
 Derivative of a bilinear form, 557  
 Derivative of matrix inversion, 561  
 Detecting linear convergence, 521  
 Detecting order of convergence, 523  
 Detecting periodicity in data, 290  
 Determining the domain of analyticity, 426  
 Determining the type of convergence in numerical experiments, 418  
 Diagonalization of local translation invariant linear grid operators, 325  
 diagonally dominant matrices from nodal analysis, 184  
 Different choices for consistent fixed point iterations (II), 531  
 Different choices for consistent iteration functions (III), 535  
 Different meanings of “convergence”, 418  
 Discrete evolutions for non-autonomous ODEs, 702  
 Discretization, 702  
 Distribution of machine numbers, 80  
 Divided differences and derivatives, 365  
 Efficiency of `fft`, 314  
 Efficiency of `fft` for different backend implementations, 322  
 Efficiency of FFT-based solver, 328  
 Efficiency of iterative methods, 552  
 Efficiency of `fft`, 314  
 Efficiency of `fft` for different backend implementations, 322  
 Efficient associative matrix multiplication, 71  
 Efficient evaluation of trigonometric interpolation polynomials, 397  
 Eigenvectors of circulant matrices, 282  
 Eigenvectors of commuting matrices, 283  
 Embedded Runge-Kutta methods, 728  
 Empiric Convergence of collocation single step methods, 757  
 Empiric convergence of equidistant trapezoidal rule, 503  
 Envelope of a matrix, 177  
 Envelope oriented matrix storage, 181  
 Error of polynomial interpolation, 422  
 Error representation for generalized Lagrangian interpolation, 421  
 Estimation of “wrong quadrature error”? , 512  
 Estimation of “wrong” error?, 723  
 Euler methods for stiff decay IVP, 753  
 Evolution operator for Lotka-Volterra ODE, 694  
 Ex. 2.3.3.4 cnt'd, 137  
 Explicit adaptive RK-SSM for stiff IVP, 733  
 Explicit Euler method as difference scheme, 698  
 Explicit Euler method for damped oscillations, 745  
 Explicit representation of error of polynomial interpolation, 421  
 Explicit trapzoidal rule for decay equation, 738  
 Exploiting trigonometric identities to avoid cancellation, 93  
 Extremal numbers in  $\mathbb{M}$ , 80  
 Failure of damped Newton method, 571  
 Failure of Krylov iterative solvers, 681  
 Fast Toeplitz solvers, 336  
 Feasibility of implicit Euler timestepping, 700  
 FFT algorithm by matrix factorization, 318  
 FFT based on general factorization, 320  
 FFT for prime vector length, 320  
 Filtering in Fourier domain, 311  
 Finite-time blow-up, 693  
 Fit of hyperplanes, 245  
 Fixed points in 1D, 531  
 Fractional order of convergence of secant method, 548  
 Frequency filtering by DFT, 293  
 Frequency identification with DFT, 289  
 From higher order ODEs to first order systems, 691  
 Full-rank condition, 202  
 Gain through adaptivity, 724  
 Gauss-Newton versus Newton, 582  
 Gauss-Radau collocation SSM for stiff IVP, 766

- Gaussian elimination, 117
- Gaussian elimination and LU-factorization, 126
- Gaussian elimination for non-square matrices, 122
- Gaussian elimination via rank-1 modifications, 123
- Gaussian elimination with pivoting for  $3 \times 3$ -matrix, 134
- Generalized bisection methods, 538
- Generalized eigenvalue problems and Cholesky factorization, 593
- Generalized Lagrange polynomials for Hermite Interpolation, 350
- Generalized polynomial interpolation, 349
- Gibbs phenomenon, 459
- Gradient method in 2D, 659
- Gram-Schmidt orthogonalization of polynomials, 489
- Gram-Schmidt orthonormalization based on **MyVector** implementation, 34
- Group property of autonomous evolutions, 695
- Growth with limited resources, 686
- Halley's iteration, 542
- Heartbeat model, 687
- Heating production in electrical circuits, 477
- Hesse matrix of least squares functional, 202
- Hidden summation, 72
- Horner scheme, 346
- Image compression, 250
- Image segmentation, 607
- Impact of choice of norm, 520
- Impact of matrix data access patterns on runtime, 52
- Impact of roundoff errors on CG, 668
- Implicit differentiation of  $F$ , 541
- Implicit nature of collocation single step methods, 756
- Implicit RK-SSMs for stiff IVP, 763
- Importance of numerical quadrature, 476
- In-situ LU-decomposition, 130
- Inequalities between vector norms, 102
- Initial guess for power iteration, 607
- Initialization of sparse matrices in Eigen, 160
- Inner products on spaces  $\mathcal{P}_m$  of polynomials, 446
- Input errors and roundoff errors, 82
- Instability of multiplication with inverse, 145
- interpolation
  - piecewise cubic monotonicity preserving, 378
  - shape preserving quadratic spline, 391
- Interpolation and approximation: enabling technologies, 409
- Interpolation error estimates and the Lebesgue constant, 423
- Interpolation error: trigonometric interpolation, 458
- Interpolation of vector-valued data, 339
- Intersection of lines in 2D, 115
- Justification of Ritz projection by min-max theorem, 635
- Kinetics of chemical reactions, 748
- Krylov methods for complex s.p.d. system matrices, 655
- Krylov subspace methods for generalized EVP, 651
- Least squares data fitting, 578
- Lebesgue constant for equidistant nodes, 368
- Linear filtering of periodic signals, 277
- Linear parameter estimation = linear data fitting, 401
- Linear parameter estimation in 1D, 193
- linear regression, 196
- Linear regression for stationary Markov chains, 332
- Linear systems with arrow matrices, 151
- Linearly convergent iteration, 521
- Local approximation by piecewise polynomials, 467
- Local convergence of Newton's method, 568
- local convergence of the secant method, 549
- Loss of sparsity when forming normal equations, 208
- LSE: key components of mathematical models in many fields, 109
- LU-decomposition of flipped "arrow matrix", 173
- Machine precision for IEEE standard, 84
- Magnetization curves, 370
- Many choices for consistent fixed point iterations, 529
- Many sequential solutions of LSE, 149
- Mathematical functions in a numerical code, 341
- Matrix inversion by means of Newton's method, 562
- Matrix norm associated with  $\infty$ -norm and 1-norm, 103
- Matrix representation of interpolation operator, 349
- Meaning of full-rank condition for linear models, 202

- Meaningful “ $O$ -bounds” for complexity, 68  
 Measuring the angles of a triangle, 194  
 Midpoint rule, 482  
 Min-max theorem, 613  
 Minimality property of Broyden’s rank-1-modification, 573  
 Model reduction in circuit simulation, 408  
 Monitoring convergence for Broyden’s quasi-Newton method, 574  
 Monomial representation, 346  
 Multi-dimensional data interpolation, 339  
 Multidimensional fixed point iteration, 534  
 Multiplication of Kronecker product with vector, 74  
 Multiplication of polynomials, 275, 276  
 Multiplying matrices in EIGEN, 61  
 Multiplying triangular matrices, 58  
 Necessary condition for L-stability, 764  
 Necessity of iterative approximation, 518  
 Newton method and minimization of quadratic functional, 580  
 Newton method in 1D, 540  
 Newton’s iteration; computational effort and termination, 566  
 Newton-Cotes formulas, 482  
 Nodal analysis of linear electric circuit, 109  
 Non-linear cubic Hermite interpolation, 379  
 Non-linear data fitting (II), 582  
 Non-linear electric circuit, 516  
 Normal equations for some examples from Section 3.0.1, 200  
 Normal equations from gradient, 200  
 Notation for single step methods, 703  
 Numerical Differentiation for computation of Jacobian, 564  
 Numerical stability and sensitive dependence on data, 105  
 Numerical summation of Fourier series, 309  
 NumPy command reshape, 52  
 Orders of simple polynomial quadrature formulas, 486  
 Oregonator reaction, 719  
 Origin of the term “Spline”, 385  
 Oscillating polynomial interpolant, 366  
 Output of explicit Euler method, 698  
 Overflow and underflow, 85  
 Parameter estimation for a linear model, 193  
 Parameter identification for linear time-invariant filters, 331  
 PCA for data classification, 255  
 PCA of stock prices, 254  
 Piecewise cubic Hermite interpolation, 376  
 Piecewise cubic interpolation schemes, 384  
 Piecewise linear interpolation, 342  
 Piecewise polynomial interpolation, 469  
 Piecewise quadratic interpolation, 373  
 Pivoting and numerical stability, 133  
 Pivoting destroys sparsity, 175  
 Polybomial interpolation vs. polynomial fitting, 404  
 Polynomial best approximation on general intervals, 416  
 Polynomial fitting, 403  
 Power iteration, 604  
 Predator-prey model, 686  
 Predicting stiffness of non-linear IVPs, 752  
 Principal axis of a point cloud, 259  
 Pseudoinverse and SVD, 244  
 QR-Algorithm, 593  
 QR-based solution of banded LSE, 227  
 QR-based solution of linear systems of equations, 226  
 QR-decomposition of “fat” matrices, 217  
 QR-decomposition of banded matrices, 221  
 Quadratic convergence, 524  
 Quadratic functional in 2D, 656  
 Quadratur  
     Gauss-Legendre Ordnung 4, 487  
 Quadrature errors for composite quadrature rules, 502  
 Radiative heat transfer, 278  
 Rank defect in linear least squares problems, 202  
 Rationale for adaptive quadrature, 507  
 Rationale for high-order single step methods, 711  
 Rationale for partial pivoting policy, 136  
 Rationale for using LU-decomposition in algorithms, 131  
 Recursive LU-factorization, 130  
 Reducing fill-in by reordering, 182  
 Reducing bandwidth by row/column permutations, 182  
 Reduction to finite signals, 275  
 Reduction to periodic convolution, 281  
 Refined local stepsize control, 726  
 Regions of stability for simple implicit RK-SSM, 761  
 Regions of stability of some explicit RK-SSM, 747  
 Relative error and number of correct digits, 83  
 Relevance of asymptotic complexity, 69  
 Removing a singularity by transformation, 497

- Reshaping matrices in EIGEN, 51  
 Residual based termination of Newton's method, 567  
 Resistance to currents map, 156  
 Restarted GMRES, 681  
 RK-SSM and quadrature rules, 715  
 Roundoff effects in normal equations, 208  
 Row and column transformations, 60  
 Row swapping commutes with forward elimination, 138  
 Row-wise & column-wise view of matrix product, 55  
 Runge's example, 419, 422  
 Runtime comparison for computation of coefficient of trigonometric interpolation polynomials, 397  
 Runtime of Gaussian elimination, 121  
 Runtimes of elementary linear algebra operations in EIGEN, 70  
 S.p.d. Hessians, 41  
 Sacrificing numerical stability for efficiency, 153  
 Scaling a matrix, 58  
 Semi-implicit Euler single-step method, 767  
 Sensitivity of linear mappings, 113  
 Shape preservation of cubic spline interpolation, 386  
 Shifted inverse iteration, 618  
 Significance of smoothness of interpoland, 422  
 Simple adaptive stepsize control, 724  
 Simple adaptive timestepping for fast decay, 737  
 Simple composite polynomial quadrature rules, 499  
 Simple preconditioners, 676  
 Simple Runge-Kutta methods by quadrature & bootstrapping, 712  
 Simplified Newton method, 563  
 Sine transform via DFT of half length, 324  
 Small residuals by Gaussian elimination, 144  
 Smoothing of a triangulation, 161  
 Solving the increment equations for implicit RK-SSMs, 760  
 Sound filtering by DFT, 294  
 Sparse *LU*-factors, 172  
 Sparse elimination for arrow matrix, 167  
 Sparse LSE in circuit modelling, 157  
 Sparse matrices from the discretization of linear partial differential equations, 157  
 Special cases in IEEE standard, 81  
 Spectrum of Fourier matrix, 286  
 Speed of convergence of Euler methods, 705  
 spline  
 interpolants, approx. complete cubic, 473  
 shape preserving quadratic interpolation, 391  
 Splitting linear and local terms, 773  
 Splitting off stiff components, 772  
 Square root iteration as Newton's method, 540  
 Square root of a s.p.d. matrix, 673  
 Stability by small random perturbations, 142  
 Stability function and exponential function, 740  
 Stability functions of explicit Runge-Kutta single step methods, 739  
 Stable discriminant formula, 92  
 Stable implementation of Householder reflections, 217  
 Stable orthonormalization by QR-decomposition, 78  
 Stable solution of LSE by means of QR-decomposition, 228  
 Stage form equations for increments, 759  
 Standard EIGEN `lu()` operator versus triangularView(), 148  
 Stepsize control detects instability, 741  
 Stepsize control in `ode45`, 728  
 Storing orthogonal transformations, 220  
 Strongly attractive limit cycle, 749  
 Subspace power methods, 638  
 Summation of exponential series, 97  
 SVD and additive rank-1 decomposition, 238  
 SVD-based computation of the rank of a matrix, 240  
 Switching to equivalent formulas to avoid cancellation, 93  
 Tables of quadrature rules, 479  
 Tangent field and solution curves, 696  
 Taylor approximation, 410  
 Termination criterion for contractive fixed point iteration, 536  
 Termination criterion for direct power iteration, 607  
 Termination criterion in `pcg`, 679  
 Termination of PCG, 678  
 Testing equality with zero, 85  
 Testing stability of matrix  $\times$  vector multiplication, 105  
 The “matrix  $\times$  vector-multiplication problem”, 101  
 The inverse matrix and solution of a LSE, 112  
 The message of asymptotic estimates, 498  
 The zoo of sparse matrix formats, 158  
 Timing polynomial evaluations, 356  
 Timing sparse elimination for the combinatorial graph Laplacian, 168

Trade cancellation for approximation, 98  
Transformation of polynomial approximation schemes, 414  
Transformation of quadrature rules, 478  
Transforming approximation error estimates, 415  
Transient circuit simulation, 688  
Transient simulation of RLC-circuit, 742  
Trend analysis, 251  
Tridiagonal preconditioning, 677  
Trigonometric interpolation of analytic functions, 463  
Two-dimensional DFT in MATLAB, 299  
Understanding the structure of product matrices, 56  
Uniqueness of SVD, 238  
Unitary similarity transformation to tridiagonal form, 594  
Unstable Gram-Schmidt orthonormalization, 77  
Using Intel Math Kernel Library (Intel MKL) from EIGEN, 66  
Vandermonde matrix, 348  
Vectorisation of a matrix, 51  
Visualization of explicit Euler method, 697  
Visualization: superposition of impulse responses, 274  
Weak locality of the natural cubic spline interpolation, 387  
Why using  $\mathbb{K} = \mathbb{C}$ ?, 284  
Wilkinson's counterexample, 140

# Abbreviations and Acronyms

CHI  $\triangleq$  cubic Hermite interpolation, 374

CSI  $\triangleq$  cubic spline interpolation, 381

FRC  $\triangleq$  full-rank condition, 202

GE  $\triangleq$  Gaussian elimination, 117

IRK-SSM  $\triangleq$  implicit Runge-Kutta single-step method, 760

IVP  $\triangleq$  initial-value problem, 684

KCL  $\triangleq$  Kirchhoff Current Law, 110

LPI  $\triangleq$  Lagrangian (interpolation polynomial) approximation scheme, 417

LSE  $\triangleq$  Linear System of Equations, 109

NEQ  $\triangleq$  normal equations, 198

NMT  $\triangleq$  natural monotonicity test, 569

ODE  $\triangleq$  ordinary differential equation, 684

POD  $\triangleq$  principal orthogonal decomposition, 257

PP  $\triangleq$  piecewise polynomial, 468

PPLIP  $\triangleq$  piecewise polynomial Lagrange interpolation, 467

QR  $\triangleq$  quadrature rule, 478

SF  $\triangleq$  stability function, 739

SSM  $\triangleq$  single-step method, 702

# German terms

$L^2$ -inner product =  $L^2$ -Skalarprodukt, 488  
back substitution = Rücksubstitution, 118  
bandwidth = Bandbreite, 177  
  
cancellation = Auslöschung, 86  
capacitance = Kapazität, 110  
capacitor = Kondensator, 110  
circulant = zirkulant, 280  
circulant matrix = zirkulante Matrix, 280  
coil = Spule, 110  
column (of a matrix) = (Matrix)spalte, 39  
column transformation = Spaltenumformung, 60  
column vector = Spaltenvektor, 38  
composite quadrature formulas = zusammenge-  
setzte Quadraturformeln, 499  
computational effort = Rechenaufwand, 67  
consistency = Konsistenz, 519  
constitutive relation = Kennlinie, 340  
constitutive relations = Bauelementgleichungen,  
110  
constrained least squares = Ausgleichsproblem  
mit Nebenbedingungen, 265  
convergence = Konvergenz, 519  
convolution = Faltung, 271  
  
damped Newton method = gedämpftes Newton-  
Verfahren, 567  
deblurring = Entrauschen, 301  
dense matrix = vollbesetzte Matrix, 156  
descent methods = Abstiegsverfahren, 655  
discrete convolution = diskrete Faltung, 276  
divided differences = dividierte Differenzen, 362  
dot product = (Euklidisches) Skalarprodukt, 54  
  
eigenspace = Eigenraum, 591  
eigenvalue = Eigenwert, 591  
electric circuit = elektrischer Schaltkreis/Netzwerk, 109  
energy norm = Energienorm, 655  
envelope = Hülle, 177  
extended normal equations = erweiterte Normalengleichungen, 209  
  
field of numbers = Zahlenkörper, 38

fill-in = "fill-in", 172  
fixed point iteration = Fixpunktiteration, 528  
floating point number = Gleitpunktzahl, 80  
floating point numbers = Gleitkommazahlen, 80  
forward elimination = Vorwärtselimination, 118  
Fourier series = Fourierreihe, 308  
frequency domain = Frequenzbereich, 289  
  
Gaussian elimination = Gaußelimination, 117  
  
Hessian = Hesse-Matrix, 41  
high pass filter = Hochpass, 293  
  
identity matrix = Einheitsmatrix, 40  
image segmentation = Bildsegmentierung, 607  
impulse response = Impulsantwort, 271  
in situ = am Ort, 130  
in situ = an Ort und Stelle, 124  
initial guess = Anfangsnäherung, 519  
initial value problem (IVP) = Anfangswertproblem,  
690  
intermediate value theorem = Zwischenwertsatz,  
537  
inverse iteration = inverse Iteration, 607  
  
Kirchhoff (current) law = Kirchhoffsche Knoten-  
regel, 110  
knot = Knoten, 380  
Krylov space = Krylovraum, 664  
  
least squares = Methode der kleinsten Quadrate,  
192  
line search = Minimierung in eine Richtung, 656  
linear system of equations = lineares Glei-  
chungssystem, 109  
low pass filter = Tiefpassfilter, 293  
lower triangular matrix = untere Dreiecksmatrix,  
40  
LU-factorization = LR-Zerlegung, 125  
  
machine number = Maschinenzahl, 80  
machine numbers = Maschinenzahlen, 79  
mass matrix = Massenmatrix, 626  
matrix factorization = Matrixzerlegung, 125  
mesh width = Gitterweite, 467

- mesh/grid = Gitter, 467  
 multiplicity= Vielfachheit, 591  
 nested = geschachtelt, 410  
 nodal analysis = Knotenanalyse, 109  
 normal equations = Normalengleichungen, 198  
 order of convergence = Konvergenzordnung, 523  
 ordinary differential equation = gewöhnliche Differentialgleichung, 690  
 overdetermined = überbestimmt, 192  
 overflow = Überlauf, 85  
 partial pivoting = Spaltenpivotsuche, 135  
 pattern (of a matrix) = Besetzungsmuster, 58  
 power method = Potenzmethode, 597  
 preconditioning = Vorkonditionierung, 673  
 predator-prey model = Räuber-Beute-Modell, 686  
 principal axis transformation = Hauptachsentransformation, 592  
 principal component analysis = Hauptkomponentenanalyse, 251  
 quadratic functional = quadratisches Funktional, 655  
 quadrature node = Quadraturknoten, 478  
 quadrature weight = Quadraturgewicht, 478  
 resistor = Widerstand, 110  
 right hand side vector = rechte-Seite-Vektor, 109  
 rounding = Rundung, 83  
 roundoff error = Rundungsfehler, 79  
 row (of a matrix) = (Matrix)zeile, 39  
 row transformation = Zeilenumformung, 60  
 row vector = Zeilenvektor, 38  
 saddle point problem = Sattelpunktproblem, 265  
 Sampling = Abtasten, 270  
 scaling = Skalierung, 58  
 single step method = Einschrittverfahren, 702  
 singular value decomposition = Singlärwertzerlegung, 235  
 sparse matrix = dünnbesetzte Matrix, 156  
 speed of convergence = Konvergenzgeschwindigkeit, 520  
 State space = Zustandsraum, 690  
 steepest descent = steilster Abstieg, 656  
 stiffness matrix = Steifigkeitsmatrix, 626  
 storage format = Speicherformat, 49  
 subspace correction = Unterraumkorrektur, 663  
 system matrix = Koeffizientenmatrix, 109  
 tensor product = Tensorprodukt, 54  
 tent/hat function = Hutfunktion, 342
- termination criterion = Abbruchbedingung, 524  
 time domain = Zeitbereich, 289  
 timestepping = Zeitdiskretisierung, 697  
 total least squares = totales Ausgleichsproblem, 263  
 transpose = transponieren/Transponierte, 38  
 truss = Stabwerk, 624  
 underflow = Unterlauf, 85  
 unit vector = Einheitsvektor, 39  
 upper triangular matrix = obere Dreiecksmatrix, 40  
 variational calculus = Variationsrechnung, 385  
 zero padding = Ergänzen durch Null, 398