

Numerical Methods

for Computational Science and Engineering

(“Numerische Methoden für CSE”, 401-0663-00 V)

Prof. Ralf Hiptmair and Dr. Vasile Gradinaru

Draft version November 5, 2011, Subversion 38355

(C) Seminar für Angewandte Mathematik, ETH Zürich

<http://www.sam.math.ethz.ch/~hiptmair/tmp/NumCSE11.pdf>

Always under construction



Contents

| | | |
|----------|---|-----------|
| 1 | Computing with Matrices and Vectors [48, Sect. 2.1&2.2], [27, Sect. 2.1&2.2] | 28 |
| 1.1 | Notations | 29 |
| 1.1.1 | Vectors | 29 |
| 1.1.2 | Matrices | 31 |
| 1.2 | Elementary operations | 40 |
| 1.3 | Complexity/computational effort | 54 |
| 1.4 | BLAS | 63 |
| 2 | Direct Methods for Linear Systems of Equations | 85 |
| 2.1 | Gaussian Elimination | 89 |
| 2.2 | LU-Decomposition/LU-Factorization | 106 |
| 2.3 | Pivoting | 124 |
| 2.4 | Supplement: Machine Arithmetic | 138 |

| | | |
|---------|--|-----|
| 2.5 | Stability of Gaussian Elimination | 148 |
| 2.5.1 | Vector norms and matrix norms [13, Sect. 2.1.2], [35, Sect. 1.2], [51, Sect. 1.11] | 148 |
| 2.5.2 | Numerical Stability [13, Sect. 2.3] | 152 |
| 2.5.3 | Roundoff analysis of Gaussian elimination | 155 |
| 2.5.4 | Conditioning | 166 |
| 2.5.5 | Sensitivity of linear systems | 172 |
| 2.6 | Sparse Matrices | 175 |
| 2.6.1 | Sparse matrix storage formats | 181 |
| 2.6.2 | Sparse matrices in MATLAB | 183 |
| 2.6.3 | LU-factorization of sparse matrices | 195 |
| 2.6.4 | Banded matrices [13, Sect. 3.7] | 209 |
| 2.7 | Stable Gaussian elimination without pivoting | 223 |
| 2.8 | QR-factorization/decomposition [35, Sect. 13], [27, Sect. 7.3] | 240 |
| 2.9 | Modification Techniques | 267 |
| 2.9.0.1 | Rank-1-modifications | 270 |
| 2.9.0.2 | Adding a column | 281 |
| 2.9.0.3 | Adding a row | 286 |

| | | | |
|----------|---|------------|-----------------------------------|
| 3 | Data Interpolation in 1D | 289 | |
| 3.1 | Abstract interpolation | 290 | NumCSE, autumn 2010 |
| 3.2 | Polynomials | 300 | |
| 3.3 | Polynomial Interpolation: Theory | 302 | |
| 3.4 | Polynomial Interpolation: Algorithms | 310 | |
| 3.4.1 | Multiple evaluations | 311 | |
| 3.4.2 | Single evaluation [13, Sect. 8.2.2] | 314 | |
| 3.4.3 | Extrapolation to zero | 319 | |
| 3.4.4 | Newton basis and divided differences [13, Sect. 8.2.4], [51, Sect. 8.2] | 329 | |
| 3.5 | Polynomial Interpolation: Sensitivity [51, Sect. 8.1.3] | 337 | |
| 3.6 | Shape preserving interpolation | 345 | R. Hiptmair |
| 3.6.1 | Piecewise linear interpolation | 351 | rev 38355, November 5, 2011 |
| 3.7 | Cubic Hermite Interpolation | 355 | |
| 3.7.1 | Definition and algorithms | 356 | |
| 3.7.2 | Shape preserving Hermite interpolation | 361 | |
| 3.8 | Splines [13, Ch. 9] | 367 | |
| 3.8.1 | Cubic spline interpolation [35, XIII, 46], [51, Sect. 8.6.1] | 369 | |
| 3.8.2 | Shape Preserving Spline Interpolation | 385 | |

| | | | |
|----------|--|------------|-----------------------------------|
| 4 | Iterative Methods for Non-Linear Systems of Equations | 401 | |
| 4.1 | Iterative methods | 404 | NumCSE, autumn 2010 |
| 4.1.1 | Speed of convergence | 409 | |
| 4.1.2 | Termination criteria | 418 | |
| 4.2 | Fixed Point Iterations [13, Sect. 5.3], [51, Sect. 6.3] | 425 | |
| 4.2.1 | Consistent fixed point iterations | 426 | |
| 4.2.2 | Convergence of fixed point iterations | 430 | |
| 4.3 | Zero Finding | 443 | |
| 4.3.1 | Bisection [13, Sect. 5.5.1] | 444 | |
| 4.3.2 | Model function methods | 446 | |
| 4.3.2.1 | Newton method in scalar case [35, Sect. 18.1], [13, Sect. 5.5.2] | 447 | R. Hiptmair |
| 4.3.2.2 | Special one-point methods | 449 | rev 38355, November 5, 2011 |
| 4.3.2.3 | Multi-point methods | 458 | |
| 4.3.3 | Note on Efficiency | 466 | |
| 4.4 | Newton's Method [35, Sect. 19], [13, Sect. 5.6] | 470 | |
| 4.4.1 | The Newton iteration | 470 | |
| 4.4.2 | Convergence of Newton's method | 477 | |
| 4.4.3 | Termination of Newton iteration | 482 | |
| 4.4.4 | Damped Newton method [13, pp. 200] | 485 | 0.0 |
| 4.4.5 | Quasi-Newton Method [51, Sect. 7.1.4] | 492 | p. 5 |

| | | | |
|----------|---|------------|-----------------------------------|
| 5 | Krylov Methods for Linear Systems of Equations | 501 | |
| 5.1 | Descent Methods [51, Sect. 4.3.3] | 503 | NumCSE, autumn 2010 |
| 5.1.1 | Quadratic minimization context | 504 | |
| 5.1.2 | Abstract steepest descent | 507 | |
| 5.1.3 | Gradient method for s.p.d. linear system of equations | 510 | |
| 5.1.4 | Convergence of the gradient method | 513 | |
| 5.2 | Conjugate gradient method (CG) [35, Ch. 9], [13, Sect. 13.4], [51, Sect. 4.3.4] | 522 | |
| 5.2.1 | Krylov spaces | 526 | |
| 5.2.2 | Implementation of CG | 528 | |
| 5.2.3 | Convergence of CG | 536 | |
| 5.3 | Preconditioning [13, Sect. 13.5], [35, Ch. 10], [51, Sect. 4.3.5] | 550 | |
| 5.4 | Survey of Krylov Subspace Methods | 566 | |
| 5.4.1 | Minimal residual methods | 566 | |
| 5.4.2 | Iterations with short recursions [51, Sect. 4.5] | 569 | |
| 6 | Eigenvalues | 574 | R. Hiptmair |
| 6.1 | Theory of eigenvalue problems [48, Ch. 7], [27, Ch. 9], [51, Sect. 1.7] | 583 | rev 38355, November 5, 2011 |
| 6.2 | “Direct” Eigensolvers | 589 | |
| 6.3 | Power Methods | 599 | |
| 6.3.1 | Direct power method [13, Sect. 7.5], [51, Sect. 5.3.1], [51, Sect. 5.3] | 599 | |
| 6.3.2 | Inverse Iteration [13, Sect. 7.6], [51, Sect. 5.3.2] | 623 | |
| 6.3.3 | Preconditioned inverse iteration (PINVIT) | 652 | |
| 6.3.4 | Subspace iterations | 658 | |
| 6.3.4.1 | Orthogonalization | 671 | |
| 6.3.4.2 | Ritz projection | 683 | |
| 6.4 | Krylov Subspace Methods [35, Sect. 30] | 694 | 0.0 |
| 6.5 | Singular Value Decomposition | 719 | p. 6 |

| | | | |
|----------|--|------------|--|
| 7 | Least Squares | 764 | |
| 7.1 | Normal Equations [13, Sect. 4.2], [35, Ch. 11] | 778 | |
| 7.2 | Orthogonal Transformation Methods [13, Sect. 4.4.2] | 784 | |
| 7.3 | Total Least Squares | 801 | |
| 7.4 | Constrained Least Squares | 804 | |
| 7.4.1 | Solution via normal equations | 805 | |
| 7.4.2 | Solution via SVD | 809 | |
| 7.5 | Non-linear Least Squares [13, Ch. 6] | 810 | |
| 7.5.1 | (Damped) Newton method | 812 | |
| 7.5.2 | Gauss-Newton method | 814 | |
| 7.5.3 | Trust region method (Levenberg-Marquardt method) | 822 | |
| 8 | Filtering Algorithms | 824 | |
| 8.1 | Discrete convolutions | 826 | |
| 8.2 | Discrete Fourier Transform (DFT) | 845 | |
| 8.2.1 | Discrete convolution via DFT | 857 | |
| 8.2.2 | Frequency filtering via DFT | 859 | |
| 8.2.3 | Real DFT | 875 | |
| 8.2.4 | Two-dimensional DFT | 877 | |
| 8.2.5 | Semi-discrete Fourier Transform [51, Sect. 10.11] | 885 | |
| 8.3 | Fast Fourier Transform (FFT) [13, Sect. 8.7.3], [35, Sect. 53], [51, Sect. 10.9.2] | 904 | |
| 8.4 | Trigonometric transformations [35, Sect. 55] | 916 | |
| 8.4.1 | Sine transform | 917 | |
| 8.4.2 | Cosine transform | 926 | |
| 8.5 | Toeplitz matrix techniques | 929 | |
| 8.5.1 | Toeplitz matrix arithmetic | 933 | |
| 8.5.2 | The Levinson algorithm | 935 | |

| | | | |
|-----------|--|-------------|-----------------------------------|
| 9 | Approximation of Functions in 1D | 940 | |
| 9.1 | Error estimates for polynomial interpolation | 943 | NumCSE, autumn 2010 |
| 9.2 | Chebyshev Interpolation | 957 | |
| 9.2.1 | Motivation and definition | 957 | |
| 9.2.2 | Chebyshev interpolation error estimates | 966 | |
| 9.2.3 | Chebyshev interpolation: computational aspects | 976 | |
| 9.3 | Trigonometric interpolation [13, Sect. 8.5] | 983 | |
| 9.4 | Approximation by piecewise polynomials | 1001 | |
| 9.4.1 | Piecewise Lagrange interpolation | 1003 | |
| 9.4.2 | Cubic Hermite interpolation: error estimates | 1012 | |
| 9.4.3 | Cubic spline interpolation: error estimates [35, Ch. 47] | 1019 | |
| 10 | Numerical Quadrature [35, VII], [13, Ch. 10] | 1023 | R. Hiptmair |
| 10.1 | Quadrature Formulas | 1026 | rev 38355, November 5, 2011 |
| 10.2 | Polynomial Quadrature Formulas [13, Sect. 10.2] | 1030 | |
| 10.3 | Composite Quadrature | 1037 | |
| 10.4 | Gauss Quadrature [35, Ch. 40-41], [13, Sect.10.3] | 1063 | |
| 10.5 | Adaptive Quadrature | 1085 | |
| 11 | Clustering Techniques | 1093 | |
| 11.1 | Kernel Matrices | 1093 | |
| 11.2 | Local Separable Approximation | 1096 | |
| 11.3 | Cluster Trees | 1109 | 0.0 |
| 11.4 | Algorithm | 1118 | p. 8 |

| | | |
|---|-------------|--|
| 12 Single Step Methods | 1136 | |
| 12.1 Initial value problems (IVP) for ODEs | 1136 | NumCSE, autumn 2010 |
| 12.1.1 Examples | 1137 | |
| 12.1.2 Theory [51, Sect. 11.1], [13, Sect. 11.3] | 1147 | |
| 12.2 Euler methods | 1156 | |
| 12.2.1 Explicit Euler method | 1157 | |
| 12.2.2 Implicit Euler method | 1160 | |
| 12.2.3 Abstract single step methods | 1164 | |
| 12.3 Convergence of single step methods [13, Sect. 11.5] [51, Sect. 11.3] | 1168 | |
| 12.4 Runge-Kutta methods [13, Sect. 11.6], [35, Ch. 76], [51, Sect. 11.8] | 1177 | |
| 12.5 Stepsize control [13, Sect. 11.7], [51, Sect. 11.8.2] | 1189 | |
| 13 Stiff Integrators [13, Sect. 11.9] | 1220 | R. Hiptmair rev 38355, October 21, 2011 |
| 13.1 Model problem analysis [35, Ch. 77], [51, Sect. 11.3.3] | 1222 | |
| 13.2 Stiff problems [51, Sect. 11.10] | 1231 | |
| 13.3 Implicit Runge-Kutta methods [13, Sect. 11.6.2], [51, Sect. 11.8.3] | 1248 | |
| 13.4 Semi-implicit Runge-Kutta methods [35, Ch. 80] | 1256 | |
| 14 Structure Preservation [31] | 1261 | |
| 14.1 Dissipative Evolutions | 1261 | |
| 14.2 Quadratic Invariants | 1261 | |
| 14.3 Reversible Integrators | 1261 | 0.0 |
| 14.4 Symplectic Integrators | 1261 | p. 9 |

Focus

- ▷ on **algorithms** (principles, scope, and limitations),
- ▷ on (efficient, stable) **implementation** in MATLAB,
- ▷ on **numerical experiments** (design and interpretation).

no emphasis on

- theory and proofs (unless essential for understanding of algorithms)
- hardware-related issues (e.g. parallelization, vectorization, memory access)
(These aspect will be covered in the new course “Parallel Computing for Scientific Simulation” offered by D-INFK from autumn term 2012)

Linear systems of
equations

Eigenvalue
problems

Least squares
problems

Interpolation

Quadrature

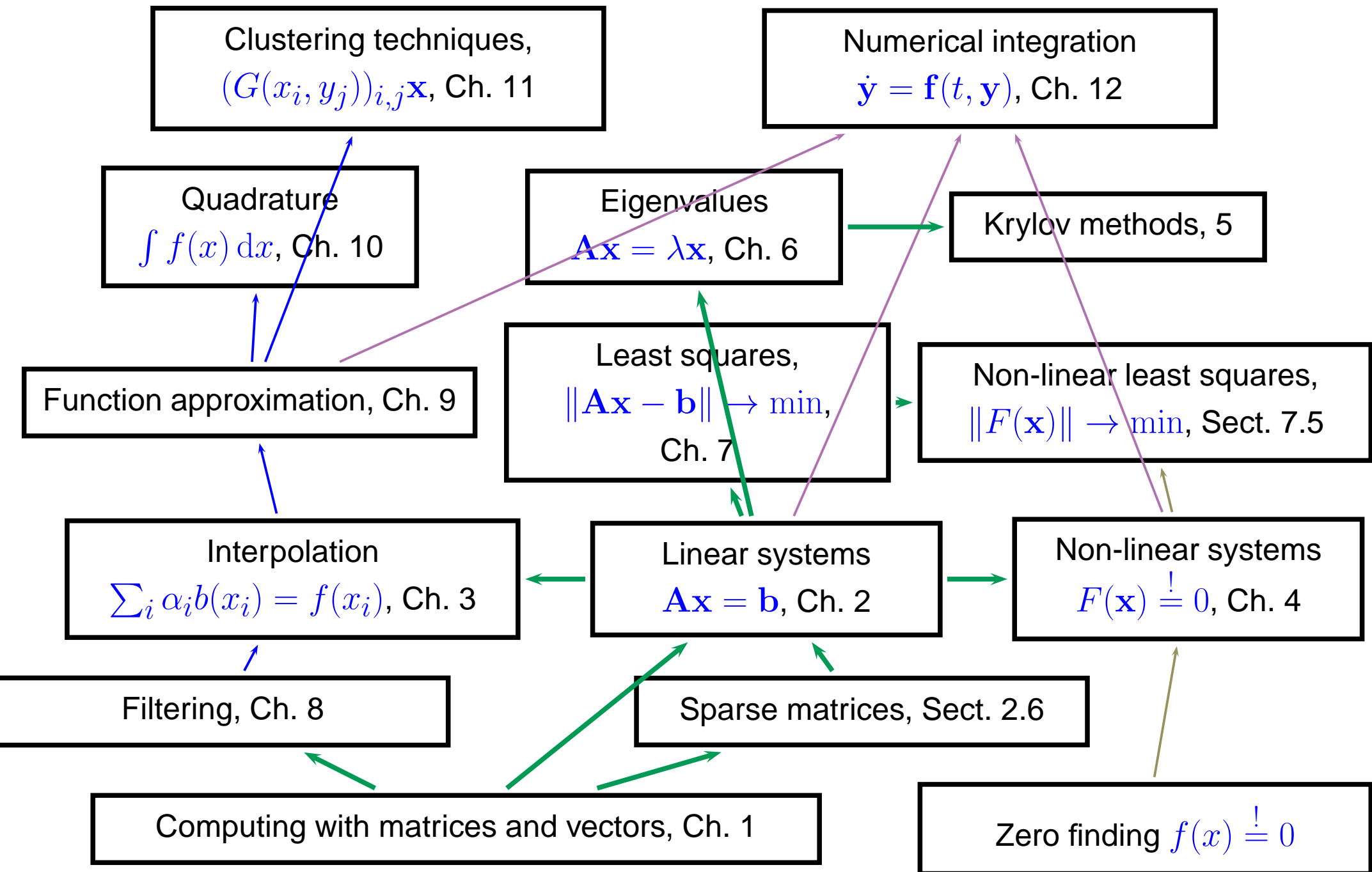
Numerical
integration of ODEs



Analysis

Linear algebra

Programming



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

Goals

- Knowledge of the fundamental algorithms in numerical mathematics
- Knowledge of the essential terms in numerical mathematics and the techniques used for the analysis of numerical algorithms
- Ability to choose the appropriate numerical method for concrete problems
- Ability to interpret numerical results
- Ability to implement numerical algorithms efficiently in MATLAB

Indispensable:

Learning by doing (→ exercises)

Reading instructions

These course materials are neither a textbook nor comprehensive lecture notes.
They are meant to be supplemented by explanations given in class.

Some pieces of advice:

- these lecture slides are not designed to be self-contained, but to supplement explanations in class.
- this document 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,
- study referenced literature to refresh prerequisite knowledge and for alternative presentation of the material (from a different angle, maybe).

What to expect

- The course is difficult and demanding (*ie.* ETH level)
- Do **not** expect to understand everything in class. The average student will
 - understand about one third of the material when attending the lectures,
 - understand another third when making a *serious effort* to solve the homework problems,
 - hopefully understand the remaining third when studying for the examination after the end of the course.

Perseverance will be rewarded!

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.

- [\[\[13\]\]](#) 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)

- [\[\[35\]\]](#) 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)

- [\[\[51\]\]](#) 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)

- [[17]] 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).

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

- [[48]] K. NIPP AND D. STOFFER, *Lineare Algebra*, vdf Hochschulverlag, Zürich, 5 ed., 2002.
- [[27]] M. GUTKNECHT, *Lineare algebra*, lecture notes, SAM, ETH Zürich, 2009.
<http://www.sam.math.ethz.ch/~mhg/unt/LA/HS07/>.
- [[63]] M. STRUWE, *Analysis für Informatiker*. Lecture notes, ETH Zürich, 2009.
<https://moodle-app1.net.ethz.ch/lms/mod/resource/index.php?id=145>.

General information

Lecturer: Prof. Ralf Hiptmair HG G 58.2, ☎ 044 632 3404, hiptmair@sam.math.ethz.ch

Assistants: Jonas Šukys, HG E 62.1, ☎ 044 632 0401, jonas.sukys@sam.math.ethz.ch
Rafaela Guberovic, HG J 47, ☎ 044 632 4320, rafaela.guberovic@sam.math.ethz.ch
Konstantin Grella, HG J 59, ☎ 044 632 6362, konstantin.grella@sam.math.ethz.ch
Robert Carnecky, CAB G 66.2, ☎ 044 632 5527, crobi@inf.ethz.ch
Sebastian Lienert, lienerts@student.ethz.ch
Simon Härdi, shaerdi@student.ethz.ch

Classes: Mon, 10.15-12.00 (HG F3), Thu, 10.15-12.00 (HG G5)

Tutorials: Mon, 13.15-15.00 (LFO C 13, ML H 34.3)

Thu, 8.15-10.00 (HG G 26.3, HG G 26.5, LFW E 13)

Thu, 13.15-15.00 (ML F 40)

Please register (in course website) for tutorial groups until September 25th:

http://www.math.ethz.ch/education/bachelor/lectures/hs2011/math/nummath_cse

Consulting hours: Tue 12-13 Rafaela Guberovic, Jonas Šukys (HG E 62.1)

Wed 12-13 Sebastian Lienert, Simon Härdi (in front of HG J 54)

Fri 12-13 Konstantin Grella, Robert Carnecky (in front of HG J 54)

To reach the desks in front of HG J 54, please take an elevator near Polyterrasse.

Assignments

- The assignment sheets will be uploaded on the course webpage on Monday every week.
- The attempted solutions have to be handed in until the following tutorial class (give them to the assistant during the tutorial class itself or put them into the plexiglass trays in HG G 53/5).
- To earn a “Testat” you are expected to
 - present your solutions of **two** homework problems at the **blackboard** in your tutorial class. In the beginning of the semester you are supposed to sign up for two homework problems. Sign up links are available on the course website. Deadline: September 25th.
 - make a serious effort to solve the exercises and hand in the attempted solutions. In order to obtain the “Testat” at least **80%** of the *core homework problems* have to be solved.

Self-grading

For each exercise sheet the students have to state which (sub)problems they have solved or partially solved by ticking the fields *solved* (100%), *partly solved* (50%) or *not solved* (0%) in a list for each sub-problem. Random checks will be conducted to ensure honest self-grading and cheaters will be

punished. In case a student is caught cheating (i.e., stating as solved exercises that he/she didn't solve) the stricter requirements for the "Testat" will apply to the perpetrator.

However, you may request thorough corrections for up to five individual homework problems during the term. Indicate your requests clearly in the self-grading sheet.

Website:

http://www.math.ethz.ch/education/bachelor/lectures/hs2011/math/nummath_cse

Online homework submission (optional)

You can use the online homework submission tool:

<http://www.math.ethz.ch/~grsam/submit/?VL=02>

- Three-hour written examination involving coding problems to be done at the computer on

TBA

- Dry-run for computer based examination:

TBA, registration via course website

- Pre-exam question session:

TBA

- Subjects of examination:

- All topics, which have been addressed in class or in a homework problem.

- Lecture documents will be available as PDF during the examination. The corresponding final version of the lecture documents will be made available on TBA

- The exam questions will be asked both in German and in English.

Changes in course documentation during HS11

These course documents are being altered all the time in response to new ideas and topics and experiences in class . The latter sometimes suggest significant changes after the material has been taught. This is not completely desirable, but may be better than keeping poorly devised parts.

The following significant changes were made to the lecture material after it had been made available in HS 2011:

- The policy for partial pivoting in Sect. 2.3, Formula (2.3.7) was changed to match that of [48, Sect. 2.5].
- The previous subsection 3.3.2 has been moved and promoted to the new Section 3.5. Only slight alterations have been made in the process.

Reporting errors

Please report errors in the electronic lecture notes via a [wiki page](#) !

(Password: CSE, please choose **EDIT** menu to enter information)

Please supply the following information:

- (sub)section where the error has been found,
- precise location (e.g, after Equation (4), Thm. 2.3.3, etc.). Refrain from giving page numbers,
- brief description of the error.

Alternative (for people not savvy with wikis): E-mail an hiptmair@sam.math.ethz.ch, **Sub-
ject:** NUMCSE

Extra questions for course evaluation

Course number (LV-ID): 401-0663-00

Date of evaluation: TBA

- D1:** I try to do all programming exercises. (HS 10: 2-9-25-39-26, 3.8)
- D2:** The programming exercises help understand the numerical methods. (HS 10: 2-9-33-42-12, 3.6)
- D3:** The programming exercises offer too little benefit for the effort spent on them. (HS 10: 7-16-21-25-29, 3.5)
- D4:** Scope and limitations of methods are properly addressed in the course. (HS 10: 9-12-19-30-19, 3.4)
- D5:** Numerical examples in class provide useful insights and motivation. (HS 10: 13-13-22-25-25, 3.4)
- D6:** There should be more examples presented and discussed in class. (HS 10: 2-15-20-28-31, 3.8)
- D7:** Too much information is crammed onto the lecture slides (HS 10: 2-7-13-16-62, 4.3)
- D8:** The course requires too much prior knowledge in linear algebra (HS 10: 18-24-24-22-13, 2.9)
- D9:** The course requires too much prior knowledge in analysis (HS 10: 25-30-21-16-5, 2.5)
- D10:** My prior knowledge of MATLAB was insufficient for the course (HS 10: 31-29-13-22-5, 2.4)
- D11:** More formal proofs would be desirable (HS 10: 56-19-15-7-4, 1.9)
- D12:** The explanations on the blackboard promote understanding (HS 10: 13-36-25-20-5, 2.7)
- D13:** The codes included in the lecture material convey useful information (HS 10: 5-20-18-25-29, 3.5)
- D14:** The master solutions for exercise problems offer too little guidance. (HS 10: 5-13-22-15-13, 3.2)
- D15:** The relevance of the numerical methods taught in the course is convincingly conveyed. (HS 10: 9-15-27-22-11, 3.1)

16: I would not mind the course being taught in English. (HS 10: 25-14-23-23-13, 2.8)

- Scoring:
- 6: I agree fully
 - 5: I agree to a large extent
 - 4: I agree partly
 - 3: I do not quite agree
 - 2: I disagree
 - 1: I disagree strongly

Evaluation of assistants:

| Assistant | shortcut |
|-----------|----------|
| | |
| | |
| | |
| | |
| | |
| | |

Please enter the shortcut code after the LV-ID in the three separate boxes.

MATLAB

- MATLAB** ("Matrix Laboratory"):
- full fledged high level *programming language* (for numerical algorithms)
 - integrated *programming environment*
 - versatile collection of *numerical libraries*

- in this course used for
- ▷ demonstrating (implementation of) algorithms
 - ▷ numerical experiments
 - ▷ programming assignments

This course assumes *familiarity with MATLAB*
as acquired in the introductory linear algebra courses.

Proficiency in MATLAB through "Learning on the job"
(using the very detailed online help facilities)

```

Terminal
File Edit View Terminal Tabs Help

>> help surf
SURF 3-D colored surface.
SURF(X,Y,Z,C) plots the colored parametric surface defined by
four matrix arguments. The view point is specified by VIEW.
The axis labels are determined by the range of X, Y and Z,
or by the current setting of AXIS. The color scaling is determined
by the range of C, or by the current setting of CAXIS. The scaled
color values are used as indices into the current COLORMAP.
The shading model is set by SHADING.

SURF(X,Y,Z) uses C = Z, so color is proportional to surface height.

SURF(x,y,Z) and SURF(x,y,Z,C), with two vector arguments replacing
the first two matrix arguments, must have length(x) = n and
length(y) = m where [m,n] = size(Z). In this case, the vertices
of the surface patches are the triples (x(j), y(i), Z(i,j)).
Note that x corresponds to the columns of Z and y corresponds to
the rows.

SURF(Z) and SURF(Z,C) use x = 1:n and y = 1:m. In this case,
the height, Z, is a single-valued function, defined over a
geometrically rectangular grid.

SURF(...,'PropertyName',PropertyValue,...) sets the value of the
specified surface property. Multiple property values can be set
with a single statement.

SURF(AX,...) plots into AX instead of GCA.

SURF returns a handle to a surface plot object.

AXIS, CAXIS, COLORMAP, HOLD, SHADING and VIEW set figure, axes, and
surface properties which affect the display of the surface.

Backwards compatibility
SURF('v6',...) creates a surface object instead of a surface plot
object for compatibility with MATLAB 6.5 and earlier.

See also SURFC, SURFL, MESH, SHADING.

>> 

```

Useful links:

[Matlab Online Documentation](#)

[MATLAB guide](#)

[MATLAB Primer](#)

Miscellaneous internet resources:

- [MATLAB classroom resources](#)

-

[Mathematical Methods for Simulation with MATLAB - Course material developed by Professor Gilbert](#)

- [Numerical algorithms](#)

- [NUMAWWW](#)

Computing with Matrices and Vectors

[48, Sect. 2.1&2.2], [27, Sect. 2.1&2.2]

1

1.1 Notations

1.1.1 Vectors

1.1.2 Matrices

Remark 1.1.1 (Matrix storage formats). (for dense/full matrices, *cf.* Sect. 2.6)

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}$$

Row major (C-arrays, Python):

| | | | | | | | | | |
|-------|---|---|---|---|---|---|---|---|---|
| A_arr | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|-------|---|---|---|---|---|---|---|---|---|

Column major (Fortran, MATLAB, OpenGL):

| | | | | | | | | | |
|-------|---|---|---|---|---|---|---|---|---|
| A_arr | 1 | 4 | 7 | 2 | 5 | 8 | 3 | 6 | 9 |
|-------|---|---|---|---|---|---|---|---|---|

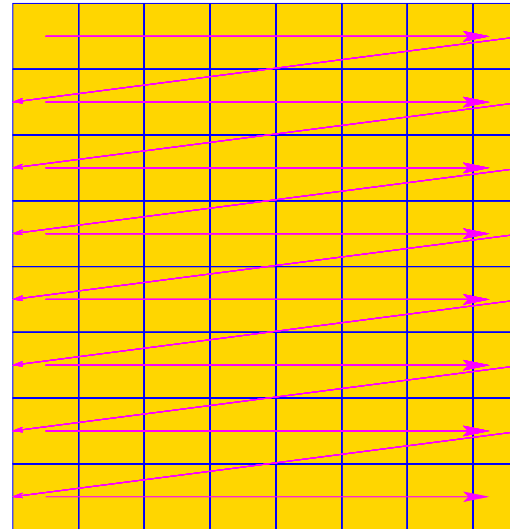
Access to entry a_{ij} of $A \in \mathbb{K}^{n,m}$, $i = 1, \dots, n$,
 $j = 1, \dots, m$:

row major:

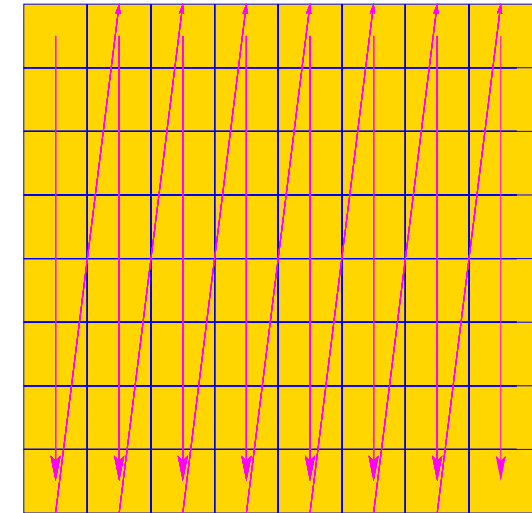
$$a_{ij} \leftrightarrow A_arr(m*(i-1)+(j-1))$$

column major:

$$a_{ij} \leftrightarrow A_arr(n*(j-1)+(i-1))$$



row major



column major



Example 1.1.2 (Impact of data access patterns on runtime).

column oriented

```
A = randn(n,n);
for j = 1:n-1,
    A(:,j+1) = A(:,j+1) - A(:,j);
end
```

access

row oriented

```
A = randn(n,n);
for i = 1:n-1,
    A(i+1,:) = A(i+1,:) - A(i,:);
end
```

access

Code 1.1.3: timing for row and column oriented matrix access in MATLAB

NumCSE,
autumn
2010

```

1 % Timing for row/column operations on matrices
2 % We conduct K runs in order to reduce the risk of skewed measurements
3 % due to OS activity during MATLAB run.
4 K = 3; res = [];
5 for n=2.^(4:13)
6     A = randn(n,n);
7
8     t1 = realmax;
9     for k=1:K, tic;
10        for j = 1:n-1, A(:,j+1) = A(:,j+1) - A(:,j); end;
11        t1 = min(toc,t1);
12    end
13    t2 = realmax;
14    for k=1:K, tic;
15        for i = 1:n-1, A(i+1,:) = A(i+1,:) - A(i,:); end;
16        t2 = min(toc,t2);
17    end
18    res = [res; n, t1 , t2];
19 end
20
21 % Plot runtimes versus matrix sizes
22 figure; plot(res(:,1),res(:,2),'r+', res(:,1),res(:,3),'m*');
23 xlabel('{\bf n}','fontsize',14);

```

R. Hiptmair
rev 38355,
February
24, 2009

```
24 ylabel( '\bf runtime [s]', 'fontsize', 14);
25 legend( 'A(:,j+1) = A(:,j+1) - A(:,j)', 'A(i+1,:) = A(i+1,:) -
      A(i,:)', ...
26         'location', 'northwest' );
27 print -depsc2 '../PICTURES/accessrtlin.eps';
28
29 figure; loglog( res(:,1), res(:,2), 'r+', res(:,1), res(:,3), 'm*' );
30 xlabel( '\bf n', 'fontsize', 14);
31 ylabel( '\bf runtime [s]', 'fontsize', 14);
32 legend( 'A(:,j+1) = A(:,j+1) - A(:,j)', 'A(i+1,:) = A(i+1,:) -
      A(i,:)', ...
33         'location', 'northwest' );
34 print -depsc2 '../PICTURES/accessrtlog.eps';
```

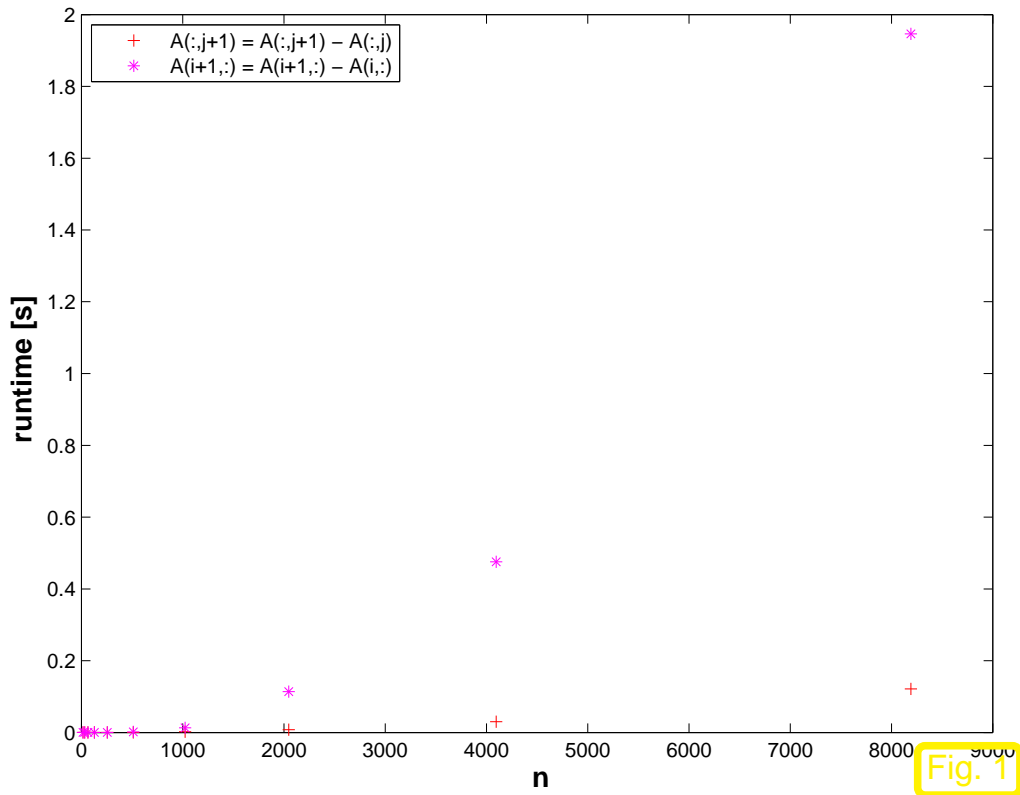


Fig. 1

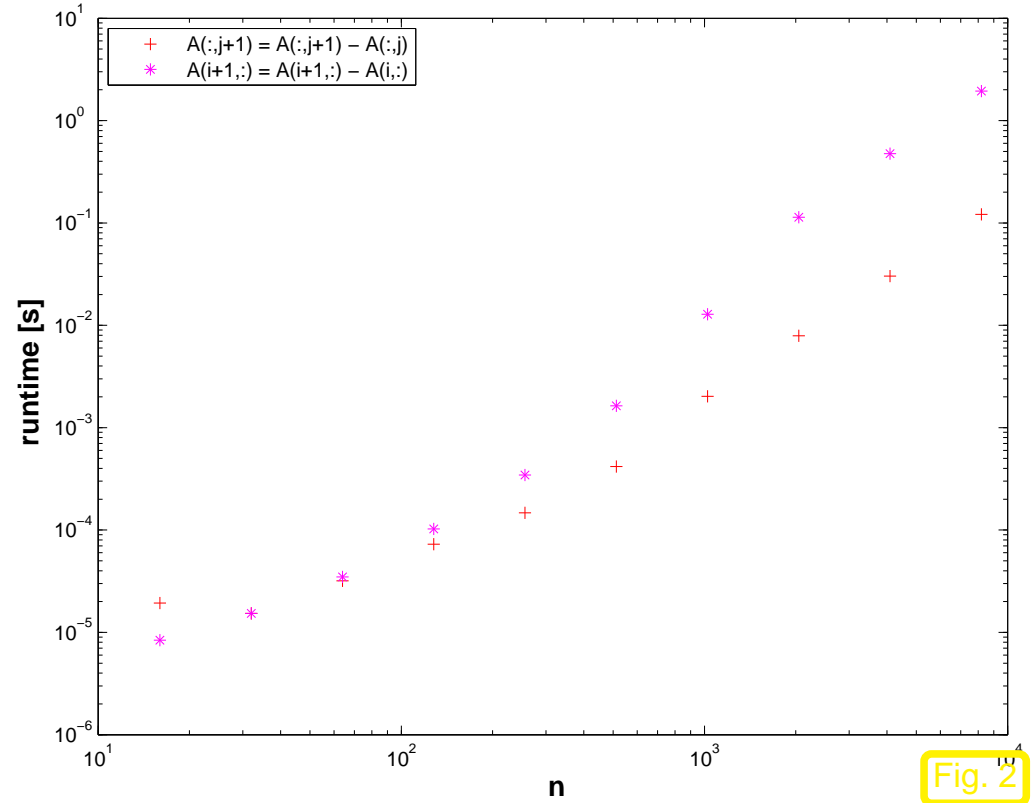


Fig. 2

(Intel Core i7, 2.66 GHz, MacOS X 1.6, MATLAB 7.12.0.635 (R2011a))



1.2 Elementary operations

Remark 1.2.10 (Block matrix product).

Given matrix dimensions $M, N, K \in \mathbb{N}$ block sizes $1 \leq n < N$ ($n' := N - n$), $1 \leq m < M$ ($m' := M - m$), $1 \leq k < K$ ($k' := K - k$) assume

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

▶
$$\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{pmatrix} = \begin{pmatrix} \mathbf{A}_{11}\mathbf{B}_{11} + \mathbf{A}_{12}\mathbf{B}_{21} & \mathbf{A}_{11}\mathbf{B}_{12} + \mathbf{A}_{12}\mathbf{B}_{22} \\ \mathbf{A}_{21}\mathbf{B}_{11} + \mathbf{A}_{22}\mathbf{B}_{21} & \mathbf{A}_{21}\mathbf{B}_{12} + \mathbf{A}_{22}\mathbf{B}_{22} \end{pmatrix} \quad (1.2.11)$$

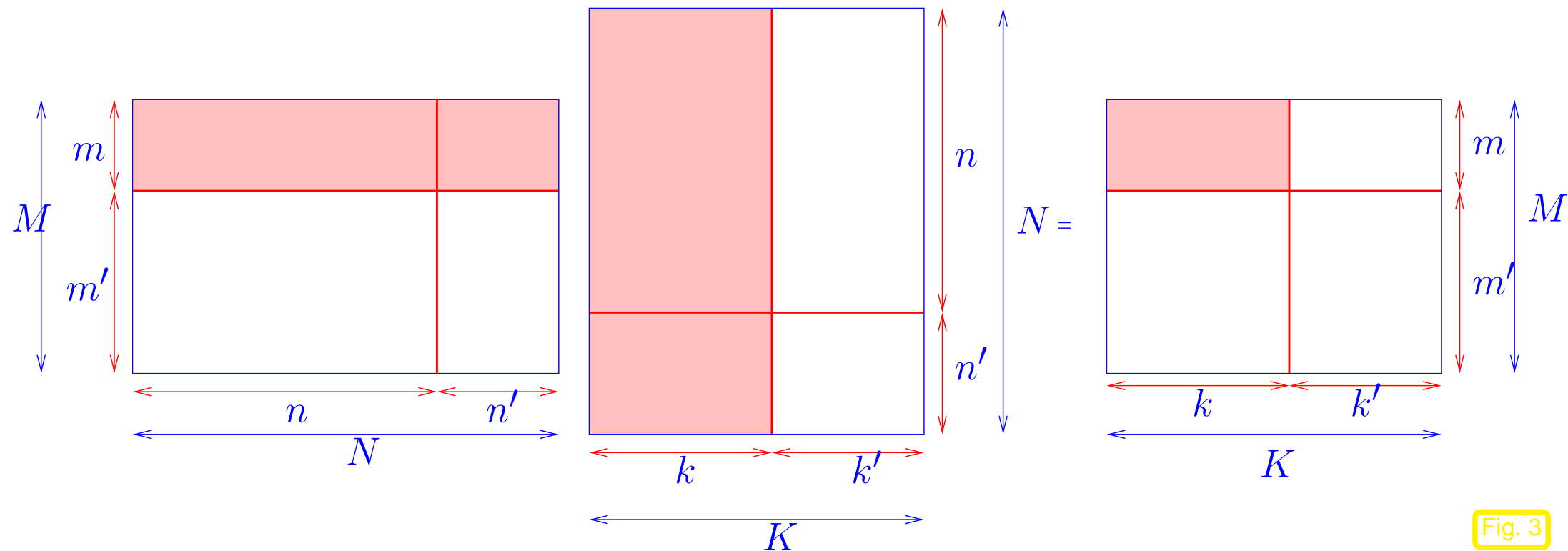


Fig. 3

1.3 Complexity/computational effort

complexity/computational effort of an algorithm $:\Leftrightarrow$ number of elementary operators

additions/multiplications

Crucial: dependence of (worst case) complexity of an algorithm on (integer) **problem size parameters** (worst case \leftrightarrow maximum for all possible data)

Usually studied: **asymptotic complexity** $\hat{=}$ “leading order term” of complexity w.r.t *large* problem size parameters

| 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 | $(\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n) \mapsto \mathbf{xy}^H$ | nm | 0 | $O(mn)$ |
| matrix product ^(*) | $(\mathbf{A} \in \mathbb{R}^{m,n}, \mathbf{B} \in \mathbb{R}^{n,k}) \mapsto \mathbf{AB}$ | mnk | $mk(n - 1)$ | $O(mnk)$ |

notation (“Landau-O”): $f(n) = O(g(n)) \Leftrightarrow \exists C > 0, N > 0: |f(n)| \leq Cg(n)$ for all $n > N$.

Example 1.3.3 (Efficient associative matrix multiplication).

$\mathbf{a} \in \mathbb{K}^m, \mathbf{b} \in \mathbb{K}^n, \mathbf{x} \in \mathbb{K}^n$:

$$\mathbf{y} = (\mathbf{a}\mathbf{b}^\top)\mathbf{x}. \quad (1.3.4)$$

$$\mathbf{T} = \mathbf{a}\mathbf{b}^\top; \mathbf{y} = \mathbf{T}\mathbf{x};$$

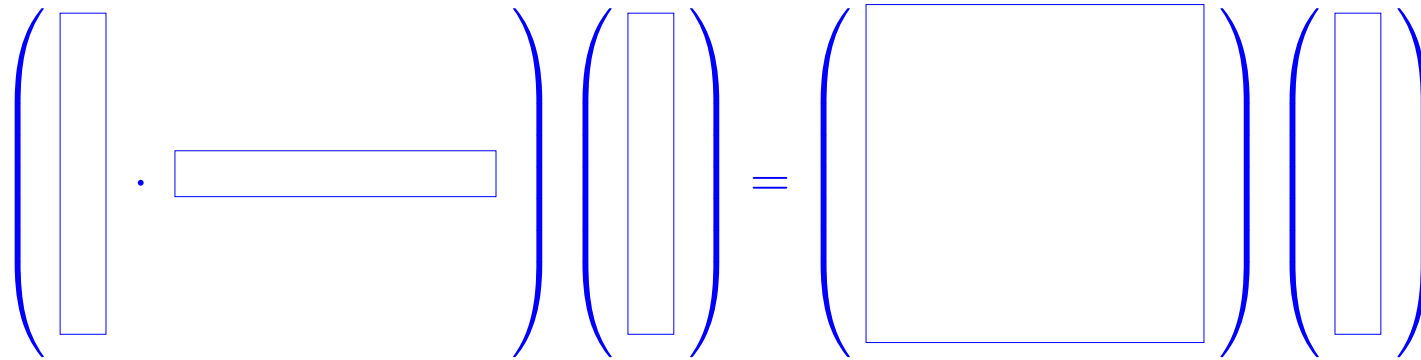
➤ complexity $O(mn)$

$$\mathbf{y} = \mathbf{a}(\mathbf{b}^\top\mathbf{x}). \quad (1.3.5)$$

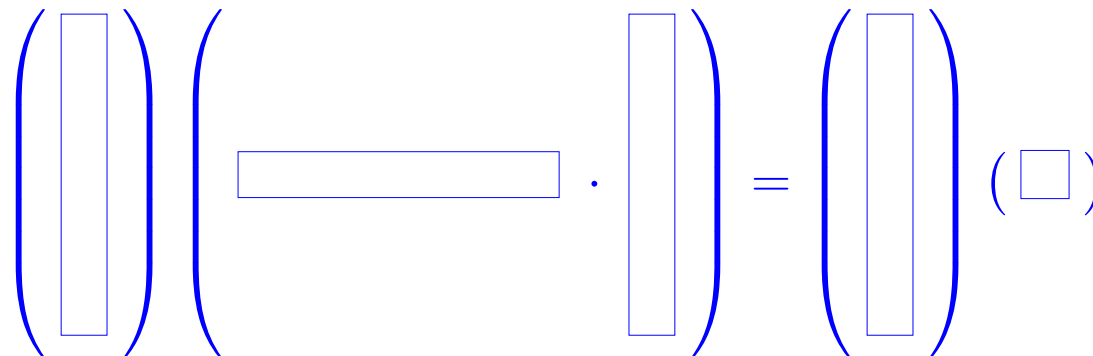
$$\mathbf{t} = \mathbf{b}^\top\mathbf{x}; \mathbf{y} = \mathbf{a}\mathbf{t};$$

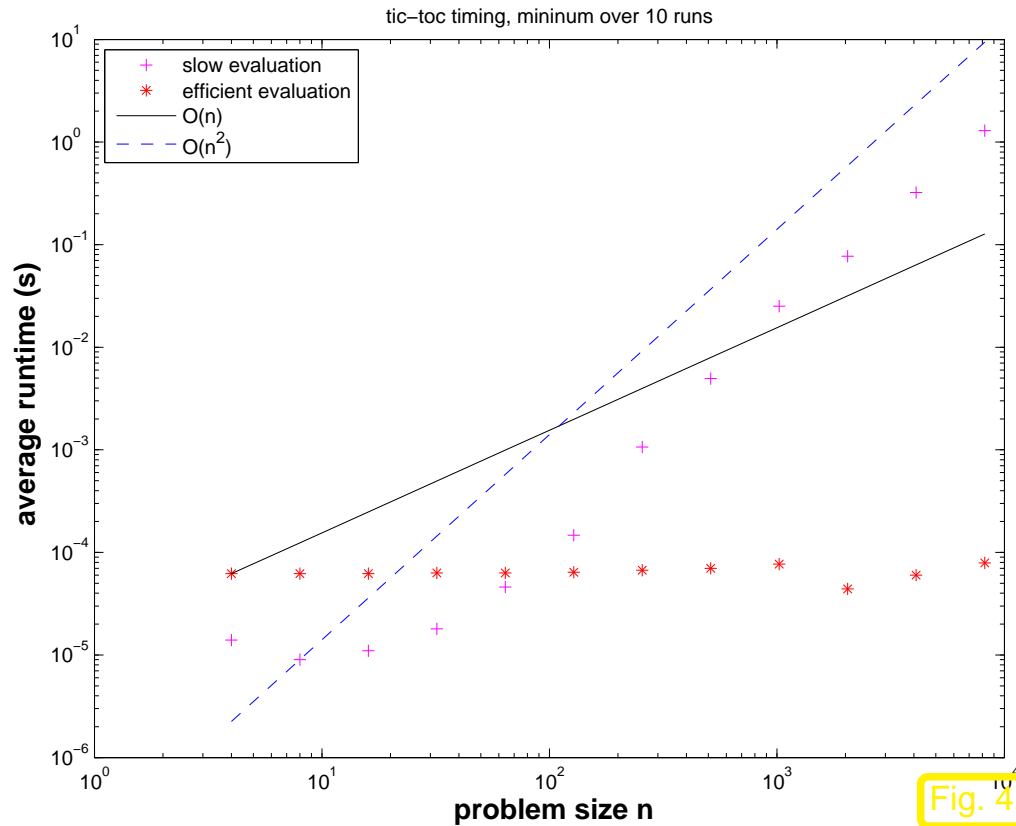
➤ complexity $O(n + m)$ (“linear complexity”)

Visualization of evaluation according to (1.3.4):



Visualization of evaluation according to (1.3.5):





◁ average runtimes for efficient/inefficient matrix×vector multiplication with rank-1 matrices (MATLABtic-toc timing)

Platform:

- MATLAB7.4.0.336 (R2007a)
- Genuine Intel(R) CPU T2500 @ 2.00GHz
- Linux 2.6.16.27-0.9-smp



1.4 BLAS

2

Direct Methods for Linear Systems of Equations

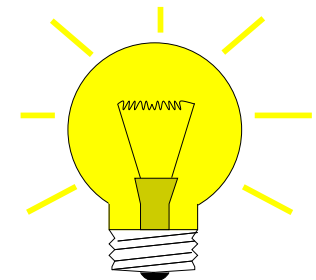
2.1 Gaussian Elimination

!

Exceptional feature of linear systems of equations (LSE):

☞ “exact” solution computable with finitely many elementary operations

Algorithm: **Gaussian elimination** (→ secondary school, linear algebra,)



Idea: transformation to “simpler”, but equivalent LSE by means of successive (invertible) *row transformations*

Code 2.1.4: Solving LSE $Ax = b$ with Gaussian elimination

```

1 function x = gausselimsolve(A,b)
2 % Gauss elimination without pivoting, x = A\b
3 % A must be an n x n-matrix, b an n-vector
4 n = size(A,1); A = [A,b]; %
5 % Forward elimination (cf. step ① in Ex. 2.1.1)
6 for i=1:n-1, pivot = A(i,i);
7   for k=i+1:n, fac = A(k,i)/pivot;
8     A(k,i+1:n+1) = A(k,i+1:n+1) -
9       fac*A(i,i+1:n+1); %
9   end
10 end
11 % Back substitution (cf. step ② in Ex. 2.1.1)
12 A(n,n+1) = A(n,n+1) /A(n,n);
13 for i=n-1:-1:1
14   for l=i+1:n
15     A(i,n+1) = A(i,n+1) -
16       A(l,n+1)*A(i,l);
16   end
17   A(i,n+1) = A(i,n+1)/A(i,i);
18 end
19 x = A(:,n+1); %

```

Algorithm 2.1.3.

Direct MATLAB implementation of Gaussian elimination for LSE $Ax = b$: grossly inefficient!

computational cost (\leftrightarrow number of elementary operations) of Gaussian elimination [48, Sect. 1.3]:

$$\text{elimination : } \sum_{i=1}^{n-1} (n-i)(2(n-i) + 3) = n(n-1)\left(\frac{2}{3}n + \frac{7}{6}\right) \text{ Ops. ,} \tag{2.1.5}$$

$$\text{back substitution : } \sum_{i=1}^n 2(n-i) + 1 = n^2 \text{ Ops. .}$$

asymptotic complexity (\rightarrow Sect. 1.3) of Gaussian elimination (without pivoting) for generic LSE $\mathbf{Ax} = \mathbf{b}$, $\mathbf{A} \in \mathbb{R}^{n,n}$ $= \frac{2}{3}n^3 + O(n^2) = O(n^3)$

Example 2.1.6 (Runtime of Gaussian elimination).

MATLAB \

▷ based on LAPACK

▷ based on BLAS (\rightarrow Sect. 1.4)

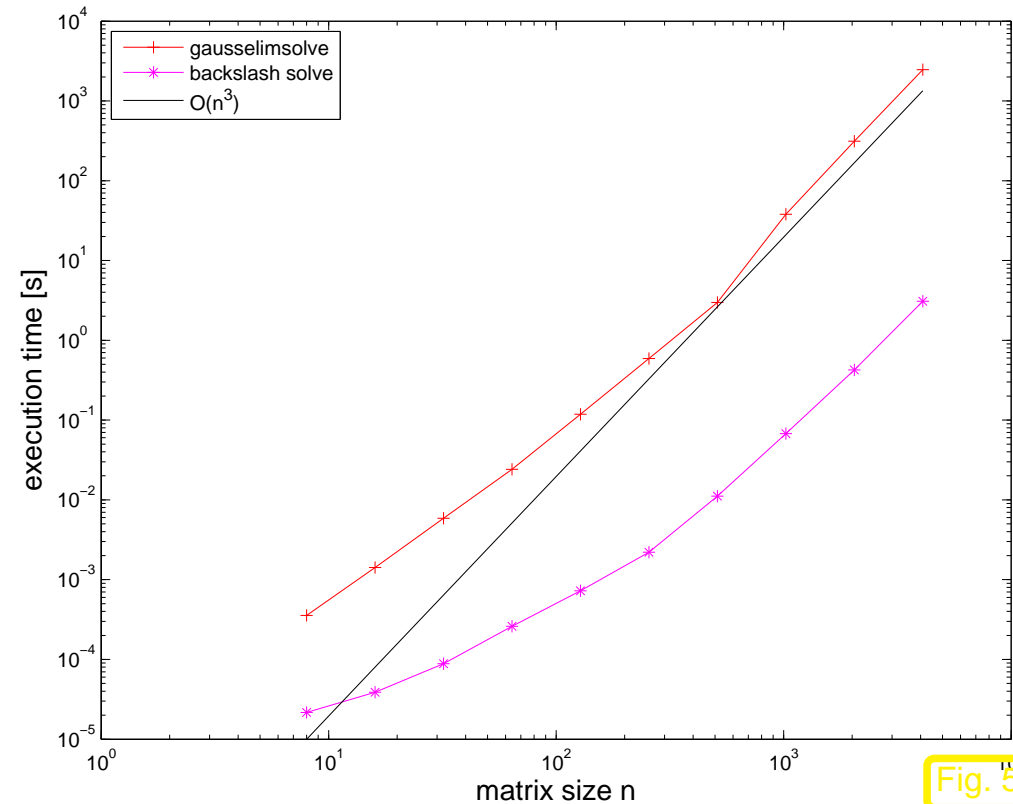


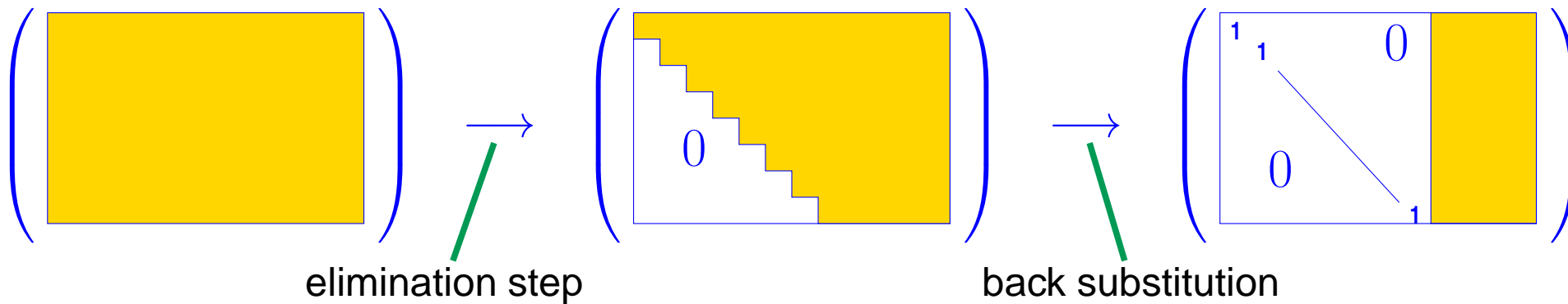
Fig. 5

Never implement Gaussian elimination yourself !

use numerical libraries (LAPACK) or MATLAB ! (MATLAB operator: \backslash)

Remark 2.1.8 (Gaussian elimination for non-square matrices).

“fat matrix”: $A \in \mathbb{K}^{n,m}$, $m > n$:



Simultaneous solving of
LSE with multiple right hand sides

Given regular $\mathbf{A} \in \mathbb{K}^{n,n}$, $\mathbf{B} \in \mathbb{K}^{n,k}$,
seek $\mathbf{X} \in \mathbb{K}^{n,k}$

$$\mathbf{AX} = \mathbf{B} \Leftrightarrow \mathbf{X} = \mathbf{A}^{-1}\mathbf{B}$$

MATLAB:

$$\mathbf{X} = \mathbf{A} \backslash \mathbf{B};$$

asymptotic complexity: $O(n^2(n+k))$

Code 2.1.9: Gaussian elimination with multiple r.h.s.

```

1 function X = gausselimmult(A,B)
2 % Gauss elimination without pivoting, X = A\B
3 n = size(A,1); m = n + size(B,2); A =
   [A,B];
4 for i=1:n-1, pivot = A(i,i);
5   for k=i+1:n, fac = A(k,i)/pivot;
6     A(k,i+1:m) = A(k,i+1:m) -
       fac*A(i,i+1:m);
7   end
8 end
9 A(n,n+1:m) = A(n,n+1:m) /A(n,n);
10 for i=n-1:-1:1
11   for l=i+1:n
12     A(i,n+1:m) = A(i,n+1:m) -
       A(l,n+1:m)*A(i,l);
13   end
14   A(i,n+1:m) = A(i,n+1:m)/A(i,i);
15 end
16 X = A(:,n+1:m);

```

2.2 LU-Decomposition/LU-Factorization

Example 2.2.1 (Gaussian elimination and LU-factorization). → [48, Sect. 2.4], [35, II.4], [27, Sect. 3.1]

LSE from Ex. 2.1.1: consider (forward) Gaussian elimination:

$$\begin{pmatrix} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 3 & -1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 4 \\ 1 \\ -3 \end{pmatrix} \iff \begin{array}{rcl} x_1 + x_2 & = & 4 \\ 2x_1 + x_2 - x_3 & = & 1 \\ 3x_1 - x_2 - x_3 & = & -3 \end{array}$$

$$\begin{pmatrix} 1 & & \\ & 1 & \\ & & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 3 & -1 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \\ -3 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & & \\ 2 & 1 & \\ & 0 & 1 \end{pmatrix} \begin{pmatrix} \mathbf{1} & \mathbf{1} & \mathbf{0} \\ 0 & -1 & -1 \\ 3 & -1 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ -7 \\ -3 \end{pmatrix} \rightarrow$$

$$\begin{pmatrix} 1 & & \\ 2 & 1 & \\ 3 & 0 & 1 \end{pmatrix} \begin{pmatrix} \mathbf{1} & \mathbf{1} & \mathbf{0} \\ 0 & -1 & -1 \\ 0 & -4 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ -7 \\ -15 \end{pmatrix} \rightarrow \underbrace{\begin{pmatrix} 1 & & \\ 2 & 1 & \\ 3 & 4 & 1 \end{pmatrix}}_{=L} \underbrace{\begin{pmatrix} 1 & 1 & 0 \\ 0 & \mathbf{-1} & \mathbf{-1} \\ 0 & 0 & 3 \end{pmatrix}}_{=U} \begin{pmatrix} 4 \\ -7 \\ 13 \end{pmatrix}$$

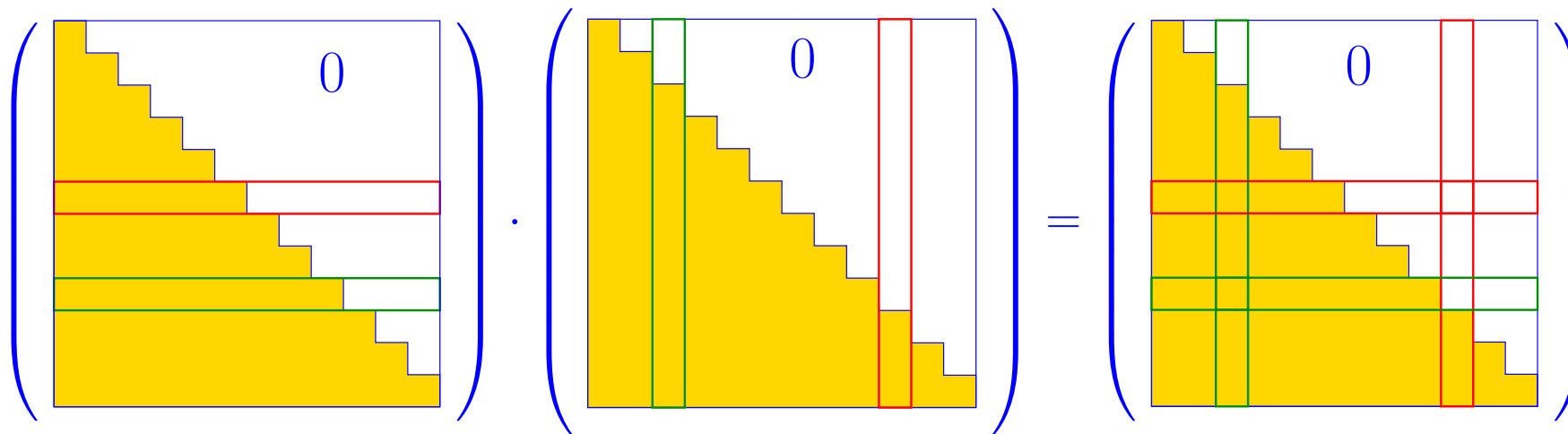
 = pivot row, pivot element bold, negative multipliers red

Lemma 2.2.4 (Group of regular diagonal/triangular matrices).

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

(assumes that \mathbf{A} is regular)

“Proof by visualization” → Rem. 1.2.4

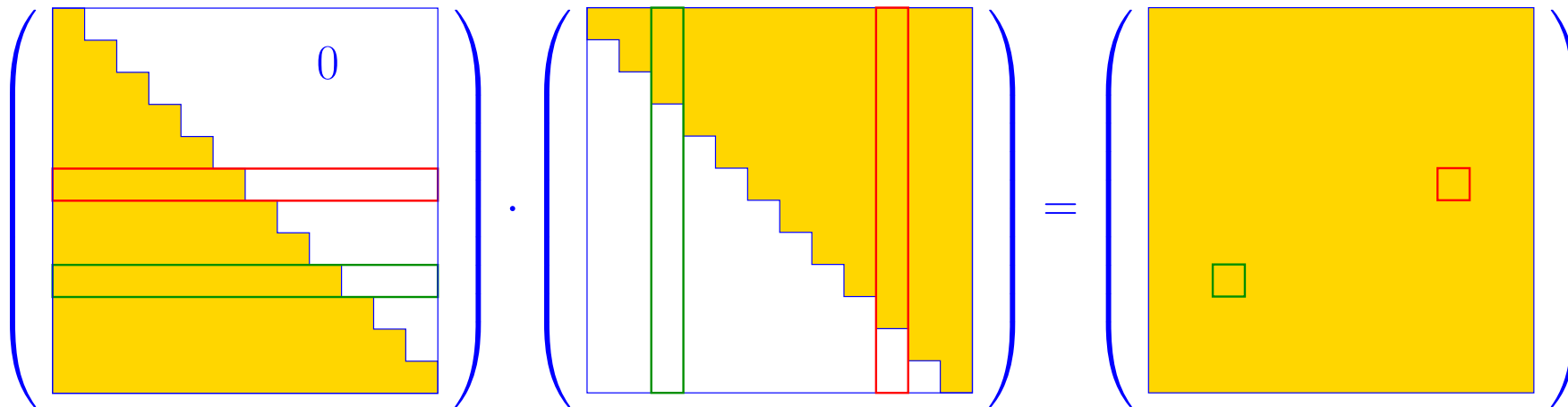


The (forward) Gaussian elimination (without pivoting), for $\mathbf{Ax} = \mathbf{b}$, $\mathbf{A} \in \mathbb{R}^{n,n}$, if possible, is algebraically equivalent to an **LU-factorization**/LU-decomposition $\mathbf{A} = \mathbf{LU}$ of \mathbf{A} into a normalized lower triangular matrix \mathbf{L} and an upper triangular matrix \mathbf{U} , [13, Thm. 3.2.1], [48, Thm. 2.10], [27, Sect. 3.1].

Lemma 2.2.5 (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.

A direct way to LU -decomposition [27, Sect. 3.1], [51, Sect. 3.3.3]:

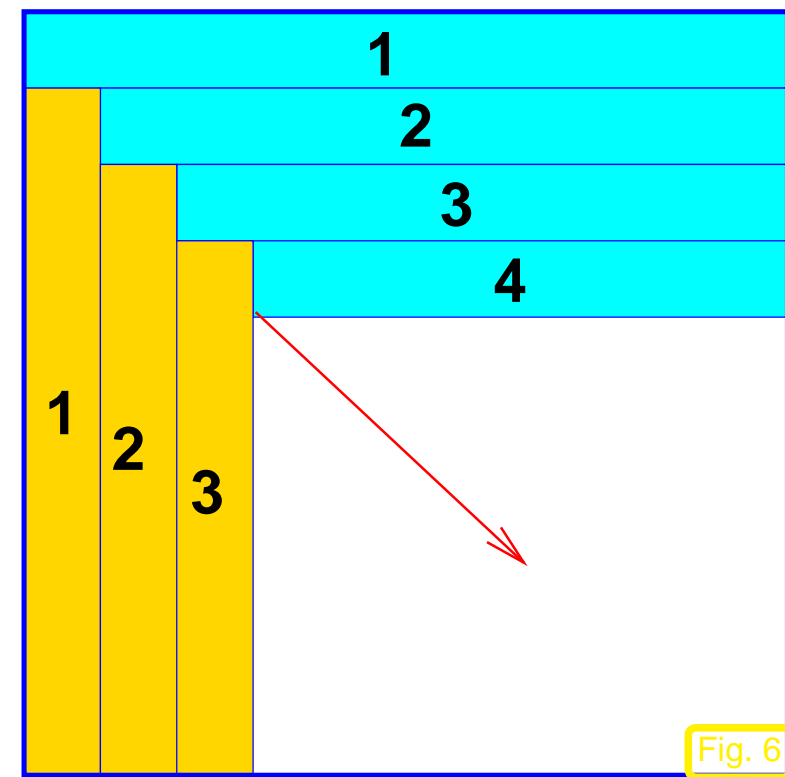
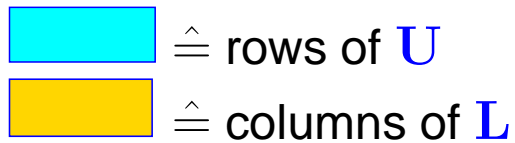


$$\mathbf{LU} = \mathbf{A} \Rightarrow a_{ik} = \sum_{j=1}^{\min\{i,k\}} l_{ij}u_{jk} = \begin{cases} \sum_{j=1}^{i-1} l_{ij}u_{jk} + 1 \cdot u_{ik} & , \text{ if } i \leq k , \\ \sum_{j=1}^{k-1} l_{ij}u_{jk} + l_{ik}u_{kk} & , \text{ if } i > k . \end{cases} \quad (2.2.6)$$

- • row by row computation of \mathbf{U}
- column by column computation of \mathbf{L}

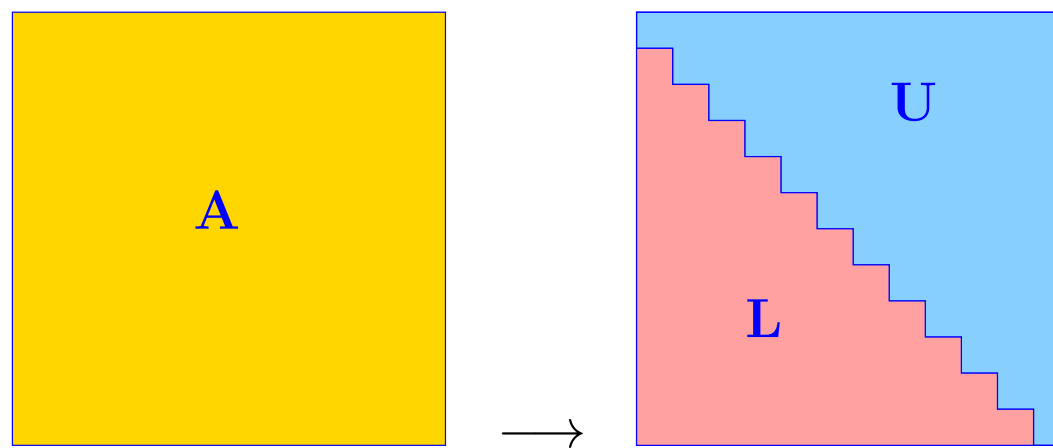
Entries of \mathbf{A} can be replaced with those of \mathbf{L} , \mathbf{U} !
(so-called *in situ*/in place computation)

(Crout's algorithm, [27, Alg. 3.1])



$$\text{asymptotic complexity of LU-factorization of } \mathbf{A} \in \mathbb{R}^{n,n} = \frac{1}{3}n^3 + O(n^2) = O(n^3) \quad (2.2.9)$$

Remark 2.2.10 (In-situ LU-decomposition).



Replace entries of \mathbf{A} with entries of \mathbf{L} (strict lower triangle) and \mathbf{U} (upper triangle).



Solving a linear system of equations by LU-factorization:

Algorithm 2.2.12 (Using LU-factorization to solve a linear system of equations).

- $\mathbf{Ax} = \mathbf{b}$:
- ① *LU*-decomposition $\mathbf{A} = \mathbf{LU}$, #elementary operations $\frac{1}{3}n(n-1)(n+1)$
 - 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 $\frac{1}{2}n(n+1)$

asymptotic complexity: (in leading order) the same as for Gaussian elimination

Remark 2.2.13 (Many sequential solutions of LSE).

Given: regular matrix $\mathbf{A} \in \mathbb{K}^{n,n}$, $n \in \mathbb{N}$, and $N \in \mathbb{N}$, both n, N large

foolish !

```

1 % Setting: N >> 1, large matrix A
2 for j=1:N
3     x = A\b;
4     b = some_function(x);
5 end
    
```

computational effort $O(Nn^3)$

smart !

```

1 % Setting: N >> 1, large matrix A
2 [L,U] = lu(A);
3 for j=1:N
4     x = U\ (L\b);
5     b = some_function(x);
6 end
    
```

computational effort $O(n^3 + Nn^2)$



Remark 2.2.16 (Block LU-factorization).

With $\mathbf{A}_{11} \in \mathbb{K}^{n,n}$ regular, $\mathbf{A}_{12} \in \mathbb{K}^{n,m}$, $\mathbf{A}_{21} \in \mathbb{K}^{m,n}$, $\mathbf{A}_{22} \in \mathbb{K}^{m,m}$:

$$\underbrace{\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}}_{\text{block LU-factorization}} = \begin{pmatrix} \mathbf{I} & 0 \\ \mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ 0 & \mathbf{S} \end{pmatrix},$$

Schur complement

$$\mathbf{S} := \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}.$$

(2.2.17)

block LU-factorization

→ block Gaussian elimination, see Rem. 2.1.12.

2.3 Pivoting

Known from linear algebra [48, Sect. 1.1]:

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

breakdown of Gaussian elimination
pivot element = 0

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} b_2 \\ b_1 \end{pmatrix}$$

Gaussian elimination feasible

Idea (in linear algebra): Avoid zero pivot elements by **swapping rows**

Example 2.3.1 (Pivoting and numerical stability). → [13, Bsp. 3.2.3]


```

1 % Example: numerical instability without
  pivoting
2 A = [5.0E-17 , 1; 1 , 1];
3 b = [1;2];
4 x1 = A\b,
5 x2 =gausselim(A,b), % see Code 2.1.8
6 [L,U] = lufak(A); % see Code 2.2.6
7 z = L\b; x3 = U\z,

```

Output of MATLAB run:

```

x1 = 1
      1
x2 = 0
      1
x3 = 0
      1

```

$$\mathbf{A} = \begin{pmatrix} \epsilon & 1 \\ 1 & 1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \Rightarrow \mathbf{x} = \begin{pmatrix} \frac{1}{1-\epsilon} \\ \frac{1-2\epsilon}{1-\epsilon} \end{pmatrix} \approx \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \text{for } |\epsilon| \ll 1.$$

What is wrong with MATLAB? Needed: insight into **roundoff errors** → Sect. 2.4



Suitable pivoting essential for controlling impact of roundoff errors
on Gaussian elimination (→ Sect. 2.5.2, [48, Sect. 2.5])

Code 2.3.6: Gaussian elimination with pivoting: extension of Code 2.1.6

```

1 function x = gepiv(A,b)
2 % Solving an LSE  $Ax=b$  by Gaussian elimination with partial pivoting
3 n = size(A,1); A = [A,b]; %
4 % Forward elimination by rank-1 modification, see Rem. 2.1.10
5 for k=1:n-1
6     [p,j] = max(abs(A(k:n,k))./max(abs(A(k:n,k:n))')) %
7     if (p < eps*norm(A(k:n,k:n),1)), %
8         disp('A nearly singular'); end
9     A([k,j+k-1],k:n+1) = A([j+k-1,k],k:n+1); %
10    A(k+1:n,k+1:n+1) =
11        A(k+1:n,k+1:n+1)-(A(k+1:n,k)*A(k,k+1:n+1))/A(k,k); %
12 end
13 % Back substitution (same as in Code 2.1.6)
14 A(n,n+1) = A(n,n+1) /A(n,n);
15 for i=n-1:-1:1
16     A(i,n+1) = (A(i,n+1) - A(i,i+1:n)*A(i+1:n,n+1))/A(i,i);
17 end
18 x = A(:,n+1); %

```

choice of pivot row index j (Line 6 of code): **relatively largest** pivot [48, Sect. 2.5],

$$j \in \{k, \dots, n\} \text{ such that } \frac{|a_{ji}|}{\max\{|a_{jl}|, l = k, \dots, n\}} \rightarrow \max \quad (2.3.7)$$

for $k = j, k \in \{i, \dots, n\}$: **partial pivoting**

Example 2.3.11 (Rationale for partial pivoting policy (2.3.7)). \rightarrow [48, Page 47]

```

1 % Example: importance of scale-invariant
  pivoting
2 epsilon = 5.0E-17;
3 A = [epsilon , 1; 1 , 1]; b = [1;2];
4 D = [1/epsilon, 0; 0 ,1];
5 A = D*A; b = D*b;
6 x1 = A\b, % MATLAB internal Gaussian
  elimination
7 x2 =gausselim(A,b), % see Code 2.1.8
8 [L,U] = lufak(A); % see Code 2.2.6
9 z = L\b; x3 = U\z,
```

Output of MATLAB run:

```

x1 = 1
    1
x2 = 0
    1
x3 = 0
    1
```



Lemma 2.3.13 (Existence of LU-factorization with pivoting). \rightarrow [13, Thm. 3.25], [35, Thm. 4.4]

For any regular $\mathbf{A} \in \mathbb{K}^{n,n}$ there is a permutation matrix (\rightarrow Def. 2.3.12) $\mathbf{P} \in \mathbb{K}^{n,n}$, 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. 2.2.3), such that $\mathbf{PA} = \mathbf{LU}$.

Example 2.3.14 (Ex. 2.3.4 cnt'd).

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 2 \\ 2 & -3 & 2 \\ 1 & 24 & 0 \end{pmatrix} \xrightarrow{\textcircled{1}} \begin{pmatrix} 2 & -3 & 2 \\ 1 & 2 & 2 \\ 1 & 24 & 0 \end{pmatrix} \xrightarrow{\textcircled{2}} \begin{pmatrix} 2 & -3 & 2 \\ 0 & 3.5 & 1 \\ 0 & 25.5 & -1 \end{pmatrix} \xrightarrow{\textcircled{3}} \begin{pmatrix} 2 & -7 & 2 \\ 0 & 25.5 & -1 \\ 0 & 3.5 & 1 \end{pmatrix} \xrightarrow{\textcircled{4}} \begin{pmatrix} 2 & -7 & 2 \\ 0 & 25.5 & -1 \\ 0 & 0 & 1.373 \end{pmatrix}$$

$$\mathbf{U} = \begin{pmatrix} 2 & -3 & 2 \\ 0 & 25.5 & -1 \\ 0 & 0 & 1.1373 \end{pmatrix}, \quad \mathbf{L} = \begin{pmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0.5 & 0.1373 & 1 \end{pmatrix}, \quad \mathbf{P} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}.$$



MATLAB function:

$$[\mathbf{L}, \mathbf{U}, \mathbf{P}] = \text{lu}(\mathbf{A}) \quad (\mathbf{P} = \text{permutation matrix})$$

Remark 2.3.15 (Row swapping commutes with forward elimination).

The LU-factorization of $\mathbf{A} \in \mathbb{K}^{n,n}$ with partial pivoting by Alg. 2.3.8 is *numerically equivalent* to the LU-factorization of \mathbf{PA} without pivoting (\rightarrow Code 2.2.6), when \mathbf{P} is a permutation matrix gathering the row swaps entailed by partial pivoting.



2.4 Supplement: Machine Arithmetic

R. Hiptmair
rev 38355,
October 14,
2011

Computer = finite automaton



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

machine numbers, set \mathbb{M}

Essential property:

\mathbb{M} is a **discrete** subset of \mathbb{R}

\mathbb{M} not closed under elementary arithmetic operations $+, -, \cdot, /$.

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

Example 2.4.6 (Input errors and roundoff errors).

Code 2.4.7: input errors and roundoff errors

```

1 >> format long;
2 >> a = 4/3; b = a-1; c = 3*b; e = 1-c
3 e = 2.220446049250313e-16
4 >> a = 1012/113; b = a-9; c = 113*b; e = 5+c
5 e = 6.750155989720952e-14
6 >> a = 83810206/6789; b = a-12345; c = 6789*b; e = c-1
7 e = -1.607986632734537e-09

```



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

correct rounding:

$$\text{rd}(x) = \arg \min_{\tilde{x} \in \mathbb{M}} |x - \tilde{x}|$$

(if non-unique, round to larger (in modulus) $\tilde{x} \in \mathbb{M}$: “rounding up”)

For any reasonable \mathbb{M} : small relative rounding error

$$\exists \text{eps} \ll 1: \frac{|\text{rd}(x) - x|}{|x|} \leq \text{eps} \quad \forall x \in \mathbb{R}. \quad (2.4.8)$$

► Realization of $\tilde{+}, \tilde{-}, \tilde{\cdot}, \tilde{/}$:

$$\star \in \{+, -, \cdot, /\}: \quad x \tilde{\star} y := \text{rd}(x \star y) \quad (2.4.9)$$

Assumption 2.4.10 (“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

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

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

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

Example 2.4.11 (Machine precision for MATLAB). (CPU Intel Pentium)

```

1 >> format hex; eps, format long; eps
2 ans = 3cb0000000000000
3 ans = 2.220446049250313e-16

```



Remark 2.4.13 (Adding ϵ to 1).

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

Code 2.4.14: $1 + \epsilon$ in MATLAB

```

1 >> fprintf ( '%30.25f\n' , 1+0.5*eps )
2 1.00000000000000000000000000000000
3 >> fprintf ( '%30.25f\n' , 1-0.5*eps )
4 0.999999999999999999998889776975
5 >>
6 fprintf ( '%30.25f\n' , (1+2/eps) - 2/eps );
0.00000000000000000000000000000000

```



Do we have to worry about these tiny roundoff errors ?



YES

(→ Sect. 2.3):

- accumulation of roundoff errors
- amplification of roundoff errors

▷ back to Gaussian elimination/LU-factorization with pivoting

2.5 Stability of Gaussian Elimination

Issue: Gauge impact of roundoff errors on Gaussian elimination with partial pivoting !

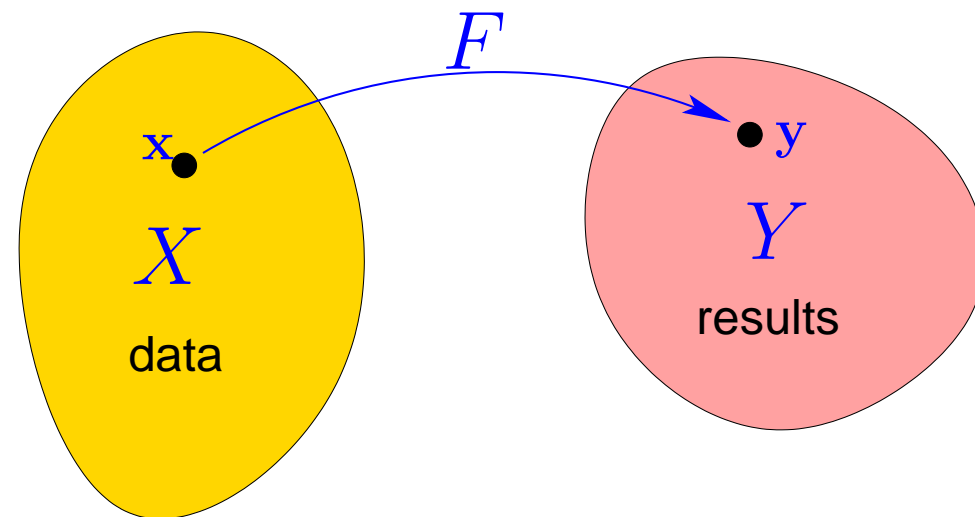
2.5.1 Vector norms and matrix norms [13, Sect. 2.1.2], [35, Sect. 1.2], [51, Sect. 1.11]

2.5.2 Numerical Stability [13, Sect. 2.3]

Abstract point of view:

Our notion of “**problem**”:

- data space X , usually $X \subset \mathbb{R}^n$
- result space Y , usually $Y \subset \mathbb{R}^m$
- mapping (problem function) $F : X \mapsto Y$



Application to linear system of equations $\mathbf{Ax} = \mathbf{b}$, $\mathbf{A} \in \mathbb{K}^{n,n}$, $\mathbf{b} \in \mathbb{K}^n$:

“The problem:”

- data $\hat{=}$ system matrix $\mathbf{A} \in \mathbb{R}^{n,n}$, right hand side vector $\mathbf{b} \in \mathbb{R}^n$
 - data space $X = \mathbb{R}^{n,n} \times \mathbb{R}^n$ with vector/matrix norms (\rightarrow Defs. 2.5.1, 2.5.5)

- problem mapping $(\mathbf{A}, \mathbf{b}) \mapsto F(\mathbf{A}, \mathbf{b}) := \mathbf{A}^{-1}\mathbf{b}$, (for regular \mathbf{A})

Numerical algorithm = Specific sequence of elementary operations
(→ programme in C++ or FORTRAN)

Below: X, Y = normed vector spaces, e.g., $X = \mathbb{R}^n, Y = \mathbb{R}^m$

Definition 2.5.11 (Stable algorithm).

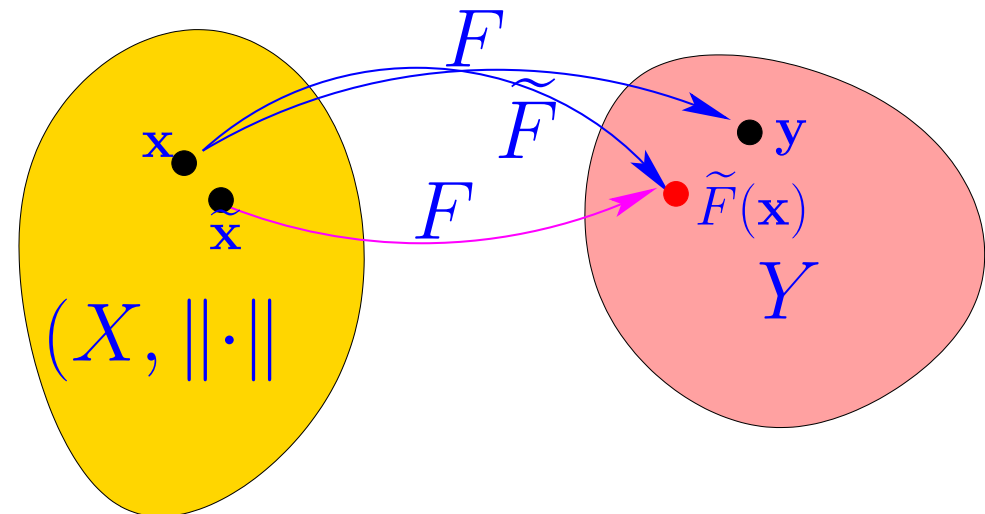
An algorithm \tilde{F} for solving a problem $F : X \mapsto Y$ is **numerically stable**, if for all $x \in X$ its result $\tilde{F}(x)$ (affected by roundoff) is the exact result for “slightly perturbed” data:

$$\exists C \approx 1: \forall x \in X: \exists \tilde{x} \in X: \|x - \tilde{x}\| \leq C \text{ eps } \|x\| \quad \wedge \quad \tilde{F}(x) = F(\tilde{x}) .$$



Terminology:

Def. 2.5.11 introduces stability in the sense of
backward error analysis



2.5.3 Roundoff analysis of Gaussian elimination

Simplification: equivalence of Gaussian elimination and LU-factorization extends to machine arithmetic, *cf.* Sect. 2.2

Lemma 2.5.12 (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:

- ➊ Gaussian elimination (forward elimination and back substitution) without pivoting, see Algorithm 2.1.3.
- ➋ LU-factorization of \mathbf{A} (\rightarrow Code 2.2.6) followed by forward and backward substitution, see Algorithm 2.2.12.

Rem. 2.3.15 \triangleright sufficient to consider LU-factorization without pivoting

Theorem 2.5.13 (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 Algorithm 2.3.8 (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 $\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{b} \in \mathbb{R}^n$, computed by Algorithm 2.3.8 (based on machine arithmetic with machine precision eps , \rightarrow Ass. 2.4.10) 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 $(\mathbf{A} + \Delta\mathbf{A})\tilde{\mathbf{x}} = \mathbf{b}$.

ρ “small” \rightarrow Gaussian elimination with partial pivoting is stable (\rightarrow Def. 2.5.11)



Bad news:

exponential growth $\rho \sim 2^n$ is possible !

Example 2.5.14 (Wilkinson's counterexample).

$$a_{ij} = \begin{cases} 1 & , \text{ if } i = j \vee j = n , \\ -1 & , \text{ if } i > j , \\ 0 & \text{ else.} \end{cases} ,$$

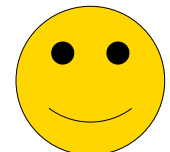
n=10:

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & 1 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & 1 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & 1 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 \end{pmatrix}$$

Partial pivoting does not trigger row permutations !

$$\blacktriangleright \quad \mathbf{A} = \mathbf{LU} , \quad l_{ij} = \begin{cases} 1 & , \text{ if } i = j , \\ -1 & , \text{ if } i > j , \\ 0 & \text{ else} \end{cases} , \quad u_{ij} = \begin{cases} 1 & , \text{ if } i = j , \\ 2^{i-1} & , \text{ if } j = n , \\ 0 & \text{ else.} \end{cases}$$

\blacktriangleright Exponential blow-up of entries of \mathbf{U} !



Observation: In practice ρ (almost) always grows only mildly (like $O(\sqrt{n})$) with n

Example 2.5.15 (Stability by small random perturbations).

```
1 % Curing Wilkinson's counterexample by random perturbation
2 % Theory: Spielman and Teng
3 res = [];
4 for n=10:10:200
5     % Build Wilkinson matrix
6     A = [ tril(-ones(n,n-1))+2*[ eye(n-1);
7         zeros(1,n-1)], ones(n,1) ];
8     % imposed solution
9     x = ((-1).^(1:n))';
10    relerr = norm(A\((A*x)-x)/norm(x);
11    % Randomly perturbed Wilkinson matrix by matrix with iid
12    % N(0,eps) distributed entries
13    Ap = A + eps*randn(size(A));
14    relerrp = norm(Ap\((A*x)-x)/norm(x);
15    res = [res; n relerr relerrp];
16 end
17 semilogy(res(:,1),res(:,2),'m-*',res(:,1),res(:,3),'r-+');
18 xlabel('matrix size n','fontsize',14);
19 ylabel('relative error','fontsize',14);
20 legend('unperturbed matrix','randn perturbed
21     matrix','location','west');
22 print -depsc2 '../PICTURES/wilkpert.eps';
```

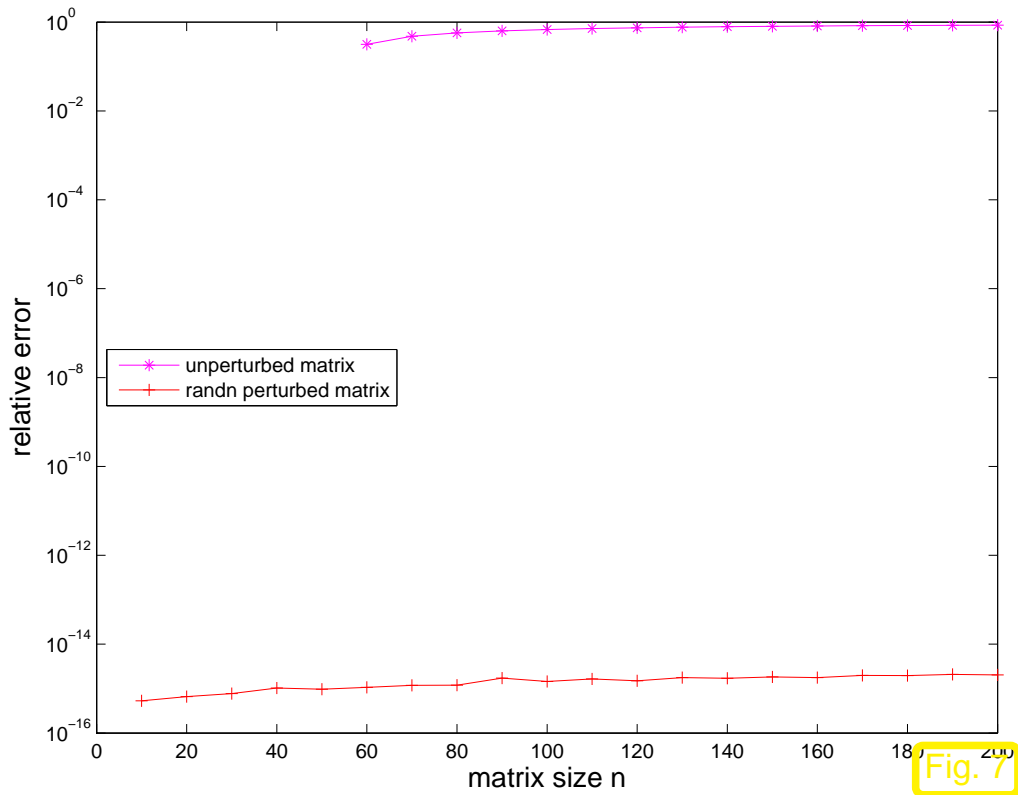


Fig. 7



Gaussian elimination/LU-factorization with partial pivoting is stable ()*
(for all practical purposes) !

(*): stability refers to maximum norm $\|\cdot\|_\infty$.

In practice *Gaussian elimination/LU-factorization with partial pivoting* produces “relatively **small residuals**”

Definition 2.5.16 (Residual).

Given an approximate solution $\tilde{\mathbf{x}} \in \mathbb{K}^n$ of the LSE $\mathbf{Ax} = \mathbf{b}$ ($\mathbf{A} \in \mathbb{K}^{n,n}$, $\mathbf{b} \in \mathbb{K}^n$), its *residual* is the vector

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}} .$$

Example 2.5.17 (Small residuals by Gaussian elimination).

Code 2.5.18: small residuals for GE

Numerical experiment with *nearly singular matrix*

$$\mathbf{A} = \mathbf{u}\mathbf{v}^T + \epsilon\mathbf{I},$$

singular rank-1 matrix

with

$$\mathbf{u} = \frac{1}{3}(1, 2, 3, \dots, 10)^T,$$

$$\mathbf{v} = \left(-1, \frac{1}{2}, -\frac{1}{3}, \frac{1}{4}, \dots, \frac{1}{10}\right)^T$$

```

1 n = 10; u = (1:n)'/3; v =
  (1./u).*(-1).^((1:n)');
2 x = ones(10,1); nx = norm(x,'inf');
3
4 result = [];
5 for epsilon = 10.^(-5:-0.5:-14)
6   A = u*v' + epsilon*eye(n);
7   b = A*x; nb = norm(b,'inf');
8   xt = A\b; % Gaussian elimination
9   r = b - A*xt; % residual
10  result = [result; epsilon,
             norm(x-xt,'inf')/nx,
             norm(r,'inf')/nb, cond(A,'inf')];
11 end

```

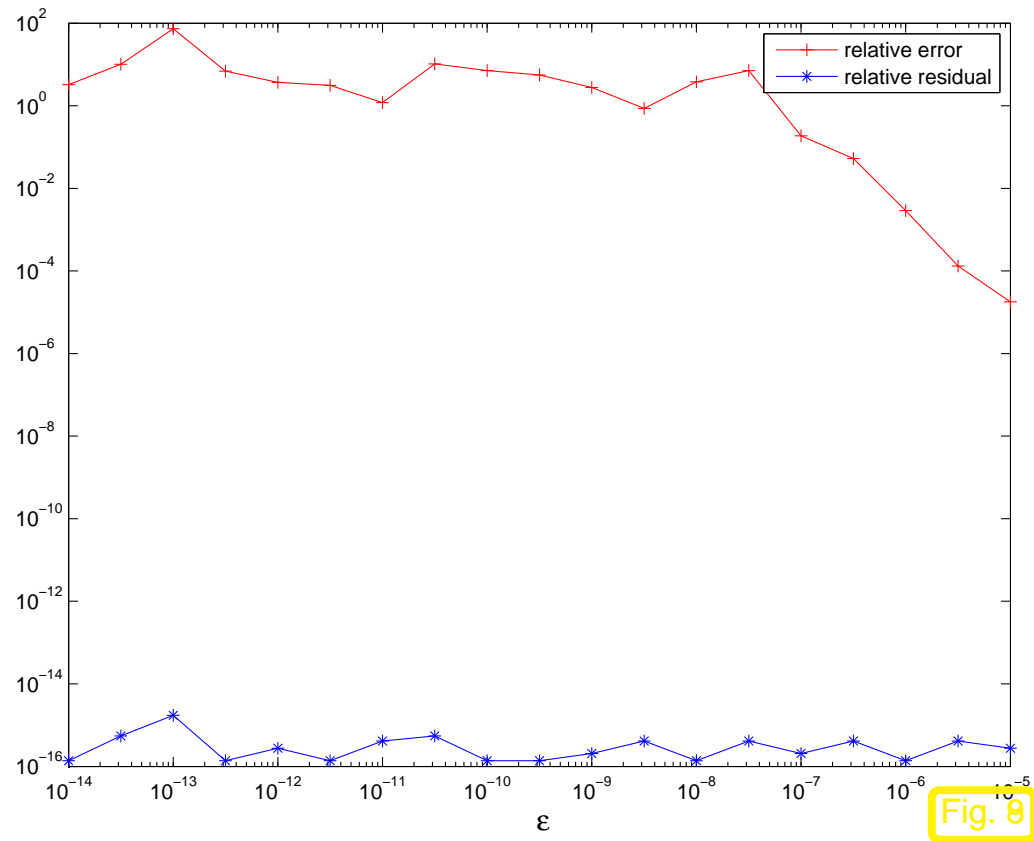


Fig. 9

Observations (w.r.t $\|\cdot\|_\infty$ -norm)

- for $\epsilon \ll 1$ large relative error in computed solution $\tilde{\mathbf{x}}$
- small residuals for any ϵ



Example 2.5.20 (Instability of multiplication with inverse).

Nearly singular matrix from Ex. 2.5.17

Code 2.5.22: instability of multiplication with inverse

```

1 n = 10; u = (1:n)'/3; v =
  (1./u).*(-1).^u;
2 x = ones(10,1); nx =
  norm(x, 'inf');
3
4 result = [];
5 for epsilon = 10.^(-5:-0.5:-14)
6   A = u*v' +
  epsilon*rand(n,n);
7   b = A*x; nb = norm(b, 'inf');
8   xt = A\b;      % Gaussian
  elimination
9   r = b - A*xt;  % residualB
10  B = inv(A); xi = B*b;
11  ri = b - A*xi; % residual
12  R = eye(n) - A*B; % residual
13  result = [result; epsilon,
  norm(r, 'inf')/nb,
  norm(ri, 'inf')/nb,
  norm(R, 'inf')/norm(B, 'inf')
  ];

```

14 end

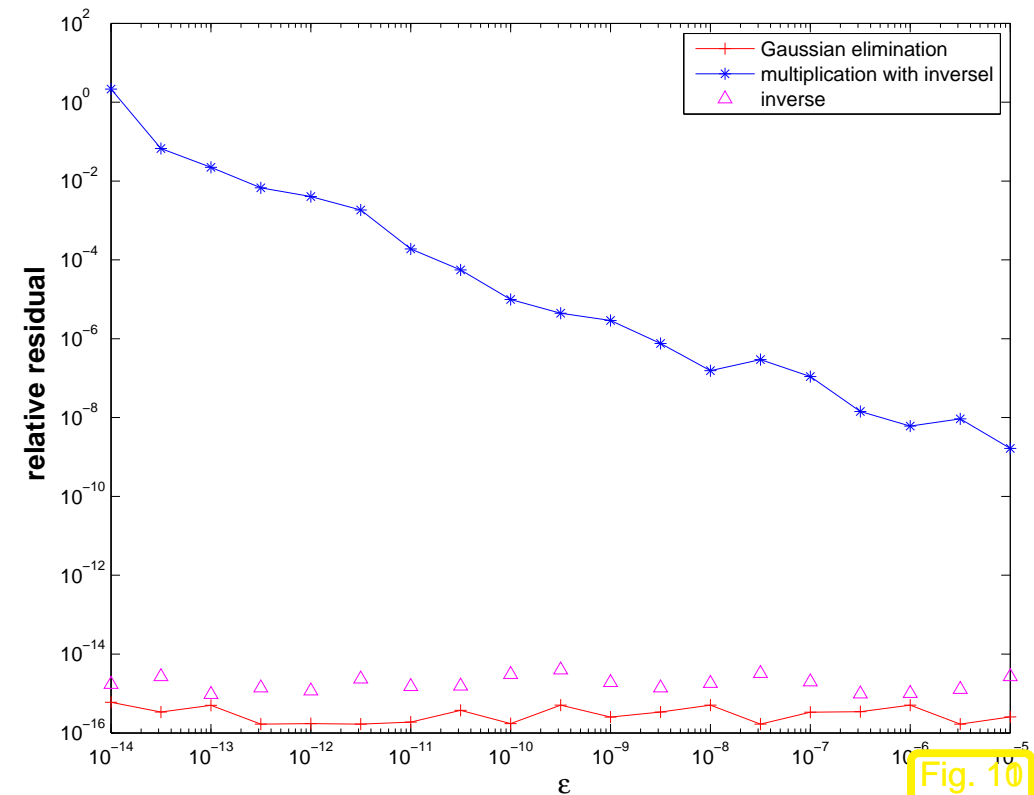


Fig. 10



2.5.4 Conditioning

Question: implications of stability results (\rightarrow previous section) for

(normwise) **relative error**: $\epsilon_r := \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|}$.

($\|\cdot\| \hat{=}$ suitable vector norm, e.g., maximum norm $\|\cdot\|_\infty$)

Perturbed linear system:

$$\mathbf{Ax} = \mathbf{b} \iff (\mathbf{A} + \Delta\mathbf{A})\tilde{\mathbf{x}} = \mathbf{b} + \Delta\mathbf{b} \quad \blacktriangleright \quad (\mathbf{A} + \Delta\mathbf{A})(\tilde{\mathbf{x}} - \mathbf{x}) = \Delta\mathbf{b} - \Delta\mathbf{Ax} . \quad (2.5.23)$$



Theorem 2.5.24 (Conditioning of LSEs). \rightarrow [51, Thm. 3.1]

If \mathbf{A} regular, $\|\Delta\mathbf{A}\| < \|\mathbf{A}^{-1}\|^{-1}$ and (2.5.23), then

$$\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{A}^{-1}\| \|\mathbf{A}\|}{1 - \|\mathbf{A}^{-1}\| \|\mathbf{A}\| \|\Delta\mathbf{A}\| / \|\mathbf{A}\|} \left(\frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} \right).$$

\uparrow
relative error
 \uparrow \uparrow
relative perturbations

Definition 2.5.26 (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 $\|\cdot\|$!

Rewriting estimate of Thm. 2.5.24 with $\Delta\mathbf{b} = 0$,

$$\epsilon_r := \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \frac{\text{cond}(\mathbf{A})\delta_A}{1 - \text{cond}(\mathbf{A})\delta_A}, \quad \delta_A := \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|}. \quad (2.5.27)$$

- (2.5.27) ➤
- 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. 2.5.11) can produce solutions with large relative error !

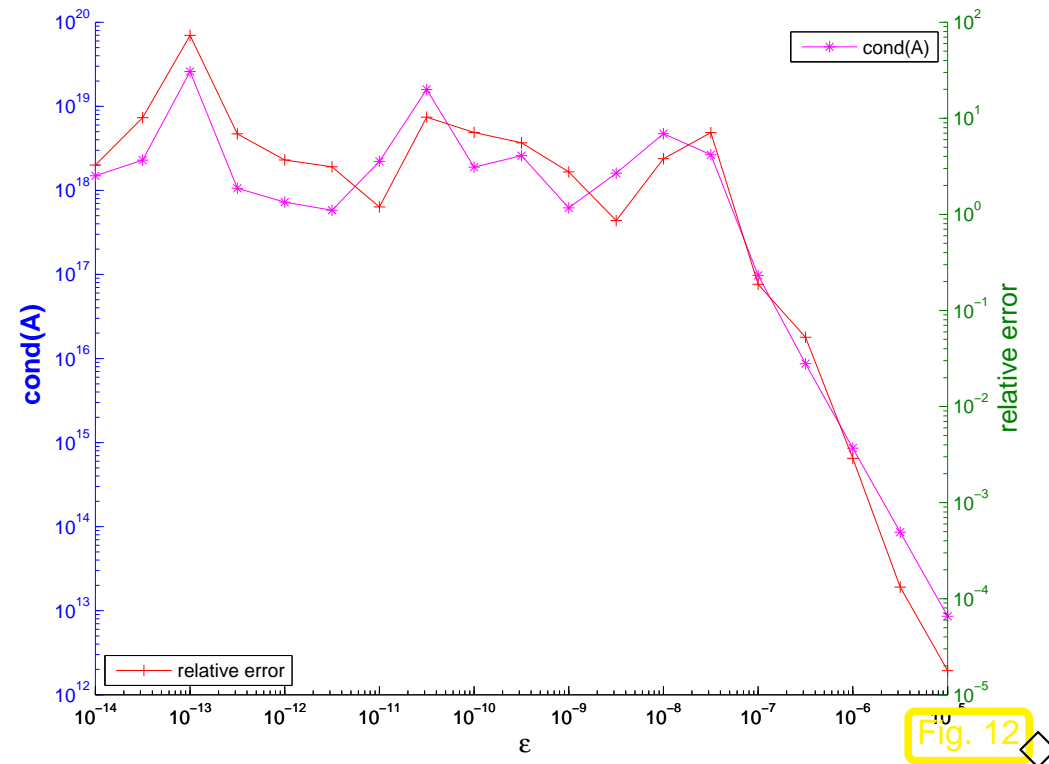
Example 2.5.28 (Conditioning and relative error). \rightarrow Ex. 2.5.17 cnt'd

Numerical experiment with *nearly singular matrix* from Ex. 2.5.17

$$\mathbf{A} = \mathbf{u}\mathbf{v}^T + \epsilon\mathbf{I},$$

$$\mathbf{u} = \frac{1}{3}(1, 2, 3, \dots, 10)^T,$$

$$\mathbf{v} = \left(-1, \frac{1}{2}, -\frac{1}{3}, \frac{1}{4}, \dots, \frac{1}{10}\right)^T$$



Example 2.5.29 (Wilkinson's counterexample cnt'd). \rightarrow Ex. 2.5.14

Code 2.5.30: GE for “Wilkinson system”

```
1 res = [];  
2 for n=10:10:1000  
3     A =  
4         [ tril(-ones(n,n-1))+2*[eye(n-1);  
5             zeros(1,n-1)], ones(n,1) ];  
6     x = ((-1).^(1:n))';  
7     relerr = norm(A\(A*x)-x)/norm(x);  
8     res = [res; n,relerr];  
9 end  
10 plot(res(:,1),res(:,2),'m-*');
```

Blow-up of entries of **U** !

↕ (*)

However, $\text{cond}_2(\mathbf{A})$ is small!

▷ **Instability** of Gaussian elimination !

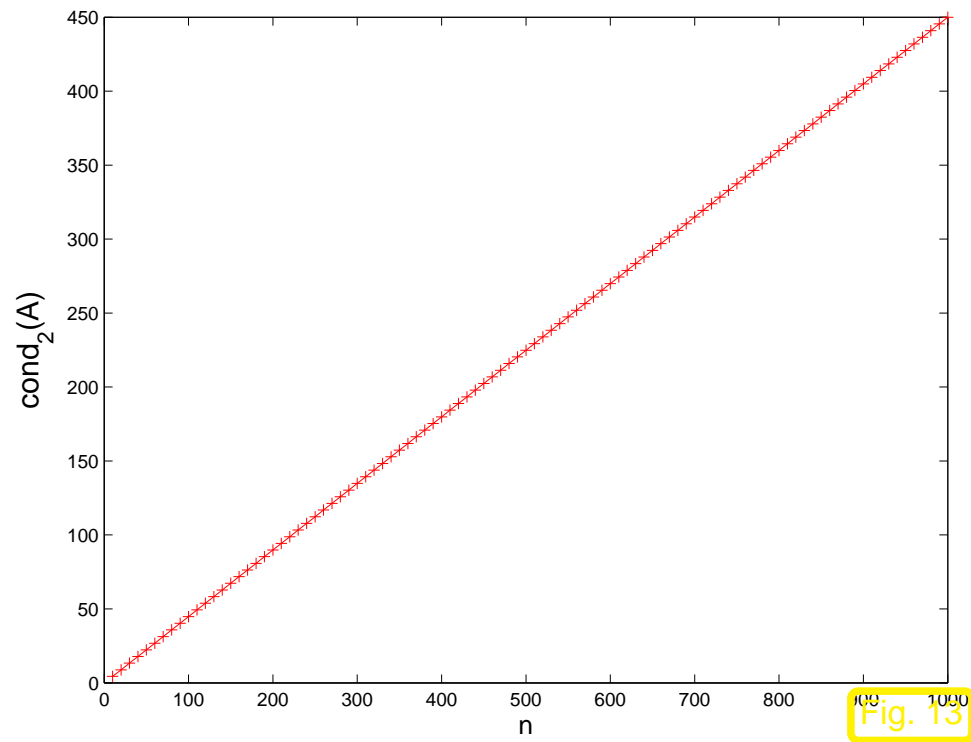


Fig. 13

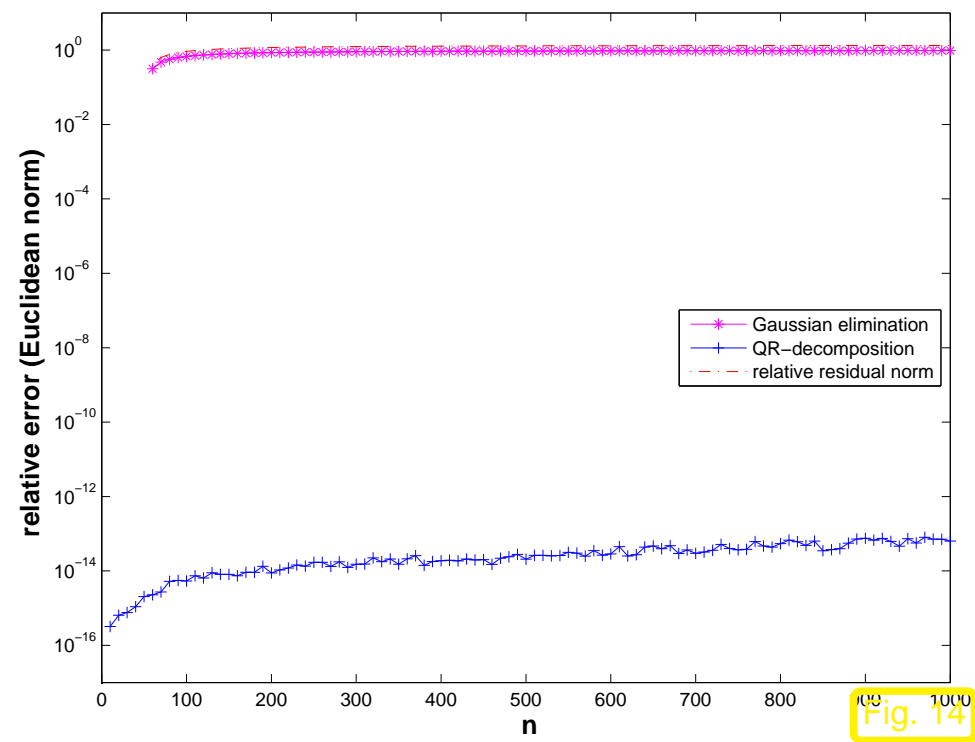


Fig. 14



2.5.5 Sensitivity of linear systems

Recall Thm. 2.5.24: for regular $\mathbf{A} \in \mathbb{K}^{n,n}$, small $\Delta\mathbf{A}$, generic vector/matrix norm $\|\cdot\|$

$$\begin{aligned} \mathbf{Ax} = \mathbf{b} \\ (\mathbf{A} + \Delta\mathbf{A})\tilde{\mathbf{x}} = \mathbf{b} + \Delta\mathbf{b} \end{aligned} \Rightarrow \frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|}{\|\mathbf{x}\|} \leq \frac{\text{cond}(\mathbf{A})}{1 - \text{cond}(\mathbf{A}) \|\Delta\mathbf{A}\| / \|\mathbf{A}\|} \left(\frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} \right). \quad (2.5.31)$$

► $\text{cond}(\mathbf{A}) \gg 1$ ➤ small relative changes of data \mathbf{A}, \mathbf{b} may effect huge relative changes in solution.

Sensitivity of a problem (for given data) gauges impact of small perturbations of the data on the result.

► $\text{cond}(\mathbf{A})$ indicates sensitivity of “LSE problem” $(\mathbf{A}, \mathbf{b}) \mapsto \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ (as “amplification factor” of relative perturbations in the data \mathbf{A}, \mathbf{b}).

Terminology:

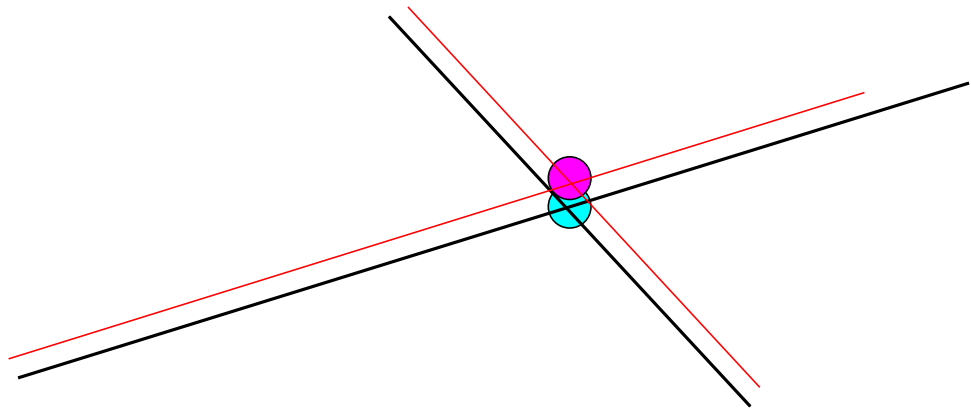
Small changes of data \Rightarrow small perturbations of result : well-conditioned problem

Small changes of data \Rightarrow large perturbations of result : ill-conditioned problem

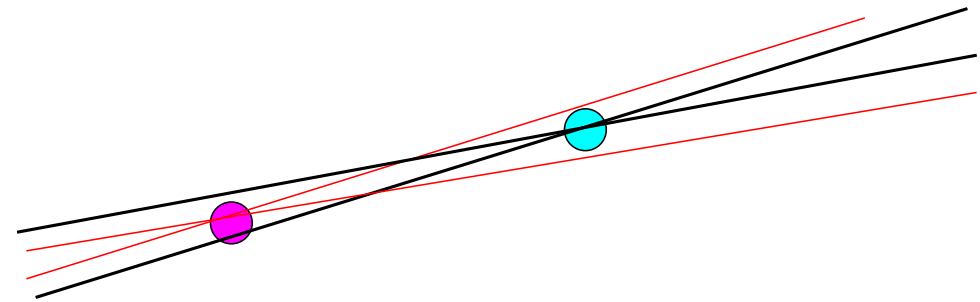
Note: sensitivity gauge depends on the chosen norm !

Example 2.5.32 (Intersection of lines in 2D).

In distance metric:



nearly orthogonal intersection: well-conditioned



glancing intersection: ill-conditioned

Hessian normal form of line # i , $i = 1, 2$:

$$L_i = \{ \mathbf{x} \in \mathbb{R}^2 : \mathbf{x}^T \mathbf{n}_i = d_i \}, \quad \mathbf{n}_i \in \mathbb{R}^2, d_i \in \mathbb{R}.$$

$$\blacktriangleright \quad \text{intersection:} \quad \underbrace{\begin{pmatrix} \mathbf{n}_1^T \\ \mathbf{n}_2^T \end{pmatrix}}_{=: \mathbf{A}} \mathbf{x} = \underbrace{\begin{pmatrix} d_1 \\ d_2 \end{pmatrix}}_{=: \mathbf{b}},$$

$\mathbf{n}_i \hat{=}$ (unit) direction vectors, $d_i \hat{=}$ distance to origin.

Code 2.5.34: condition numbers of 2×2 matrices

```

1 r = [];
2 for phi=pi/200:pi/100:pi/2
3     A = [1,cos(phi);
4         0,sin(phi)];
5     r = [r; phi,
6         cond(A),cond(A,'inf')];
7 end
8 plot(r(:,1),r(:,2),'r-',
9     r(:,1),r(:,3),'b--');
10 xlabel('{\bf angle of n_1,
11     n_2}','fontsize',14);
12 ylabel('{\bf condition
13     numbers}','fontsize',14);
14 legend('2-norm','max-norm');
15 print -depsc2
16     '../PICTURES/linesec.eps';

```

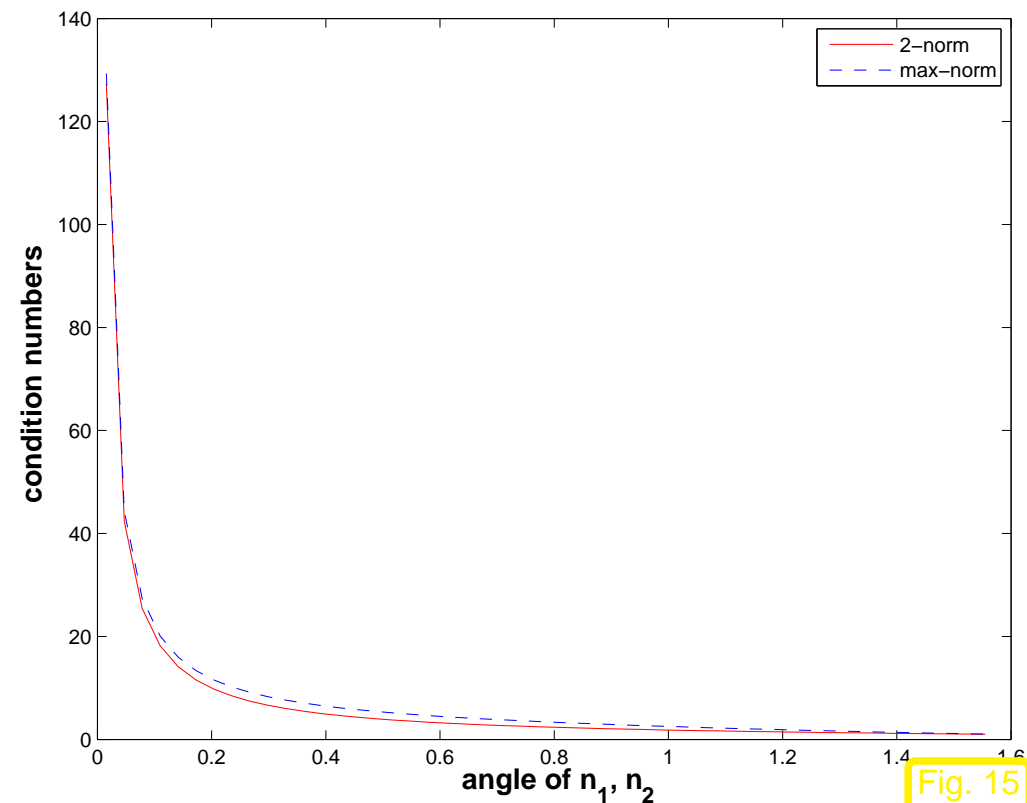


Fig. 15

Heuristics:

$$\text{cond}(\mathbf{A}) \gg 1 \iff \text{columns/rows of } \mathbf{A} \text{ "almost linearly dependent"}$$

2.6 Sparse Matrices

A classification of matrices:

Dense matrices (*ger.*: vollbesetzt)



sparse matrices (*ger.*: dünnbesetzt)

Notion 2.6.1 (Sparse matrix). $\mathbf{A} \in \mathbb{K}^{m,n}$, $m, n \in \mathbb{N}$, is *sparse*, if

$$\text{nnz}(\mathbf{A}) := \#\{(i, j) \in \{1, \dots, m\} \times \{1, \dots, n\} : a_{ij} \neq 0\} \ll mn .$$

Example 2.6.3 (Nodal analysis of (linear) electric circuit). [51, Sect. 4.7.1],

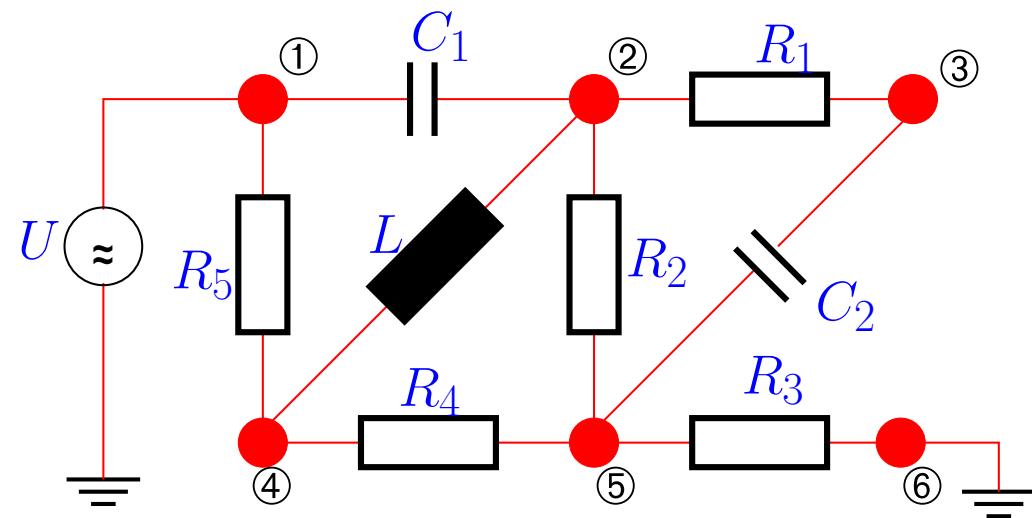
Node (*ger.*: Knoten) $\hat{=}$ junction of wires

\rightarrow number nodes $1, \dots, n$

I_{kj} : current from node $k \rightarrow$ node j , $I_{kj} = -I_{jk}$

Kirchhoff current law (KCL) : sum of node currents = 0:

$$\forall k \in \{1, \dots, n\}: \sum_{j=1}^n I_{kj} = 0 . \quad (2.6.4)$$



Unknowns: **nodal potentials** U_k , $k = 1, \dots, n$.

(some may be known: grounded nodes, voltage sources)

Constitutive relations for circuit elements: (in *frequency domain* with angular frequency $\omega > 0$):

- Ohmic resistor: $I = \frac{U}{R}$, $[R] = 1\text{VA}^{-1}$
 - capacitor: $I = \imath\omega CU$, capacitance $[C] = 1\text{AsV}^{-1}$
 - coil/inductor: $I = \frac{U}{\imath\omega L}$, inductance $[L] = 1\text{VsA}^{-1}$
- $I_{kj} = \begin{cases} R^{-1}(U_k - U_j) , \\ \imath\omega C(U_k - U_j) , \\ -\imath\omega^{-1}L^{-1}(U_k - U_j) . \end{cases}$

Constitutive relations + (2.6.4) ➤ linear system of equations:

$$\begin{aligned}
 \textcircled{2} : & \quad \imath\omega C_1(U_2 - U_1) + R_1^{-1}(U_2 - U_3) - \imath\omega^{-1}L^{-1}(U_2 - U_4) + R_2^{-1}(U_2 - U_5) = 0 , \\
 \textcircled{3} : & \quad R_1^{-1}(U_3 - U_2) + \imath\omega C_2(U_3 - U_5) = 0 , \\
 \textcircled{4} : & \quad R_5^{-1}(U_4 - U_1) - \imath\omega^{-1}L^{-1}(U_4 - U_2) + R_4^{-1}(U_4 - U_5) = 0 , \\
 \textcircled{5} : & \quad R_2^{-1}(U_5 - U_2) + \imath\omega C_2(U_5 - U_3) + R_4^{-1}(U_5 - U_4) + R_3^{-1}(U_5 - U_6) = 0 , \\
 & \quad U_1 = U \quad , \quad U_6 = 0 .
 \end{aligned}$$

$$\begin{pmatrix} \omega C_1 + \frac{1}{R_1} - \frac{i}{\omega L} + \frac{1}{R_2} & -\frac{1}{R_1} & \frac{i}{\omega L} & -\frac{1}{R_2} & 0 \\ -\frac{1}{R_1} & \frac{1}{R_1} + \omega C_2 & 0 & -\omega C_2 & 0 \\ \frac{i}{\omega L} & 0 & \frac{1}{R_5} - \frac{i}{\omega L} + \frac{1}{R_4} & -\frac{1}{R_4} & 0 \\ -\frac{1}{R_2} & -\omega C_2 & -\frac{1}{R_4} & \frac{1}{R_2} + \omega C_2 + \frac{1}{R_4} + R_3^{-1} & 0 \end{pmatrix} \begin{pmatrix} U_2 \\ U_3 \\ U_4 \\ U_5 \end{pmatrix} = \begin{pmatrix} \omega C_1 U \\ 0 \\ \frac{1}{R_5} U \\ 0 \end{pmatrix}$$



Example 2.6.7 (Sparse LSE in circuit modelling).

Modern electric circuits (VLSI chips):

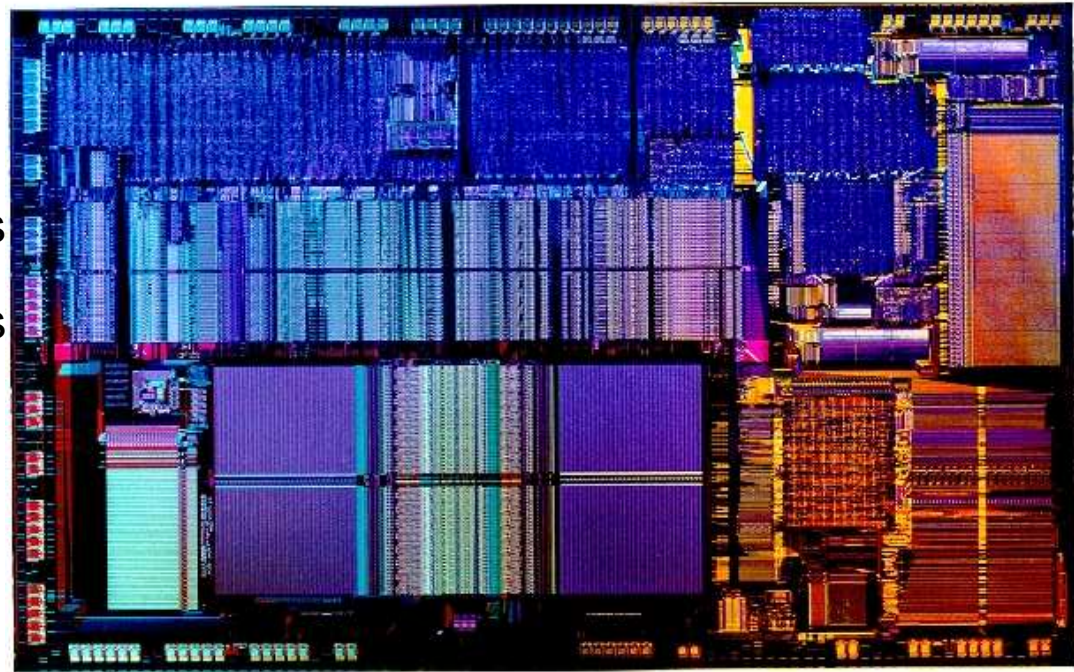
$10^5 - 10^7$ circuit elements

- 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 ➤ **sparse** circuit matrix



2.6.1 Sparse matrix storage formats

Special **sparse matrix storage formats** store *only* non-zero entries:

(\triangleright usually $O(n + m)$ storage required for sparse $n \times m$ -matrix)

- 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)

\blacktriangleright mandatory for large sparse matrices.

Example 2.6.8 (**Compressed row-storage (CRS)** format).

Data for matrix $\mathbf{A} = (a_{ij}) \in \mathbb{K}^{n,n}$ kept in three arrays

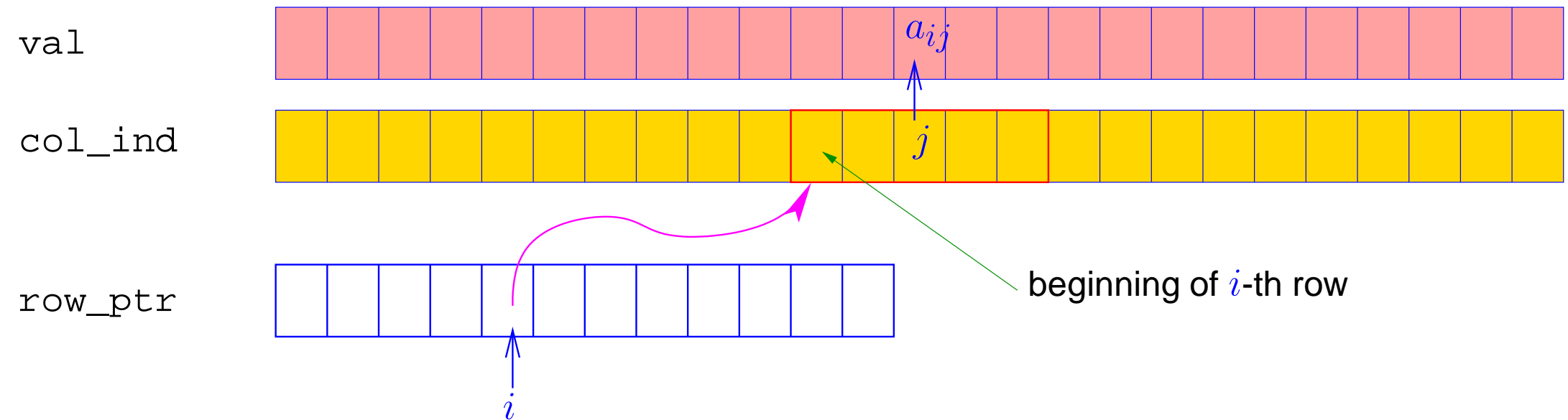

```

double * val           size nnz(A) := #{(i, j) ∈ {1, ..., n}^2, aij ≠ 0}
unsigned int * col_ind size nnz(A)
unsigned int * row_ptr  size n + 1 & row_ptr[n + 1] = nnz(A) + 1
    
```

$\text{nnz}(\mathbf{A}) \hat{=}$ (number of nonzeros) of \mathbf{A}

Access to matrix entry $a_{ij} \neq 0, 1 \leq i, j \leq n$:

$$\text{val}[k] = a_{ij} \Leftrightarrow \begin{cases} \text{col_ind}[k] = j, \\ \text{row_ptr}[i] \leq k < \text{row_ptr}[i + 1], \end{cases} \quad 1 \leq k \leq \text{nnz}(\mathbf{A}).$$



$$A = \begin{pmatrix} 10 & 0 & 0 & 0 & -2 & 0 \\ 3 & 9 & 0 & 0 & 0 & 3 \\ 0 & 7 & 8 & 7 & 0 & 0 \\ 3 & 0 & 8 & 7 & 5 & 0 \\ 0 & 8 & 0 & 9 & 9 & 13 \\ 0 & 4 & 0 & 0 & 2 & -1 \end{pmatrix}$$

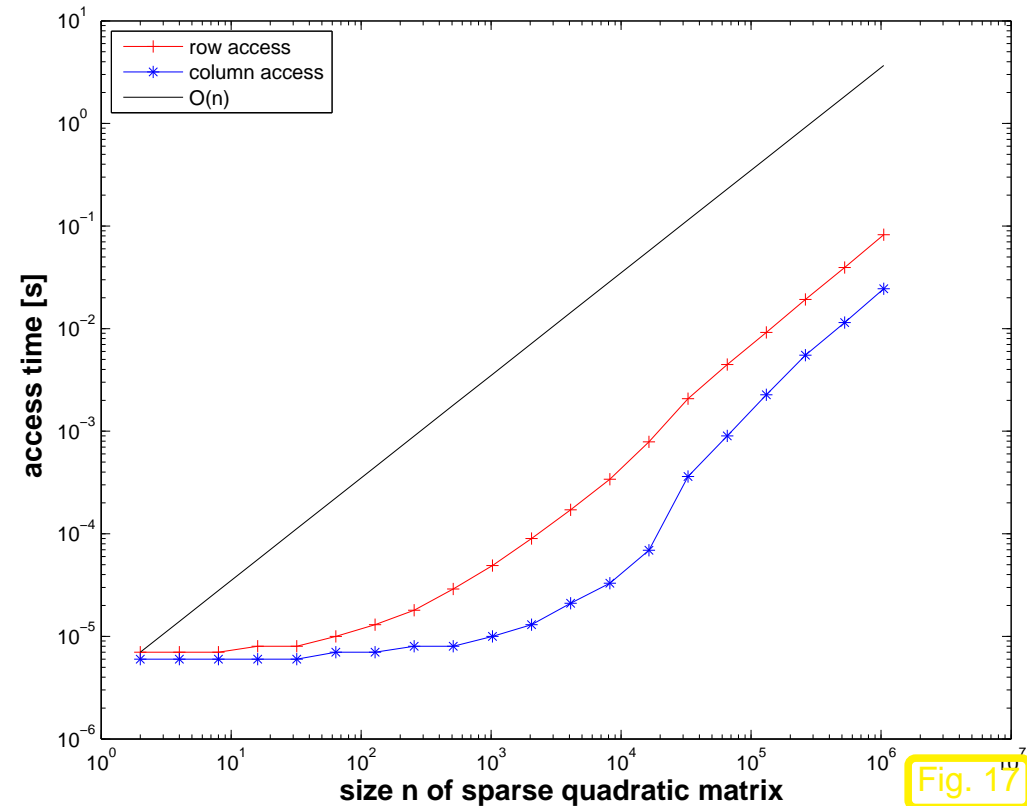
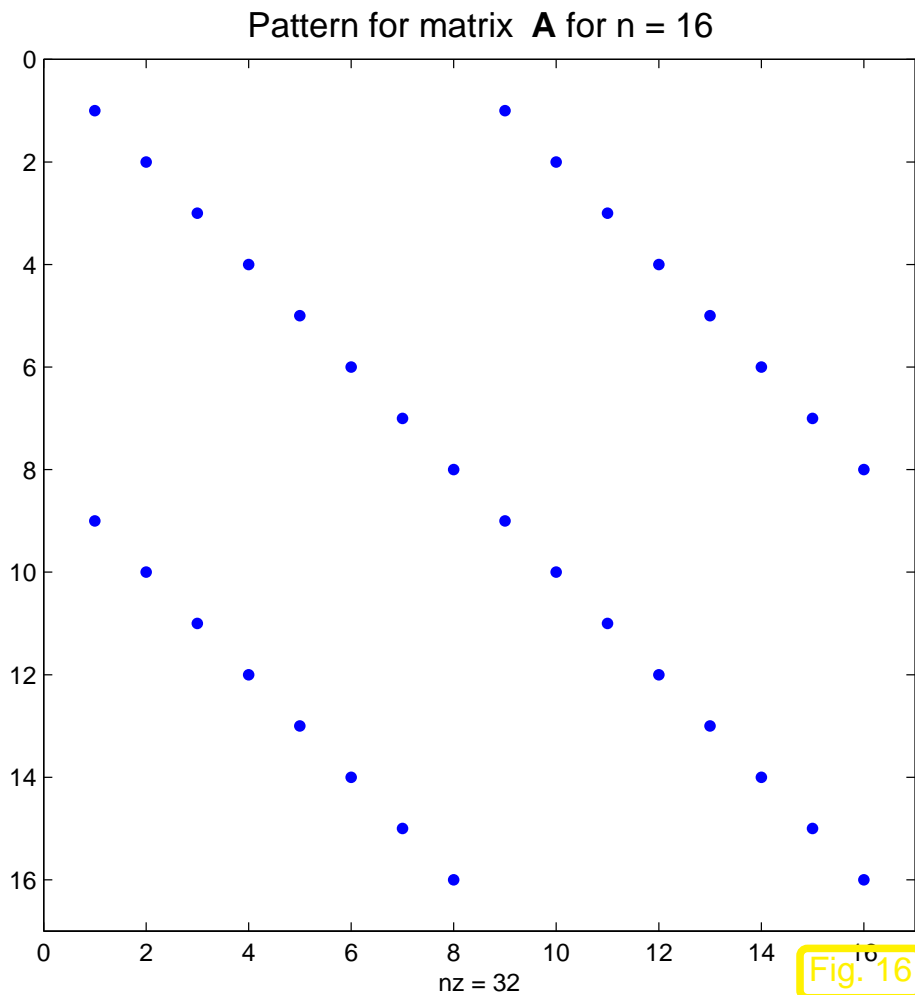
| | | | | | | | | | | | | | |
|---------|----|----|---|---|----|----|----|---|---------|----|---|---|----|
| val | 10 | -2 | 3 | 9 | 3 | 7 | 8 | 7 | 3 ... 9 | 13 | 4 | 2 | -1 |
| col_ind | 1 | 5 | 1 | 2 | 6 | 2 | 3 | 4 | 1 ... 5 | 6 | 2 | 5 | 6 |
| row_ptr | 1 | 3 | 6 | 9 | 13 | 17 | 20 | | | | | | |



2.6.2 Sparse matrices in MATLAB

Initialization: `A = sparse(m,n); A = spalloc(m,n,nnz)`
`A = sparse(i,j,s,m,n);`
`A = spdiags(B,d,m,n); A = speye(n); A = spones(S);`

Example 2.6.9 (Accessing rows and columns of sparse matrices).



Code 2.6.10: timing access to rows/columns of a sparse matrix

```

1 figure ; spy(spdiags(repmat([-1 2 5],16,1),[-8,0,8],16,16)); %
2 title('Pattern for matrix {\bf A} for n = 16','fontsize',14);
3 print -depsc2 '../PICTURES/spdiagsmatspy.eps';
4
5 t = [];
6 for i=1:20
7     n = 2^i; m = n/2;

```

```
8 A = spdiags(repmat([-1 2 5],n,1),[-n/2,0,n/2],n,n); %
9
10 t1 = inf; for j=1:5, tic; v = A(m,:)+j; t1 = min(t1,toc); end
11 t2 = inf; for j=1:5, tic; v = A(:,m)+j; t2 = min(t2,toc); end
12 t = [t; size(A,1), nnz(A), t1, t2 ];
13 end
14
15 figure;
16 loglog(t(:,1),t(:,3),'r+-', t(:,1),t(:,4),'b*-',...
17         t(:,1),t(1,3)*t(:,1)/t(1,1),'k-');
18 xlabel('{\bf size n of sparse quadratic matrix}','fontsize',14);
19 ylabel('{\bf access time [s]}','fontsize',14);
20 legend('row access','column
        access','O(n)','location','northwest');
21
22 print -depsc2 '../PICTURES/sparseaccess.eps';
```



Example 2.6.11 (Efficient Initialization of sparse matrices in MATLAB).

Code 2.6.12: Initialization of sparse matrices: version I

```
1 A1 = sparse(n,n);
```

```
2 for i=1:n
3     for j=1:n
4         if (abs(i-j) == 1), A1(i,j) = A1(i,j) + 1; end;
5         if (abs(i-j) == round(n/3)), A1(i,j) = A1(i,j) -1; end;
6     end; end
```

Code 2.6.13: Initialization of sparse matrices: version II

```
1 dat = [];
2 for i=1:n
3     for j=1:n
4         if (abs(i-j) == 1), dat = [dat; i,j,1.0]; end;
5         if (abs(i-j) == round(n/3)), dat = [dat; i,j,-1.0];
6     end; end; end;
7 A2 = sparse(dat(:,1),dat(:,2),dat(:,3),n,n);
```

Code 2.6.14: Initialization of sparse matrices: version III

```
1 dat = zeros(6*n,3); k = 0;
2 for i=1:n
3     for j=1:n
4         if (abs(i-j) == 1), k=k+1; dat(k,:) = [i,j,1.0];
5         end;
6         if (abs(i-j) == round(n/3))
7             k=k+1; dat(k,:) = [i,j,-1.0];
8         end;
```

```

9  end; end;
10 A3 = sparse(dat(1:k,1),dat(1:k,2),dat(1:k,3),n,n);

```

Code 2.6.15: Initialization of sparse matrices: driver script

```

1  % Driver routine for initialization of sparse matrices
2  K = 3; r = [];
3  for n=2.^(8:14)
4      t1= 1000; for k=1:K, fprintf('sparse1, %d, %d\n',n,k); tic;
5          sparse1; t1 = min(t1,toc); end
6          t2= 1000; for k=1:K, fprintf('sparse2, %d, %d\n',n,k); tic;
7              sparse2; t2 = min(t2,toc); end
8              t3= 1000; for k=1:K, fprintf('sparse3, %d, %d\n',n,k); tic;
9                  sparse3; t3 = min(t3,toc); end
10                 r = [r; n, t1 , t2, t3];
11             end
12
13 loglog(r(:,1),r(:,2),'r*',r(:,1),r(:,3),'m+',r(:,1),r(:,4),'b^');
14 xlabel('{\bf matrix size n}','fontsize',14);
15 ylabel('{\bf time [s]}','fontsize',14);
16 legend('Initialization I','Initialization II','Initialization
17         III',...
18         'location','northwest');
19 print -depsc2 '../PICTURES/sparseinit.eps';

```

Timings:

- Linux lions 2.6.16.27-0.9-smp #1 SMP Tue Feb 13 09:35:18 UTC 2007 i686 i686 i386 GNU/Linux
- CPU: Genuine Intel(R) CPU T2500 2.00GHz
- MATLAB 7.4.0.336 (R2007a)

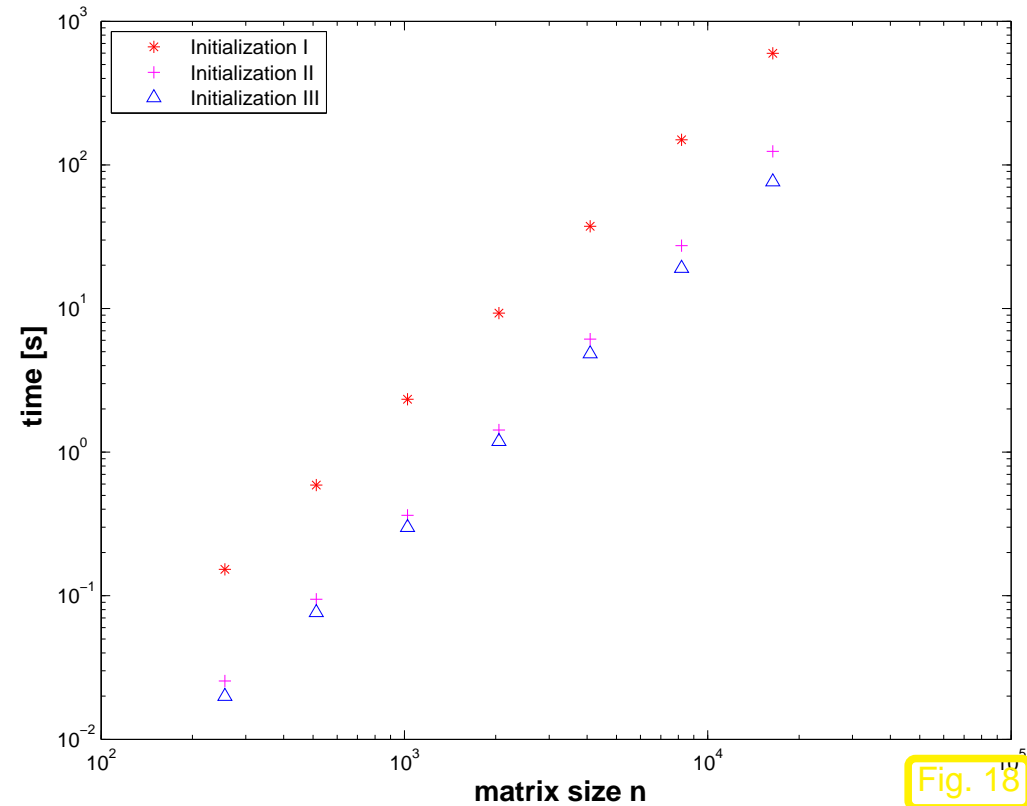


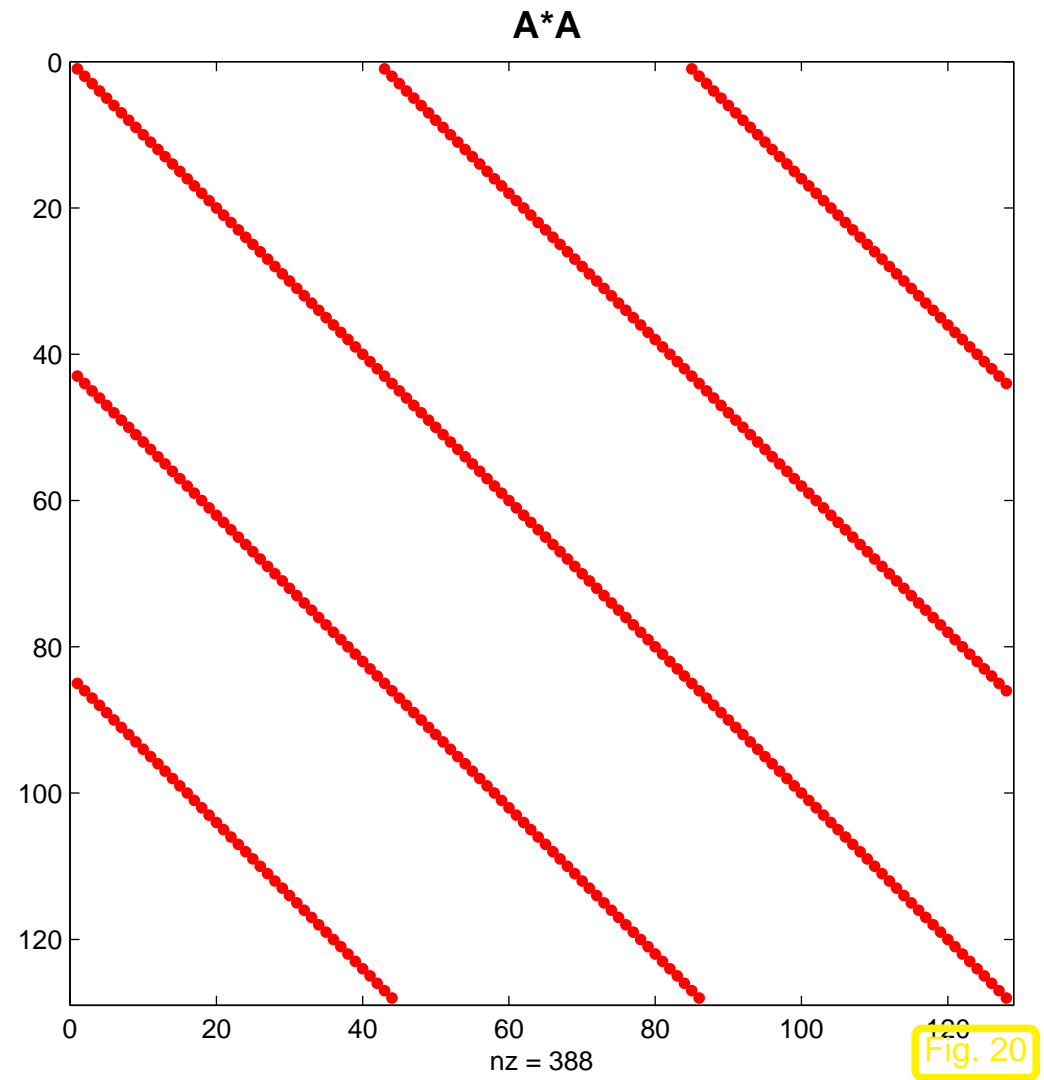
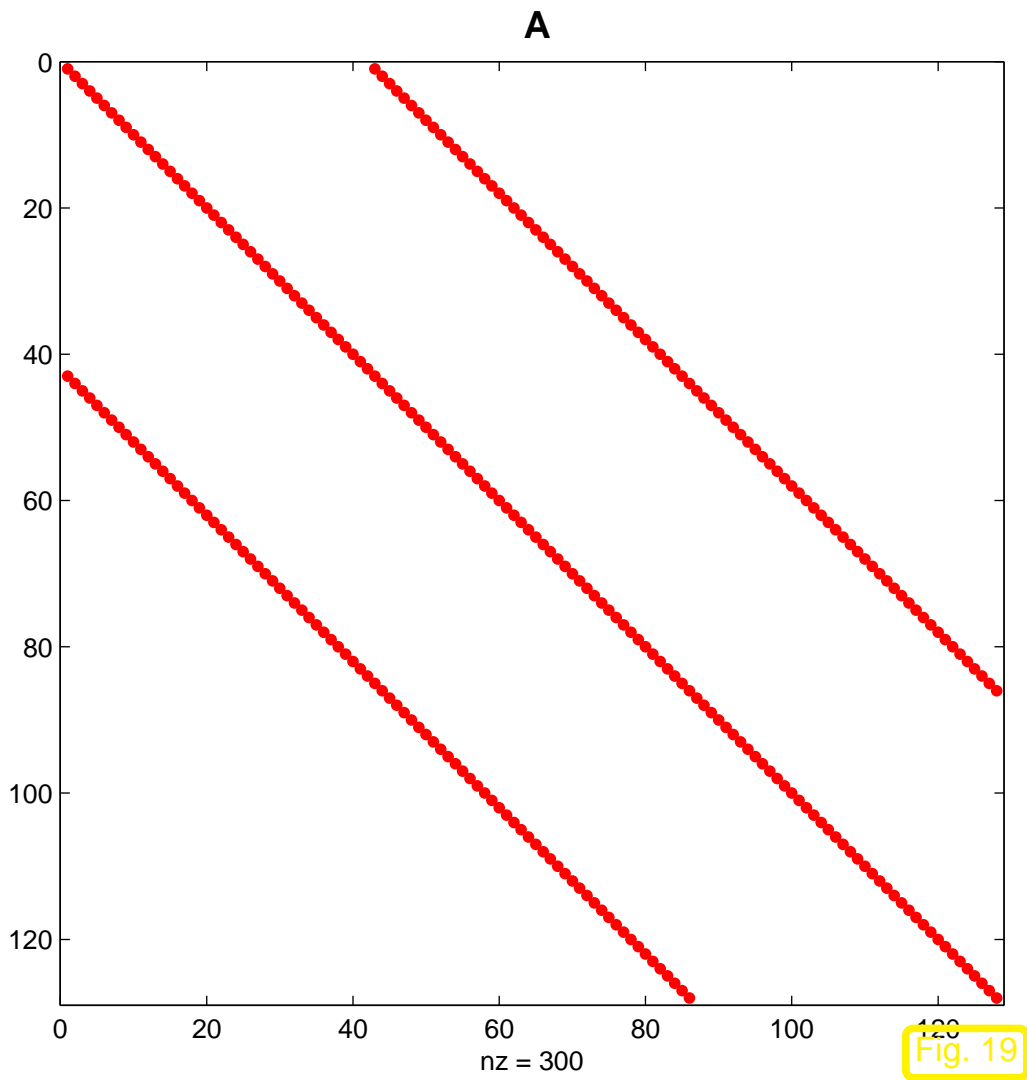
Fig. 18



Example 2.6.16 (Multiplication of sparse matrices).

Sparse matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ initialized by

```
A = spdiags([ (1:n)' , ones(n,1) , (n:-1:1)' ] , ...
            [-floor(n/3) , 0 , floor(n/3)] , n , n);
```



➤ A^2 is still a sparse matrix (\rightarrow Notion 2.6.1)

runtimes for matrix multiplication $A*A$ in MATLAB
(`tic/toc` timing) \triangleright

(same platform as in Ex. 2.6.11)

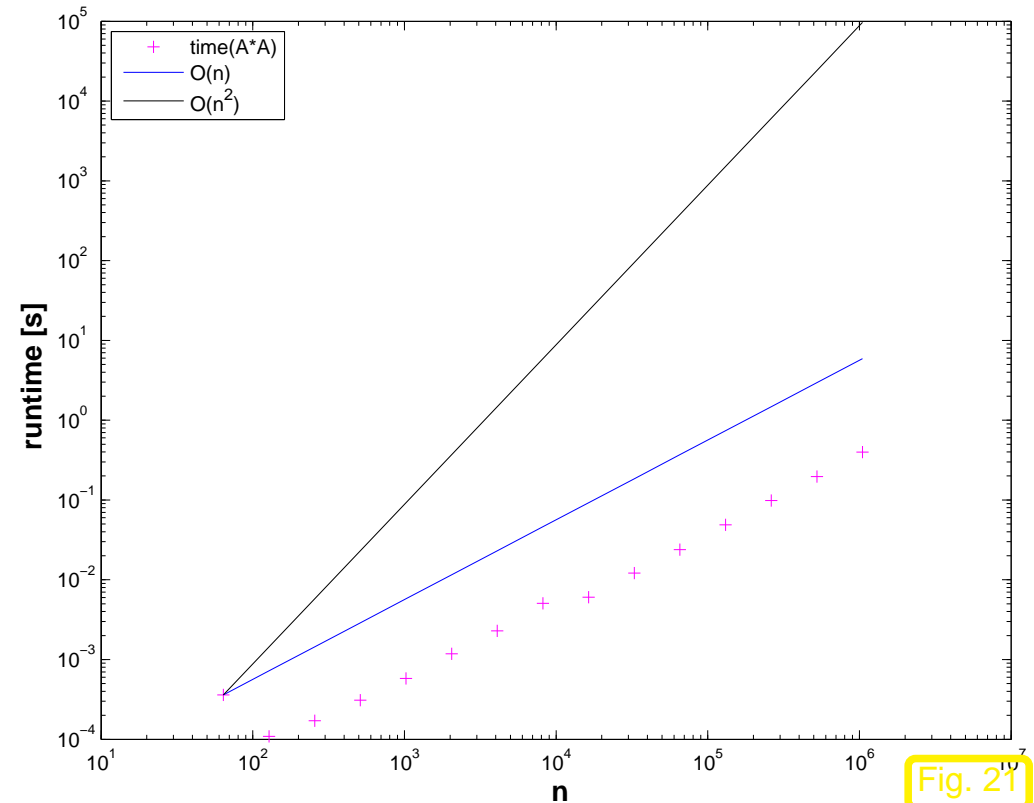


Fig. 21

\blacktriangleright $O(n)$ asymptotic complexity: “optimal” implementation ! \diamond

Remark 2.6.17 (Silly MATLAB).

A strange behavior of MATLAB from version R2010a:

```

1 % MATLAB script demonstrating the awkward effects of treating entries of
2 % sparse matrices as sparse matrices themselves.
3 A = spdiags(ones(100,3),-1:1,100,100);
4 b = A(1,1),

```

```

5 c = full(b),
6 whos('b','c');
7 sum=0; tic; for i=1:1e6, sum = sum + b; end, toc
8 sum=0; tic; for i=1:1e6, sum = sum + c; end, toc

```



2.6.3 LU-factorization of sparse matrices

Example 2.6.18 (LU-factorization of sparse matrices).

$$\mathbf{A} = \left(\begin{array}{cccc|cccc}
 3 & -1 & & & -1 & & & \\
 -1 & \ddots & \ddots & & & \ddots & & \\
 & \ddots & \ddots & -1 & & & & \\
 -1 & & -1 & 3 & & & & -1 \\
 & & & & 3 & -1 & & \\
 & & & & -1 & \ddots & \ddots & \\
 & & & & & \ddots & \ddots & -1 \\
 & & & & & & -1 & 3
 \end{array} \right) \in \mathbb{R}^{n,n}, n \in \mathbb{N}$$

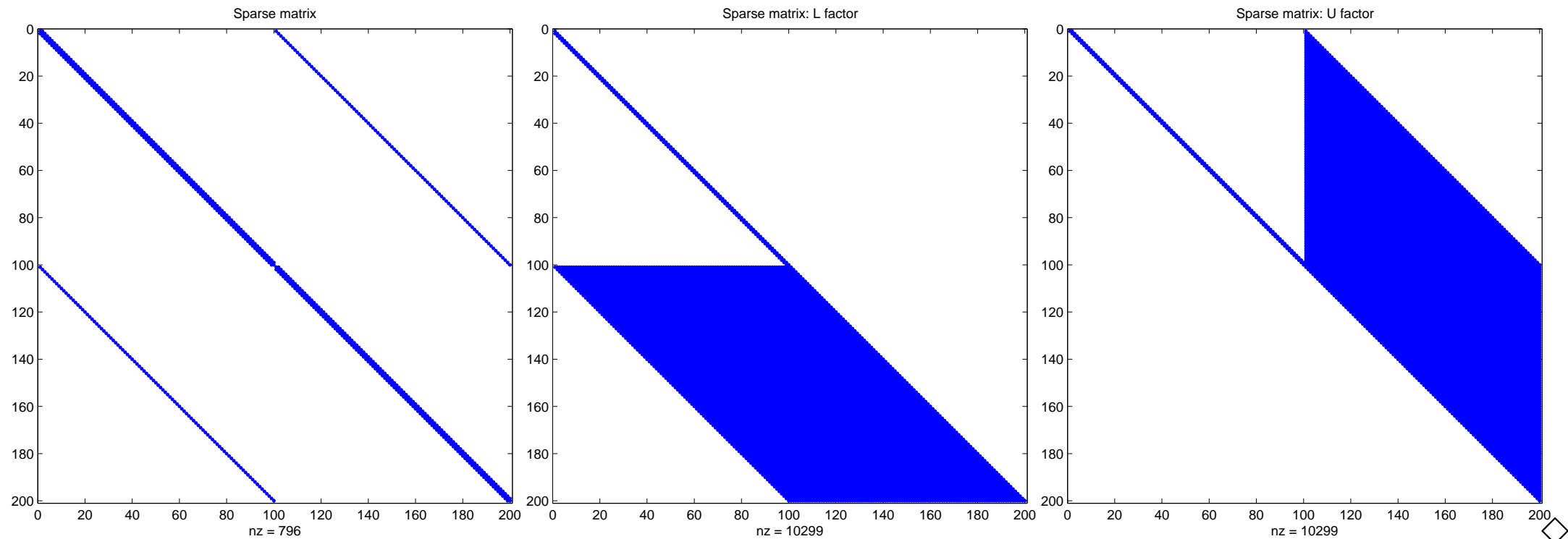
Code 2.6.19: LU-factorization of sparse matrix

NumCSE,
autumn
2010

```

1 % Demonstration of fill-in for LU-factorization of sparse matrices
2 n = 100;
3 A = [ gallery('tridiag',n,-1,3,-1), speye(n); speye(n) ,
      gallery('tridiag',n,-1,3,-1)]; [L,U,P] = lu(A);
4 figure; spy(A); title('Sparse matrix'); print -depsc2
  './PICTURES/sparseA.eps';
5 figure; spy(L); title('Sparse matrix: L factor'); print -depsc2
  './PICTURES/sparseL.eps';
6 figure; spy(U); title('Sparse matrix: U factor'); print -depsc2
  './PICTURES/sparseU.eps';

```

R. Hiptmair
rev 38355,
October 7,
2011

\mathbf{A} sparse $\not\Rightarrow$ LU -factors sparse

Definition 2.6.20 (Fill-in).

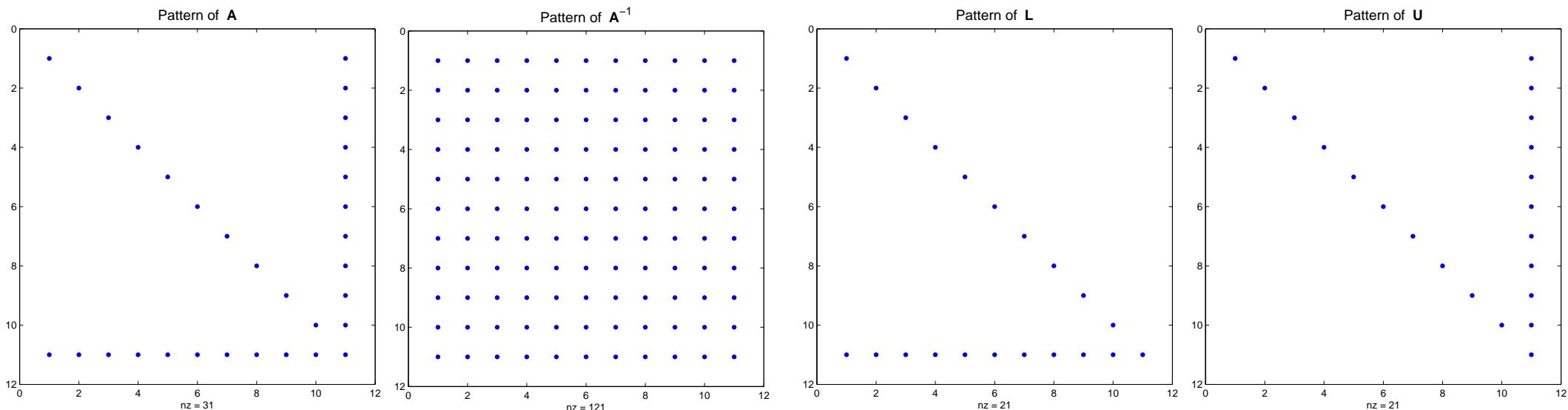
Let $\mathbf{A} = \mathbf{L}\mathbf{U}$ be an LU -factorization (\rightarrow Sect. 2.2) of $\mathbf{A} \in \mathbb{K}^{n,n}$. If $l_{ij} \neq 0$ or $u_{ij} \neq 0$ though $a_{ij} = 0$, then we encounter **fill-in** at position (i, j) .

Example 2.6.21 (Sparse LU -factors).

```

1 % Simple example for dense inverse despite sparse LU-factors
2 A = [diag(1:10), ones(10,1); ones(1,10), 2];
3 [L,U] = lu(A); spy(A); spy(L); spy(U); spy(inv(A));

```



\mathbf{L}, \mathbf{U} sparse $\not\Rightarrow \mathbf{A}^{-1}$ sparse !



Example 2.6.22 (“arrow matrix”).

$$\mathbf{A} = \left(\begin{array}{c|ccc} \alpha & & & \mathbf{b}^\top \\ \hline & & & \\ \mathbf{c} & & & \mathbf{D} \end{array} \right),$$

$\alpha \in \mathbb{R}$,
 $\mathbf{b}, \mathbf{c} \in \mathbb{R}^{n-1}$,
 $\mathbf{D} \in \mathbb{R}^{n-1, n-1}$ regular diagonal matrix, \rightarrow Def. 2.2.3

(2.6.23)

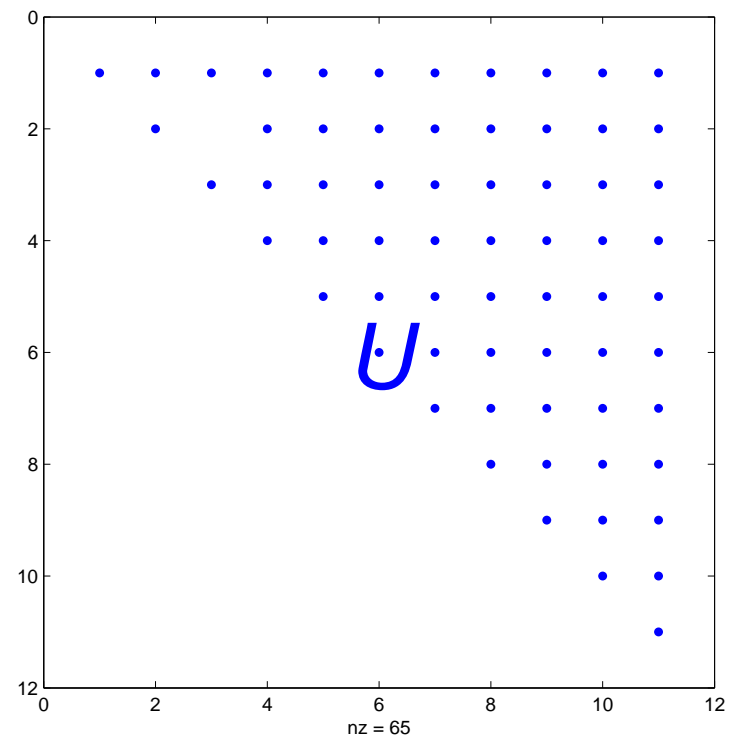
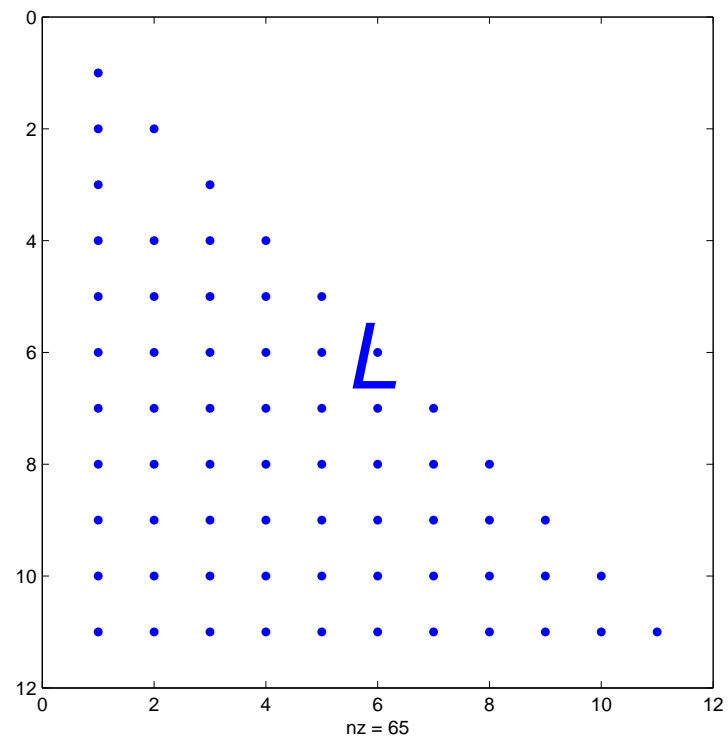
Run algorithm 2.3.8 (Gaussian elimination without pivoting):

- factor matrices with $O(n^2)$ non-zero entries.
- computational costs: $O(n^3)$

Code 2.6.24: LU-factorization of arrow matrix

```
1 n = 10; A = [ n+1, (n:-1:1) ;
               ones(n,1), eye(n,n) ];
2 [L,U,P] = lu(A); spy(L); spy(U);
```

Obvious fill-in (\rightarrow Def. 2.6.20)



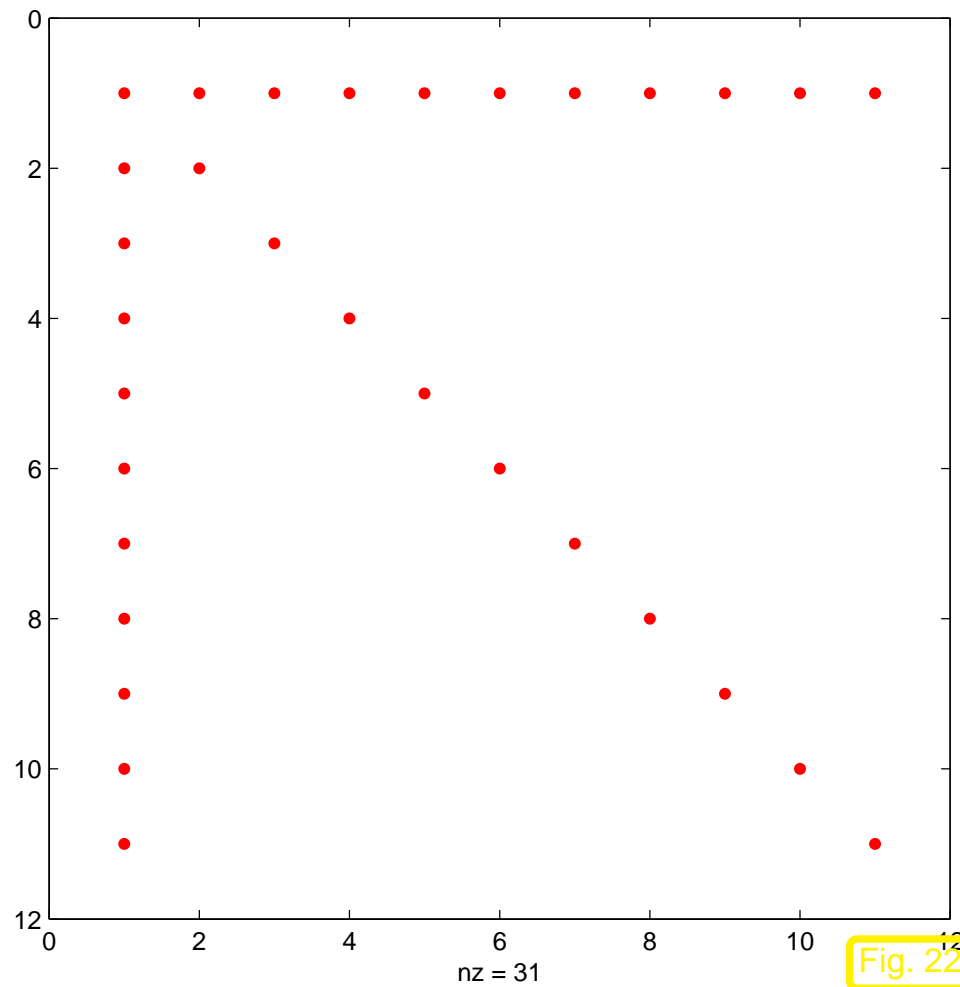


Fig. 22

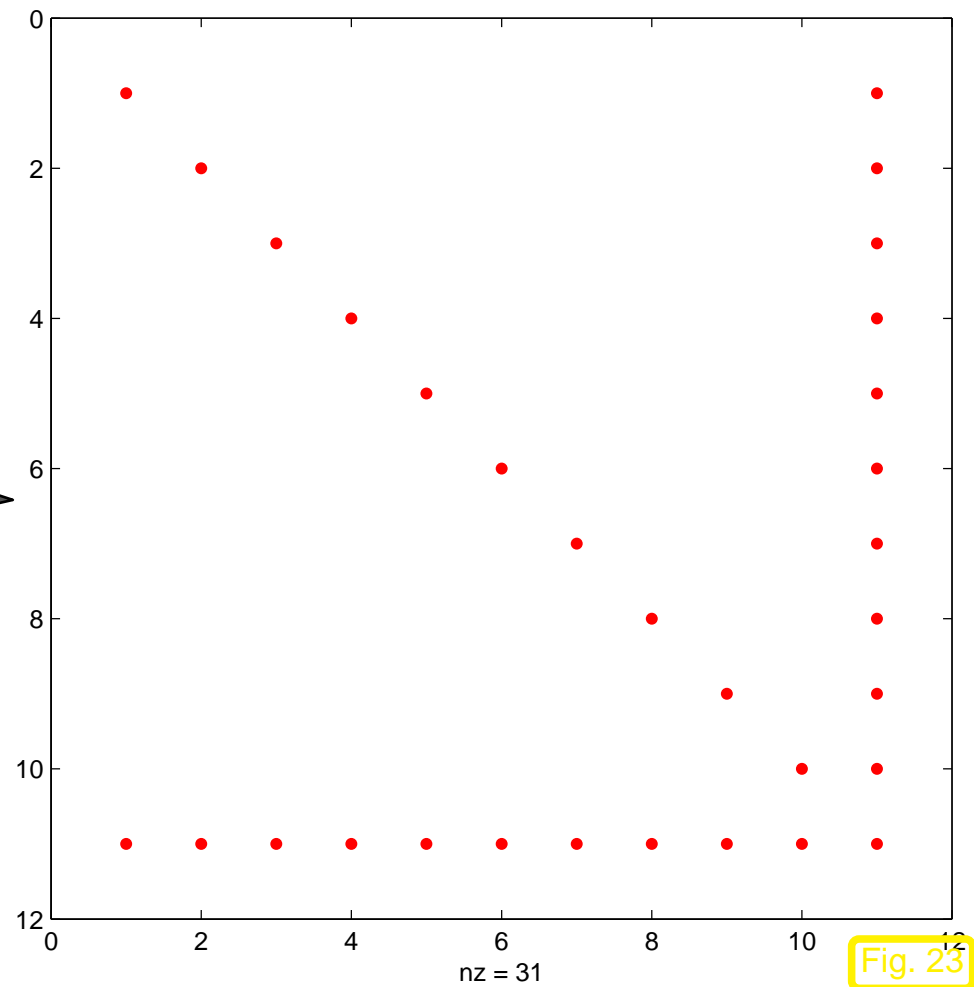


Fig. 23

Code 2.6.25: Permuting arrow matrix, see Figs. 29, 30

```

1 n = 10; A = [ n+1, (n:-1:1) ; ones(n,1), eye(n,n) ];
2 % Permutation matrix (→ Def. 2.3.12) encoding cyclic permutation
3 P = [ zeros(n,1), eye(n) ; 1, zeros(1,n) ];
4
5 figure('name','A');
6 spy(A,'r. '); print -depsc '../PICTURES/InvArrowSpy.eps';
7 figure('name','PAPT');

```

```
spy(P*A*P', 'r. '); print -depsc '../PICTURES/ArrowSpy.eps';
```

Solving LSE $Ax = y$ with A from 2.6.23: two MATLAB codes

“naive” implementation via “\”:

Code 2.6.28: LSE with arrow matrix, implementation I

```
1 function x =
   sa1(alpha,b,c,d,y)
2 A = [alpha, b'; c, diag(d)];
3 x = A\y;
```

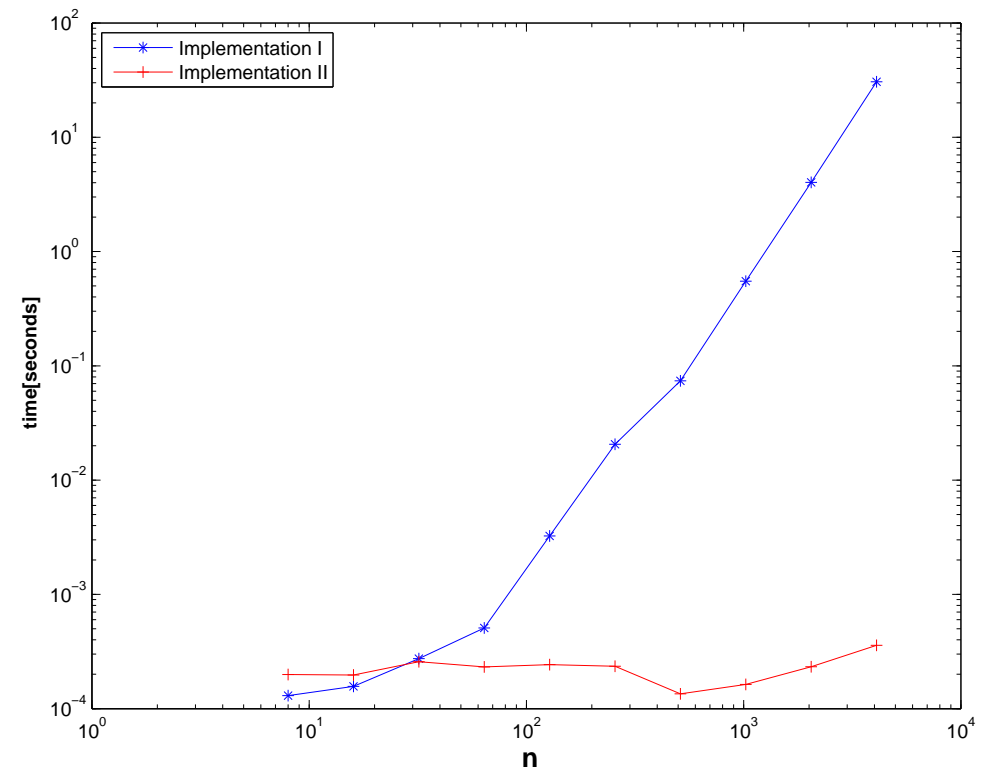
“structure aware” implementation:

Code 2.6.30: LSE with arrow matrix, implementation II

```
1 function x =
   sa2(alpha,b,c,d,y)
2 z = b./d;
3 xi = (y(1) -
   dot(z,y(2:end)))...
4       /(alpha-dot(z,c));
5 x = [xi; (y(2:end)-xi*c)./d];
```


Measuring run times:

```
t = [];
for i=3:12
    n = 2^n; alpha = 2;
    b = ones(n,1); c = (1:n)';
    d = -ones(n,1); y = (-1).^(1:(n+1))';
    tic; x1 = sa1(alpha,b,c,d,y); t1 = toc;
    tic; x2 = sa2(alpha,b,c,d,y); t2 = toc;
    t = [t; n t1 t2];
end
loglog(t(:,1),t(:,2), ...
... 'b-*',t(:,1),t(:,3),'r-+');
```



Platform as in Ex. 2.6.11

MATLAB can do much better !

R. Hiptmair

rev 38355,
October 7,
2011

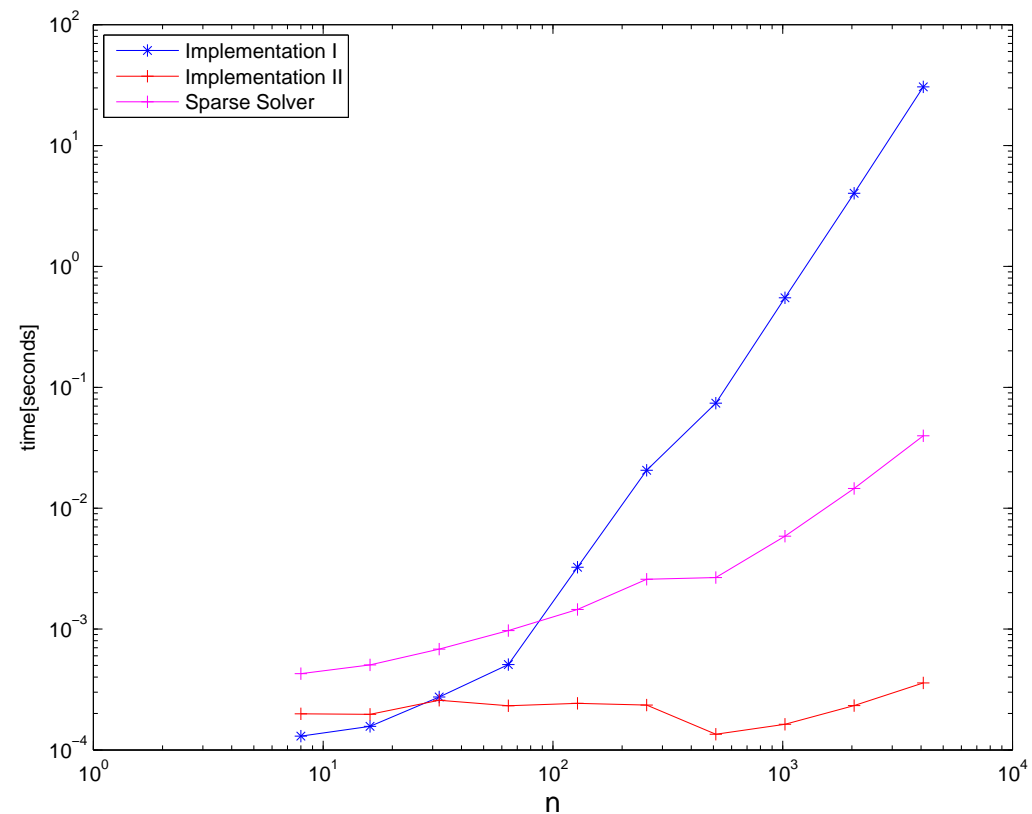
Use sparse matrix format:

Code 2.6.31: sparse solver for arrow matrix

```

1 function x =
   sa3(alpha,b,c,d,y)
2 n = length(d);
3 A = [alpha, b'; ...
4      c, spdiags(d,0,n,n)];
5 x = A\y;

```



Exploit structure of (sparse) linear systems of equations !



Caution:

stability at risk

Example 2.6.32 (Pivoting destroys sparsity).

Code 2.6.33: fill-in due to pivoting

NumCSE,
autumn
2010

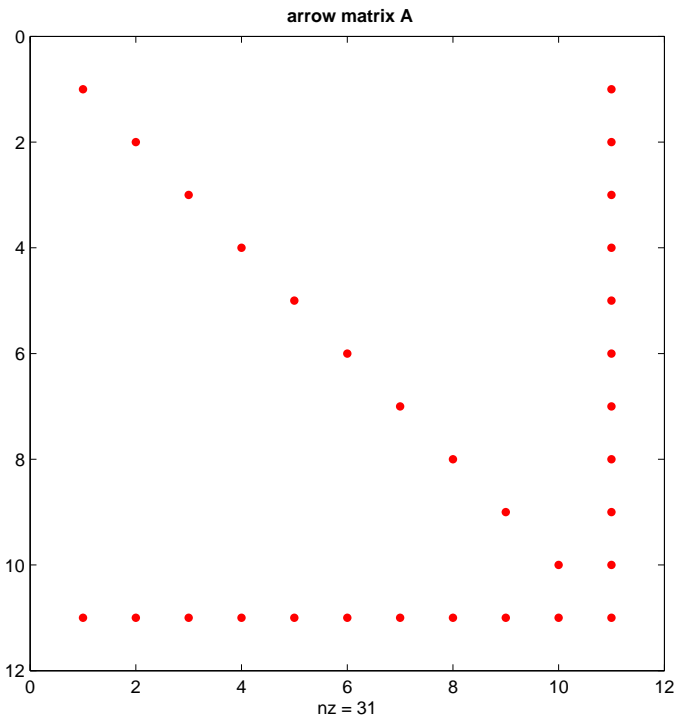
```

1 % Study of fill-in with LU-factorization due to pivoting
2 n = 10; D = diag(1./(1:n));
3 A = [ D , 2*ones(n,1); 2*ones(1,n) , 2];
4 [L,U,P] = lu(A);
5 figure; spy(A,'r'); title ('{\bf arrow matrix A}');
6 print -depsc2 '../PICTURES/fillinpivotA.eps';
7 figure; spy(L,'r'); title ('{\bf L factor}');
8 print -depsc2 '../PICTURES/fillinpivotL.eps';
9 figure; spy(U,'r'); title ('{\bf U factor}');
10 print -depsc2 '../PICTURES/fillinpivotU.eps';

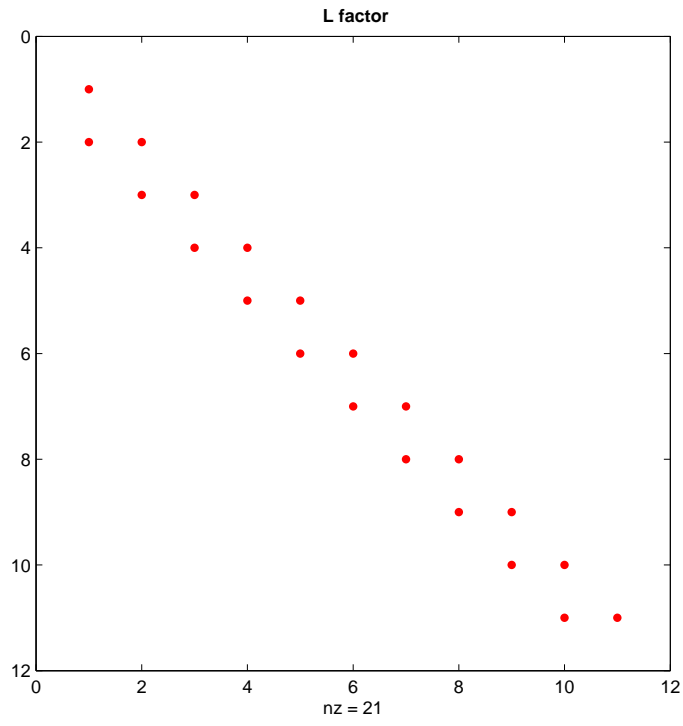
```

$$\mathbf{A} = \begin{pmatrix} 1 & & & 2 \\ & \frac{1}{2} & & 2 \\ & & \dots & \vdots \\ & & & \frac{1}{10} & 2 \\ 2 & & \dots & & 2 \end{pmatrix} \rightarrow \text{arrow matrix, Ex. 2.6.21}$$

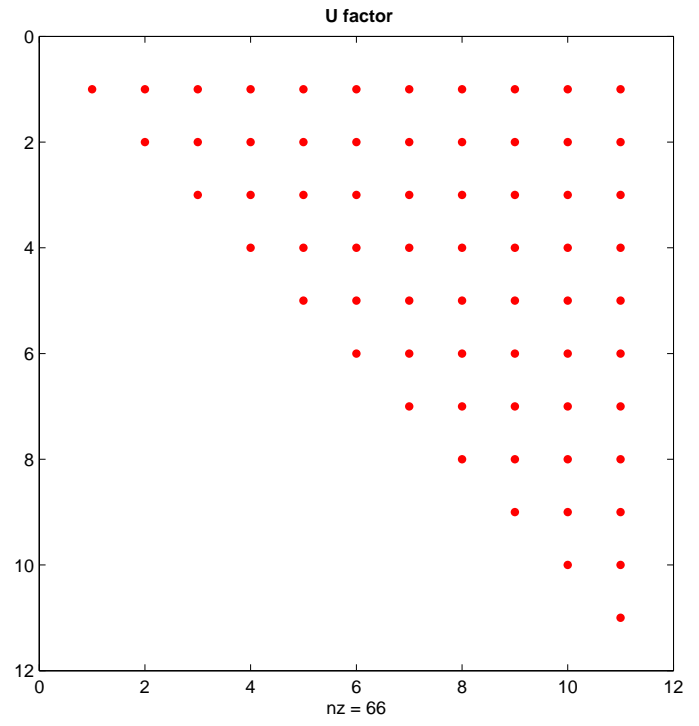
R. Hiptmair
rev 38355,
October 7,
2011



A



L



U



Definition 2.6.34 (bandwidth).

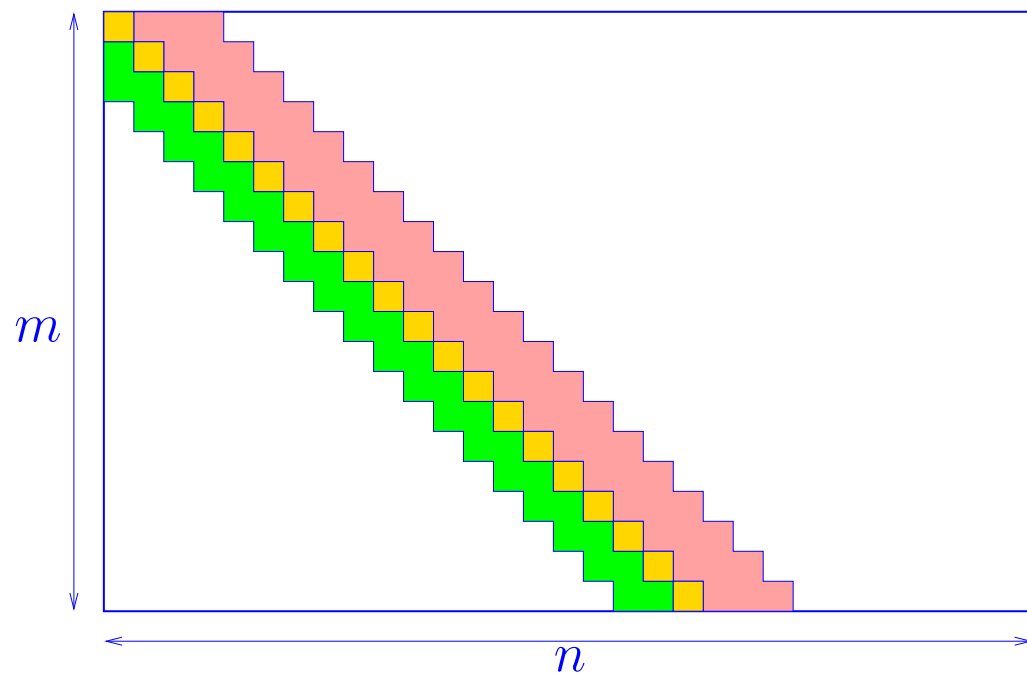
For $\mathbf{A} = (a_{ij})_{i,j} \in \mathbb{K}^{m,n}$ we call

$$\overline{m}(\mathbf{A}) := \min\{k \in \mathbb{N} : j - i > k \Rightarrow a_{ij} = 0\} \text{ upper bandwidth ,}$$

$$\underline{m}(\mathbf{A}) := \min\{k \in \mathbb{N} : i - j > k \Rightarrow a_{ij} = 0\} \text{ lower bandwidth .}$$

$$m(\mathbf{A}) := \overline{m}(\mathbf{A}) + \underline{m}(\mathbf{A}) + 1 = \text{bandwidth von } \mathbf{A} \text{ (ger.: Bandbreite)}$$

- $m(\mathbf{A}) = 1 \quad \triangleright \quad \mathbf{A}$ diagonal matrix, \rightarrow Def. 2.2.3
- $\overline{m}(\mathbf{A}) = \underline{m}(\mathbf{A}) = 1 \quad \triangleright \quad \mathbf{A}$ tridiagonal matrix
- More general: $\mathbf{A} \in \mathbb{R}^{n,n}$ with $m(\mathbf{A}) \ll n \hat{=} \text{banded matrix}$



- : diagonal
- : super-diagonals
- : sub-diagonals

$$\triangleleft \quad \overline{m}(\mathbf{A}) = 3, \underline{m}(\mathbf{A}) = 2$$

► for banded matrix $\mathbf{A} \in \mathbb{K}^{m,n}$: $\text{nnz}(\mathbf{A}) \leq \min\{m, n\}m(\mathbf{A})$

MATLAB function for creating banded matrices:

dense matrix : $\mathbf{X} = \text{diag}(\mathbf{v}) ;$

sparse matrix : $\mathbf{X} = \text{spdiags}(\mathbf{B}, \mathbf{d}, m, n) ;$ (sparse storage !)

tridiagonal matrix : $\mathbf{X} = \text{gallery}('tridiag', c, d, e) ;$ (sparse storage !)

Definition 2.6.35 (Matrix envelope (ger.: Hülle)).

For $\mathbf{A} \in \mathbb{K}^{n,n}$ define

row bandwidth $m_i^R(\mathbf{A}) := \max\{0, i - j : a_{ij} \neq 0, 1 \leq j \leq n\}, i \in \{1, \dots, n\}$

column bandwidth $m_j^C(\mathbf{A}) := \max\{0, j - i : a_{ij} \neq 0, 1 \leq i \leq n\}, j \in \{1, \dots, n\}$

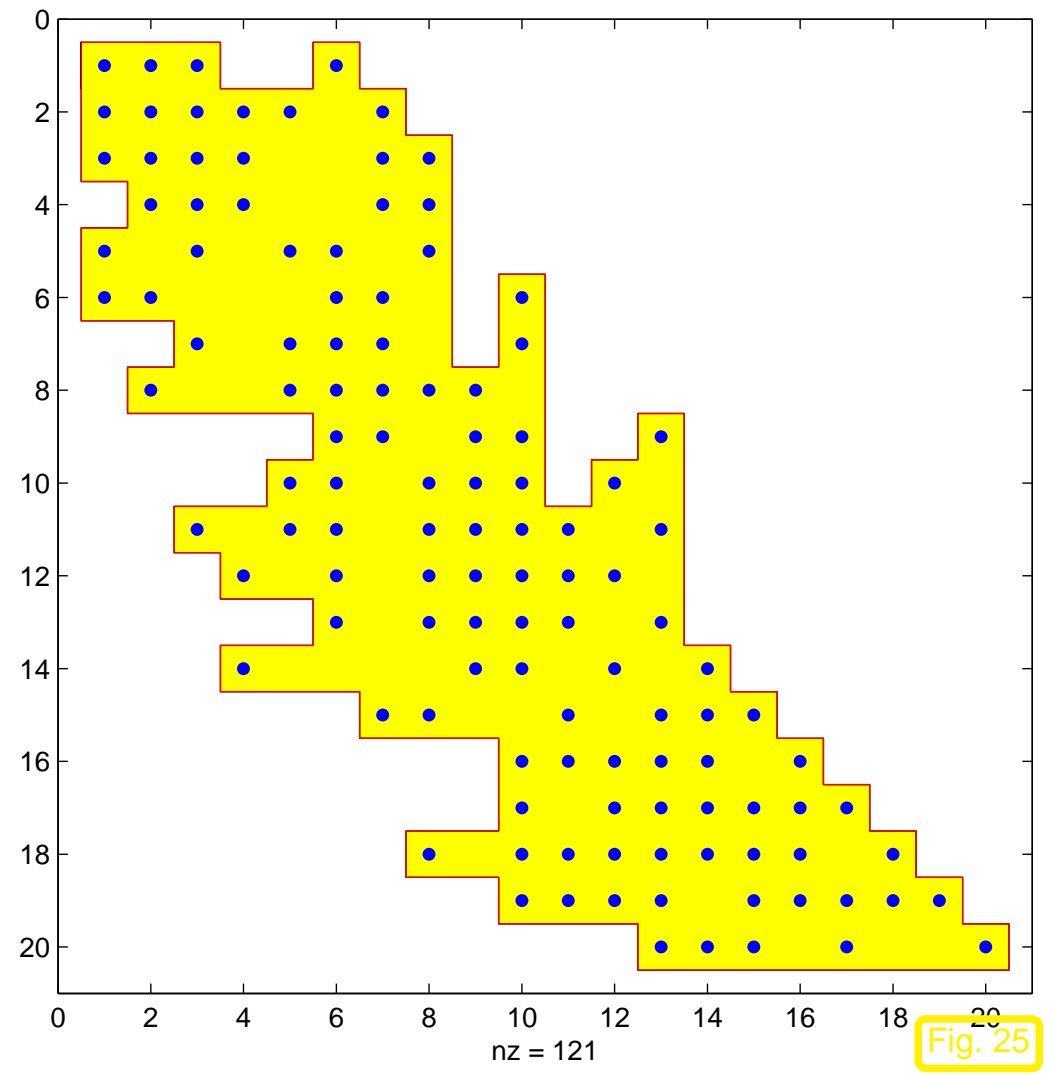
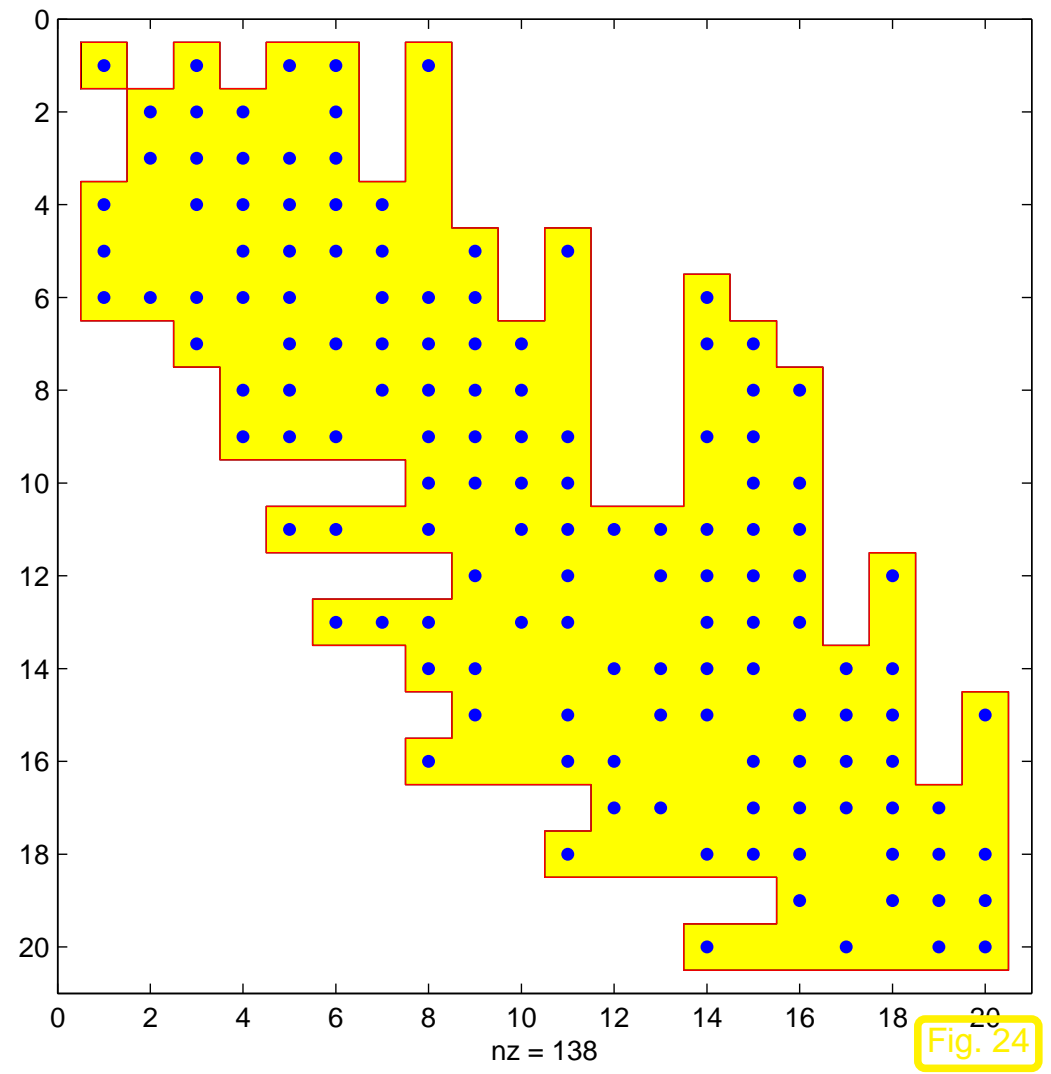
envelope $\text{env}(\mathbf{A}) := \left\{ (i, j) \in \{1, \dots, n\}^2 : \begin{array}{l} i - m_i^R(\mathbf{A}) \leq j \leq i, \\ j - m_j^C(\mathbf{A}) \leq i \leq j \end{array} \right\}$

Example 2.6.36 (Envelope of a matrix).

$$\mathbf{A} = \begin{pmatrix} * & 0 & * & 0 & 0 & 0 & 0 \\ 0 & * & 0 & 0 & * & 0 & 0 \\ * & 0 & * & 0 & 0 & 0 & * \\ 0 & 0 & 0 & * & * & 0 & * \\ 0 & * & 0 & * & * & * & 0 \\ 0 & 0 & 0 & 0 & * & * & 0 \\ 0 & 0 & * & * & 0 & 0 & * \end{pmatrix} \begin{array}{l} m_1^R(A) = 0 \\ m_2^R(A) = 0 \\ m_3^R(A) = 2 \\ m_4^R(A) = 0 \\ m_5^R(A) = 3 \\ m_6^R(A) = 1 \\ m_7^R(A) = 4 \end{array}$$

$\text{env}(A) =$ red elements

$*$ $\hat{=}$ non-zero matrix entry $a_{ij} \neq 0$



Theorem 2.6.37 (Envelope and fill-in). \rightarrow [51, Sect. 3.9]

If $\mathbf{A} \in \mathbb{K}^{n,n}$ is regular with LU-factorization $\mathbf{A} = \mathbf{L}\mathbf{U}$, then fill-in (\rightarrow Def. 2.6.20) is confined to $\text{env}(\mathbf{A})$.

Gaussian elimination *without pivoting*

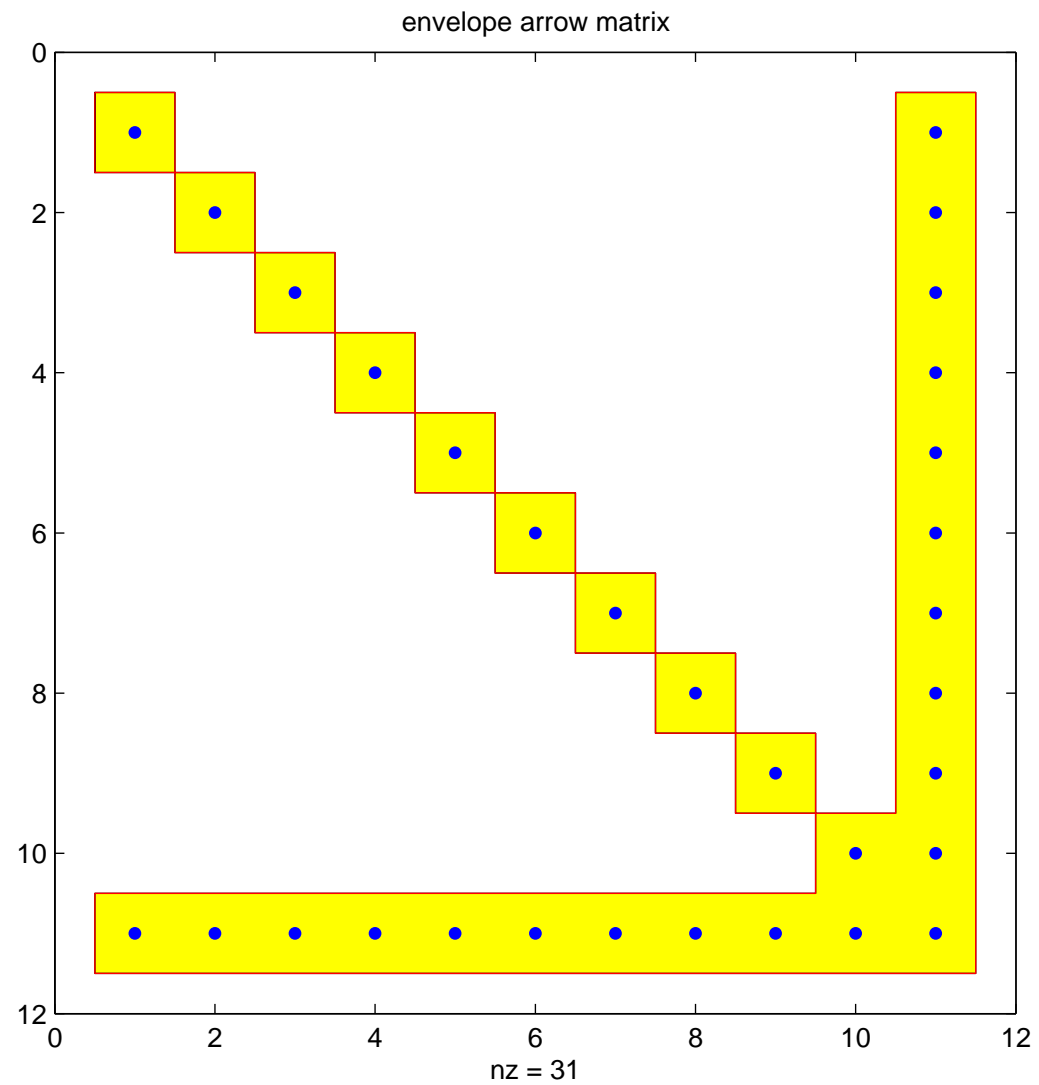
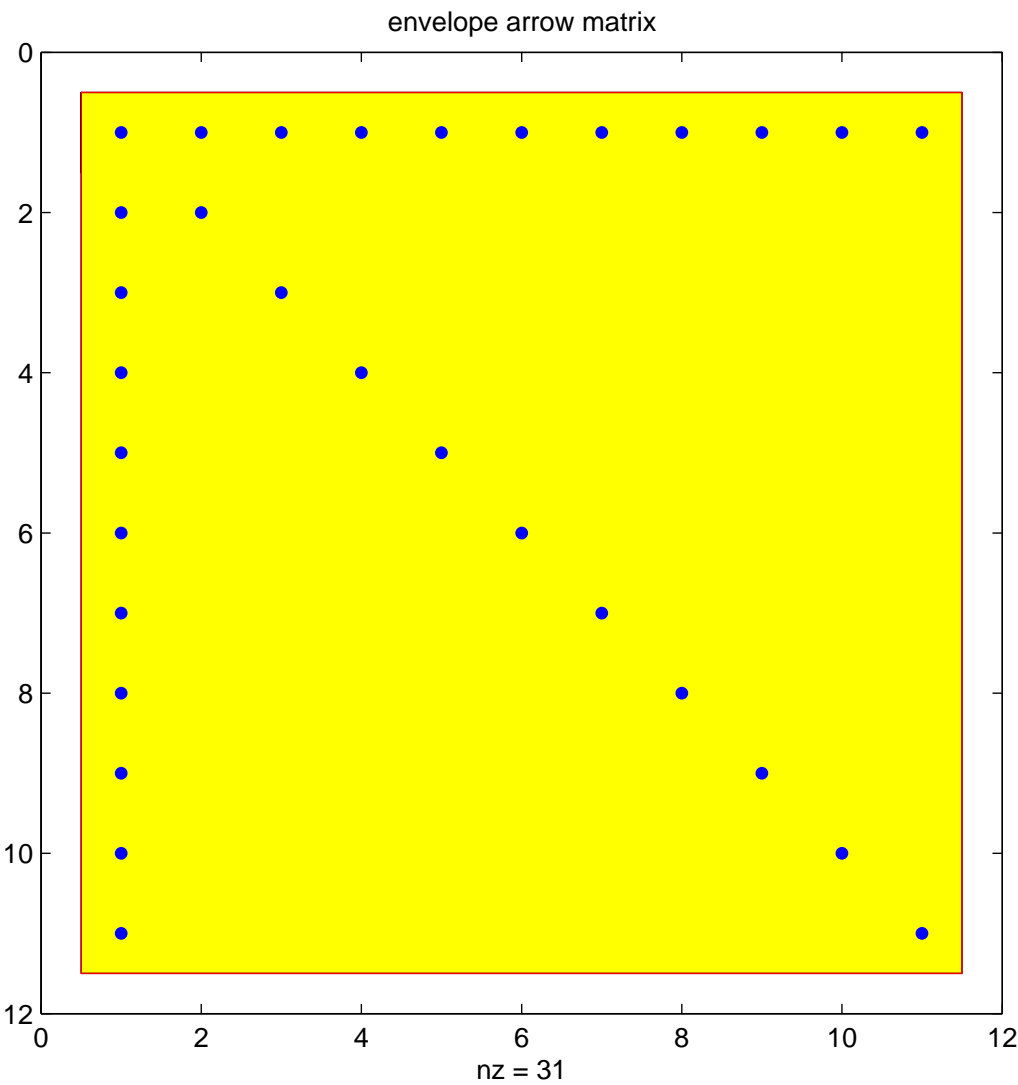
► **Envelope-aware LU-factorization:**

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 2.6.47 (Reducing bandwidth by row/column permutations).

Recall: cyclic permutation of rows/columns of arrow matrix \rightarrow Ex. 2.6.22



Another example: Reflection at cross diagonal ➤ reduction of $\# \text{env}(\mathbf{A})$

$$\begin{pmatrix} * & 0 & 0 & * & * & * \\ 0 & * & 0 & 0 & 0 & 0 \\ * & 0 & * & 0 & 0 & 0 \\ * & 0 & 0 & * & * & * \end{pmatrix} \rightarrow \begin{pmatrix} * & * & * & 0 & 0 & * \\ * & * & * & 0 & 0 & * \\ 0 & 0 & 0 & * & 0 & 0 \\ * & * & * & 0 & 0 & * \end{pmatrix}$$

$$i \leftarrow N + 1 - i$$

$$\# \text{env}(\mathbf{A}) = 30$$

$$\# \text{env}(\mathbf{A}) = 22$$



Example 2.6.48 (Reducing fill-in by reordering).

M: 114×114 symmetric matrix (from computational PDEs)

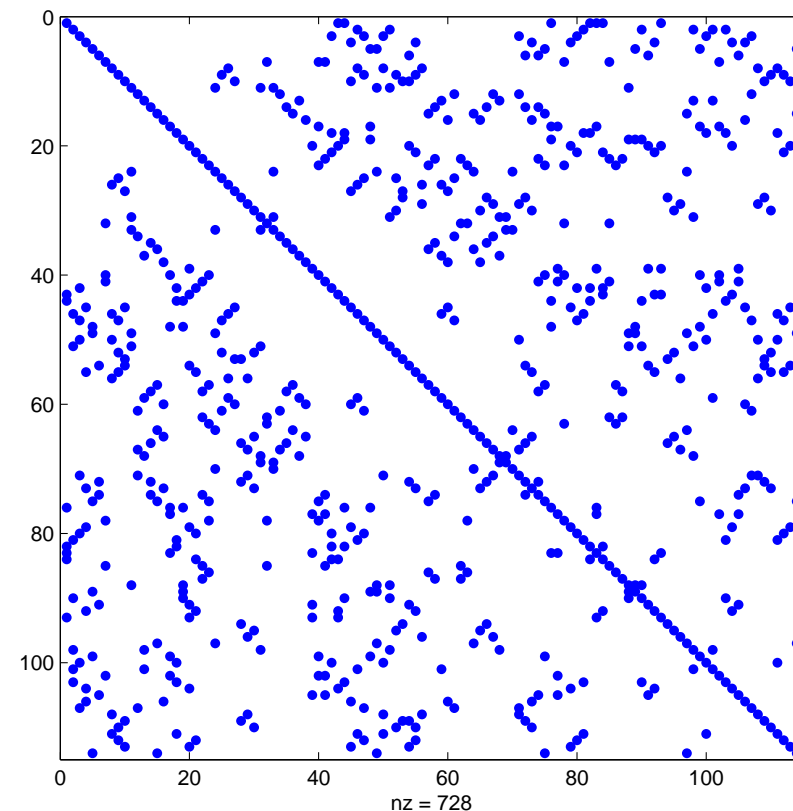
Code 2.6.49: preordering in MATLAB

```

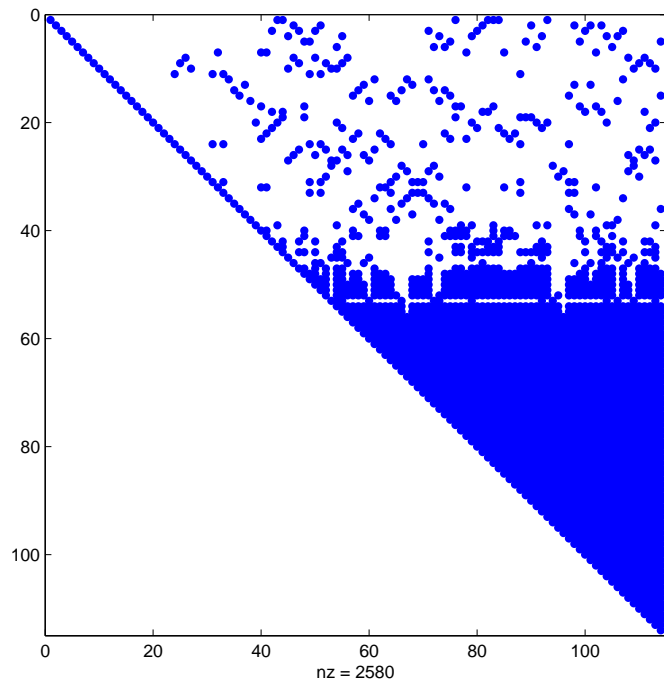
1 spy(M);
2 [L,U] = lu(M); spy(U);
3 r = symrcm(M);
4 [L,U] = lu(M(r,r)); spy(U);
5 m = symamd(M);
6 [L,U] = lu(M(m,m)); spy(U);
    
```

Pattern of **M** →

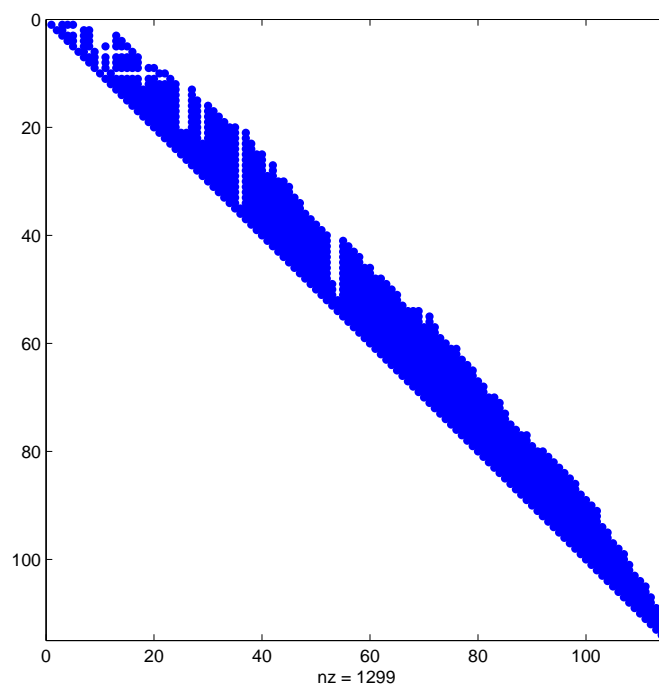
(Here: no row swaps from pivoting !)



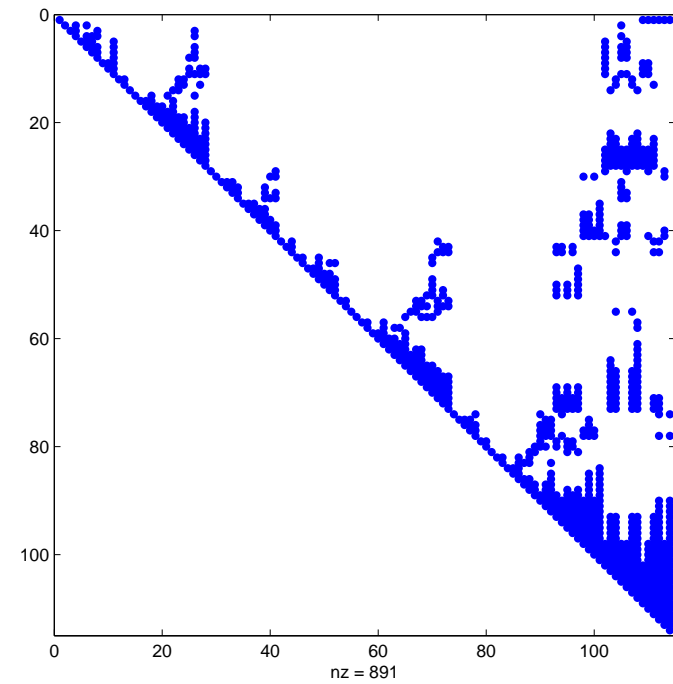
Examine patterns of LU-factors (→ Sect. 2.2) after reordering:



no reordering



reverse Cuthill-McKee



approximate minimum degree

Advice: Use numerical libraries for solving LSE with sparse system matrices !

→ SuperLU (<http://www.cs.berkeley.edu/~demmel/SuperLU.html>)

→ UMFPACK (<http://www.cise.ufl.edu/research/sparse/umfpack/>)

→ Pardiso (<http://www.pardiso-project.org/>)

→ Matlab- (on sparse storage formats)

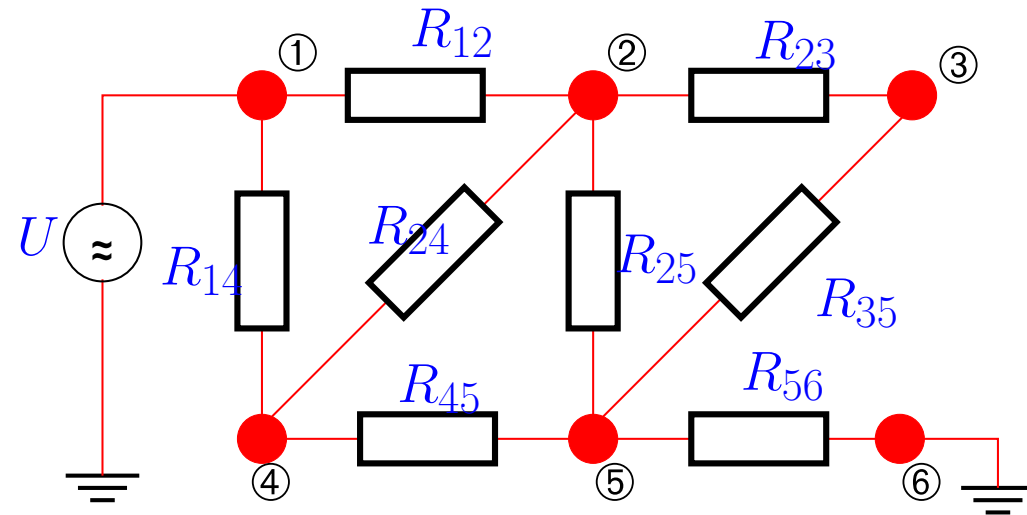
2.7 Stable Gaussian elimination without pivoting

Example 2.7.1 (Diagonally dominant matrices from nodal analysis). → Ex. 2.6.3

Consider:

electrical circuit entirely composed of Ohmic resistors.

Circuit equations from nodal analysis, see Ex. 2.6.3:



$$\begin{aligned}
 \textcircled{2} : & R_{12}^{-1}(U_2 - U_1) + R_{23}^{-1}(U_2 - U_3) - R_{24}^{-1}(U_2 - U_4) + R_{25}^{-1}(U_2 - U_5) = 0, \\
 \textcircled{3} : & R_{23}^{-1}(U_3 - U_2) + R_{35}^{-1}(U_3 - U_5) = 0, \\
 \textcircled{4} : & R_{14}^{-1}(U_4 - U_1) - R_{24}^{-1}(U_4 - U_2) + R_{45}^{-1}(U_4 - U_5) = 0, \\
 \textcircled{5} : & R_{25}^{-1}(U_5 - U_2) + R_{35}^{-1}(U_5 - U_3) + R_{45}^{-1}(U_5 - U_4) + R_{56}^{-1}(U_5 - U_6) = 0, \\
 & U_1 = U, \quad U_6 = 0.
 \end{aligned}$$



$$\begin{pmatrix} \frac{1}{R_{12}} + \frac{1}{R_{23}} + \frac{1}{R_{24}} + \frac{1}{R_{25}} & -\frac{1}{R_{23}} & -\frac{1}{R_{24}} & -\frac{1}{R_{25}} \\ & \frac{1}{R_{23}} + \frac{1}{R_{35}} & 0 & -\frac{1}{R_{35}} \\ & -\frac{1}{R_{24}} & \frac{1}{R_{24}} + \frac{1}{R_{45}} & -\frac{1}{R_{45}} \\ & -\frac{1}{R_{25}} & -\frac{1}{R_{45}} & \frac{1}{R_{22}} + \frac{1}{R_{35}} + \frac{1}{R_{45}} + \frac{1}{R_{56}} \end{pmatrix} \begin{pmatrix} U_2 \\ U_3 \\ U_4 \\ U_5 \end{pmatrix} = \begin{pmatrix} \frac{1}{R_{12}} \\ 0 \\ \frac{1}{R_{14}} \\ 0 \end{pmatrix} U$$

◇

Definition 2.7.7 (Diagonally dominant matrix). \rightarrow [51, 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 2.7.8 (LU-factorization of diagonally dominant matrices).

$$\mathbf{A} \text{ regular, diagonally dominant with positive diagonal} \Leftrightarrow \left\{ \begin{array}{l} \mathbf{A} \text{ has LU-factorization} \\ \updownarrow \\ \text{Gaussian elimination feasible without pivoting}^{(*)} \end{array} \right.$$

Proof. (2.1.11) \rightarrow induction w.r.t. n

Definition 2.7.9 (Symmetric positive definite (s.p.d.) matrices). \rightarrow [13, Def. 3.31], [51, Def. 1.22]

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

$$\mathbf{M} = \mathbf{M}^H \wedge \mathbf{x}^H \mathbf{M} \mathbf{x} > 0 \Leftrightarrow \mathbf{x} \neq \mathbf{0}.$$

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

Lemma 2.7.12. *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.

Theorem 2.7.13 (Gaussian elimination for s.p.d. matrices).

Every symmetric/Hermitian positive definite matrix (\rightarrow Def. 2.7.9) possesses an LU-decomposition (\rightarrow Sect. 2.2).

Lemma 2.7.14 (Cholesky decomposition for s.p.d. matrices). \rightarrow [27, Sect. 3.4], [35, Sect. II.5], [51, 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*).

MATLAB function:

```
R = chol(A)
```

Gaussian elimination *without pivoting* is a *numerically stable* way to solve LSEs with s.p.d. system matrix.

2.8 QR-factorization/decomposition [35, Sect. 13], [27, Sect. 7.3]

Remark 2.8.1 (Sensitivity of linear mappings).

Consider problem map (\rightarrow Sect. 2.5.2)

$$F : \begin{cases} \mathbb{K}^n \mapsto \mathbb{K}^n \\ \mathbf{x} \mapsto \mathbf{A}\mathbf{x} \end{cases}$$

for given regular $\mathbf{A} \in \mathbb{K}^{n,n}$ \triangleright $\mathbf{x} \hat{=}$ “data”

$$\Rightarrow \frac{\|\Delta\mathbf{y}\|}{\|\mathbf{y}\|} \leq \frac{\|\mathbf{A}\| \|\Delta\mathbf{x}\|}{\|\mathbf{A}^{-1}\|^{-1} \|\mathbf{x}\|} = \text{cond}(\mathbf{A}) \frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \quad (2.8.2)$$

relative perturbation in result relative perturbation in data

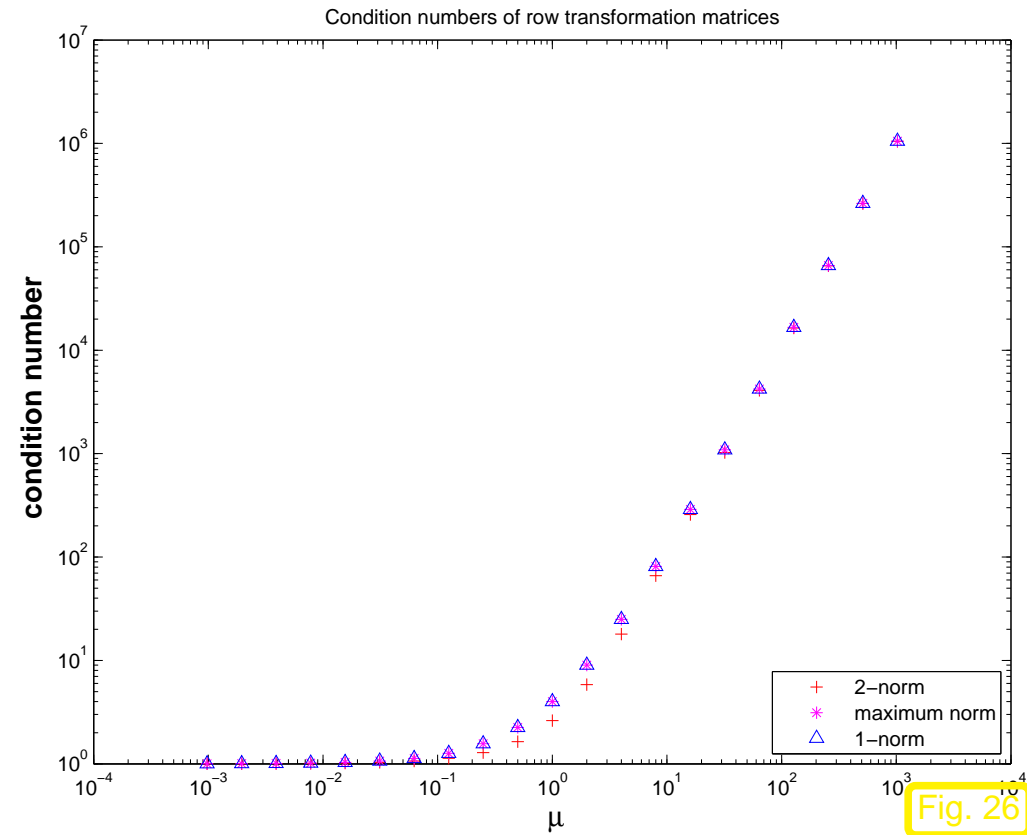
\triangleright Condition number $\text{cond}(\mathbf{A})$ (\rightarrow Def. 2.5.26) bounds amplification of relative error in argument vector in matrix \times vector-multiplication $\mathbf{x} \mapsto \mathbf{A}\mathbf{x}$. △

Example 2.8.3 (Conditioning of row transformations).

2×2 Row transformation matrix (*cf.* elimination matrices of Gaussian elimination, Sect. 2.2):

$$\mathbf{T}(\mu) = \begin{pmatrix} 1 & 0 \\ \mu & 1 \end{pmatrix}$$

Condition numbers of $\mathbf{T}(\mu)$



Goal: find unitary row transformation rendering certain matrix elements zero.

$$Q \begin{pmatrix} \text{yellow square} \end{pmatrix} = \begin{pmatrix} 0 \text{ in first column} \text{ yellow square} \end{pmatrix} \quad \text{with} \quad Q^H = Q^{-1}.$$

Example 2.8.7 (“Annihilating” orthogonal transformations in 2D).

In 2D: two possible orthogonal transformations make 2nd component of $\mathbf{a} \in \mathbb{R}^2$ vanish:

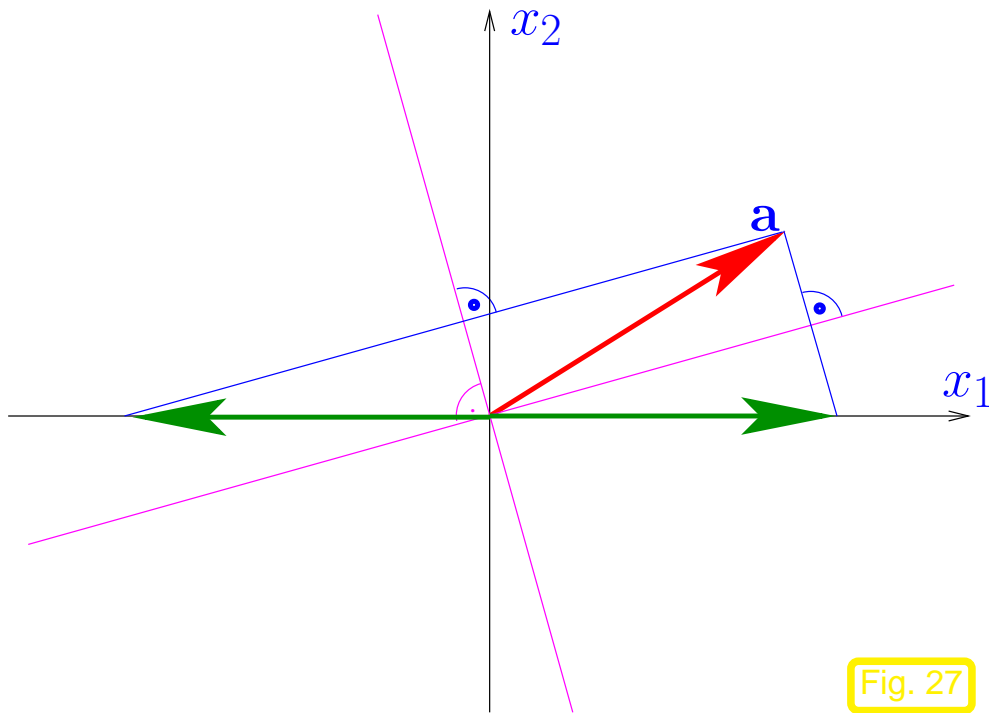


Fig. 27

reflection at angle bisector,

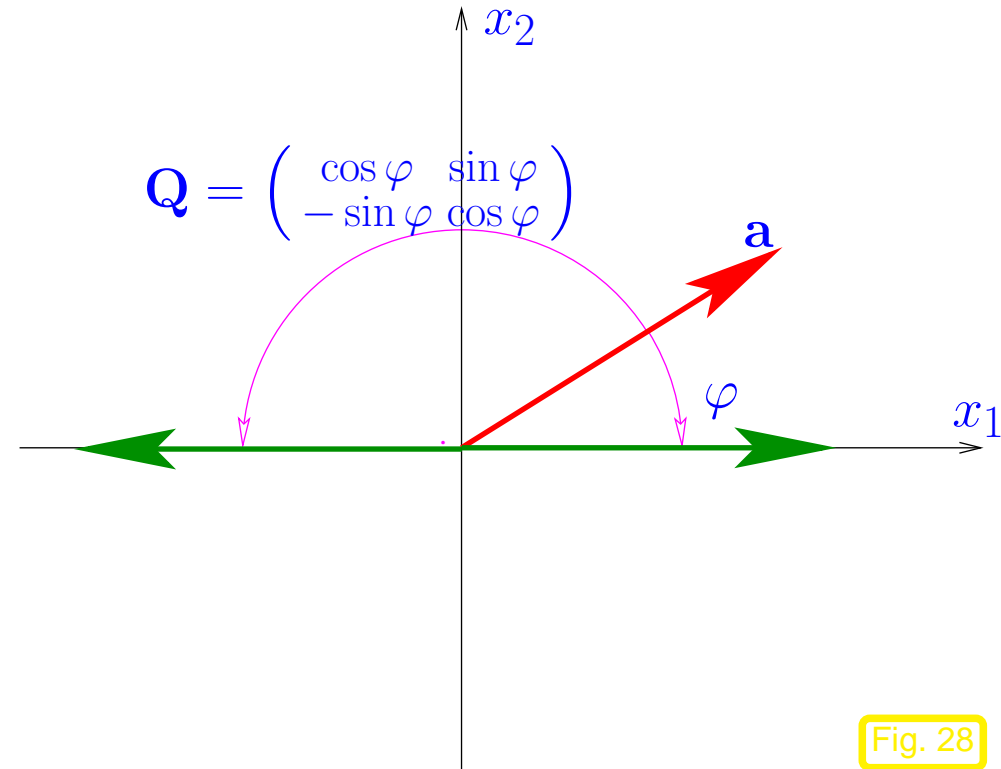


Fig. 28

rotation turning \mathbf{a} onto x_1 -axis.

$$Q = \begin{pmatrix} \cos \varphi & \sin \varphi \\ -\sin \varphi & \cos \varphi \end{pmatrix}$$



Note: two possible reflections/rotations



In n D: given $\mathbf{a} \in \mathbb{R}^n$ find orthogonal matrix $Q \in \mathbb{R}^{n,n}$ such that $Q\mathbf{a} = \|\mathbf{a}\|_2 \mathbf{e}_1$, $\mathbf{e}_1 \hat{=} 1$ st unit vector.

Choice 1: Householder reflections

$$\mathbf{Q} = \mathbf{H}(\mathbf{v}) := \mathbf{I} - 2 \frac{\mathbf{v}\mathbf{v}^H}{\mathbf{v}^H\mathbf{v}} \quad \text{with} \quad \mathbf{v} = \frac{1}{2}(\mathbf{a} \pm \|\mathbf{a}\|_2 \mathbf{e}_1). \quad (2.8.8)$$

Choice 2: successive Givens rotations [35, Sect. 14] (\rightarrow 2D case)

$$\mathbf{G}_{1k}(a_1, a_k)\mathbf{a} := \begin{pmatrix} \bar{\gamma} & \cdots & \bar{\sigma} & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots \\ -\sigma & \cdots & \gamma & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 1 \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_k \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} a_1^{(1)} \\ \vdots \\ 0 \\ \vdots \\ a_n \end{pmatrix}, \quad \text{if} \quad \begin{aligned} \gamma &= \frac{a_1}{\sqrt{|a_1|^2 + |a_k|^2}}, \\ \sigma &= \frac{a_k}{\sqrt{|a_1|^2 + |a_k|^2}}. \end{aligned} \quad (2.8.10)$$

MATLAB-Function: `[G,x] = planerot(a);`

Code 2.8.11: (plane) Givens rotation

```

1 function [G,x] = planerot(a)
2 % plane Givens rotation.
3 if (a(2) ~= 0)
4     r = norm(a); G = [a'; -a(2) a(1)]/r; x = [r; 0];
5 else , G = eye(2); end

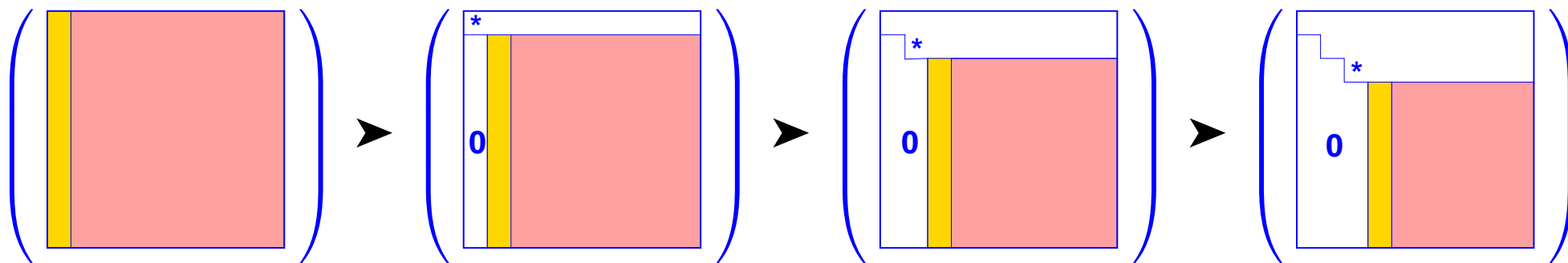
```



Mapping $\mathbf{a} \in \mathbb{K}^n$ to a multiple of \mathbf{e}_1 by $n - 1$ successive Givens rotations:

$$\begin{pmatrix} a_1 \\ \vdots \\ \vdots \\ \vdots \\ a_n \end{pmatrix} \xrightarrow{\mathbf{G}_{12}(a_1, a_2)} \begin{pmatrix} a_1^{(1)} \\ 0 \\ a_3 \\ \vdots \\ a_n \end{pmatrix} \xrightarrow{\mathbf{G}_{13}(a_1^{(1)}, a_3)} \begin{pmatrix} a_1^{(2)} \\ 0 \\ 0 \\ a_4 \\ \vdots \\ a_n \end{pmatrix} \xrightarrow{\mathbf{G}_{14}(a_1^{(2)}, a_4)} \dots \xrightarrow{\mathbf{G}_{1n}(a_1^{(n-2)}, a_n)} \begin{pmatrix} a_1^{(n-1)} \\ 0 \\ \vdots \\ \vdots \\ 0 \end{pmatrix} \quad (2.8.12)$$



Transformation to *upper triangular form* (\rightarrow Def. 2.2.3) by successive unitary transformations:



 = "target column \mathbf{a} " (determines unitary transformation),
 = modified in course of transformations.



QR-factorization
(QR-decomposition) of $\mathbf{A} \in \mathbb{C}^{n,n}$: $\mathbf{A} = \mathbf{QR}$, $\mathbf{Q} := \mathbf{Q}_1^H \cdots \mathbf{Q}_{n-1}^H$ unitary matrix, \mathbf{R} upper triangular matrix.

$$\mathbf{Q}_{n-1} \mathbf{Q}_{n-2} \cdots \mathbf{Q}_1 \mathbf{A} = \mathbf{R},$$

Generalization to $\mathbf{A} \in \mathbb{K}^{m,n}$:

$$m > n: \quad \left(\begin{array}{c} \mathbf{A} \end{array} \right) = \left(\begin{array}{c} \mathbf{Q} \end{array} \right) \left(\begin{array}{c} \mathbf{R} \end{array} \right), \quad \mathbf{A} = \mathbf{QR}, \quad \begin{array}{l} \mathbf{Q} \in \mathbb{K}^{m,n}, \\ \mathbf{R} \in \mathbb{K}^{n,n}, \end{array}$$

(2.8.15)

where $\mathbf{Q}^H \mathbf{Q} = \mathbf{I}$ (orthonormal columns), \mathbf{R} upper triangular matrix.

Lemma 2.8.16 (Uniqueness of QR-factorization).

The “economical” QR-factorization (2.8.15) of $\mathbf{A} \in \mathbb{K}^{m,n}$, $m \geq n$, with $\text{rank}(\mathbf{A}) = n$ is unique, if we demand $r_{ii} > 0$.

$$m < n: \quad \left(\begin{array}{|c|} \hline \mathbf{A} \\ \hline \end{array} \right) = \left(\begin{array}{|c|} \hline \mathbf{Q} \\ \hline \end{array} \right) \left(\begin{array}{|c|} \hline \mathbf{R} \\ \hline \end{array} \right),$$

$$\mathbf{A} = \mathbf{QR}, \quad \mathbf{Q} \in \mathbb{K}^{m,m}, \quad \mathbf{R} \in \mathbb{K}^{m,n},$$

where \mathbf{Q} unitary, \mathbf{R} upper triangular matrix.

MATLAB functions:

$$[\mathbf{Q}, \mathbf{R}] = \text{qr}(\mathbf{A}) \quad \mathbf{Q} \in \mathbb{K}^{m,m}, \mathbf{R} \in \mathbb{K}^{m,n} \text{ for } \mathbf{A} \in \mathbb{K}^{m,n}$$

$$[\mathbf{Q}, \mathbf{R}] = \text{qr}(\mathbf{A}, 0) \quad \mathbf{Q} \in \mathbb{K}^{m,n}, \mathbf{R} \in \mathbb{K}^{n,n} \text{ for } \mathbf{A} \in \mathbb{K}^{m,n}, m > n$$

(economical QR-factorization)

Computational effort for Householder QR-factorization of $\mathbf{A} \in \mathbb{K}^{m,n}, m > n$:

$$[\mathbf{Q}, \mathbf{R}] = \text{qr}(\mathbf{A}) \quad \rightarrow \text{Costs: } O(m^2n)$$

$$[\mathbf{Q}, \mathbf{R}] = \text{qr}(\mathbf{A}, 0) \quad \rightarrow \text{Costs: } O(mn^2)$$

Example 2.8.18 (Complexity of Householder QR-factorization).

```
1 % Timing QR factorizations
2
3 K = 3; r = [];
4 for n=2.^(2:6)
5     m = n*n;
6
7     A = (1:m)'*(1:n) + [eye(n);ones(m-n,n)];
8     t1 = 1000; for k=1:K, tic; [Q,R] = qr(A); t1 = min(t1,toc);
9     clear Q,R; end
10    t2 = 1000; for k=1:K, tic; [Q,R] = qr(A,0); t2 = min(t2,toc);
11    clear Q,R; end
12    t3 = 1000; for k=1:K, tic; R = qr(A); t3 = min(t3,toc); clear
13    R; end
14    r = [r; n , m , t1 , t2 , t3];
15 end
```


tic-toc-timing of different variants of QR-factorization in MATLAB

► Use $[Q,R] = \text{qr}(A,0)$, if output sufficient!

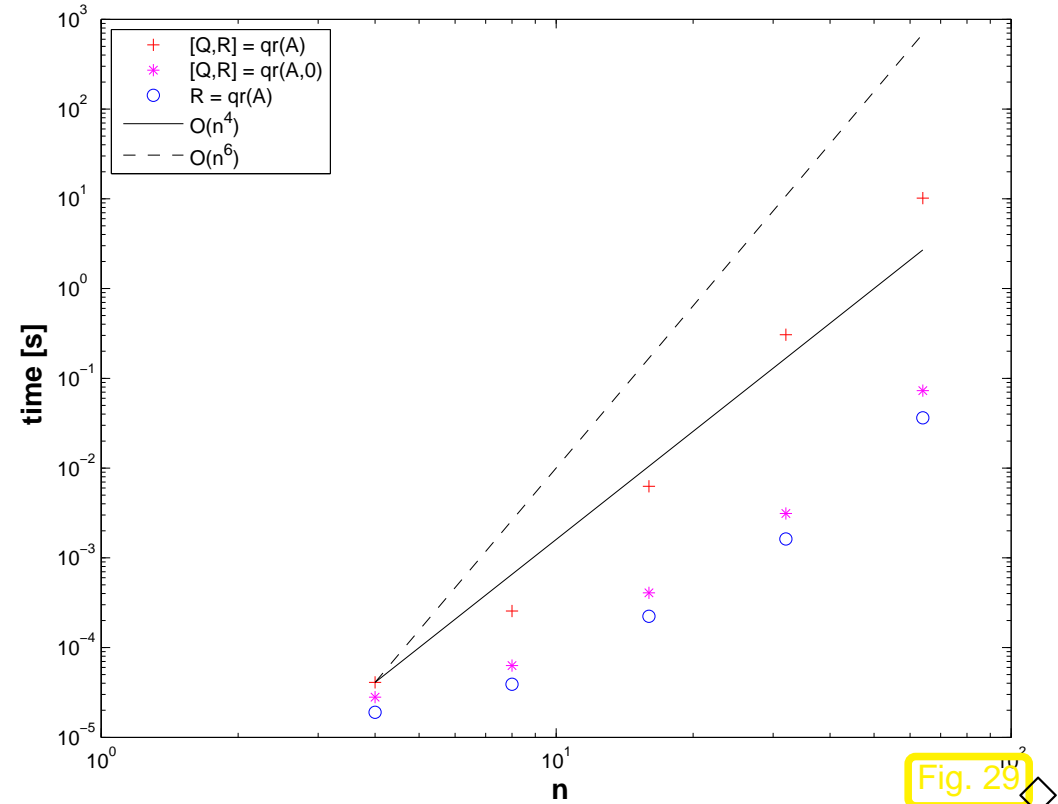


Fig. 29

Remark 2.8.20 (QR-orthogonalization).

$$\begin{pmatrix} \text{A} \end{pmatrix} = \begin{pmatrix} \text{Q} \end{pmatrix} \begin{pmatrix} \text{R} \end{pmatrix}, \quad \text{A}, \text{Q} \in \mathbb{K}^{m,n}, \text{R} \in \mathbb{K}^{n,n}.$$

If $m > n$, $\text{rank}(\mathbf{R}) = \text{rank}(\mathbf{A}) = n$ (full rank)

- $\{\mathbf{q}_{\cdot,1}, \dots, \mathbf{q}_{\cdot,n}\}$ is **orthonormal basis** of $\text{Im}(\mathbf{A})$ with
 $\text{Span}\{\mathbf{q}_{\cdot,1}, \dots, \mathbf{q}_{\cdot,k}\} = \text{Span}\{\mathbf{a}_{\cdot,1}, \dots, \mathbf{a}_{\cdot,k}\}, 1 \leq k \leq n$.



Algorithm 2.8.22 (Solving linear system of equations by means of QR-decomposition).

① QR-decomposition $\mathbf{A} = \mathbf{QR}$, computational costs $\frac{2}{3}n^3 + O(n^2)$ (about twice as expensive as *LU*-decomposition without pivoting)

$\mathbf{Ax} = \mathbf{b}$: ② orthogonal transformation $\mathbf{z} = \mathbf{Q}^H \mathbf{b}$, computational costs $4n^2 + O(n)$ (in the case of *compact storage* of reflections/rotations)

③ **Backward substitution**, solve $\mathbf{Rx} = \mathbf{z}$, computational costs $\frac{1}{2}n(n+1)$

✌ Computing the generalized QR-decomposition $\mathbf{A} = \mathbf{QR}$ by means of Householder reflections or Givens rotations is (numerically stable) for any $\mathbf{A} \in \mathbb{C}^{m,n}$.

✌ For *any* regular system matrix an LSE can be solved by means of

QR-decomposition + orthogonal transformation + backward substitution

in a stable manner.

Code 2.8.24: QR-fac. ↔ Gaussian elimination

```

1 res = [];
2 for n=10:10:1000
3     A=[ tril(-ones(n,n-1))...
4         +2*[ eye(n-1);...
5             zeros(1,n-1)],ones(n,1)];
6     x=((-1).^(1:n))';
7     b=A*x;
8     [Q,R]=qr(A);
9     errlu=norm(A\b-x)/norm(x);
10    errqr=norm(R\((Q'*b)-x)/norm(x);
11    res=[res; n,errlu,errqr];
12 end
13 semilogy(res(:,1),res(:,2),'m-*',
14           res(:,1),res(:,3),'b-+');

```

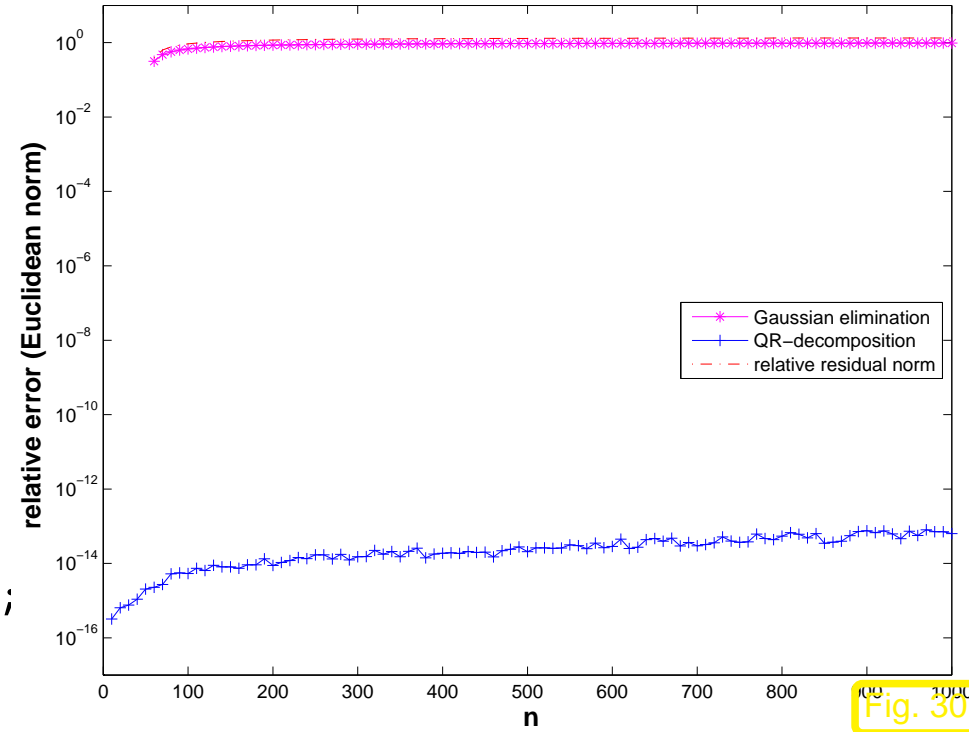


Fig. 30

superior stability of QR-decomposition !



Fill-in for QR-decomposition ?

bandwidth

$$A \in \mathbb{C}^{n,n} \text{ with QR-decomposition } A = QR \Rightarrow m(R) \leq m(A) \text{ (} \rightarrow \text{Def. 2.6.34)}$$

Example 2.8.25 (QR-based solution of tridiagonal LSE).

Elimination of Sub-diagonals by $n - 1$ successive Givens rotations:

$$\begin{pmatrix} * & * & 0 & 0 & 0 & 0 & 0 & 0 \\ * & * & * & 0 & 0 & 0 & 0 & 0 \\ 0 & * & * & * & 0 & 0 & 0 & 0 \\ 0 & 0 & * & * & * & 0 & 0 & 0 \\ 0 & 0 & 0 & * & * & * & 0 & 0 \\ 0 & 0 & 0 & 0 & * & * & * & 0 \\ 0 & 0 & 0 & 0 & 0 & * & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & * & * \end{pmatrix} \xrightarrow{G_{12}} \begin{pmatrix} * & * & * & 0 & 0 & 0 & 0 & 0 \\ 0 & * & * & 0 & 0 & 0 & 0 & 0 \\ 0 & * & * & * & 0 & 0 & 0 & 0 \\ 0 & 0 & * & * & * & 0 & 0 & 0 \\ 0 & 0 & 0 & * & * & * & 0 & 0 \\ 0 & 0 & 0 & 0 & * & * & * & 0 \\ 0 & 0 & 0 & 0 & 0 & * & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & * & * \end{pmatrix} \xrightarrow{G_{23}} \dots \xrightarrow{G_{n-1,n}} \begin{pmatrix} * & * & * & 0 & 0 & 0 & 0 & 0 \\ 0 & * & * & * & 0 & 0 & 0 & 0 \\ 0 & 0 & * & * & * & 0 & 0 & 0 \\ 0 & 0 & 0 & * & * & * & 0 & 0 \\ 0 & 0 & 0 & 0 & * & * & * & 0 \\ 0 & 0 & 0 & 0 & 0 & * & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & * \end{pmatrix}$$

MATLAB code (c, d, e, b = column vectors of length n , $n \in \mathbb{N}$, $e(n), c(n)$ not used):

$$\mathbf{A} = \begin{pmatrix} d_1 & c_1 & 0 & \dots & 0 \\ e_1 & d_2 & c_2 & & \vdots \\ 0 & e_2 & d_3 & c_3 & \\ \vdots & \ddots & \ddots & \ddots & c_{n-1} \\ 0 & \dots & 0 & e_{n-1} & d_n \end{pmatrix} \leftrightarrow \text{spdiags}([e, d, c], [-1 \ 0 \ 1], n, n)$$

Code 2.8.26: solving a tridiagonal system by means of QR-decomposition

```

1 function y = tridiagqr(c,d,e,b)
2 n = length(d); t = norm(d)+norm(e)+norm(c);
3 for k=1:n-1

```

```
4 [R,z] = planerot([d(k);e(k)]);  
5 if (abs(z(1))/t < eps), error('Matrix singular'); end;  
6 d(k) = z(1); b(k:k+1) = R*b(k:k+1);  
7 Z = R*[c(k), 0;d(k+1), c(k+1)];  
8 c(k) = Z(1,1); d(k+1) = Z(2,1);  
9 e(k) = Z(1,2); c(k+1) = Z(2,2);  
10 end  
11 A = spdiags([d,[0;c(1:end-1)]],[0;0;e(1:end-2)]],[0 1 2],n,n);  
12 y = A\b;
```

Asymptotic complexity $O(n)$



2.9 Modification Techniques

Example 2.9.1 (Resistance to currents map).

Large (linear) electric circuit (modelling
→ Ex. 2.6.3) ▷

Sought:

Dependence of (certain) branch currents on “continuously varying” resistance R_x

(▷ currents for many different values of R_x)

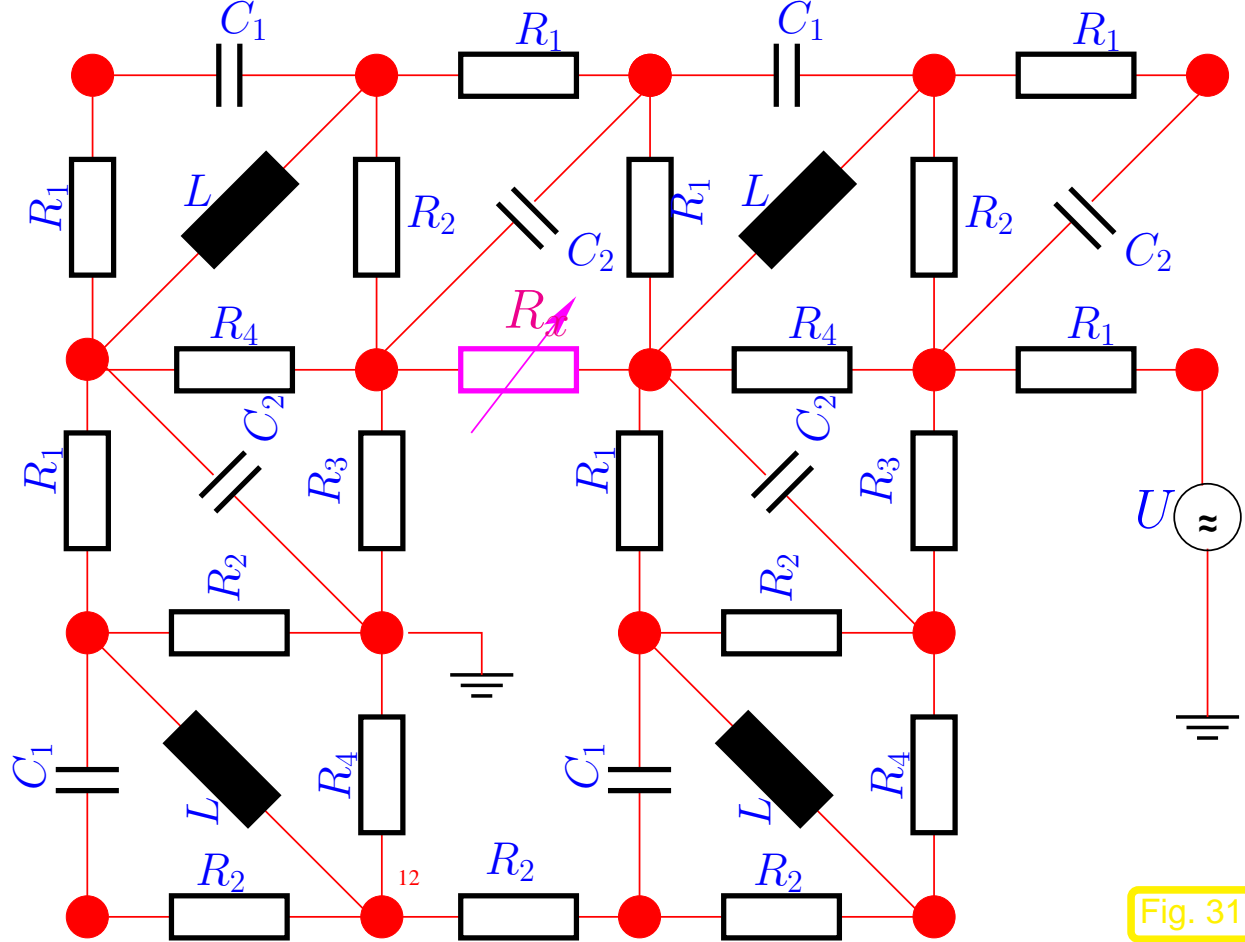


Fig. 31

Problem: Efficient *update* of matrix factorizations in the case of ‘slight’ changes of the matrix [23, Sect. 12.6], [61, Sect. 4.9].

Example 2.9.2 (Changing entries/rows/columns of a matrix).

Changing a single entry: given $x \in \mathbb{K}$

$$\mathbf{A}, \tilde{\mathbf{A}} \in \mathbb{K}^{n,n}: \tilde{a}_{ij} = \begin{cases} a_{ij} & , \text{ if } (i, j) \neq (i^*, j^*) , \\ x + a_{ij} & , \text{ if } (i, j) = (i^*, j^*) , \end{cases} \quad , \quad i^*, j^* \in \{1, \dots, n\} . \quad (2.9.3)$$

$$\blacktriangleright \quad \boxed{\tilde{\mathbf{A}} = \mathbf{A} + x \cdot \mathbf{e}_{i^*} \mathbf{e}_{j^*}^T} . \quad (2.9.4)$$

Changing a single row: given $\mathbf{x} \in \mathbb{K}^n$

$$\mathbf{A}, \tilde{\mathbf{A}} \in \mathbb{K}^{n,n}: \tilde{a}_{ij} = \begin{cases} a_{ij} & , \text{ if } i \neq i^* , \\ x_j + a_{ij} & , \text{ if } i = i^* , \end{cases} \quad , \quad i^*, j^* \in \{1, \dots, n\} .$$

$$\blacktriangleright \quad \boxed{\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{e}_{i^*} \mathbf{x}^T} . \quad (2.9.5)$$



$$\mathbf{A} \in \mathbb{K}^{n,n} \mapsto \tilde{\mathbf{A}} := \mathbf{A} + \boxed{\mathbf{u}\mathbf{v}^H}, \quad \mathbf{u}, \mathbf{v} \in \mathbb{K}^n. \quad (2.9.6)$$

general rank-1-matrix

Remark 2.9.7 (Solving LSE in the case of rank-1-modification).

Lemma 2.9.8 (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 $\mathbf{I} + \mathbf{V}^H\mathbf{A}^{-1}\mathbf{U}$ regular.

Task:

Solve $\tilde{\mathbf{A}}\mathbf{x} = \mathbf{b}$, when LU-factorization $\mathbf{A} = \mathbf{L}\mathbf{U}$ already known

Apply Lemma 2.9.8 for $k = 1$:

$$\mathbf{x} = \left(\mathbf{I} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^H}{1 + \mathbf{v}^H \mathbf{A}^{-1} \mathbf{u}} \right) \mathbf{A}^{-1} \mathbf{b}.$$

Efficient implementation !

Asymptotic complexity $O(n^2)$
(back substitutions)



Code 2.9.9: solving a rank-1 modified LSE

```

1 function x = smw(L,U,u,v,b)
2 t = L\b; z = U\t;
3 t = L\u; w = U\t;
4 alpha = 1+dot(v,w);
5 if (abs(alpha) <
      eps*norm(U,1)), error('Nearly
      singular matrix'); end;
6 x = z - w*dot(v,z)/alpha;

```



Task: Efficient computation of QR-factorization (\rightarrow Sect. 2.8) $\tilde{\mathbf{A}} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$ of $\tilde{\mathbf{A}}$ from (2.9.6), when QR-factorization $\mathbf{A} = \mathbf{Q}\mathbf{R}$ already known

① With $\mathbf{w} := \mathbf{Q}^H \mathbf{u}$: $\mathbf{A} + \mathbf{u} \mathbf{v}^H = \mathbf{Q}(\mathbf{R} + \mathbf{w} \mathbf{v}^H)$

➔ Asymptotic complexity $O(n^2)$ (depends on how \mathbf{Q} is stored \rightarrow Rem. 2.8.21)

② Objective: $\mathbf{w} \rightarrow \|\mathbf{w}\| \mathbf{e}_1$ ➤ via $n - 1$ Givens rotations, see (2.8.10).

$$\mathbf{w} = \begin{pmatrix} * \\ * \\ \vdots \\ * \\ * \\ * \\ * \\ * \end{pmatrix} \xrightarrow{\mathbf{G}_{n-1,n}} \begin{pmatrix} * \\ * \\ \vdots \\ * \\ * \\ * \\ * \\ 0 \end{pmatrix} \xrightarrow{\mathbf{G}_{n-2,n-1}} \begin{pmatrix} * \\ * \\ \vdots \\ * \\ * \\ 0 \\ 0 \end{pmatrix} \xrightarrow{\mathbf{G}_{n-3,n-2}} \dots \xrightarrow{\mathbf{G}_{1,2}} \begin{pmatrix} * \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad (2.9.10)$$

Note: rotations affect \mathbf{R} !

$$\mathbf{R} = \begin{pmatrix} * & * & \dots & * & * & * & * \\ 0 & * & \dots & * & * & * & * \\ \vdots & & \ddots & & & & \vdots \\ 0 & \dots & 0 & * & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \\ 0 & \dots & 0 & 0 & 0 & 0 & * \end{pmatrix} \xrightarrow{\mathbf{G}_{n-1,n}} \begin{pmatrix} * & * & \dots & * & * & * & * \\ 0 & * & \dots & * & * & * & * \\ \vdots & \ddots & & & & \vdots & \\ 0 & \dots & 0 & * & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{n-2,n-1}}$$

$$\rightarrow \begin{pmatrix} * & * & \dots & * & * & * & * \\ 0 & * & \dots & * & * & * & * \\ \vdots & & \ddots & & & & \vdots \\ 0 & \dots & 0 & * & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{n-3,n-2}} \dots \xrightarrow{\mathbf{G}_{1,2}} \begin{pmatrix} * & * & \dots & * & * & * & * \\ * & * & \dots & * & * & * & * \\ & & \ddots & & & & \vdots \\ 0 & \dots & * & * & * & * & * \\ 0 & \dots & 0 & * & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \end{pmatrix} =: \mathbf{R}_1$$

upper Hessenberg matrix: Entry $(i, j) = 0$, if $i > j + 1$.

► $\mathbf{A} + \mathbf{u}\mathbf{v}^H = \mathbf{Q}\mathbf{Q}_1^H \left(\underbrace{\mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^H}_{\text{upper Hessenberg matrix}} \right)$ with unitary $\mathbf{Q}_1 := \mathbf{G}_{12} \cdots \mathbf{G}_{n-1,n}$.

→ Asymptotic complexity $O(n^2)$

③ Successive Givens rotations: $\mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^H \mapsto$ upper triangular form

$$\mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^H = \begin{pmatrix} * & * & \dots & * & * & * & * \\ * & * & \dots & * & * & * & * \\ & & \ddots & & & & \\ 0 & \dots & * & * & * & * & * \\ 0 & \dots & 0 & * & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{12}} \begin{pmatrix} * & * & \dots & * & * & * & * \\ 0 & * & \dots & * & * & * & * \\ & & \ddots & & & & \\ 0 & \dots & * & * & * & * & * \\ 0 & \dots & 0 & * & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{23}} \dots$$

$$\xrightarrow{\mathbf{G}_{n-2,n-1}} \begin{pmatrix} * & * & \cdots & * & * & * & * \\ 0 & * & \cdots & * & * & * & * \\ & & \ddots & & & & \\ 0 & \cdots & 0 & * & * & * & * \\ 0 & \cdots & 0 & 0 & * & * & * \\ 0 & \cdots & 0 & 0 & 0 & * & * \\ 0 & \cdots & 0 & 0 & 0 & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{n-1,n}} \begin{pmatrix} * & * & \cdots & * & * & * & * \\ 0 & * & \cdots & * & * & * & * \\ & & \ddots & & & & \\ 0 & \cdots & 0 & * & * & * & * \\ 0 & \cdots & 0 & 0 & * & * & * \\ 0 & \cdots & 0 & 0 & 0 & * & * \\ 0 & \cdots & 0 & 0 & 0 & 0 & * \end{pmatrix} =: \tilde{\mathbf{R}}. \quad (2.9.11)$$

➔ Asymptotic complexity $O(n^2)$



$$\mathbf{A} + \mathbf{u}\mathbf{v}^H = \tilde{\mathbf{Q}}\tilde{\mathbf{R}} \quad \text{mit } \tilde{\mathbf{Q}} = \mathbf{Q}\mathbf{Q}_1^H \mathbf{G}_{n-1,n}^H \cdots \mathbf{G}_{12}^H.$$

MATLAB-function: `[Q1,R1] = qrupdate(Q,R,u,v);`

Special case: rank-1-modifications preserving symmetry & positivity (→ Def. 2.7.9):

$$\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n} \mapsto \tilde{\mathbf{A}} := \mathbf{A} + \alpha \mathbf{v}\mathbf{v}^H, \quad \mathbf{v} \in \mathbb{K}^n, \alpha > 0. \quad (2.9.12)$$

Task: Efficient computation of Cholesky factorization $\tilde{\mathbf{A}} = \tilde{\mathbf{R}}^H \tilde{\mathbf{R}}$ (→ Lemma 2.7.14) of $\tilde{\mathbf{A}}$ from (2.9.12), when Cholesky factorization $\mathbf{A} = \mathbf{R}^H \mathbf{R}$ of \mathbf{A} already known

With $\mathbf{w} := \mathbf{R}^{-H} \mathbf{v}$: $\mathbf{A} + \alpha \mathbf{v} \mathbf{v}^H = \mathbf{R}^H (\mathbf{I} + \alpha \mathbf{w} \mathbf{w}^H) \mathbf{R}$.

➔ Asymptotic complexity $O(n^2)$ (backward substitution !)

② Idea: formal Gaussian elimination: with $\tilde{\mathbf{w}} = (w_2, \dots, w_n)^T \rightarrow$ see (2.1.11)

$$\mathbf{I} + \alpha \mathbf{w} \mathbf{w}^H = \left(\begin{array}{c|c} 1 + \alpha w_1^2 & \alpha w_1 \tilde{\mathbf{w}}^H \\ \hline \alpha w_1 \tilde{\mathbf{w}} & \mathbf{I} + \alpha \tilde{\mathbf{w}} \tilde{\mathbf{w}}^H \end{array} \right) \rightarrow \left(\begin{array}{c|c} 1 + \alpha w_1^2 & \alpha w_1 \tilde{\mathbf{w}}^H \\ \hline 0 & \mathbf{I} + \alpha^{(1)} \tilde{\mathbf{w}} \tilde{\mathbf{w}}^H \end{array} \right) \quad (2.9.13)$$

where $\alpha^{(1)} := \alpha - \frac{\alpha^2 w_1^2}{1 + \alpha w_1^2}$.

same structure \triangleright recursion

► Computation of Cholesky-factorization

$$\mathbf{I} + \alpha \mathbf{w} \mathbf{w}^H = \mathbf{R}_1^H \mathbf{R}_1 .$$

Motivation: “recursion”
(2.9.13).

→ asymptotic complexity $O(n^2)$
($O(n)$, if only \mathbf{d}, \mathbf{s} computed
→ (2.9.16))

Code 2.9.15: Cholesky factorization of rank-1-modified identity matrix

```

1 function [d,s] = roid(alpha,w)
2 n = length(w);
3 d = []; s = [];
4 for i=1:n
5     t = alpha*w(i);
6     d = [d; sqrt(1+t*w(i))];
7     s = [s; t/d(i)];
8     alpha = alpha - s(i)^2;
9 end
    
```

③ Special structure of \mathbf{R}_1 :

$$\mathbf{R}_1 = \begin{pmatrix} d_1 & & & & & \\ & \ddots & & & & \\ & & \ddots & & & \\ & & & \ddots & & \\ & & & & \ddots & \\ & & & & & d_n \end{pmatrix} + \begin{pmatrix} s_1 & & & & & \\ & \ddots & & & & \\ & & \ddots & & & \\ & & & \ddots & & \\ & & & & \ddots & \\ & & & & & s_n \end{pmatrix} \begin{pmatrix} 0 & w_2 & w_3 & \cdots & \cdots & w_n \\ 0 & 0 & w_3 & \cdots & \cdots & w_n \\ \vdots & & \ddots & & & \vdots \\ \vdots & & & & 0 & w_{n-1} & w_n \\ 0 & \cdots & \cdots & & 0 & w_n \\ 0 & \cdots & & & \cdots & 0 \end{pmatrix} \quad (2.9.16)$$

$$\mathbf{R}_1 = \begin{pmatrix} d_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & d_n \end{pmatrix} + \begin{pmatrix} s_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & s_n \end{pmatrix} \begin{pmatrix} 0 & 1 & 1 & \cdots & \cdots & 1 \\ 0 & 0 & 1 & \cdots & \cdots & 1 \\ \vdots & & \ddots & & & \vdots \\ \vdots & & & 0 & 1 & 1 \\ 0 & \cdots & \cdots & 0 & 1 \\ 0 & \cdots & \cdots & \cdots & 0 \end{pmatrix} \begin{pmatrix} w_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & w_n \end{pmatrix}$$

► Smart multiplication

$$\tilde{\mathbf{R}} := \mathbf{R}_1 \mathbf{R}$$

→ Complexity $O(n^2)$

$$\mathbf{A} + \alpha \mathbf{v} \mathbf{v}^H = \tilde{\mathbf{R}}^H \tilde{\mathbf{R}}$$

Code 2.9.18: Update of Cholesky factorization in the case of s.p.d. preserving rank-1-modification

```

1  function Rt = roudchol(R,alpha,v)
2  w = R' \ v;
3  [d,s] = roid(alpha,w);
4  T = zeros(1,n);
5  for j=n-1:-1:1
6      T = [w(j+1)*R(j+1,:)+T(1,:);T];
7  end
8  Rt = spdiags(d,0,n,n)*R+spdiags(s,0,n,n)*T;

```

$$\mathbf{A} \in \mathbb{K}^{m,n} \mapsto \tilde{\mathbf{A}} = [(\mathbf{A})_{:,1}, \dots, (\mathbf{A})_{:,k-1}, \mathbf{v}, (\mathbf{A})_{:,k}, \dots, (\mathbf{A})_{:,n}] , \quad \mathbf{v} \in \mathbb{K}^m . \quad (2.9.19)$$

Known: QR-factorization $\mathbf{A} = \mathbf{QR}$, $\mathbf{Q} \in \mathbb{K}^{m,m}$ unitary $\mathbf{R} \in \mathbb{K}^{m,n}$ upper triangular matrix.

Task: Efficient computation of QR-factorization $\tilde{\mathbf{A}} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$ of $\tilde{\mathbf{A}}$ from (2.9.19), $\tilde{\mathbf{Q}} \in \mathbb{K}^{m,m}$ unitary, $\tilde{\mathbf{R}} \in \mathbb{K}^{m,n+1}$ upper triangular

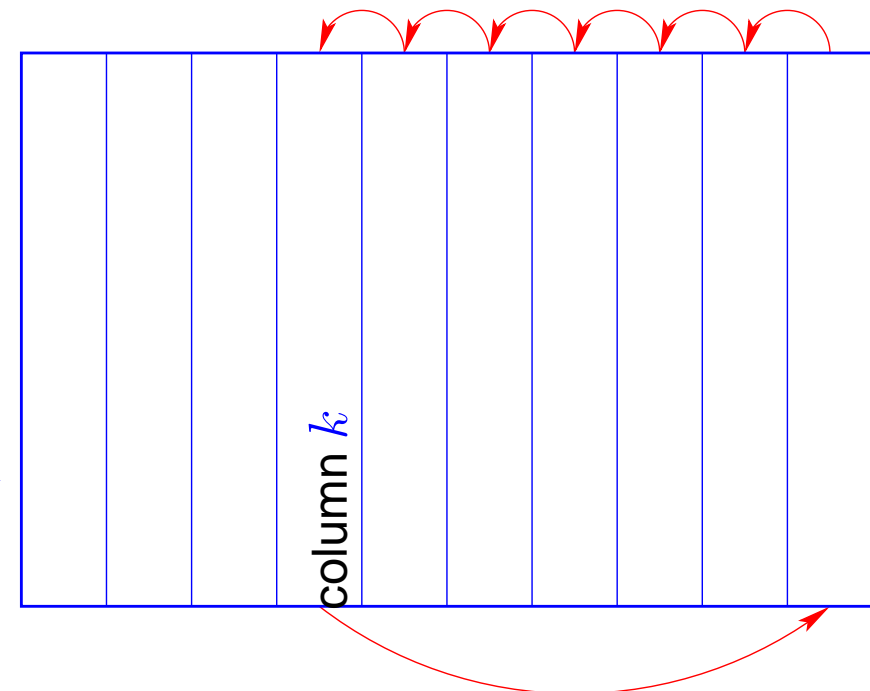
Idea: Easy, if $k = n + 1$ (adding last column)

► \exists column permutation

$$k \mapsto n + 1 , i \mapsto i - 1 , i = k + 1, \dots, n + 1$$

\sim permutation matrix

$$\mathbf{P} = \begin{pmatrix} 1 & 0 & \cdots & \cdots & 0 \\ 0 & \ddots & & & \\ & & 1 & 0 & \\ \vdots & & 0 & 1 & \\ \vdots & & & & \ddots \\ 0 & \cdots & 1 & 0 & 1 \end{pmatrix} \in \mathbb{R}^{n+1,n+1}$$



$$\tilde{\mathbf{A}} \longrightarrow \mathbf{A}_1 = \tilde{\mathbf{A}}\mathbf{P} = [\mathbf{a}_1, \dots, \mathbf{a}_n, \mathbf{v}] = \mathbf{Q} \begin{bmatrix} \mathbf{R} & \mathbf{Q}^H \mathbf{v} \end{bmatrix} = \mathbf{Q} \begin{pmatrix} \text{yellow columns } \mathbf{R} & \text{pink column } \mathbf{Q}^H \mathbf{v} \end{pmatrix}$$

column $\mathbf{Q}^H \mathbf{v}$
 case $m > n + 1$

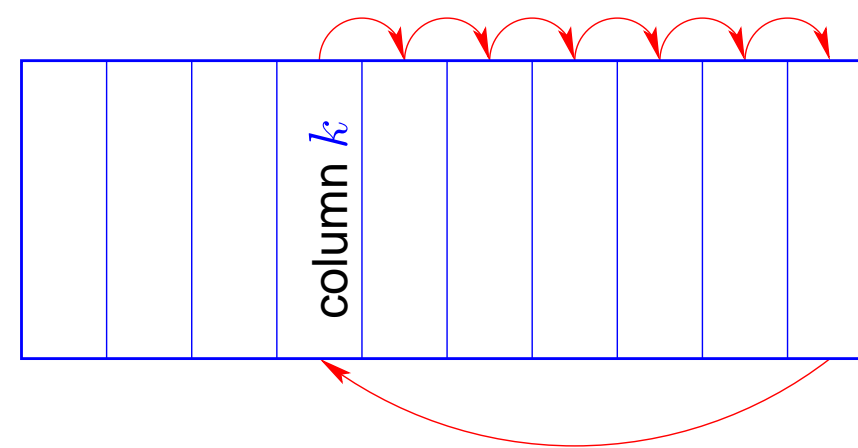
① If $m > n + 1$: \exists orthogonal transformation $\mathbf{Q}_1 \in \mathbb{K}^{m,m}$ (Householder reflection) with

$$\mathbf{Q}_1 \mathbf{Q}^H \mathbf{v} = \begin{pmatrix} * \\ \vdots \\ * \\ * \\ 0 \\ \vdots \\ 0 \end{pmatrix} \left\{ \begin{array}{l} n+1 \\ \\ \\ \\ m-n-1 \end{array} \right. \blacktriangleright \mathbf{Q}_1 \mathbf{Q}^H \mathbf{A}_1 = \begin{pmatrix} * & \dots & & \dots & * \\ 0 & * & & & * \\ \vdots & & \ddots & & \vdots \\ & & & * & * \\ \vdots & & & & * \\ 0 & \dots & & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & \dots & & \dots & 0 \end{pmatrix} \left\{ \begin{array}{l} n+1 \\ \\ \\ \\ \\ m-n-1 \end{array} \right.$$

➔ Computational effort $O(m-n)$ (a single reflection)

inverse permutation:

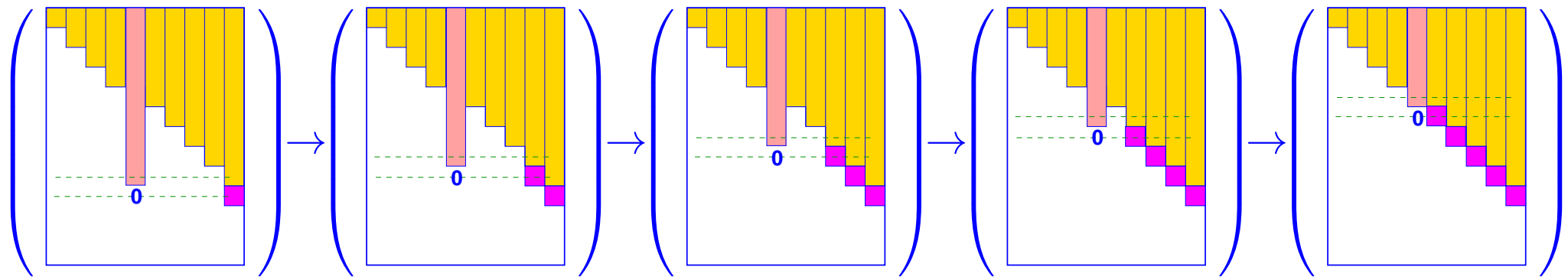
~ right multiplication with \mathbf{P}^H



$$\mathbf{Q}_1 \mathbf{Q}^H \tilde{\mathbf{A}} = \mathbf{Q}_1 \mathbf{Q}^H \mathbf{A}_1 \mathbf{P}^H = \begin{pmatrix} * & \dots & * & \dots & * \\ 0 & * & * & & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ & & \vdots & \ddots & \vdots \\ \vdots & & * & * & * \\ 0 & \dots & 0 & \dots & 0 \\ \vdots & & \vdots & & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{pmatrix} = \begin{pmatrix} \text{yellow trapezoidal blocks} \\ \text{pink vertical bar} \\ \text{yellow trapezoidal blocks} \end{pmatrix}$$

② $n + 1 - k$ successive Givens rotations \Rightarrow upper triangular matrix $\tilde{\mathbf{R}}$

$$\mathbf{Q}_1 \mathbf{Q}^H \mathbf{A}_1 = \begin{pmatrix} * & \dots & * & \dots & * \\ 0 & * & * & & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ & & & * & * \\ & & & * & * \\ \vdots & & * & & 0 \\ 0 & \dots & 0 & \dots & \vdots \\ \vdots & & \vdots & & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{pmatrix} \xrightarrow{\mathbf{G}_{n,n+1}} \dots \xrightarrow{\mathbf{G}_{k,k+1}} \begin{pmatrix} * & \dots & * & \dots & * \\ 0 & * & * & & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ & & & * & * \\ & & & 0 & \dots \\ \vdots & & 0 & & * \\ 0 & \dots & 0 & \dots & 0 \\ \vdots & & \vdots & & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{pmatrix}$$



--- $\hat{=}$ rows targeted by plane rotations, ■ $\hat{=}$ new entries $\neq 0$

➔ Asymptotic complexity $O((n - k)^2)$

$$\mathbf{A} \in \mathbb{K}^{m,n} \mapsto \tilde{\mathbf{A}} = \begin{bmatrix} (\mathbf{A})_{1,:} \\ \vdots \\ (\mathbf{A})_{k-1,:} \\ \mathbf{v}^T \\ (\mathbf{A})_{k,:} \\ \vdots \\ (\mathbf{A})_{m,:} \end{bmatrix}, \quad \mathbf{v} \in \mathbb{K}^n. \quad (2.9.20)$$

Given: QR-factorization $\mathbf{A} = \mathbf{QR}$, $\mathbf{Q} \in \mathbb{K}^{m+1,m+1}$ unitary, $\mathbf{R} \in \mathbb{K}^{m,n}$ upper triangular matrix.

Task: efficient computation of QR-factorization $\tilde{\mathbf{A}} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$ of $\tilde{\mathbf{A}}$ from (2.9.20), $\tilde{\mathbf{Q}} \in \mathbb{K}^{m+1,m+1}$ unitary, $\tilde{\mathbf{R}} \in \mathbb{K}^{m+1,n+1}$ upper triangular matrix.

- ① \exists (partial) cyclic row permutation $m + 1 \leftarrow k, i \leftarrow i + 1, i = k, \dots, m$:
 \rightarrow **unitary** permutation matrix (\rightarrow Def. 2.3.12) $\mathbf{P} \in \{0, 1\}^{m+1, m+1}$

$$\mathbf{P}\tilde{\mathbf{A}} = \begin{pmatrix} \mathbf{A} \\ \mathbf{v}^T \end{pmatrix} \blacktriangleright \begin{pmatrix} \mathbf{Q}^H & 0 \\ 0 & 1 \end{pmatrix} \mathbf{P}\tilde{\mathbf{A}} = \begin{pmatrix} \mathbf{R} \\ \mathbf{v}^T \end{pmatrix} = \begin{pmatrix} \text{[Upper triangular matrix with yellow diagonal]} \\ \mathbf{v}^T \end{pmatrix}$$

case $m = n$

- ② Transform into upper triangular form by m successive Givens rotations:

$$\begin{pmatrix} * & \dots & & \dots & * \\ 0 & * & & & \vdots \\ \vdots & 0 & \dots & & \\ \vdots & \vdots & & * & \vdots \\ 0 & 0 & & 0 & * & * \\ 0 & 0 & & 0 & 0 & * \\ * & \dots & \dots & * & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{1, m+1}} \begin{pmatrix} * & \dots & & \dots & * \\ 0 & * & & & \vdots \\ \vdots & 0 & \dots & & \\ \vdots & \vdots & & * & \vdots \\ 0 & 0 & & 0 & * & * \\ 0 & 0 & & 0 & 0 & * \\ 0 & * & \dots & * & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{2, m+1}} \dots$$

$$\dots \xrightarrow{\mathbf{G}_{m-1,m+1}} \begin{pmatrix} * & \dots & & \dots & * \\ 0 & * & & & \vdots \\ \vdots & 0 & \ddots & & \\ \vdots & \vdots & & * & \vdots \\ 0 & 0 & & 0 & * \\ 0 & 0 & & 0 & * \\ 0 & \dots & & \dots & 0 \\ & & & & * \end{pmatrix} \xrightarrow{\mathbf{G}_{m,m+1}} \begin{pmatrix} * & \dots & & \dots & * \\ 0 & * & & & \vdots \\ \vdots & 0 & \ddots & & \\ \vdots & \vdots & & * & \vdots \\ 0 & 0 & & 0 & * \\ 0 & 0 & & 0 & * \\ 0 & \dots & & \dots & 0 \\ & & & & 0 \end{pmatrix} := \tilde{\mathbf{R}} \quad (2.9.21)$$

③ With $\mathbf{Q}_1 = \mathbf{G}_{m,m+1} \cdots \mathbf{G}_{1,m+1}$

$$\tilde{\mathbf{A}} = \mathbf{P}^T \begin{pmatrix} \mathbf{Q} & 0 \\ 0 & 1 \end{pmatrix} \mathbf{Q}_1^H \tilde{\mathbf{R}} = \tilde{\mathbf{Q}} \tilde{\mathbf{R}} \quad \text{with unitary } \tilde{\mathbf{Q}} \in \mathbb{K}^{m+1,m+1} .$$

3

Data Interpolation in 1D

3.1 Abstract interpolation

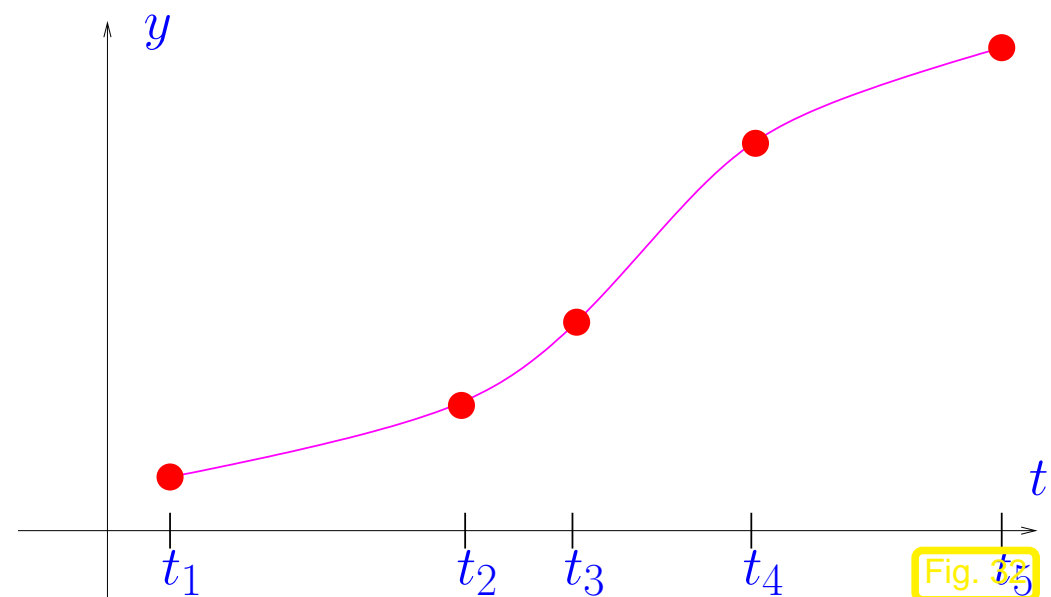
Example 3.1.2 (Constitutive relations (*ger.* Kennlinien) from measurements).

Known: several *accurate* measurements

$$(t_i, y_i), \quad i = 1, \dots, m$$

Examples:

| t | y |
|--------------------|-------------------|
| voltage U | current I |
| pressure p | density ρ |
| magnetic field H | magnetic flux B |
| ... | ... |



R. Hiptmair
rev 38355,
September
7, 2011



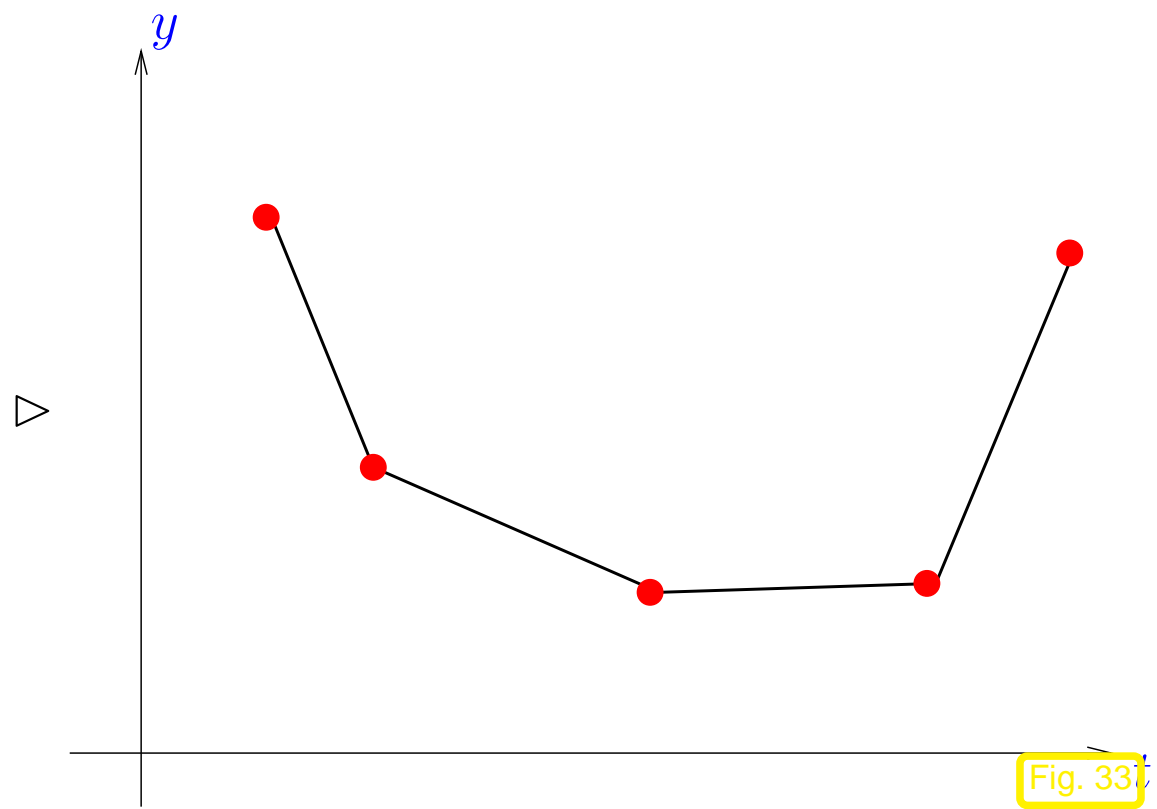
Remark 3.1.3 (Function representation).

```
1 class Function {  
2   private :  
3     // various internal data describing f  
4   public :  
5     // Constructor: accepts information necessary for specifying the  
6     // function  
7     Function(/* .... */);  
8     // Evaluation operator  
9     double operator () (double t) const;  
};
```

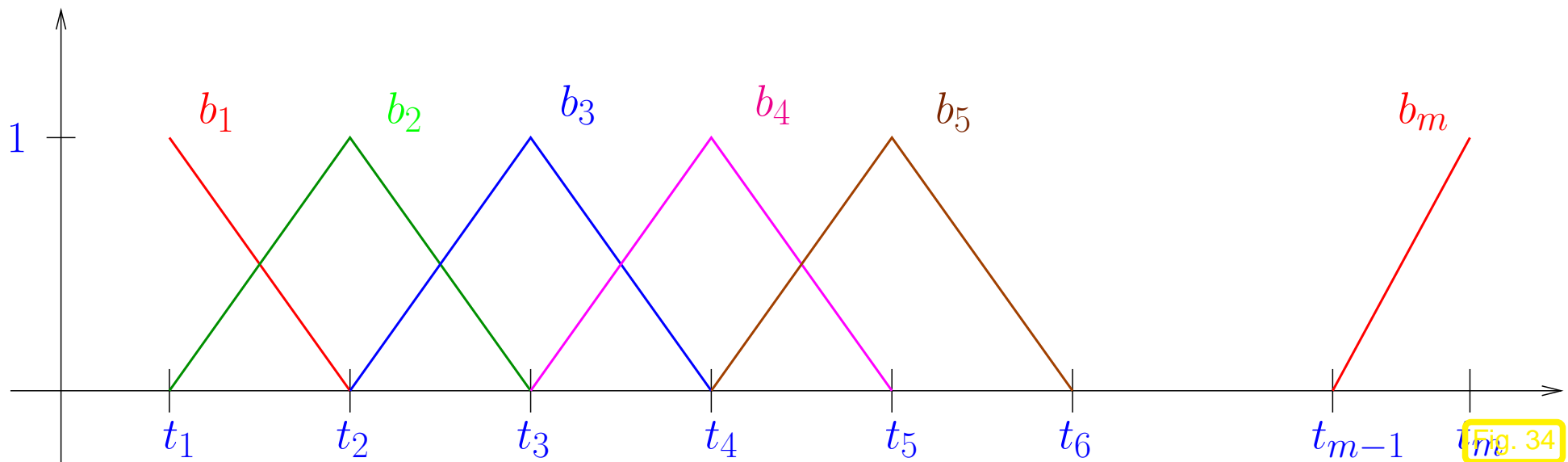


Example 3.1.5 (Piecewise linear interpolation). See also Sect. 3.6.1

Piecewise linear interpolant of data



“Tent function” (“hat function”) basis:



Remark 3.1.8 (Interpolation as linear mapping).



Remark 3.1.11 (“Software solution” of interpolation problem).

```
1 class Interpolant {  
2   private :  
3   // various internal data describing f  
4   public :  
5   // Constructor: computation of coefficients  $c_j$  of representation (3.1.4)  
6   Interpolant(const vector<double> &t, const vector<double> &y) ;  
7   // Evaluation operator for interpolant f  
8   double operator () (double t) const ;  
9 } ;
```



3.2 Polynomials

Notation: Vector space of the polynomials of degree $\leq k$, $k \in \mathbb{N}$:

$$\mathcal{P}_k := \{t \mapsto \alpha_k t^k + \alpha_{k-1} t^{k-1} + \dots + \alpha_1 t + \alpha_0, \alpha_j \in \mathbb{K}\} . \quad (3.2.1)$$

leading coefficient

Theorem 3.2.2 (Dimension of space of polynomials).

$$\dim \mathcal{P}_k = k + 1 \quad \text{and} \quad \mathcal{P}_k \subset C^\infty(\mathbb{R}).$$

Remark 3.2.3 (Polynomials in Matlab).

MATLAB: $\alpha_k t^k + \alpha_{k-1} t^{k-1} + \dots + \alpha_0 \rightarrow$ Vector $(\alpha_k, \alpha_{k-1}, \dots, \alpha_0)$ (ordered!).



Remark 3.2.4 (Horner scheme). \rightarrow [13, Bem. 8.11]

Evaluation of a polynomial in monomial representation:

Horner scheme

$$p(t) = t(\dots t(t(\alpha_n t + \alpha_{n-1}) + \alpha_{n-2}) + \dots + \alpha_1) + \alpha_0 . \quad (3.2.5)$$

```

1 function y = polyval(p,x)
2 y = p(1); for i=2:length(p), y = x*y+p(i); end

```

Asymptotic complexity: $O(n)$ Use: MATLAB “built-in”-function `polyval(p,x)`;

3.3 Polynomial Interpolation: Theory

Interpolation problem (\rightarrow Sect. 3.1): (re-)construction of a polynomial through points (t_i, y_i) .

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 $y_0, \dots, y_n \in \mathbb{K}$ compute $p \in \mathcal{P}_n$ such that

$$p(t_j) = y_j \quad \text{for } j = 0, \dots, n. \quad (3.3.1)$$

For nodes $t_0 < t_1 < \dots < t_n$ (\rightarrow Lagrange interpolation) consider

Lagrange polynomials
$$L_i(t) := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j}, \quad i = 0, \dots, n. \quad (3.3.2)$$

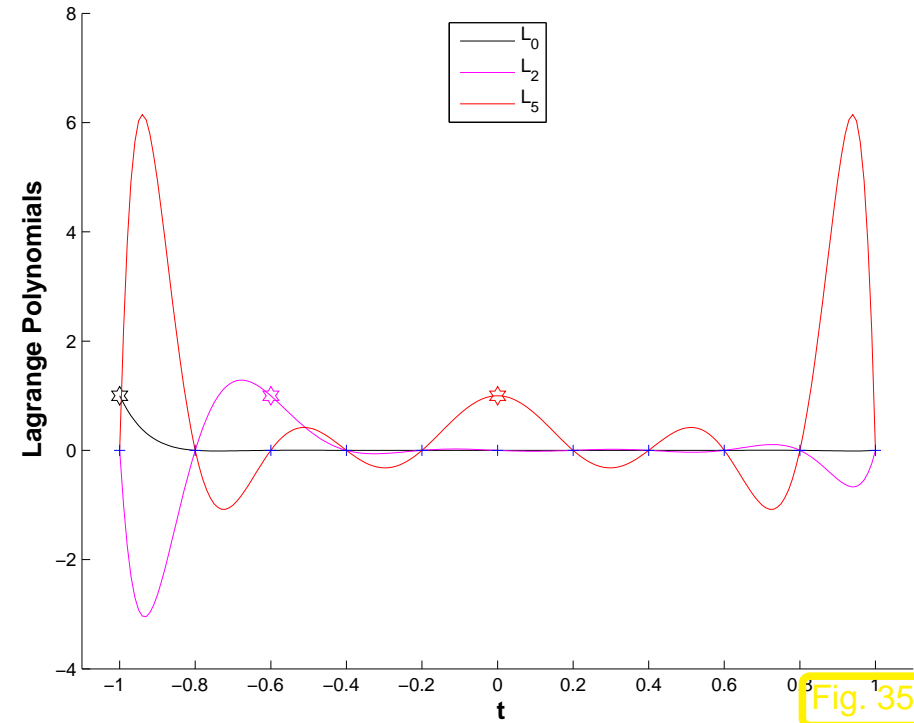


$L_i \in \mathcal{P}_n$ and $L_i(t_j) = \delta_{ij}$

Example 3.3.3. Lagrange polynomials for uniformly spaced nodes

$$\mathcal{T} := \left\{ t_j = -1 + \frac{2}{n} j \right\}, \quad j = 0, \dots, n.$$

Plot $n = 10, j = 0, 2, 5 \rightarrow$



Theorem 3.3.5 (Existence & uniqueness of Lagrange interpolation polynomial). \rightarrow [51, Thm. 8.1], [13, Satz 8.3]

The general polynomial interpolation problem (3.3.1) admits a unique solution $p \in \mathcal{P}_n$.

Theorem 3.3.7 (Lagrange interpolation as linear mapping). \rightarrow Rem. 3.1.8

The polynomial interpolation in the nodes $\mathcal{T} := \{t_j\}_{j=0}^n$ defines a linear operator

$$I_{\mathcal{T}} : \begin{cases} \mathbb{K}^{n+1} & \rightarrow \mathcal{P}_n, \\ (y_0, \dots, y_n)^T & \mapsto \text{interpolating polynomial } p. \end{cases} \quad (3.3.8)$$

3.4 Polynomial Interpolation: Algorithms

3.4.1 Multiple evaluations

```

1 class PolyInterp {
2   private :
3     // various internal data describing p
4   public :
5     // Constructor taking node vector (t0, ..., tn) as argument
6     PolyInterp(const vector<double> &t);
7     // Evaluation operator at x for data (y0, ..., yn)
8     double eval(double x, const vector<double> &y) const;
9 };

```

Barycentric interpolation formula

$$p(t) = \frac{\sum_{i=0}^n \frac{\lambda_i}{t - t_i} y_i}{\sum_{i=0}^n \frac{\lambda_i}{t - t_i}}. \quad (3.4.1)$$

with $\lambda_i = \frac{1}{(t_i - t_0) \cdots (t_i - t_{i-1})(t_i - t_{i+1}) \cdots (t_i - t_n)}$, $i = 0, \dots, n$ → precompute !

Computational effort: • computation of λ_i : $O(n^2)$ (only once),

• every subsequent evaluation of p : $O(n)$,

⇒ total effort $O(Nn) + O(n^2)$

Code 3.4.2: Evaluation of the interpolation polynomials with barycentric formula

NumCSE,
autumn
2010

```

1 function p = intpolyval(t,y,x)
2 % t: row vector of nodes  $t_0, \dots, t_n$ 
3 % y: row vector of data  $y_0, \dots, y_n$ 
4 % x: row vector of evaluation points  $x_1, \dots, x_N$ 
5 n = length(t); % number of interpolation nodes = degree of polynomial - 1
6 N = length(x); % Number of evaluation points stored in x
7 % Precompute the weights  $\lambda_i$  with effort  $O(n^2)$ 
8 for k = 1:n
9     lambda(k) = 1 / prod(t(k) - t([1:k-1,k+1:n])); end;
10 for i = 1:N
11 % Compute quotient of weighted sums of  $\frac{\lambda_i}{t-t_i}$ , effort  $O(n)$ 
12 z = (x(i)-t); j = find(z == 0);
13 if (~isempty(j)), p(i) = y(j); % avoid division by zero
14 else
15     mu = lambda./z; p(i) = sum(mu.*y)/sum(mu);
16 end
17 end

```

R. Hiptmair
rev 38355,
October 26,
2011

3.4.2 Single evaluation [13, Sect. 8.2.2]

For $\{i_0, \dots, i_m\} \subset \{0, \dots, n\}$, $0 \leq m \leq n$:

p_{i_0, \dots, i_m} = interpolation polynomial of degree m through $(t_{i_0}, y_{i_0}), \dots, (t_{i_m}, y_{i_m})$,

recursive definition:

$$\begin{aligned}
 p_i(t) &\equiv y_i, & i = 0, \dots, n, \\
 p_{i_0, \dots, i_m}(t) &= \frac{(t - t_{i_0})p_{i_1, \dots, i_m}(t) - (t - t_{i_m})p_{i_0, \dots, i_{m-1}}(t)}{t_{i_m} - t_{i_0}} \\
 &= \text{TODO!} .
 \end{aligned} \tag{3.4.3}$$

Aitken-Neville algorithm:

| $n =$ | 0 | 1 | 2 | 3 |
|-------|-----------------|----------------------|-----------------------|------------------------|
| t_0 | $y_0 =: p_0(x)$ | $\nearrow p_{01}(x)$ | $\nearrow p_{012}(x)$ | $\nearrow p_{0123}(x)$ |
| t_1 | $y_1 =: p_1(x)$ | $\nearrow p_{12}(x)$ | $\nearrow p_{123}(x)$ | |
| t_2 | $y_2 =: p_2(x)$ | $\nearrow p_{23}(x)$ | | |
| t_3 | $y_3 =: p_3(x)$ | | | |

Code 3.4.4: Aitken-Neville algorithm

```

1 function v = ANipoleval(t,y,x)
2 for i=1:length(y)
3     for k=i-1:-1:1
4         y(k) =
5             y(k+1)+(y(k+1)-y(k))*...
6                 (x-t(i))/(t(i)-t(k))
7     end
8 end
v = y(1);
    
```

```
1 class PolyEval {
2     private :
3         // evaluation point and various internal data describing the 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    };
```

Example 3.4.5 (Timing polynomial evaluations).

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 point $x \in [0, n]$.

Minimum `tic-toc`-computational time
over 100 runs →

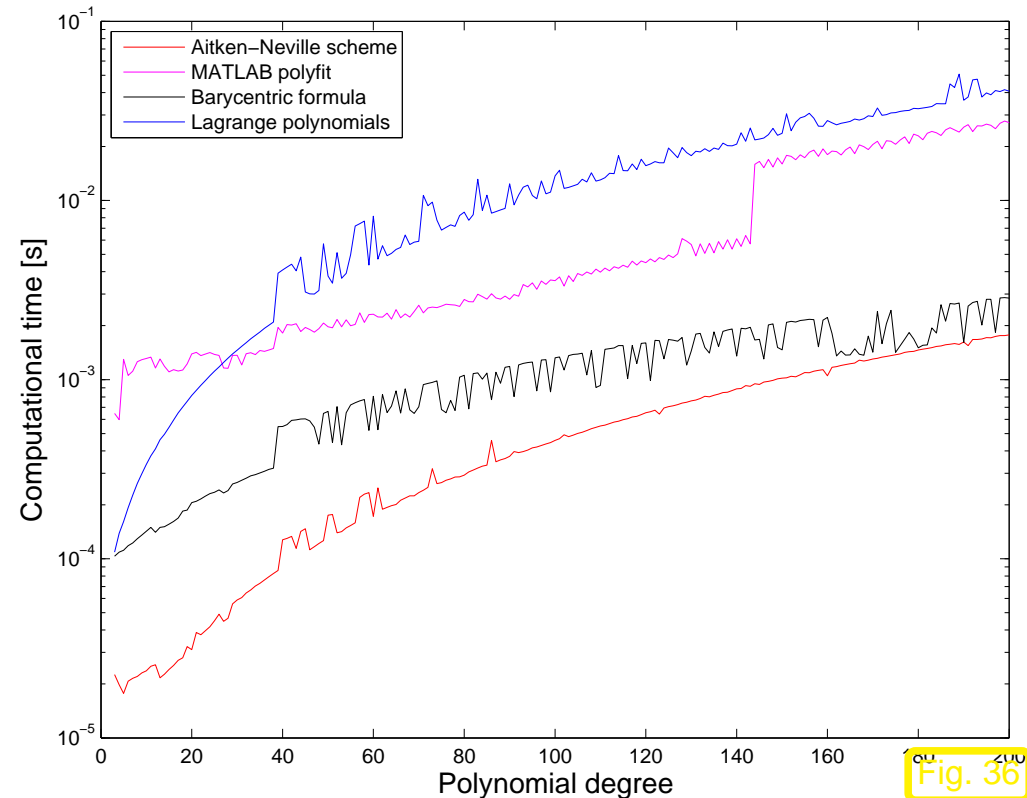


Fig. 36

Code 3.4.6: MATLAB polynomial evaluation using built-in function `polyfit`

```

1 function v=ipoleval(t,y,x)
2   p = polyfit(t,y,length(y)-1);
3   v=polyval(p,x);

```

Code 3.4.7: Lagrange polynomial interpolation and evaluation

```
1 function p = intpolyval_lag(t,y,x)
2 p=zeros(size(x));
3 for k=1:length(t); p=p + y(k)*lagrangepoly(x, k-1, t); end
4
5 function L=lagrangepoly(x, index, nodes)
6 L=1;
7 for j=[0:index-1, index+1:length(nodes)-1];
8     L = L .* (x-nodes(j+1)) ./ ( nodes(index+1)-nodes(j+1) );
9 end
```



3.4.3 Extrapolation to zero

Problem: compute $\lim_{t \rightarrow 0} f(t)$ with prescribed precision, when the evaluation of the function $y=f(t)$ is numerically unstable (\rightarrow Sect. 2.5.2) for $|t| \ll 1$.

Idea:



- ① evaluation of $f(t_i)$ for different t_i , $i = 0, \dots, n$, $|t_i| > 0$.
- ② $f(0) \approx p(0)$ with interpolation polynomial $p \in \mathcal{P}_n$, $p(t_i) = f(t_i)$.

Example 3.4.8 (Numeric differentiation through extrapolation).

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**

$$f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}. \quad (3.4.9)$$

► straightforward implementation:

MATLAB-CODE : Numeric differentiation through finite differences & relative errors.

```
x=1.1; h=2.[-1:-5:-36];
atanerr = abs(dirnumdiff(atan,x,h)-1/(1+x^2))*(1+x^2);
sqrterr = abs(dirnumdiff(sqrt,x,h)-1/(2*sqrt(x)))*(2*sqrt(x));
experr = abs(dirnumdiff(exp,x,h)-exp(x))/exp(x);

function[df]=dirnumdiff(f,x,h)
df=(f(x+h)-f(x))./h;
end
```

$$f(x) = \arctan(x)$$

| h | Relative error |
|-----------|------------------|
| 2^{-1} | 0.20786640808609 |
| 2^{-6} | 0.00773341103991 |
| 2^{-11} | 0.00024299312415 |
| 2^{-16} | 0.00000759482296 |
| 2^{-21} | 0.00000023712637 |
| 2^{-26} | 0.00000001020730 |
| 2^{-31} | 0.00000005960464 |
| 2^{-36} | 0.00000679016113 |

$$f(x) = \sqrt{x}$$

| h | Relative error |
|-----------|------------------|
| 2^{-1} | 0.09340033543136 |
| 2^{-6} | 0.00352613693103 |
| 2^{-11} | 0.00011094838842 |
| 2^{-16} | 0.00000346787667 |
| 2^{-21} | 0.00000010812198 |
| 2^{-26} | 0.00000001923506 |
| 2^{-31} | 0.00000001202188 |
| 2^{-36} | 0.00000198842224 |

$$f(x) = \exp(x)$$

| h | Relative error |
|-----------|------------------|
| 2^{-1} | 0.29744254140026 |
| 2^{-6} | 0.00785334954789 |
| 2^{-11} | 0.00024418036620 |
| 2^{-16} | 0.00000762943394 |
| 2^{-21} | 0.00000023835113 |
| 2^{-26} | 0.00000000429331 |
| 2^{-31} | 0.00000012467100 |
| 2^{-36} | 0.00000495453865 |

Code 3.4.10: Numerical differentiation by extrapolation to zero

```

1 function d = diffex(f,x,h0,tol)
2 % f: handle of to a function defined in a neighborhood of  $x \in \mathbb{R}$ ,
3 % x: point at which approximate derivative is desired,
4 % h0: initial distance from x,
5 % tol: relative target tolerance
6 h = h0;
7 % Aitken-Neville scheme, see Code 3.4.3 ( $x = 0!$ )
8 y(1) = (f(x+h0)-f(x-h0))/(2*h0);
9 for i=2:10
10     h(i) = h(i-1)/2;
11     y(i) = f(x+h(i))-f(x-h(i))/h(i-1);
12     for k=i-1:-1:1

```

```

13     y(k) = y(k+1)-(y(k+1)-y(k))*h(i)/(h(i)-h(k));
14 end
15 % termination of extrapolation, when desired tolerance is achieved
16 if (abs(y(2)-y(1)) < tol*abs(y(1))), break; end %
17 end
18 d = y(1);

```

diffex2(@atan,1.1,0.5) diffex2(@sqrt,1.1,0.5) diffex2(@exp,1.1,0.5)

| 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 | 0.00000000004816 | 6 | 0.00000000000002 |
| | | 7 | 0.00000000000021 | | |



3.4.4 Newton basis and divided differences [13, Sect. 8.2.4], [51, Sect. 8.2]

```
1 class PolyEval {  
2   private :  
3   // evaluation point and various internal data describing the 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    double operator () (double x) const;  
11 };
```

Tool: “update friendly” representation: **Newton basis** for \mathcal{P}_n

$$N_0(t) := 1, \quad N_1(t) := (t - t_0), \quad \dots, \quad N_n(t) := \prod_{i=0}^{n-1} (t - t_i). \quad (3.4.13)$$

➤ LSE for polynomial interpolation problem in Newton basis:

$$a_j \in \mathbb{R}: \quad a_0 N_0(t_j) + a_1 N_1(t_j) + \dots + a_n N_n(t_j) = y_j, \quad j = 0, \dots, n.$$

⇔ triangular linear system

$$\begin{pmatrix} 1 & 0 & \cdots & 0 \\ 1 & (t_1 - t_0) & \ddots & \vdots \\ \vdots & \vdots & \ddots & 0 \\ 1 & (t_n - t_0) & \cdots & \prod_{i=0}^{n-1} (t_n - t_i) \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{pmatrix} .$$

Simpler and more efficient algorithm using **divided differences**:

$$\begin{aligned} y[t_i] &= y_i \\ y[t_i, \dots, t_{i+k}] &= \frac{y[t_{i+1}, \dots, t_{i+k}] - y[t_i, \dots, t_{i+k-1}]}{t_{i+k} - t_i} \end{aligned} \quad \text{(recursion)} \quad (3.4.14)$$

R. Hiptmair
rev 38355,
October 26,
2011

Recursive calculation by **divided differences scheme**, cf. Aitken-Neville scheme, Code 3.4.3:

$$\begin{array}{l|l} t_0 & y[t_0] \\ & > y[t_0, t_1] \\ t_1 & y[t_1] & > y[t_0, t_1, t_2] \\ & > y[t_1, t_2] & > y[t_0, t_1, t_2, t_3], \\ t_2 & y[t_2] & > y[t_1, t_2, t_3] \\ & > y[t_2, t_3] \\ t_3 & y[t_3] \end{array} \quad (3.4.15)$$

Code 3.4.16: Divided differences, recursive implementation, in situ computation

```

1 function y = divdiff(t,y)
2 n = length(y)-1;
3 if (n > 0)
4     y(1:n) = divdiff(t(1:n),y(1:n));
5     for j=0:n-1
6         y(n+1) = (y(n+1)-y(j+1))/(t(n+1)-t(j+1));
7     end
8 end

```

Code 3.4.18: Divided differences evaluation by modified Horner scheme

```

1 function p = evaldivdiff(t,y,x)
2 dd=divdiff(t,y); % Compute divided differences, see Code 3.4.16
3 n = length(y)-1;
4 p=dd(n+1);
5 for j=n:-1:1, p = (x-t(j)).*p+dd(j); end

```

- Computational effort:
- $O(n^2)$ for computation of divided differences,
 - $O(n)$ for every single evaluation of $p(t)$.

Example 3.4.19 (Class PolyEval).

Code 3.4.20: "Update friendly" polynomial interpolant

NumCSE,
autumn
2010

```

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 };
14
15 PolyEval::PolyEval() {}
16
17 void PolyEval::addPoint(double td, double yd)
18 { t.push_back(td); y.push_back(yd); divdiff(); }
19
20 // Update coefficients in Newton basis representation, cf. Code 3.4.16
21 void PolyEval::divdiff() {
22     int n = t.size();
23     for (int j=0; j<n-1; j++) y[n-1] =

```

R. Hiptmair
rev 38355,
October 26,
2011

```
24     ((y[n-1]-y[j])/(t[n-1]-t[j]));  
25 }  
26 double PolyEval::operator() (double x) const {  
27     double s = y.back();  
28     for(int i = y.size()-2; i >= 0; --i) s = s*(x-t[i])+y[i];  
29     return s;  
30 }
```



3.5 Polynomial Interpolation: Sensitivity [51, Sect. 8.1.3]

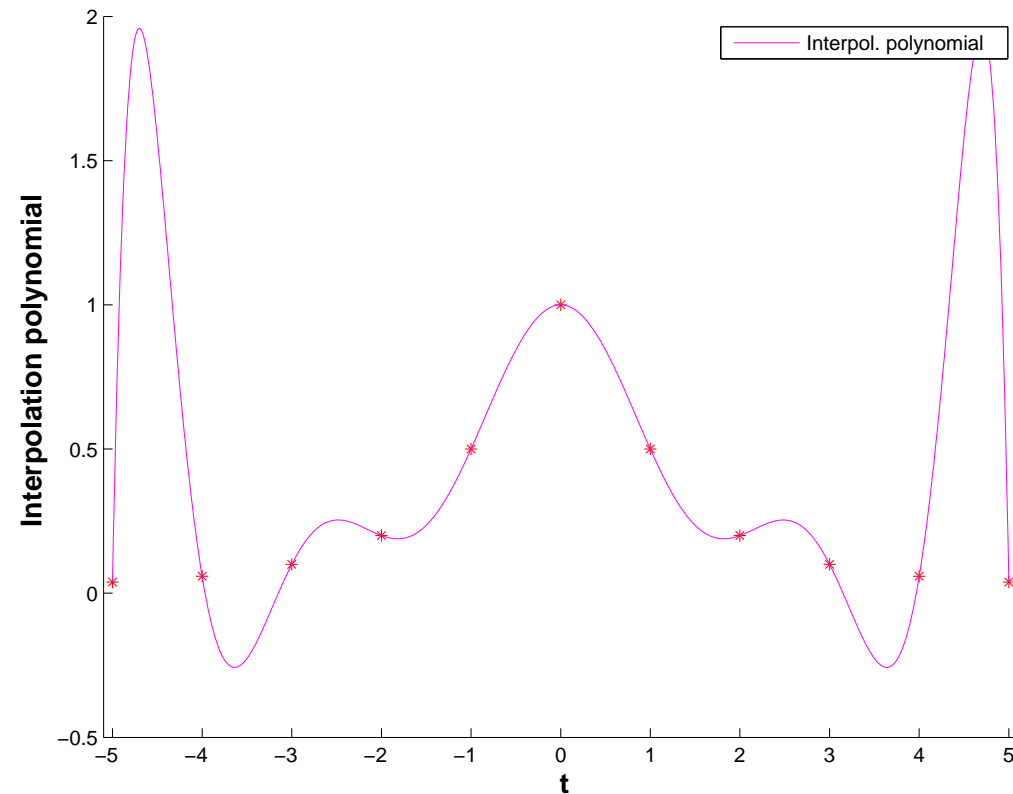
Example 3.5.1 (Oscillating interpolation polynomial: Runge's counterexample). → [13, Sect. 8.3], [51, Ex. 8.1]

Interpolation polynomial with uniformly spaced nodes:

$$\mathcal{T} := \left\{ -5 + \frac{10}{n} j \right\}_{j=0}^n,$$
$$y_j = \frac{1}{1 + t_j^2}, \quad j = 0, \dots, n.$$

Plot $n = 10 \rightarrow$

See example 9.1.5.



possible strong oscillations of interpolating polynomials
of *high degree* on *uniformly spaced* nodes!



Lemma 3.5.5 (Absolute conditioning of polynomial interpolation).

Given a mesh $\mathcal{T} \subset \mathbb{R}$ with generalized Lagrange polynomials $L_i, i = 0, \dots, n$, and fixed $I \subset \mathbb{R}$, the norm of the interpolation operator satisfies

$$\|\mathcal{I}_{\mathcal{T}}\|_{\infty \rightarrow \infty} := \sup_{\mathbf{y} \in \mathbb{K}^{n+1} \setminus \{0\}} \frac{\|\mathcal{I}_{\mathcal{T}}(\mathbf{y})\|_{L^\infty(I)}}{\|\mathbf{y}\|_\infty} = \left\| \sum_{i=0}^n |L_i| \right\|_{L^\infty(I)}, \quad (3.5.6)$$

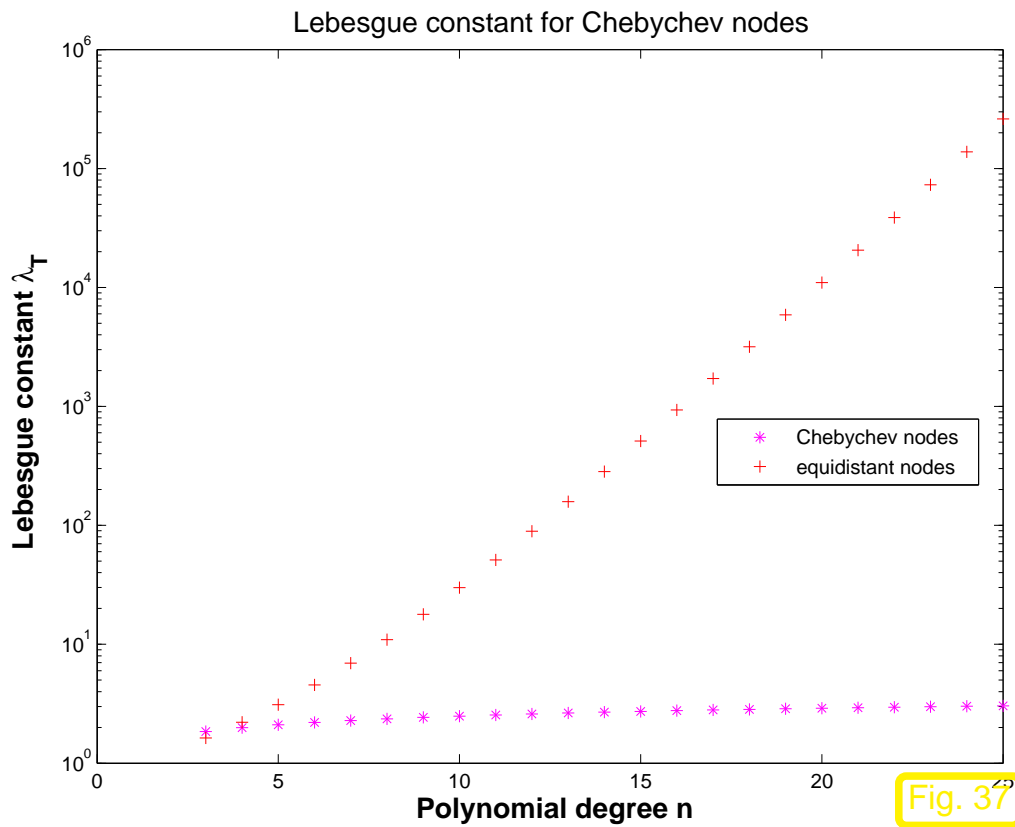
$$\|\mathcal{I}_{\mathcal{T}}\|_{2 \rightarrow 2} := \sup_{\mathbf{y} \in \mathbb{K}^{n+1} \setminus \{0\}} \frac{\|\mathcal{I}_{\mathcal{T}}(\mathbf{y})\|_{L^2(I)}}{\|\mathbf{y}\|_2} \leq \left(\sum_{i=0}^n \|L_i\|_{L^2(I)}^2 \right)^{\frac{1}{2}}. \quad (3.5.7)$$

Terminology:

Lebesgue constant of \mathcal{T} : $\lambda_{\mathcal{T}} := \left\| \sum_{i=0}^n |L_i| \right\|_{L^\infty(I)}$

Remark 3.5.8 (Lebesgue constant for equidistant nodes).

$$I = [-1, 1], \quad \mathcal{T} = \left\{ -1 + \frac{2k}{n} \right\}_{k=0}^n \quad (\text{uniformly spaced nodes})$$



Theory [11]: for uniformly spaced nodes

$$\lambda_{\mathcal{T}} \geq C e^{n/2}$$

for $C > 0$ independent of n .



3.6 Shape preserving interpolation

Example 3.6.1 (Magnetization curves).

For many materials physics stipulates properties of the functional dependence of magnetic flux B from magnetic field strength H :

- $H \mapsto B(H)$ monotone (increasing),
- $H \mapsto B(H)$ concave

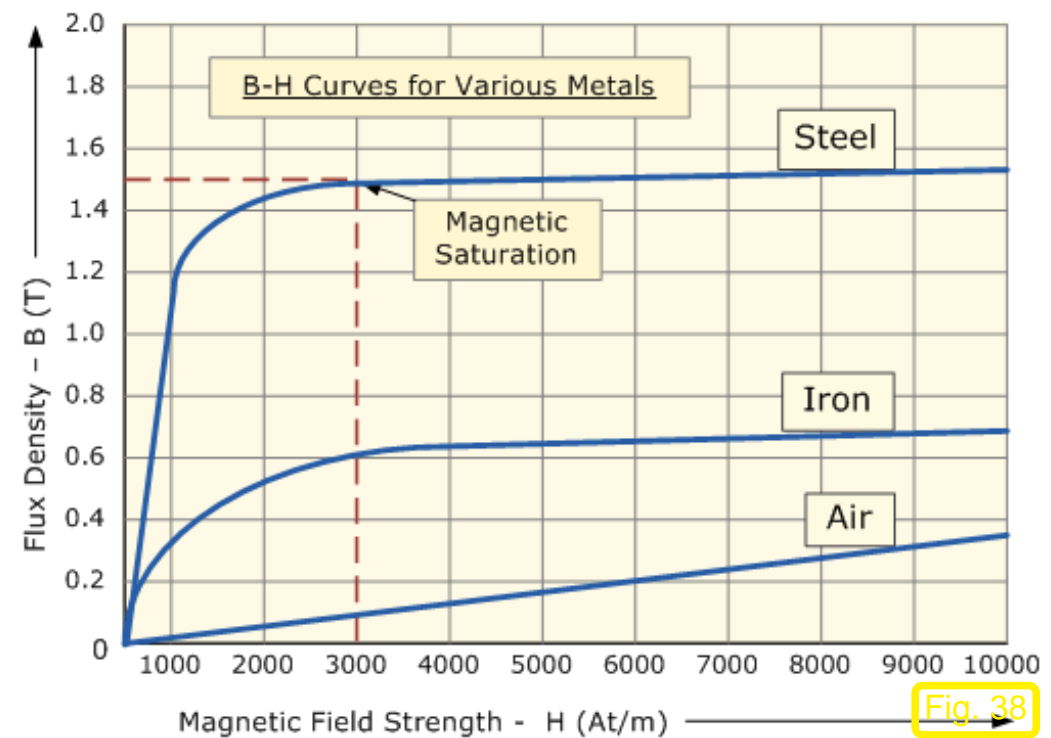


Fig. 38



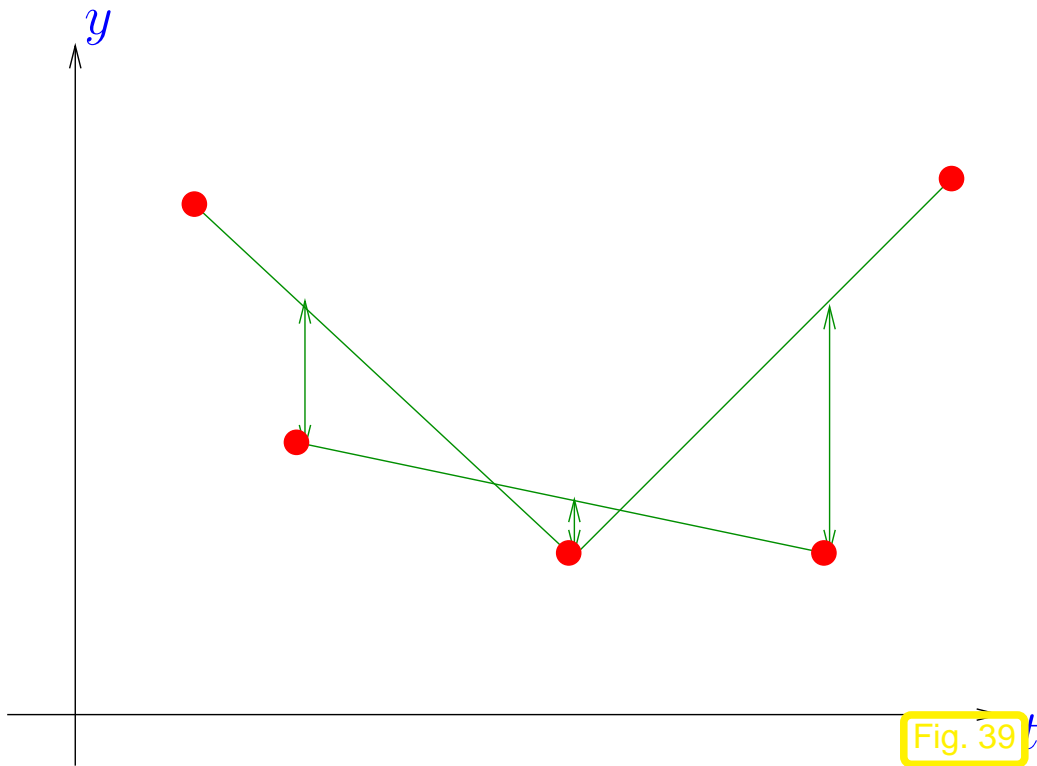
Definition 3.6.2 (monotonic data).

The data (t_i, y_i) are called *monotonic* when $y_i \geq y_{i-1}$ or $y_i \leq y_{i-1}$, $i = 1, \dots, n$.

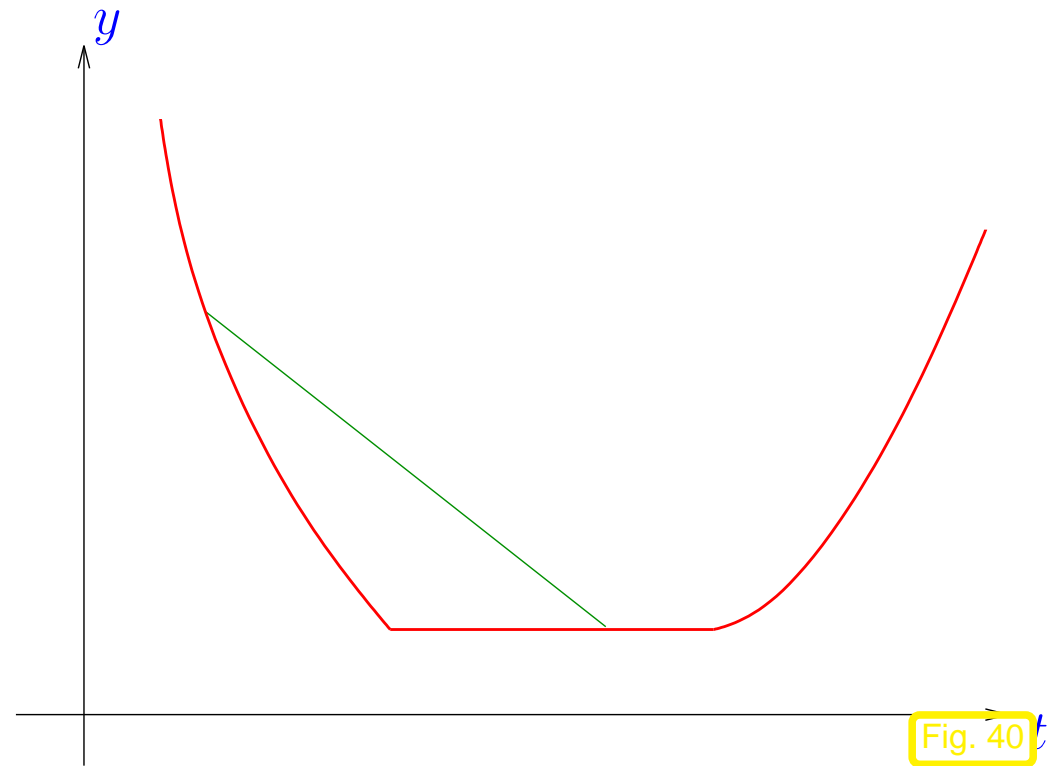
Definition 3.6.3 (Convex/concave data).

The data $\{(t_i, y_i)\}_{i=0}^n$ are called **convex** (**concave**) if

$$\Delta_j \stackrel{(\geq)}{\leq} \Delta_{j+1}, \quad j = 1, \dots, n-1, \quad \Delta_j := \frac{y_j - y_{j-1}}{t_j - t_{j-1}}, \quad j = 1, \dots, n.$$



Convex data



Convex function

Now consider interpolation: data $(t_i, y_i), i = 0, \dots, n \longrightarrow$ interpolant f

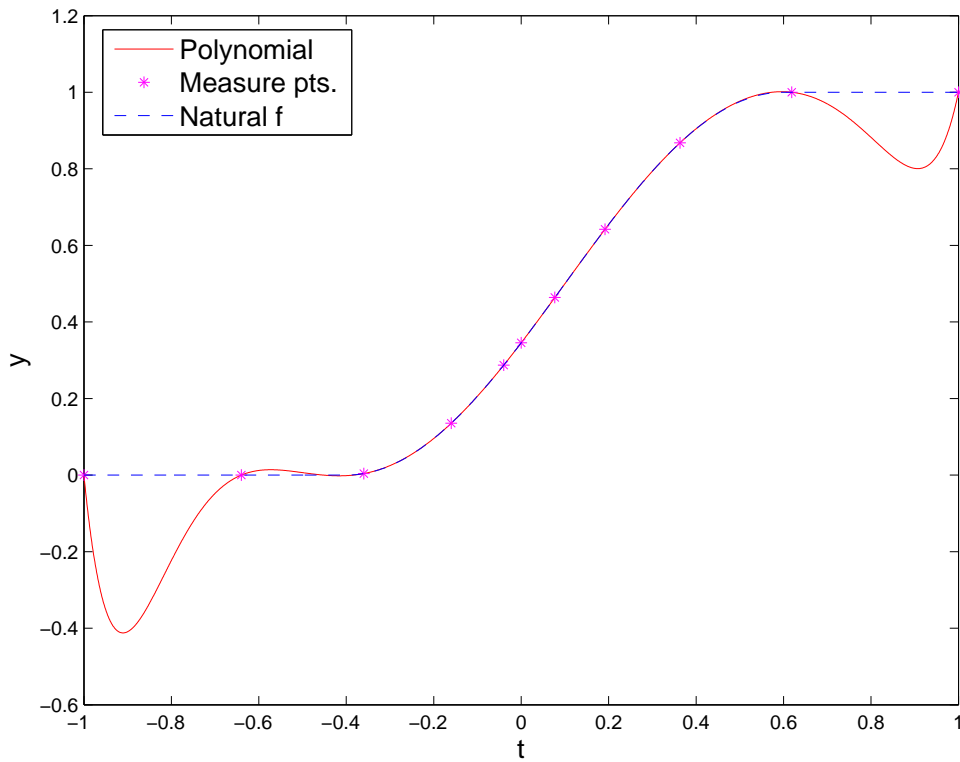
Goal: **shape preserving interpolation:**

| | | |
|----------------|---|-----------------------------|
| 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 $f _I$. |

Example 3.6.5 (Bad behavior of global polynomial interpolants).



← Interpolating polynomial, degree = 10



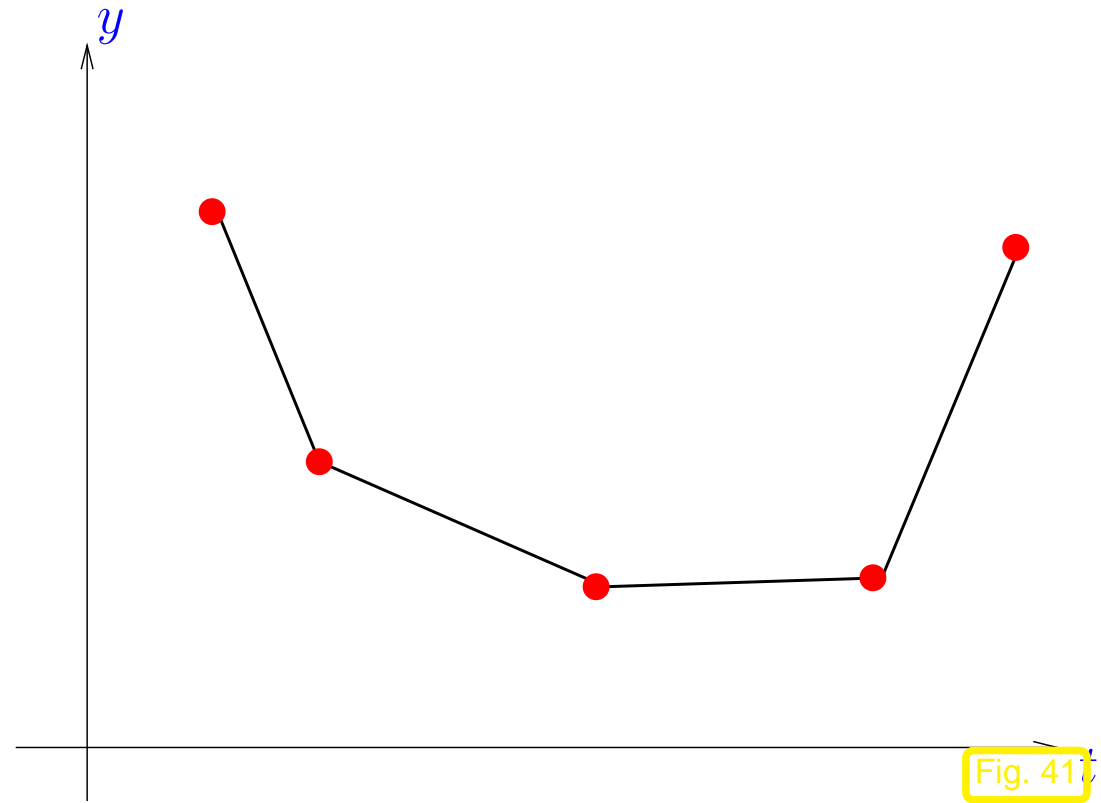
3.6.1 Piecewise linear interpolation

Data: $(t_i, y_i) \in \mathbb{R}^2, i = 0, \dots, n, n \in \mathbb{N}, t_0 < t_1 < \dots < t_n.$

Piecewise linear interpolant, cf. Ex. 3.1.5:

$$s(x) = \frac{(t_{i+1} - t)y_i + (t - t_i)y_{i+1}}{t_{i+1} - t_i} \quad t \in [t_i, t_{i+1}].$$

Piecewise linear interpolant of data in Fig. 49:



Obvious: linear interpolation is **linear** (as mapping $\mathbf{y} \mapsto s$) and **local**:

$$y_j = \delta_{ij}, \quad i, j = 0, \dots, n \quad \Rightarrow \quad \text{supp}(s) \subset [t_{i-1}, t_{i+1}].$$

Theorem 3.6.6 (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 $\Rightarrow s|_I$ is positive/negative,
- if $(t_i, y_i)|_I$ are monotonic (increasing/decreasing) $\Rightarrow s|_I$ is monotonic (increasing/decreasing),
- if $(t_i, y_i)|_I$ are convex/concave $\Rightarrow s|_I$ is convex/concave.

Example 3.6.7 (Piecewise quadratic interpolation).

Nodes as in Ex. 3.6.5

Piecewise linear/quadratic interpolation \triangleright

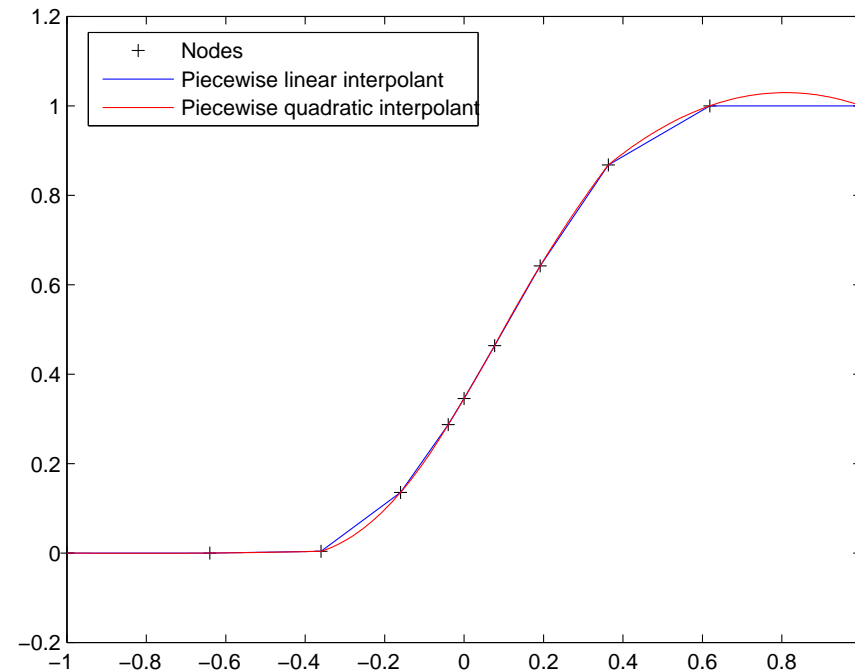


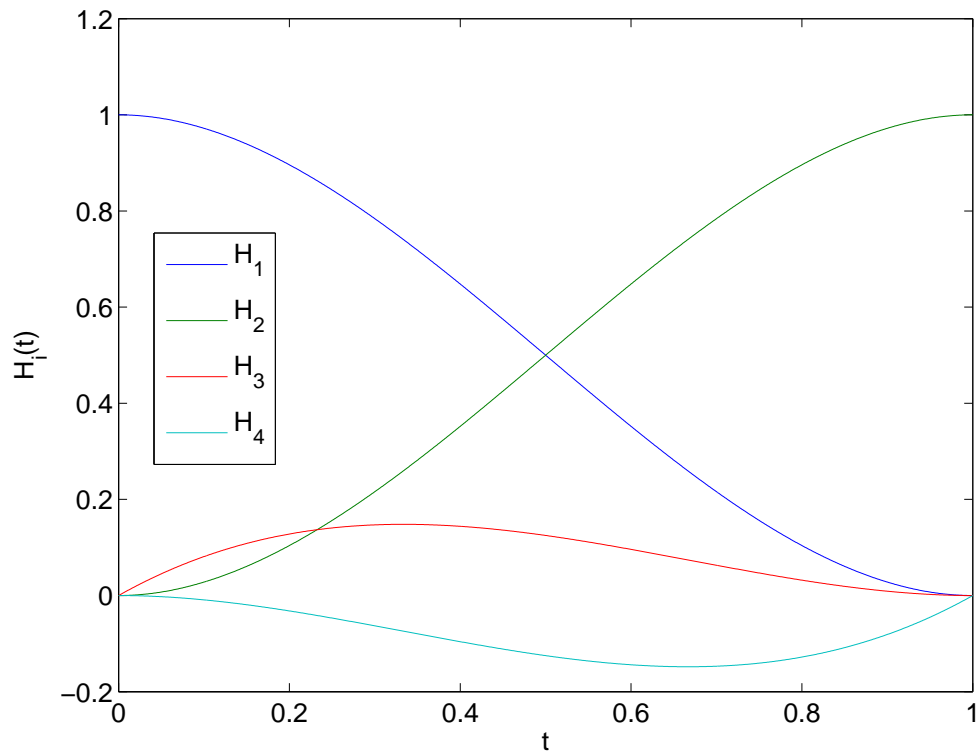
Fig. 42

3.7 Cubic Hermite Interpolation

3.7.1 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: function $f \in C^1([t_0, t_n])$, $f(t_i) = y_i$, $i = 0, \dots, n$



$$\begin{aligned}
 H_1(t) &:= \phi\left(\frac{t_i-t}{h_i}\right), & H_2(t) &:= \phi\left(\frac{t-t_{i-1}}{h_i}\right), \\
 H_3(t) &:= -h_i\psi\left(\frac{t_i-t}{h_i}\right), & H_4(t) &:= h_i\psi\left(\frac{t-t_{i-1}}{h_i}\right), \\
 h_i &:= t_i - t_{i-1}, \\
 \phi(\tau) &:= 3\tau^2 - 2\tau^3, \\
 \psi(\tau) &:= \tau^3 - \tau^2.
 \end{aligned}
 \tag{3.7.2}$$

◁ Local basis polynomial on $[0, 1]$

| | $H(t_{i-1})$ | $H(t_i)$ | $H'(t_{i-1})$ | $H'(t_i)$ |
|-------|--------------|----------|---------------|-----------|
| H_1 | 1 | 0 | 0 | 0 |
| H_2 | 0 | 1 | 0 | 0 |
| H_3 | 0 | 0 | 1 | 0 |
| H_4 | 0 | 0 | 0 | 1 |

How to choose the slopes c_i ?

Natural attempt: (weighted) average of local slopes:

$$c_i = \begin{cases} \Delta_1 & , \text{ for } i = 0 , \\ \Delta_n & , \text{ for } i = n , \\ \frac{t_{i+1}-t_i}{t_{i+1}-t_{i-1}}\Delta_i + \frac{t_i-t_{i-1}}{t_{i+1}-t_{i-1}}\Delta_{i+1} & , \text{ if } 1 \leq i < n . \end{cases} , \quad \Delta_j := \frac{y_j - y_{j-1}}{t_j - t_{j-1}} , j = 1, \dots, n .$$

(3.7.4)

► **Linear** local interpolation operator

Example 3.7.5 (Piecewise cubic 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

$$f(x) := \sin(5x) e^x.$$

Use of weighted averages of slopes as
in (3.7.4).

See Code 3.7.5.

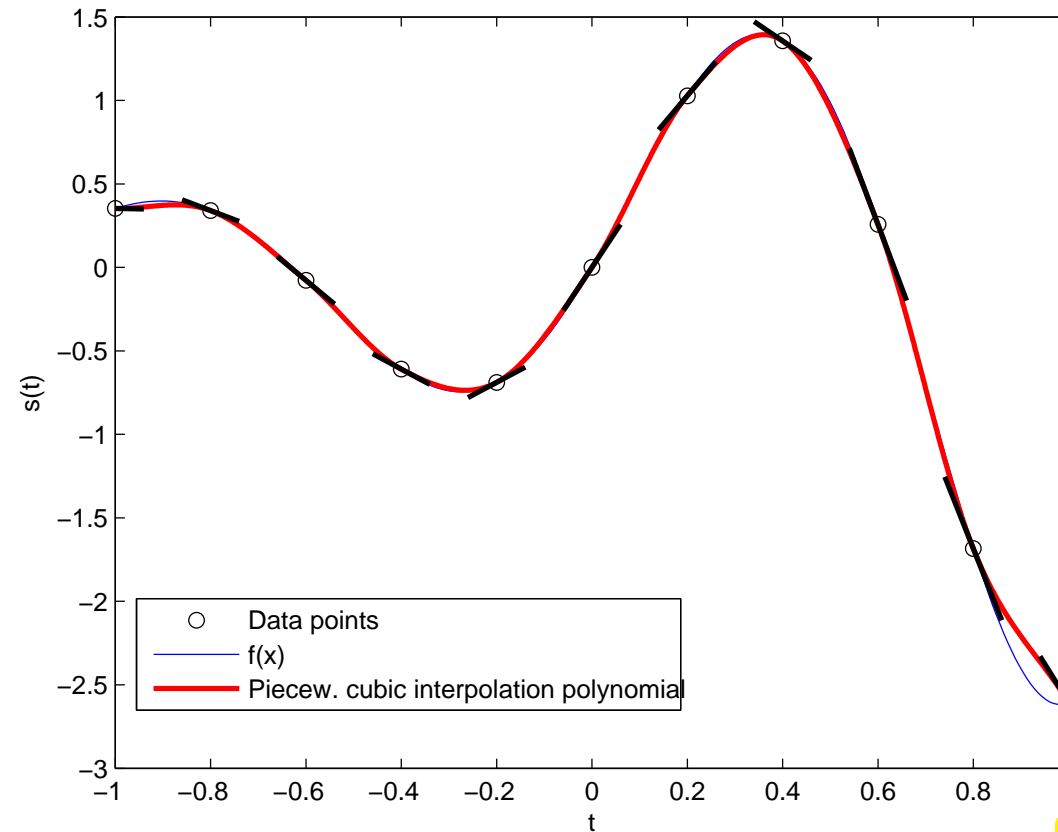


Fig. 43

No preservation of monotonicity!



3.7.2 Shape preserving Hermite interpolation

$$c_i = \begin{cases} 0 & , \text{ if } \operatorname{sgn}(\Delta_i) \neq \operatorname{sgn}(\Delta_{i+1}) , \\ \text{weighted average of } \Delta_i, \Delta_{i+1} & \text{otherwise} \end{cases} , \quad i = 1, \dots, n-1 . \quad (3.7.7)$$

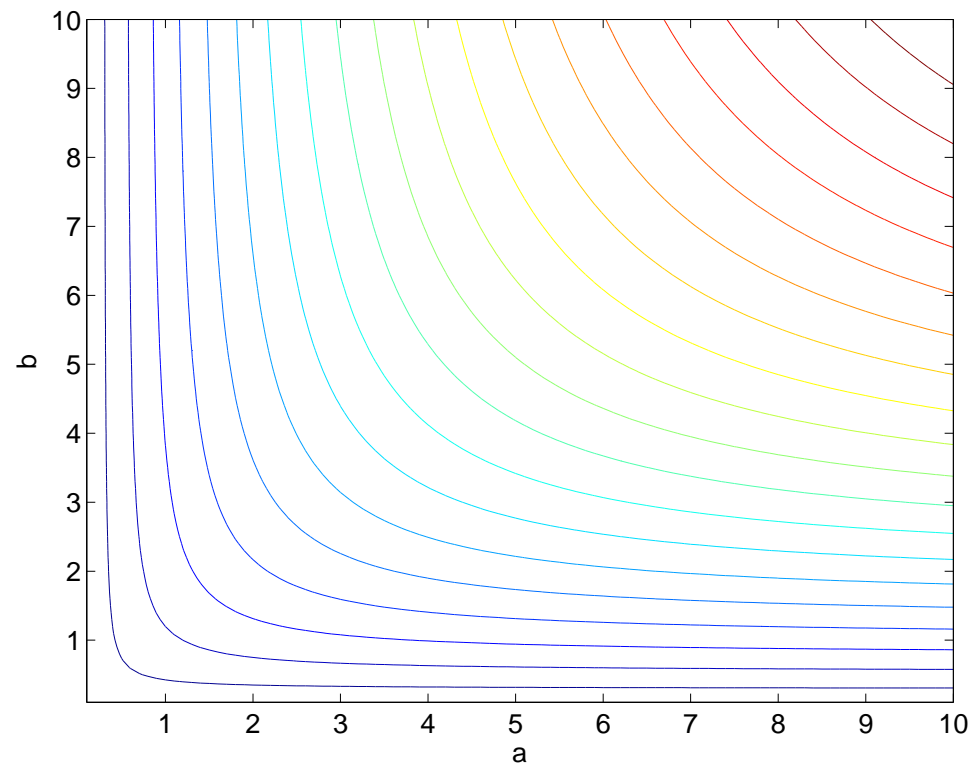
Which kind of average ?

$$c_i = \frac{1}{\frac{w_a}{\Delta_i} + \frac{w_b}{\Delta_{i+1}}} \quad (3.7.8)$$

= weighted **harmonic mean** of the slopes with weights w_a, w_b , ($w_a + w_b = 1$).

Harmonic mean = “smoothed $\min(\cdot, \cdot)$ -function”.

Contour plot of the harmonic mean of a and b \rightarrow
($w_a = w_b = 1/2$).



Concrete choice of the weights:

$$w_a = \frac{2h_{i+1} + h_i}{3(h_{i+1} + h_i)}, \quad w_b = \frac{h_{i+1} + 2h_i}{3(h_{i+1} + h_i)},$$

$$\begin{array}{l} \text{sgn}(\Delta_1) = \text{sgn}(\Delta_2) \\ \rightarrow \end{array} \quad c_i = \begin{cases} \Delta_1 & , \text{ if } i = 0 , \\ \frac{3(h_{i+1} + h_i)}{\frac{2h_{i+1} + h_i}{\Delta_i} + \frac{2h_i + h_{i+1}}{\Delta_{i+1}}} & , \text{ for } i \in \{1, \dots, n-1\} , \\ \Delta_n & , \text{ if } i = n , \end{cases} \quad h_i := t_i - t_{i-1} .$$

(3.7.9)

Example 3.7.10 (Monotonicity preserving piecewise cubic polynomial interpolation).

Data from ex. 3.6.5

MATLAB-function:

```
v = pchip(t,y,x);
```

t: Sampling points

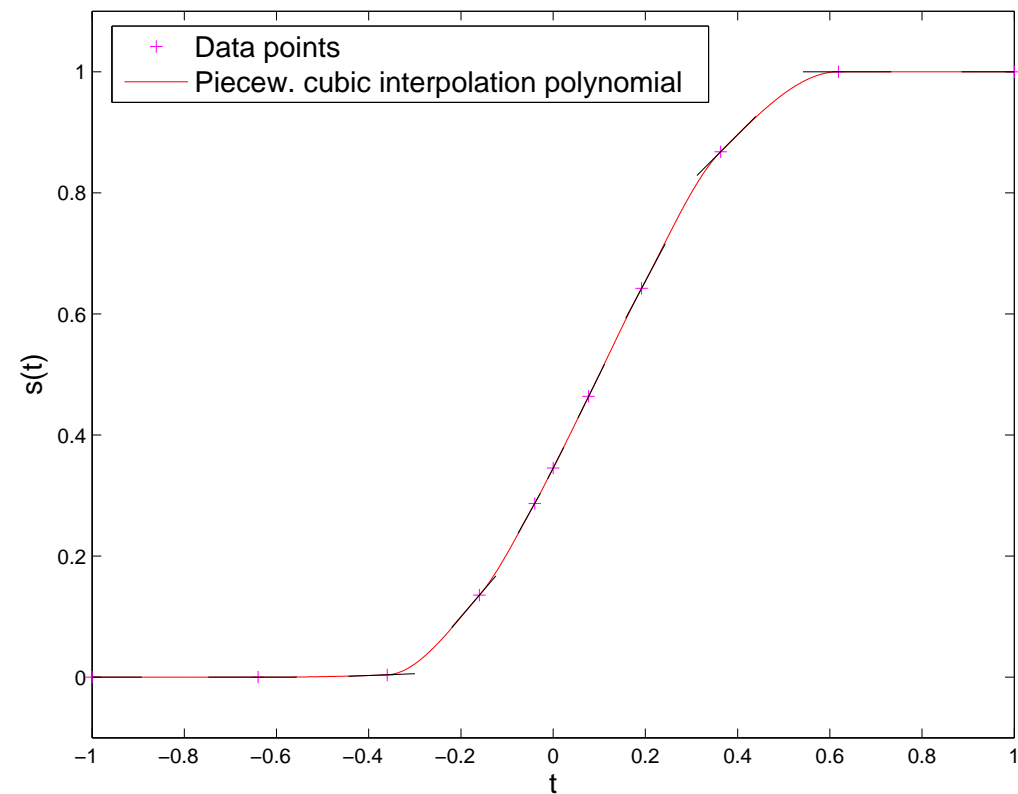
y: Sampling values

x: Evaluation points

v: Vector $s(x_i)$

Local interpolation operator

! Non linear interpolation operator



Theorem 3.7.11 (Monotonicity preservation of limited cubic Hermite interpolation).

The cubic Hermite interpolation polynomial with slopes as in (3.7.9) provides a local monotonicity-preserving C^1 -interpolant.

Definition 3.8.1 (Spline space). \rightarrow [51, 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

$$\mathcal{S}_{d,\mathcal{M}} := \{s \in C^{d-1}(I) : s_j := s|_{[t_{j-1}, t_j]} \in \mathcal{P}_d \forall j = 1, \dots, n\} .$$

$d - 1$ -times continuously differentiable locally polynomial of degree d

- $d = 0$: \mathcal{M} -piecewise constant *discontinuous* functions
- $d = 1$: \mathcal{M} -piecewise linear *continuous* functions
- $d = 2$: *continuously differentiable* \mathcal{M} -piecewise quadratic functions

3.8.1 Cubic spline interpolation [35, XIII, 46], [51, Sect. 8.6.1]

Another special case: **cubic spline interpolation**, $d = 3$ (related to Hermite interpolation, Sect. 3.7)

$$(3.7.1) \quad \blacktriangleright \quad s_{|[t_{j-1}, t_j]}(t) = \begin{aligned} & s(t_{j-1}) \cdot (1 - 3\tau^2 + 2\tau^3) + \\ & s(t_j) \cdot (3\tau^2 - 2\tau^3) + \\ & h_j s'(t_{j-1}) \cdot (\tau - 2\tau^2 + \tau^3) + \\ & h_j s'(t_j) \cdot (-\tau^2 + \tau^3), \end{aligned} \quad (3.8.3)$$

with $h_j := t_j - t_{j-1}$, $\tau := (t - t_{j-1})/h_j$.

$s \in C^2([t_0, t_n]) \Rightarrow n - 1$ continuity constraints for $s''(t)$ at the internal nodes

$$s''_{|[t_{j-1}, t_j]}(t_j) = s''_{|[t_j, t_{j+1}]}(t_j), \quad j = 1, \dots, n - 1. \quad (3.8.4)$$

(3.8.4) $\rightarrow n - 1$ **linear** equations for n slopes $c_j := s'(t_j)$

$$\frac{1}{h_j} c_{j-1} + \left(\frac{2}{h_j} + \frac{2}{h_{j+1}} \right) c_j + \frac{1}{h_{j+1}} c_{j+1} = 3 \left(\frac{y_j - y_{j-1}}{h_j^2} + \frac{y_{j+1} - y_j}{h_{j+1}^2} \right), \quad (3.8.6)$$

for $j = 1, \dots, n - 1$.

(3.8.6) \Leftrightarrow *underdetermined* $(n - 1) \times (n + 1)$ linear system of equations

$$\begin{pmatrix} b_0 & a_1 & b_1 & 0 & \cdots & \cdots & 0 \\ 0 & b_1 & a_2 & b_2 & & & \\ & 0 & \ddots & \ddots & \ddots & & \vdots \\ \vdots & & & \ddots & \ddots & \ddots & \\ & & & \ddots & a_{n-2} & b_{n-2} & 0 \\ 0 & \cdots & \cdots & 0 & b_{n-2} & a_{n-1} & b_{n-1} \end{pmatrix} \begin{pmatrix} c_0 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} 3 \left(\frac{y_1 - y_0}{h_1^2} + \frac{y_2 - y_1}{h_2^2} \right) \\ \vdots \\ 3 \left(\frac{y_{n-1} - y_{n-2}}{h_{n-1}^2} + \frac{y_n - y_{n-1}}{h_n^2} \right) \end{pmatrix}. \quad (3.8.7)$$

with

$$\begin{aligned} b_i &:= \frac{1}{h_{i+1}}, \quad i = 0, 1, \dots, n - 1, \\ a_i &:= \frac{2}{h_i} + \frac{2}{h_{i+1}}, \quad i = 0, 1, \dots, n - 1. \\ &[\quad b_i, a_i > 0 \quad , \quad a_i = 2(b_i + b_{i-1}) \quad] \end{aligned}$$

→ two additional constraints are required, (at least) three different choices are possible:

① **Complete cubic spline interpolation:** $s'(t_0) = c_0, s'(t_n) = c_n$ prescribed.

② **Natural cubic spline interpolation:** $s''(t_0) = s''(t_n) = 0$

$$\frac{2}{h_1}c_0 + \frac{1}{h_1}c_1 = 3 \frac{y_1 - y_0}{h_1^2}, \quad \frac{1}{h_n}c_{n-1} + \frac{2}{h_n}c_n = 3 \frac{y_n - y_{n-1}}{h_n^2}.$$

③ **Periodic cubic spline interpolation:** $s'(t_0) = s'(t_n), s''(t_0) = s''(t_n)$

$n \times n$ -linear system with s.p.d. coefficient matrix

$$\mathbf{A} := \begin{pmatrix} a_1 & b_1 & 0 & \cdots & 0 & b_0 \\ b_1 & a_2 & b_2 & & & 0 \\ 0 & \cdots & \cdots & \cdots & & \vdots \\ \vdots & & \cdots & \cdots & \cdots & 0 \\ 0 & & & \cdots & a_{n-1} & b_{n-1} \\ b_0 & 0 & \cdots & 0 & b_{n-1} & a_0 \end{pmatrix}, \quad \begin{aligned} b_i &:= \frac{1}{h_{i+1}}, \quad i = 0, 1, \dots, n-1, \\ a_i &:= \frac{2}{h_i} + \frac{2}{h_{i+1}}, \quad i = 0, 1, \dots, n-1. \end{aligned}$$

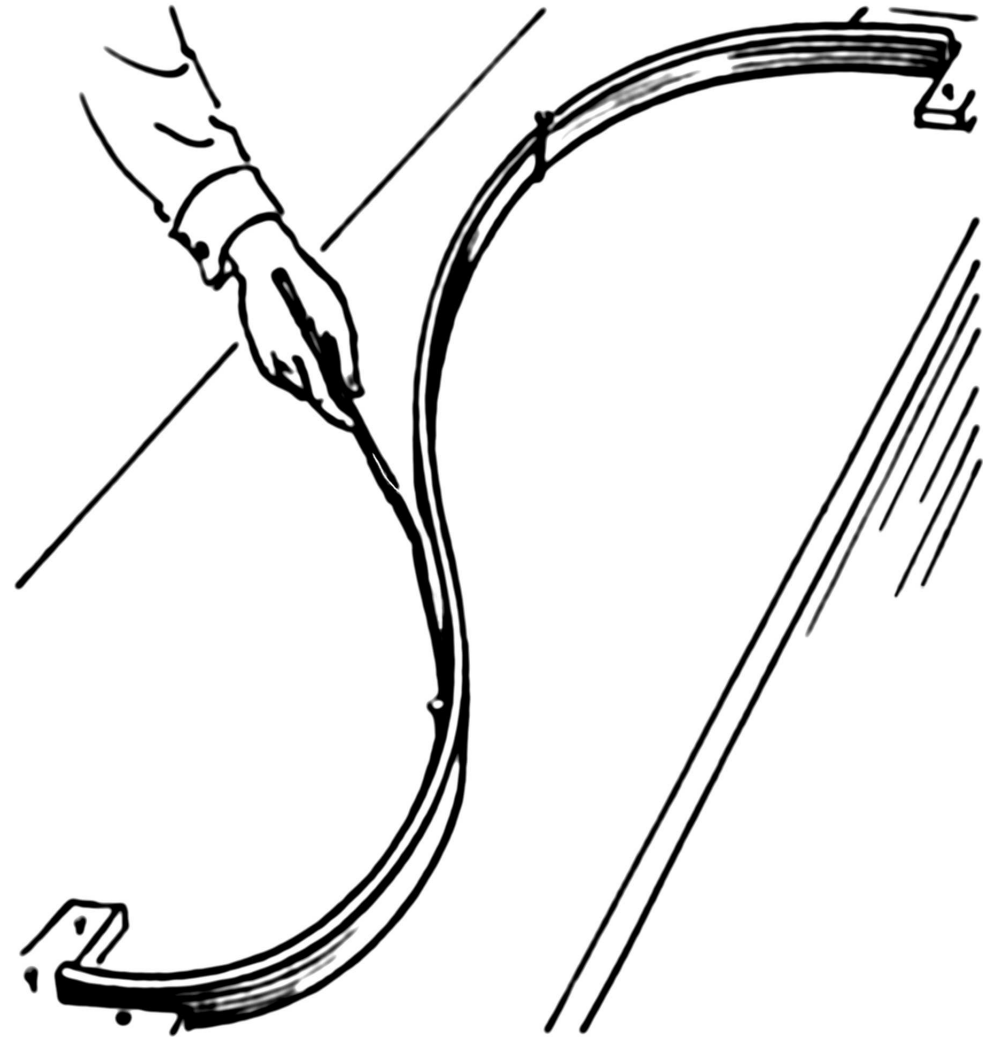
MATLAB-function: $\mathbf{v} = \text{spline}(\mathbf{t}, \mathbf{y}, \mathbf{x})$: natural / complete spline interpolation
(see spline-toolbox in MATLAB)

Structural properties of cubic spline interpolants

Remark 3.8.9 (Extremal properties of natural cubic spline interpolants). \rightarrow [51, Sect. 8.6.1, Property 8.2]

Remark 3.8.12 (Origin of the term “Spline”).

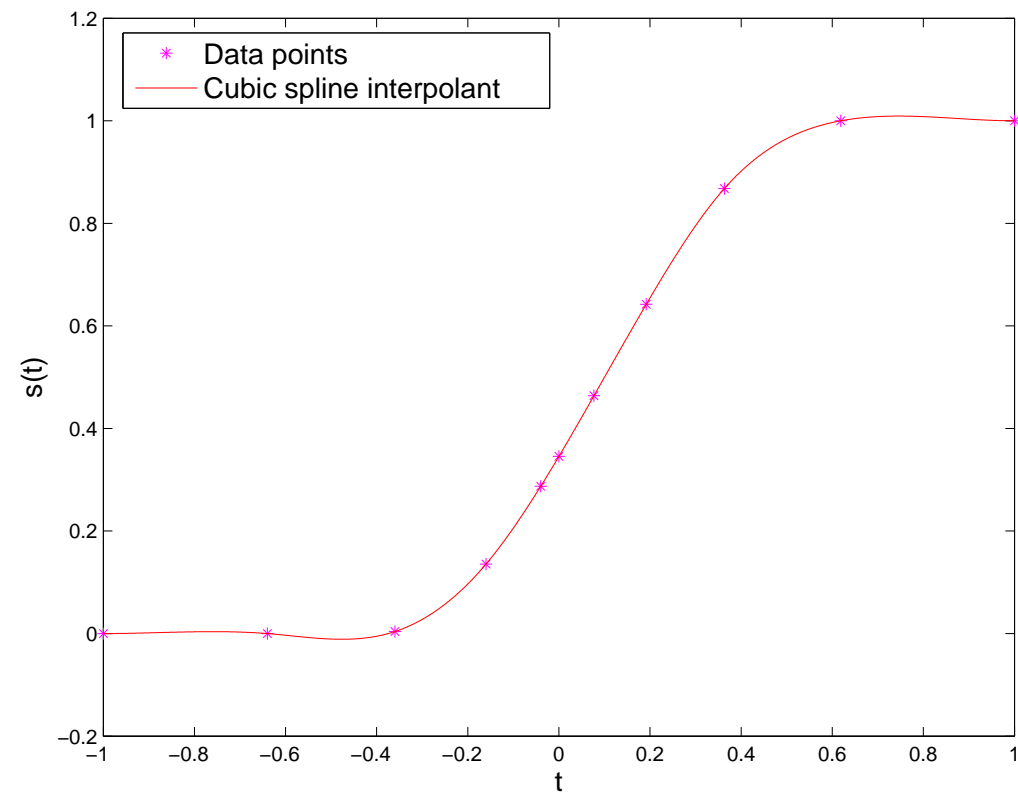
Cubic spline interpolation
before MATLAB



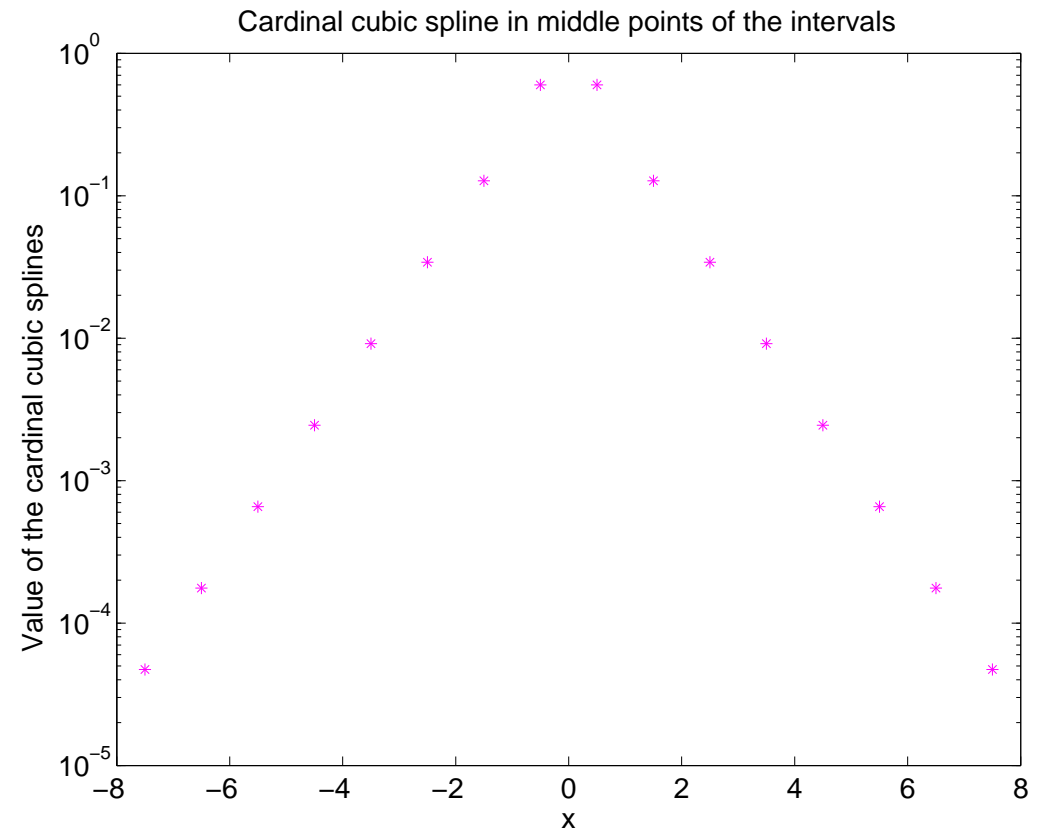
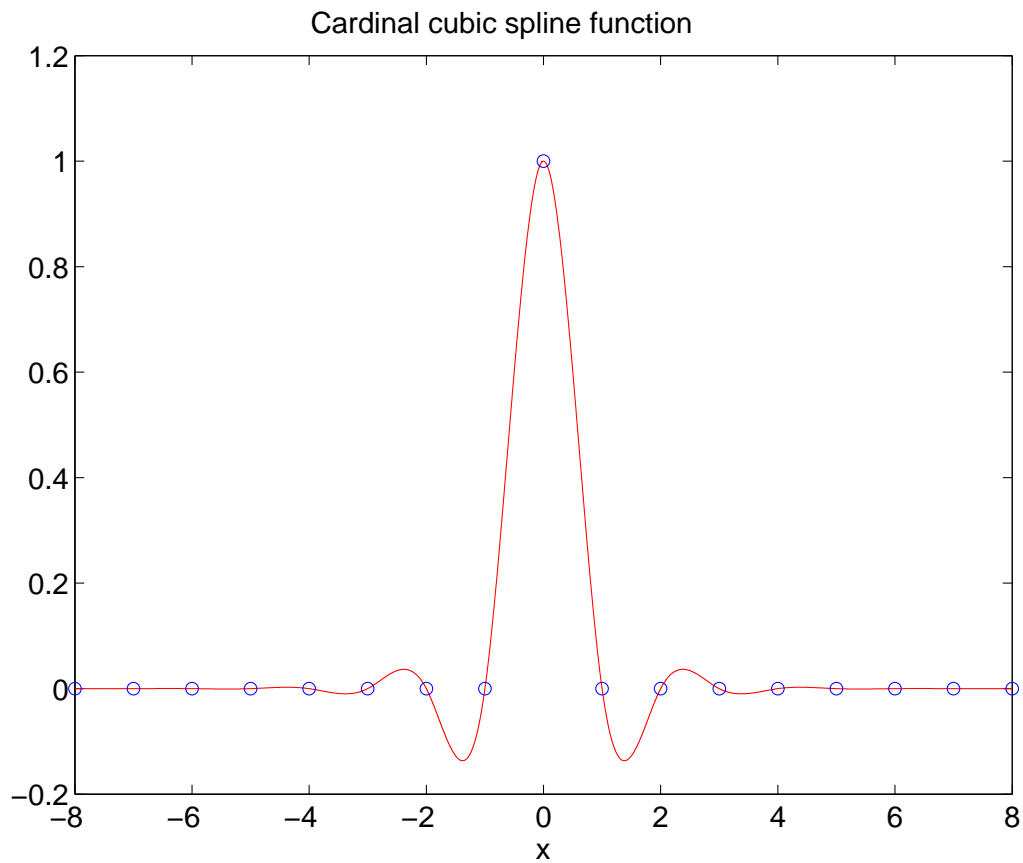
Remark 3.8.13 (Shape preservation).

Data $s(t_j) = y_j$ from Ex. 3.6.5 and

$$c_0 := \frac{y_1 - y_0}{t_1 - t_0},$$
$$c_n := \frac{y_n - y_{n-1}}{t_n - t_{n-1}}.$$



Example 3.8.14 (Locality of the natural cubic spline interpolation).

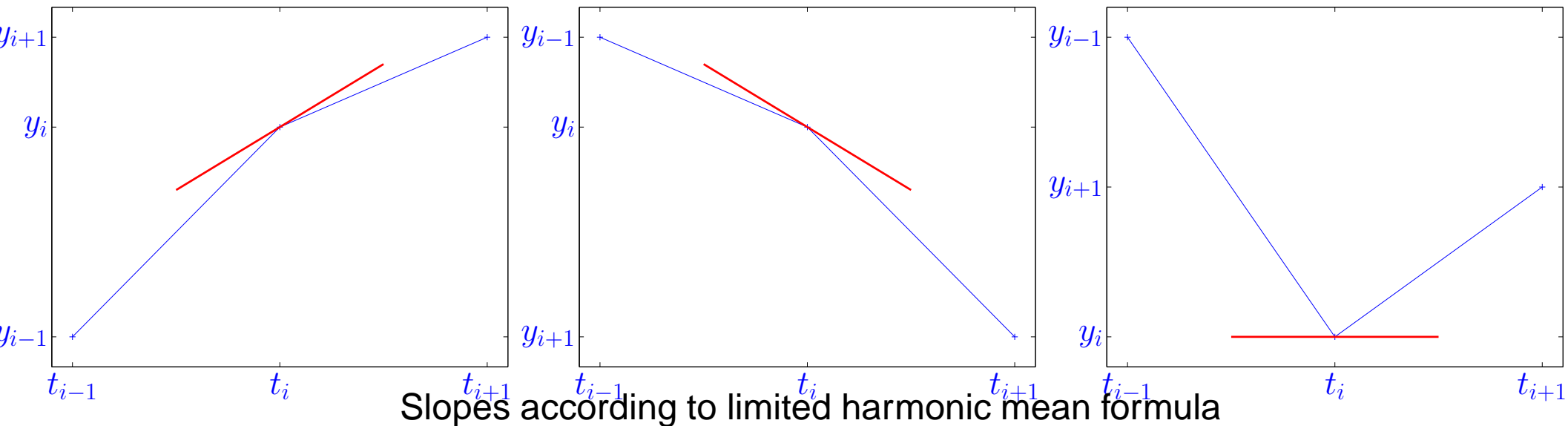


3.8.2 Shape Preserving Spline Interpolation

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:
- knot set $\mathcal{M} \subset [t_0, t_n]$ (\rightarrow Def 3.8.1),
 - 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 (\rightarrow Sect. 3.6)

① Shape preserving choice of slopes c_i , $i = 0, \dots, n$ [43, 50], analogous to Sect. 3.7.2



② Choice of “extra knots” $p_i \in]t_{i-1}, t_i]$, $i = 1, \dots, n$:

These points will be used to build the knot set for the final **quadratic spline**:

$$\mathcal{M} = \{t_0 < p_1 \leq t_1 < p_2 \leq \dots < p_n \leq t_n\} .$$

```

p = (t(1)-1)*ones(1,length(t)-1);
for j=1:n-1
    if (c(j) ~= c(j+1))
        p(j)=(y(j+1)-y(j)+...
            t(j)*c(j)-t(j+1)*c(j+1))/...
            (c(j)-c(j+1));
    end
    if ((p(j)<t(j))|(p(j)>t(j+1)))
        p(j) = 0.5*(t(j)+t(j+1));
    end;
end
end

```

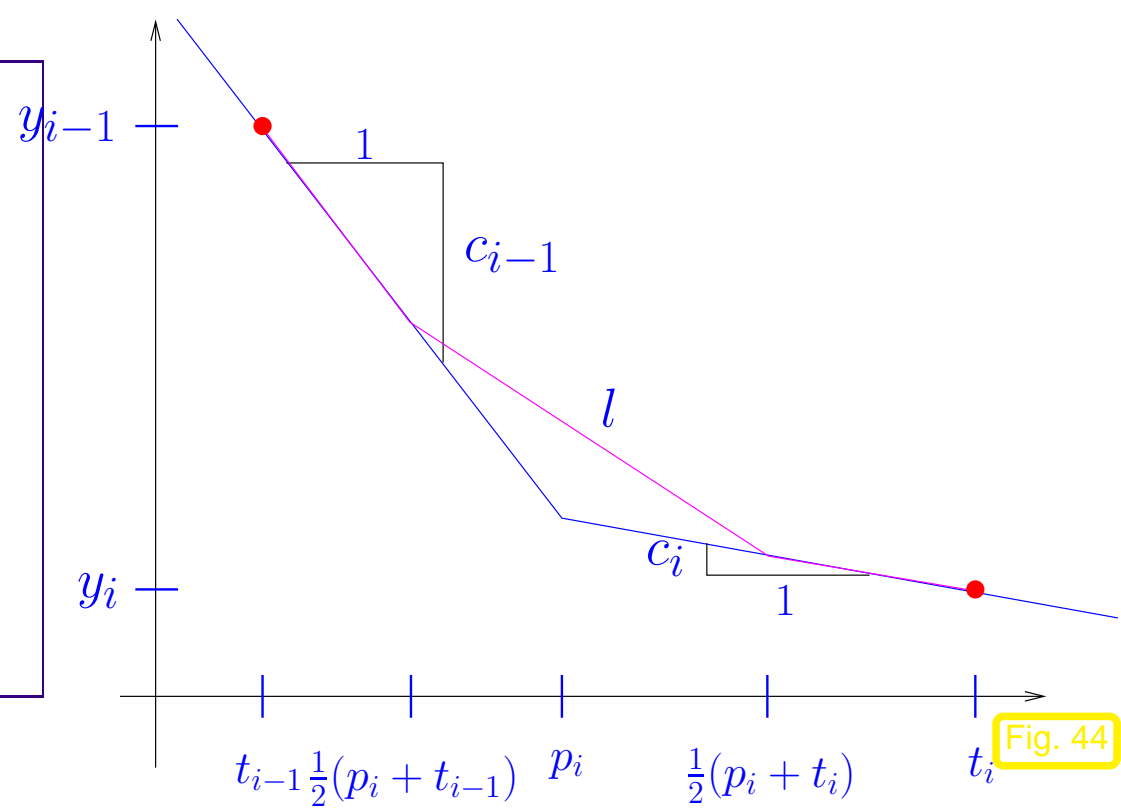


Fig. 44

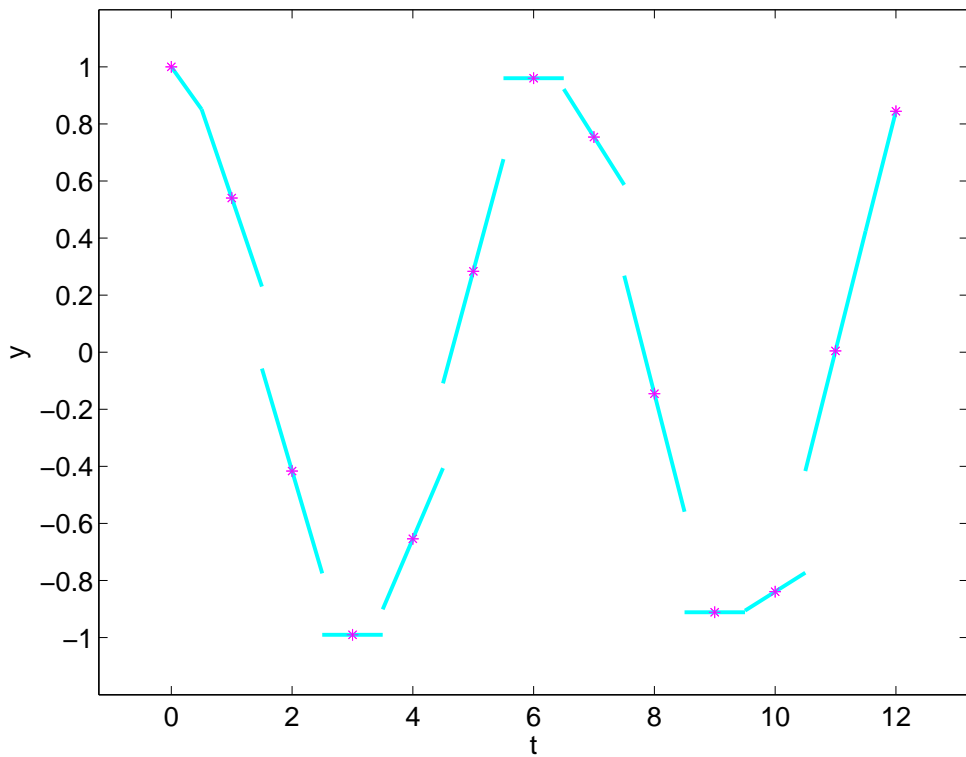
③ Set $l =$ linear spline (polygon) on the knot set \mathcal{M}' (middle points of \mathcal{M})

$$\mathcal{M}' = \{t_0 < \frac{1}{2}(t_0 + p_1) < \frac{1}{2}(p_1 + t_1) < \frac{1}{2}(t_1 + p_2) < \dots < \frac{1}{2}(t_{n-1} + p_n) < \frac{1}{2}(p_n + t_n) < t_n\};$$

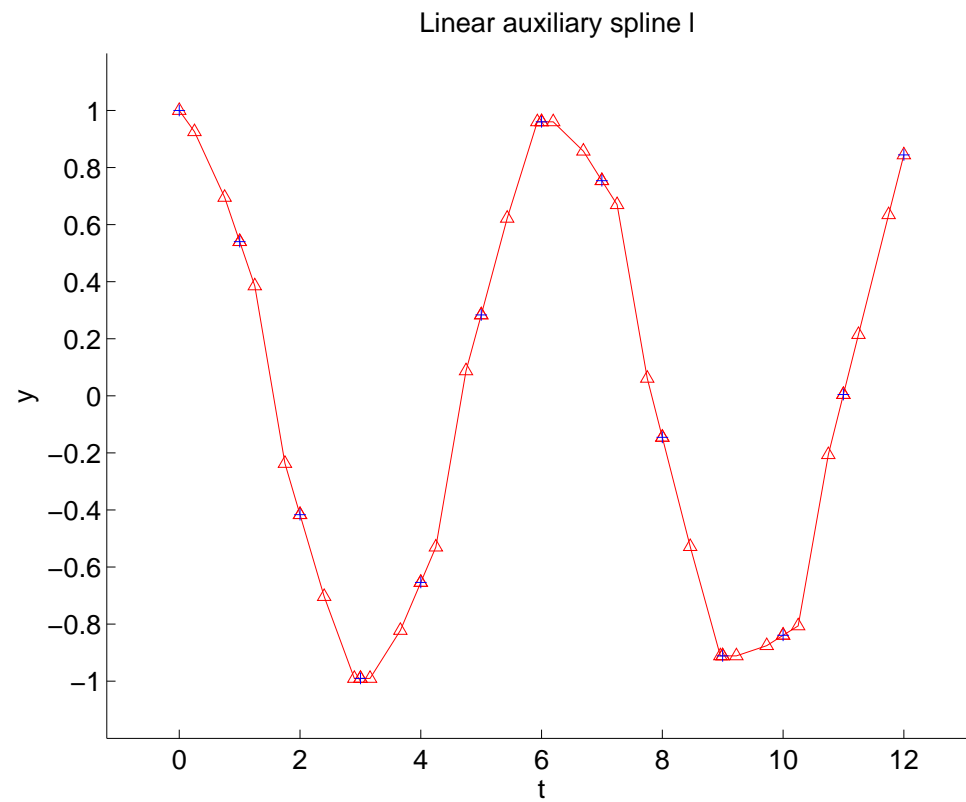
with $l(t_i) = y_i$, $l'(t_i) = c_i$.

Example 3.8.16 (Auxiliary construction for shape preserving quadratic spline interpolation).

Data points: $t=(0:12); y = \cos(t);$



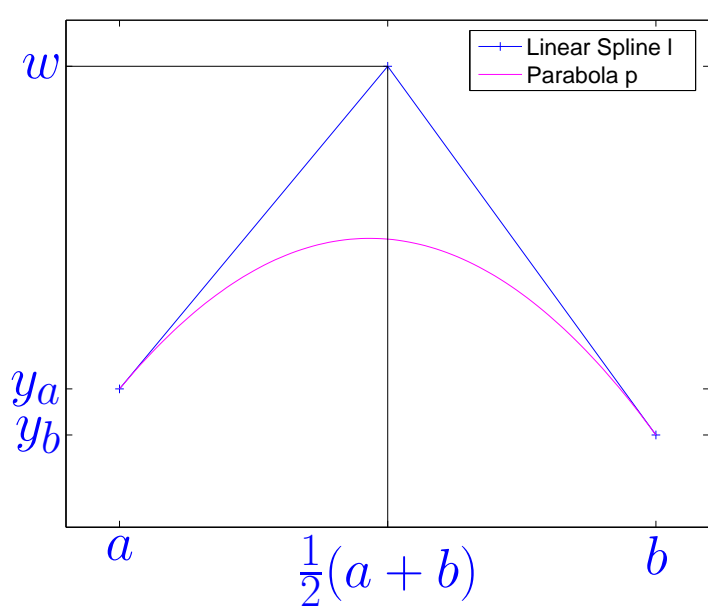
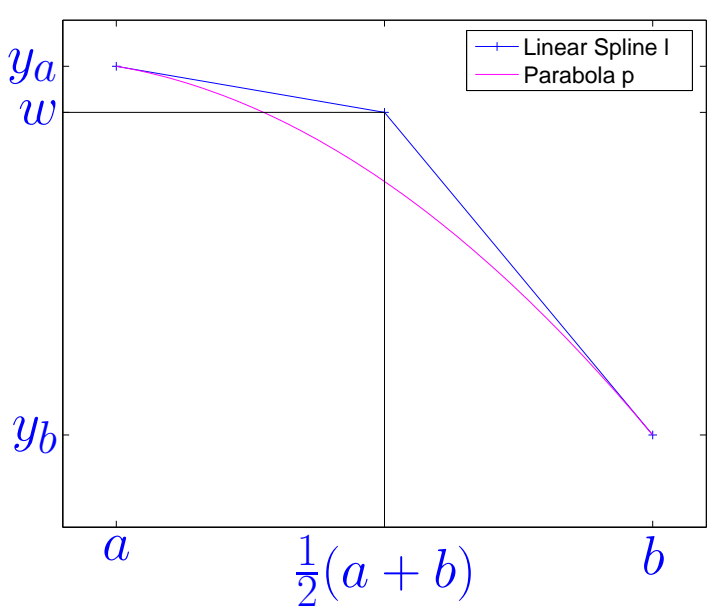
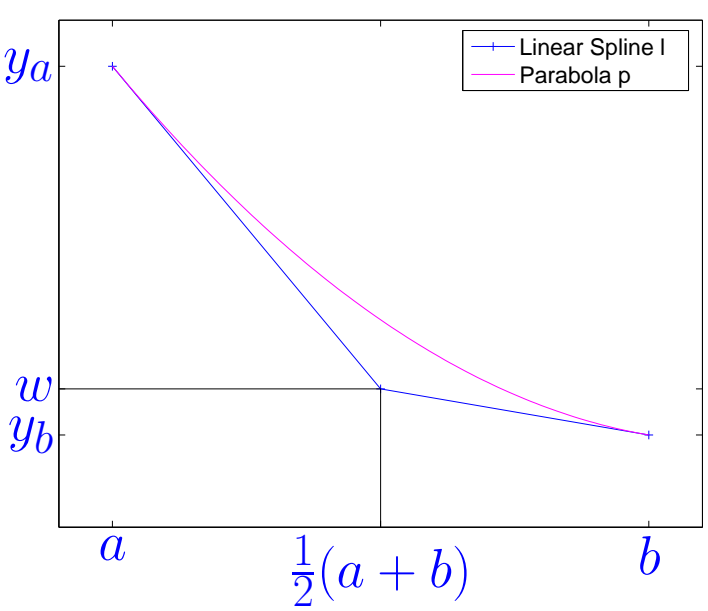
Local slopes $c_i, i = 0, \dots, n$



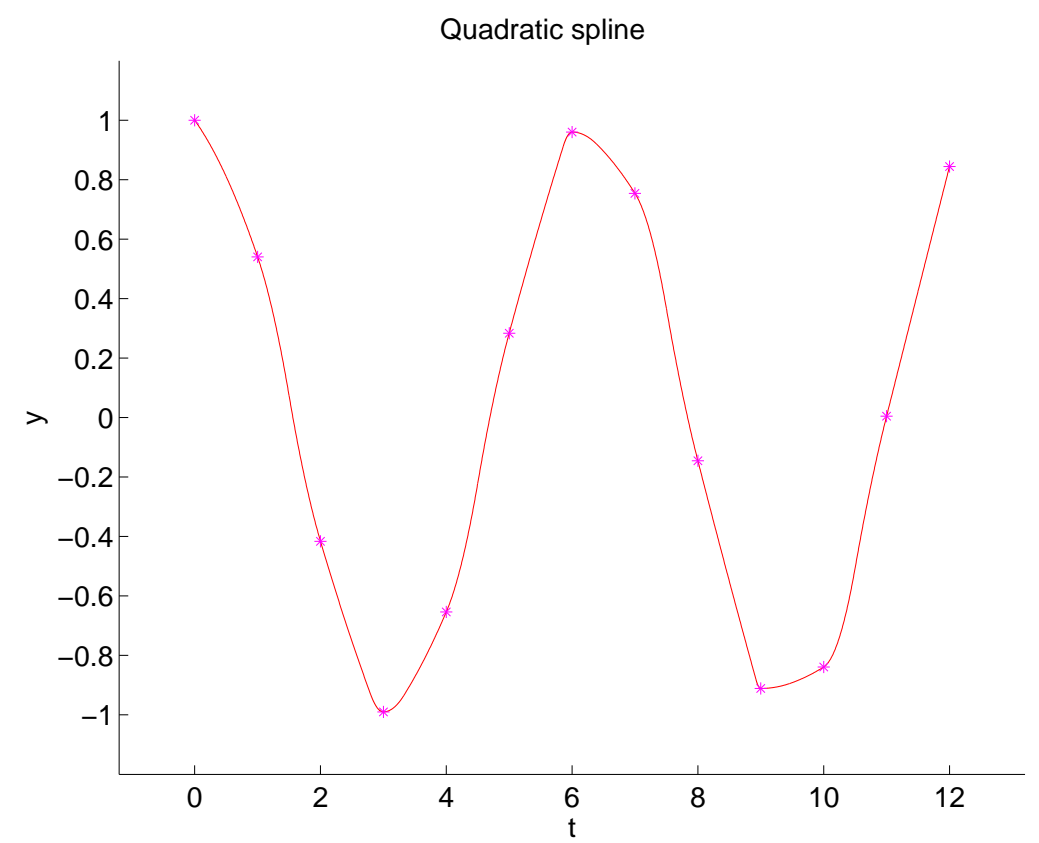
Linear auxiliary spline l

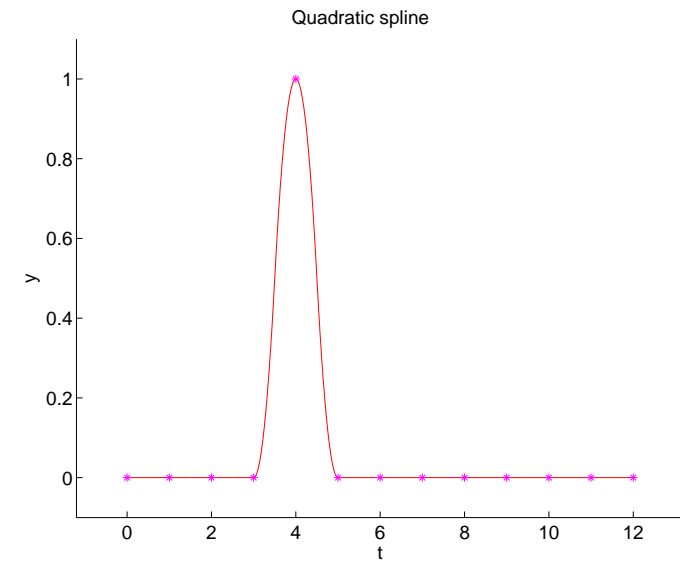
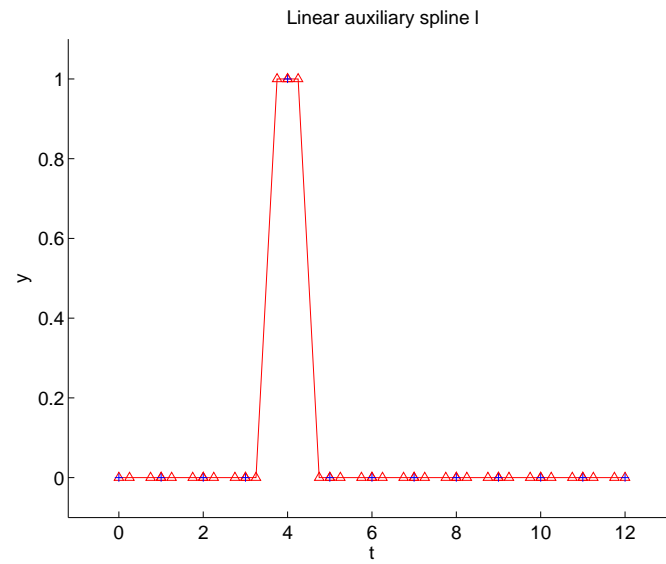
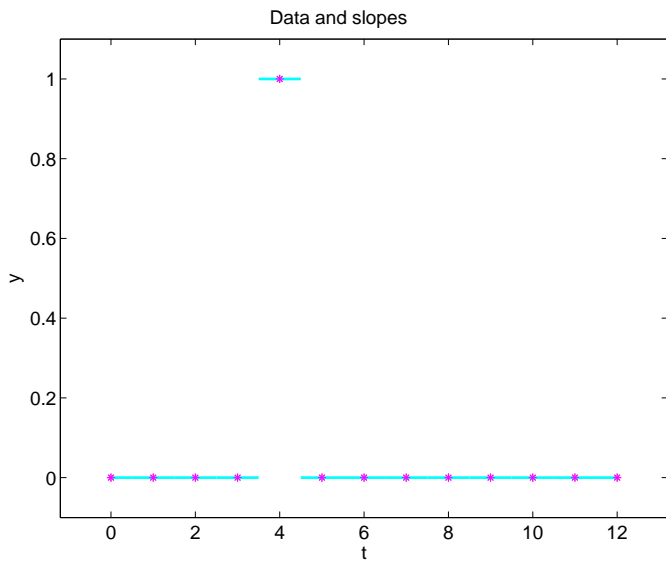


④ Local quadratic approximation / interpolation of l :



Continuation of Ex. 3.8.16:
Interpolating quadratic spline

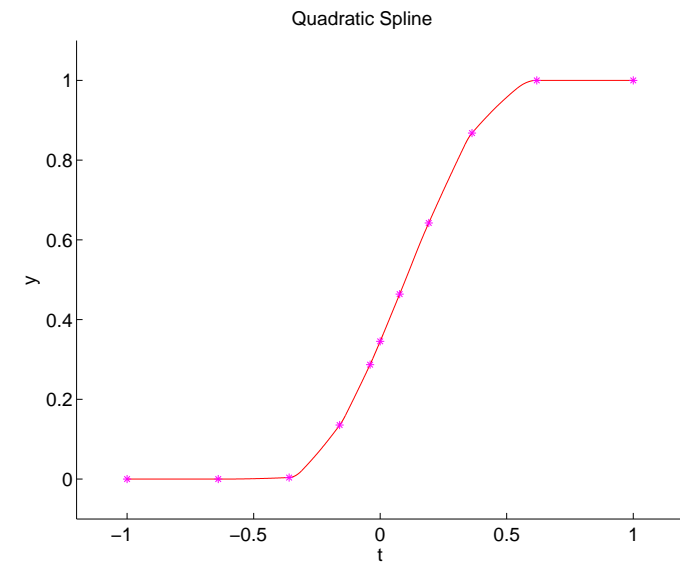
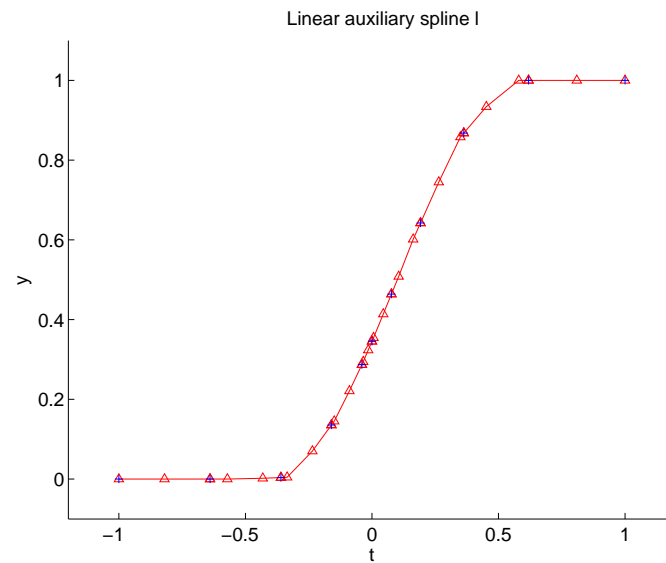
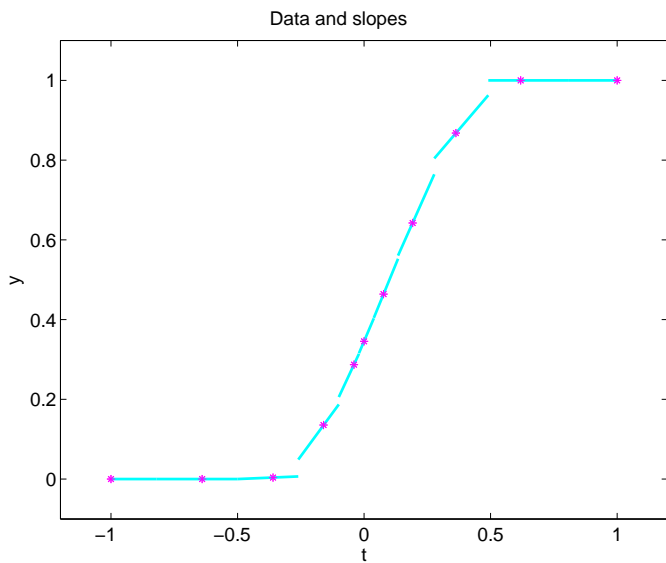




Shape preserving quadratic spline interpolation = **local** **but** **not linear**

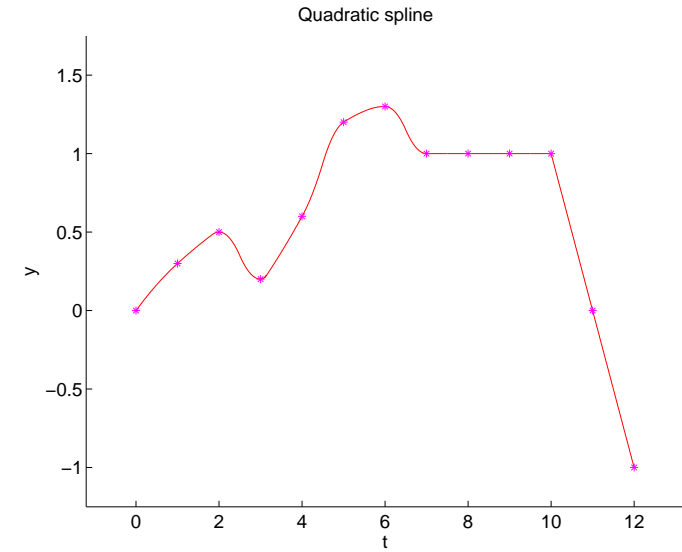
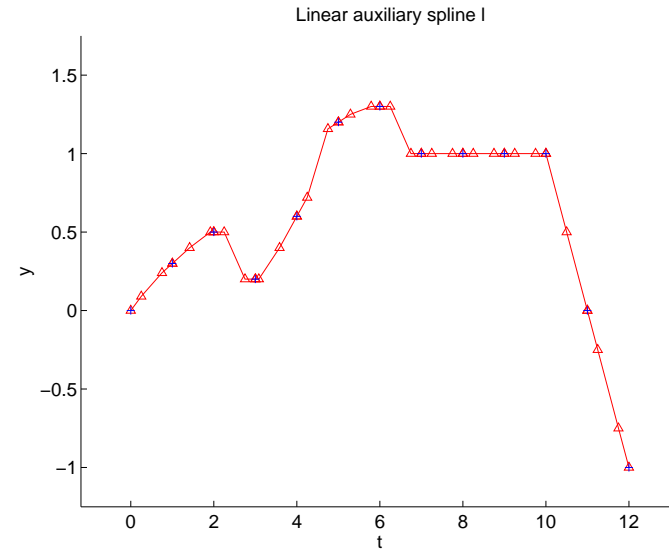
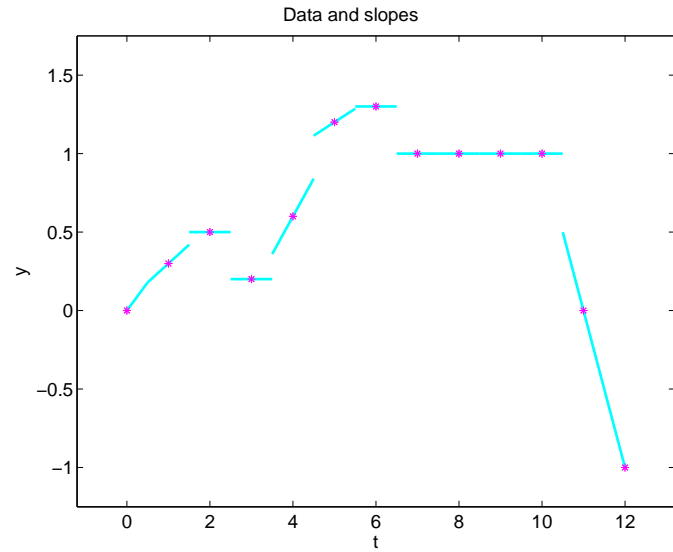
Example 3.8.17 (Shape preserving quadratic spline interpolation).

Data from Ex. 3.6.5:



Data from [43]:

| | | | | | | | | | | | | | |
|-------|---|-----|-----|-----|-----|-----|-----|---|---|---|----|----|----|
| 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 | 1 | 0 | -1 |



Iterative Methods for Non-Linear Systems of Equations

4

Example 4.0.1 (Non-linear electric circuit).

Schmitt trigger circuit

NPN bipolar junction transistor:

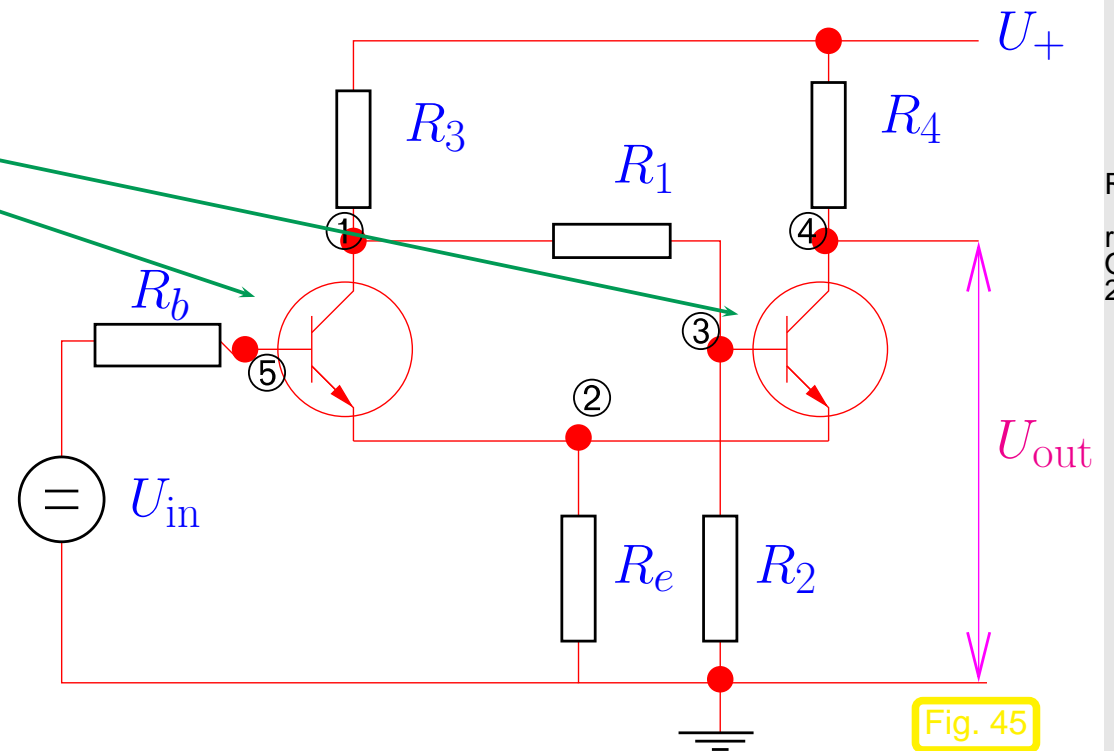
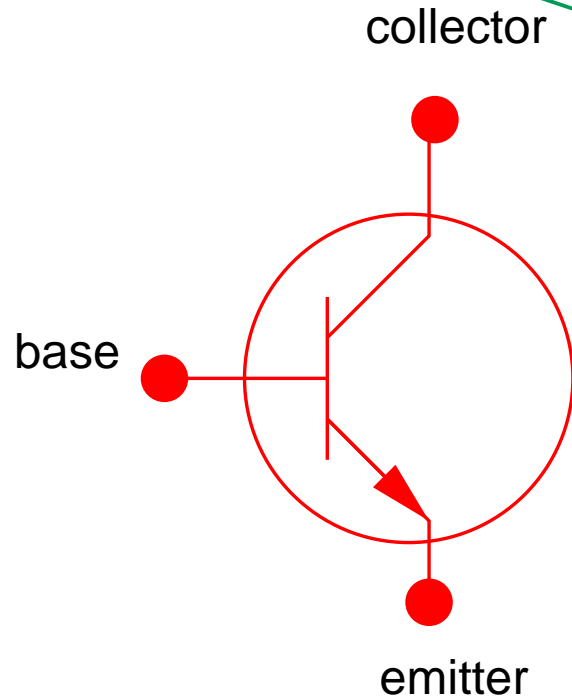


Fig. 45

Ebers-Moll model (large signal approximation):

$$\begin{aligned}
 I_C &= I_S \left(e^{\frac{U_{BE}}{U_T}} - e^{\frac{U_{BC}}{U_T}} \right) - \frac{I_S}{\beta_R} \left(e^{\frac{U_{BC}}{U_T}} - 1 \right) = I_C(U_{BE}, U_{BC}) , \\
 I_B &= \frac{I_S}{\beta_F} \left(e^{\frac{U_{BE}}{U_T}} - 1 \right) + \frac{I_S}{\beta_R} \left(e^{\frac{U_{BC}}{U_T}} - 1 \right) = I_B(U_{BE}, U_{BC}) , \\
 I_E &= I_S \left(e^{\frac{U_{BE}}{U_T}} - e^{\frac{U_{BC}}{U_T}} \right) + \frac{I_S}{\beta_F} \left(e^{\frac{U_{BE}}{U_T}} - 1 \right) = I_E(U_{BE}, U_{BC}) .
 \end{aligned}
 \tag{4.0.2}$$

I_C, I_B, I_E : current in collector/base/emitter,

U_{BE}, U_{BC} : potential drop between base-emitter, base-collector.

Non-linear system of equations from nodal analysis (\rightarrow Ex. 2.6.3):

$$\begin{aligned}
 \textcircled{1} : & R_3(U_1 - U_+) + R_1(U_1 - U_3) + I_B(U_5 - U_1, U_5 - U_2) = 0 , \\
 \textcircled{2} : & R_e U_2 + I_E(U_5 - U_1, U_5 - U_2) + I_E(U_3 - U_4, U_3 - U_2) = 0 , \\
 \textcircled{3} : & R_1(U_3 - U_1) + I_B(U_3 - U_4, U_3 - U_2) = 0 , \\
 \textcircled{4} : & R_4(U_4 - U_+) + I_C(U_3 - U_4, U_3 - U_2) = 0 , \\
 \textcircled{5} : & R_b(U_5 - U_{in}) + I_B(U_5 - U_1, U_5 - U_2) = 0 .
 \end{aligned}
 \tag{4.0.3}$$

5 equations \leftrightarrow 5 unknowns U_1, U_2, U_3, U_4, U_5

Formally:

$$(4.0.3) \longleftrightarrow F(\mathbf{u}) = 0$$



Given: function $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n, n \in \mathbb{N}$



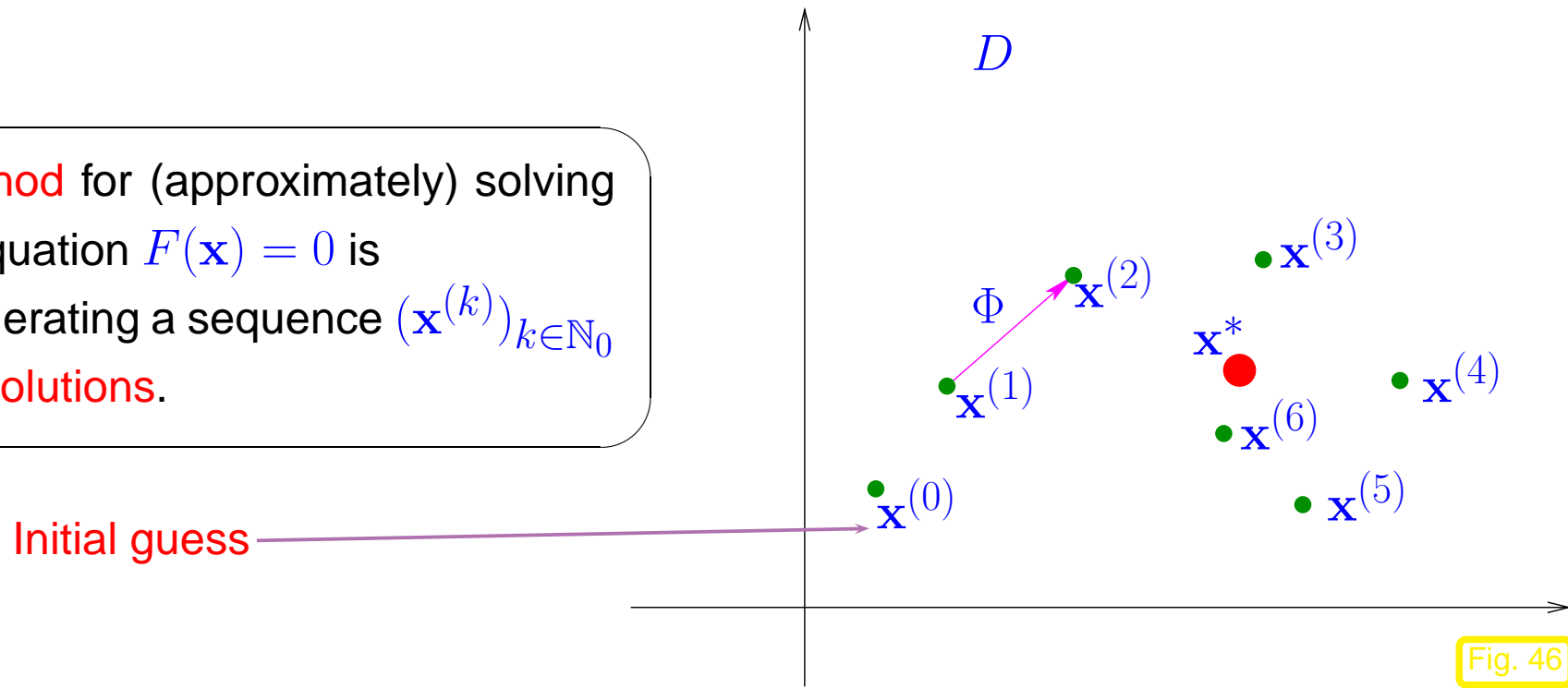
Possible meaning: Availability of a **procedure** function $y=F(x)$ evaluating F

Sought: solution of **non-linear equation** $F(\mathbf{x}) = 0$

In general no existence & uniqueness of solutions

4.1 Iterative methods

An **iterative method** for (approximately) solving the non-linear equation $F(\mathbf{x}) = 0$ is an algorithm generating a sequence $(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$ of **approximate solutions**.



Fundamental concepts: **convergence** \rightarrow **speed of convergence**
consistency

local convergence



(Only initial guesses “sufficiently close” to \mathbf{x}^* guarantee convergence.)

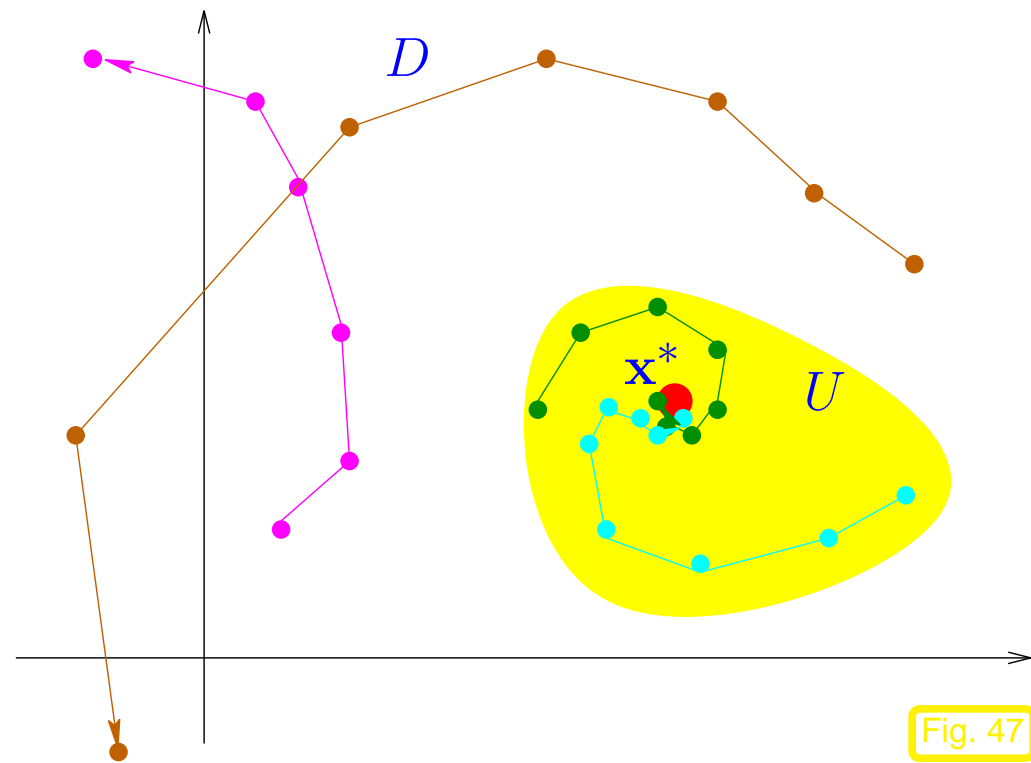


Fig. 47

Goal: Find iterative methods that converge (locally) to a solution of $F(\mathbf{x}) = 0$.

4.1.1 Speed of convergence

“Speed of convergence” \leftrightarrow decrease of norm (see Def. 2.5.1) of iteration error

Definition 4.1.6 (Linear convergence).

A sequence $\mathbf{x}^{(k)}$, $k = 0, 1, 2, \dots$, in \mathbb{R}^n *converges linearly* to $\mathbf{x}^* \in \mathbb{R}^n$, if

$$\exists L < 1: \quad \left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\| \leq L \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| \quad \forall k \in \mathbb{N}_0 .$$

Terminology: least upper bound for L gives the *rate of convergence*

Remark 4.1.7 (Impact of choice of norm).

| | | |
|--|--------------------|-------------------|
| <i>Fact of convergence</i> of iteration is | independent | of choice of norm |
| <i>Fact of linear convergence</i> | depends | on choice of norm |
| <i>Rate of linear convergence</i> | depends | on choice of norm |



Example 4.1.12 (Linearly convergent iteration).

Iteration ($n = 1$):

$$x^{(k+1)} = x^{(k)} + \frac{\cos x^{(k)} + 1}{\sin x^{(k)}}.$$

Code 4.1.13: simple fixed point iteration

```

1 y = [ ];
2 for i = 1:15
3     x = x + (cos(x)+1)/sin(x);
4     y = [y,x];
5 end
6 err = y - x;
7 rate = err(2:15)./err(1:14);

```

Note: $x^{(15)}$ replaces the exact solution x^* in the computation of the rate of convergence.

| 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 |

Linear convergence as in Def. 4.1.6



error graphs = straight lines in lin-log scale

→ Rem. 4.1.10

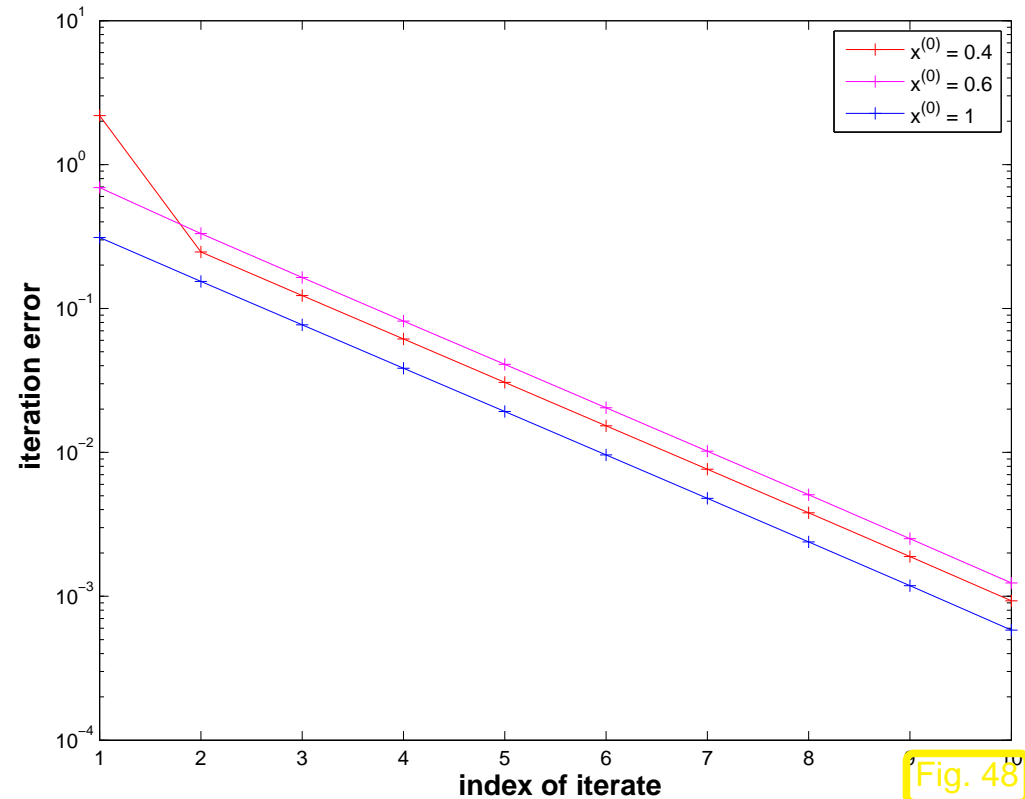


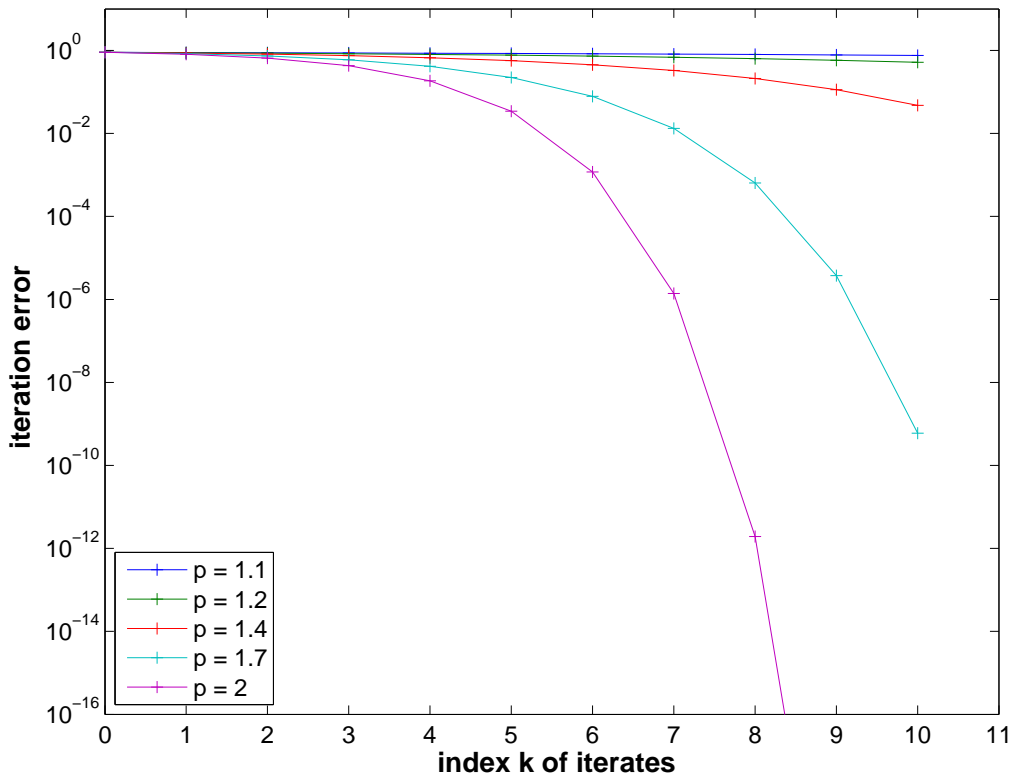
Fig. 48

Definition 4.1.14 (Order of convergence). → [35, Sect. 17.2], [13, Def. 5.14], [51, Def. 6.1]

A **convergent** sequence $\mathbf{x}^{(k)}$, $k = 0, 1, 2, \dots$, in \mathbb{R}^n converges with **order** p to $\mathbf{x}^* \in \mathbb{R}^n$, if

$$\exists C > 0: \quad \left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\| \leq C \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|^p \quad \forall k \in \mathbb{N}_0,$$

and, in addition, $C < 1$ in the case $p = 1$ (linear convergence → Def. 4.1.6).



◁ Qualitative error graphs for convergence of order p (lin-log scale)

Example 4.1.16 (quadratic convergence). (= convergence of order 2)

Iteration for computing \sqrt{a} , $a > 0$:

$$x^{(k+1)} = \frac{1}{2} \left(x^{(k)} + \frac{a}{x^{(k)}} \right) \Rightarrow |x^{(k+1)} - \sqrt{a}| = \frac{1}{2x^{(k)}} |x^{(k)} - \sqrt{a}|^2. \quad (4.1.17)$$

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 | 2.000000000000000000 | 0.58578643762690485 | |
| 1 | 1.500000000000000000 | 0.08578643762690485 | |
| 2 | 1.416666666666666652 | 0.00245310429357137 | 1.850 |
| 3 | 1.41421568627450966 | 0.00000212390141452 | 1.984 |
| 4 | 1.41421356237468987 | 0.00000000000159472 | 2.000 |
| 5 | 1.41421356237309492 | 0.00000000000000022 | 0.630 |

Note the **doubling** of the number of significant digits in each step !

[impact of roundoff !]



4.1.2 Termination criteria

④ **Residual based** termination: STOP convergent iteration $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}_0}$, when

$$\|F(\mathbf{x}^{(k)})\| \leq \tau, \quad \tau \hat{=} \text{prescribed tolerance} > 0.$$



no guaranteed accuracy

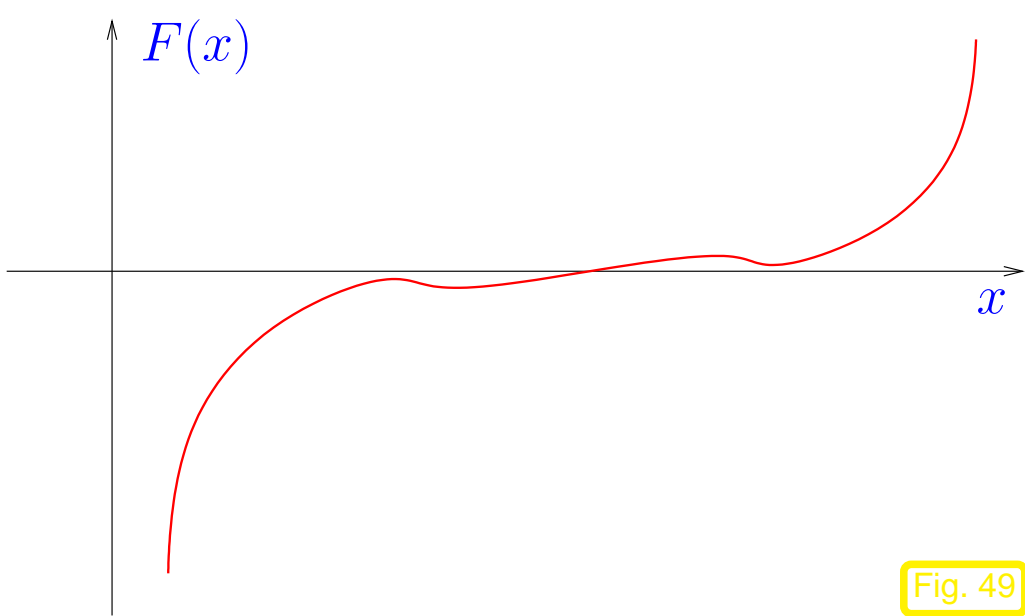


Fig. 49

$$\|F(\mathbf{x}^{(k)})\| \text{ small } \not\Rightarrow |x - x^*| \text{ small}$$

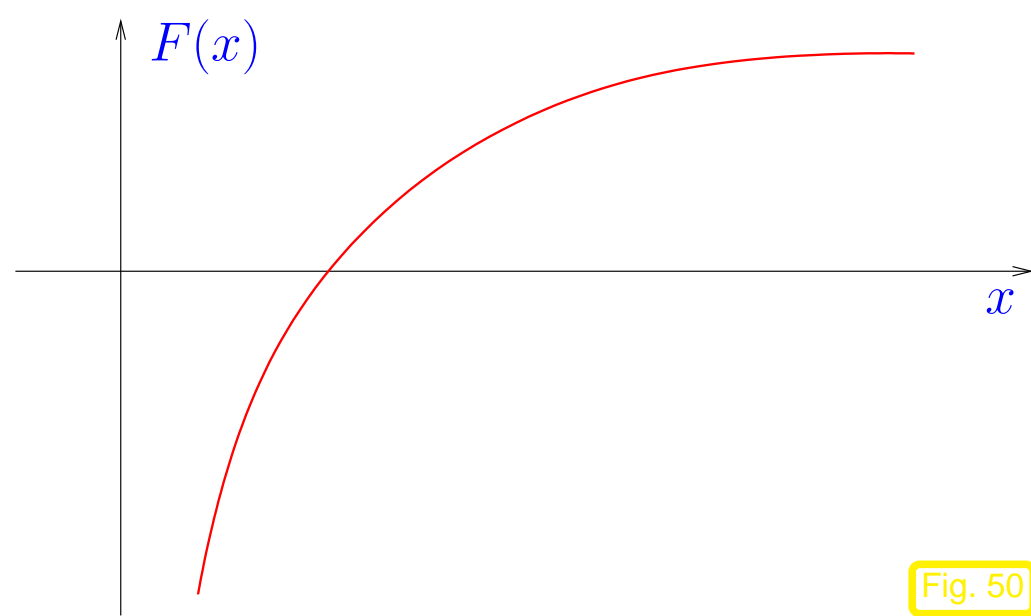


Fig. 50

$$\|F(\mathbf{x}^{(k)})\| \text{ small } \Rightarrow |x - x^*| \text{ small}$$

⑤

Correction based termination: STOP convergent iteration $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}_0}$, when

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq \begin{cases} \tau_{\text{abs}} \\ \text{or} \\ \tau_{\text{rel}} \|\mathbf{x}^*\| \end{cases}, \quad \begin{matrix} \tau_{\text{abs}} \\ \tau_{\text{rel}} \end{matrix} \text{ prescribed absolute relative } \text{tolerances} > 0.$$

Remark 4.1.23 (A posteriori termination criterion for linearly convergent iterations). \rightarrow [13, Lemma 5.

Known: iteration linearly convergent with rate of convergence $0 < L < 1$:

Iterates satisfy:

$$\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\| \leq \frac{L}{1-L} \left\| \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \right\| . \quad (4.1.24)$$



Example 4.1.25 (A posteriori error bound for linearly convergent iteration).

Iteration of Example 4.1.12:

$$x^{(k+1)} = x^{(k)} + \frac{\cos x^{(k)} + 1}{\sin x^{(k)}} \Rightarrow x^{(k)} \rightarrow \pi \quad \text{for } x^{(0)} \text{ close to } \pi .$$

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 | 0.000000000009257 |
| 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!



4.2 Fixed Point Iterations [13, Sect. 5.3], [51, Sect. 6.3]

A **fixed point iteration** is defined by **iteration function** $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$:

$$\begin{array}{l}
 \text{iteration function} \quad \Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n \\
 \text{initial guess} \quad \mathbf{x}^{(0)} \in U
 \end{array}
 \begin{array}{l}
 \triangleright \\
 \underbrace{\text{iterates } (\mathbf{x}^{(k)})_{k \in \mathbb{N}_0} : \mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)})}_{\rightarrow \text{1-point method, cf. (4.1.2)}}
 \end{array}
 .$$

4.2.1 Consistent fixed point iterations

Definition 4.2.1 (Consistency of fixed point iterations, *c.f.* Def. 4.1.4).

A fixed point iteration $\mathbf{x}^{(k+1)} = \Phi(\mathbf{x}^{(k)})$ is **consistent** with $F(\mathbf{x}) = 0$, if

$$F(\mathbf{x}) = 0 \quad \text{and} \quad \mathbf{x} \in U \cap D \quad \Leftrightarrow \quad \Phi(\mathbf{x}) = \mathbf{x} .$$

General construction of fixed point iterations that is consistent with $F(\mathbf{x}) = 0$:

rewrite $F(\mathbf{x}) = 0 \Leftrightarrow \Phi(\mathbf{x}) = \mathbf{x}$ and then

use the **fixed point iteration** $\mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)})$. (4.2.2)

Note: there are *many* ways to transform $F(\mathbf{x}) = 0$ into a fixed point form !

Example 4.2.3 (Options for fixed point iterations).

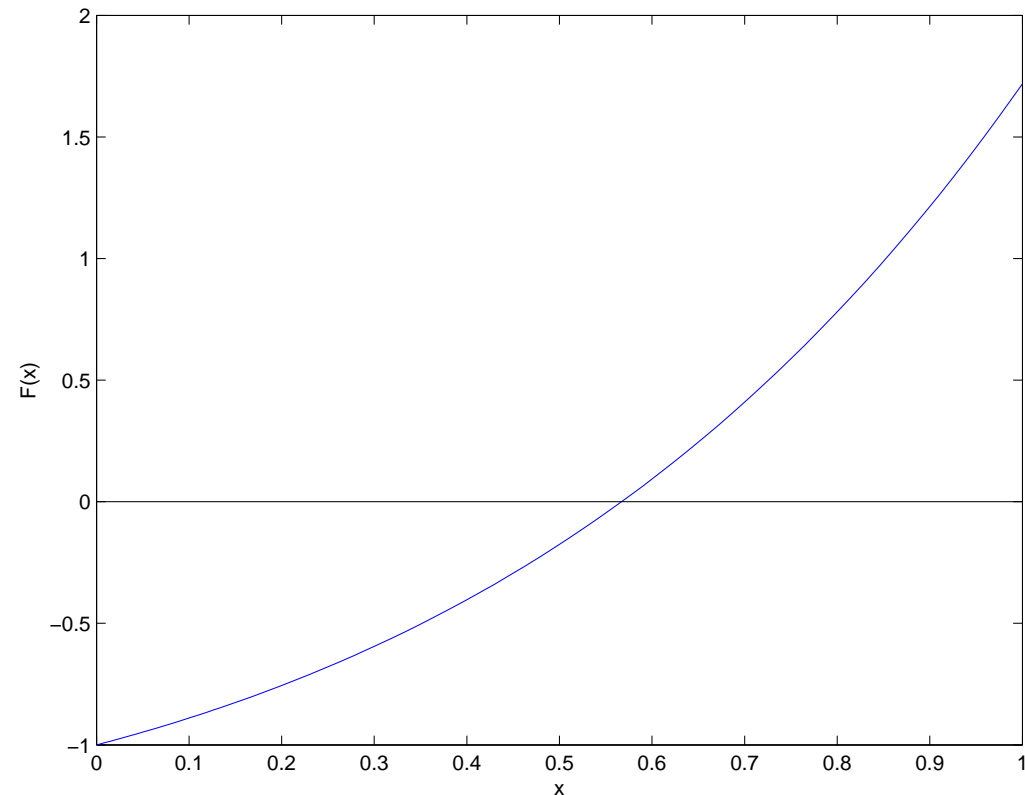
$$F(x) = xe^x - 1, \quad x \in [0, 1].$$

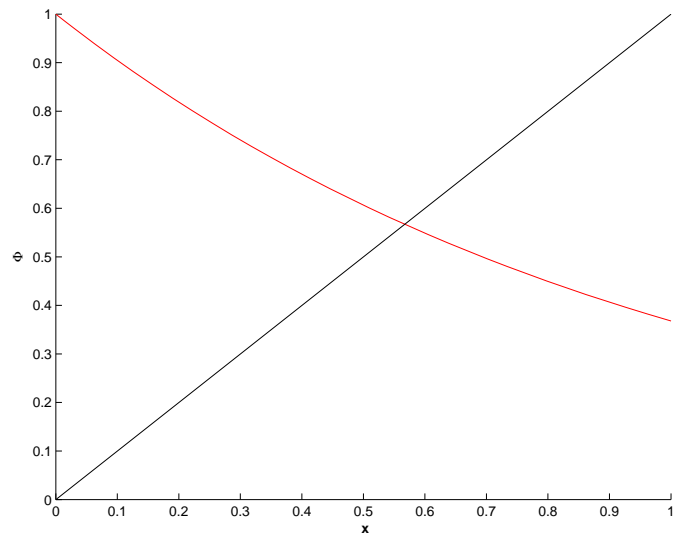
Different fixed point forms:

$$\Phi_1(x) = e^{-x},$$

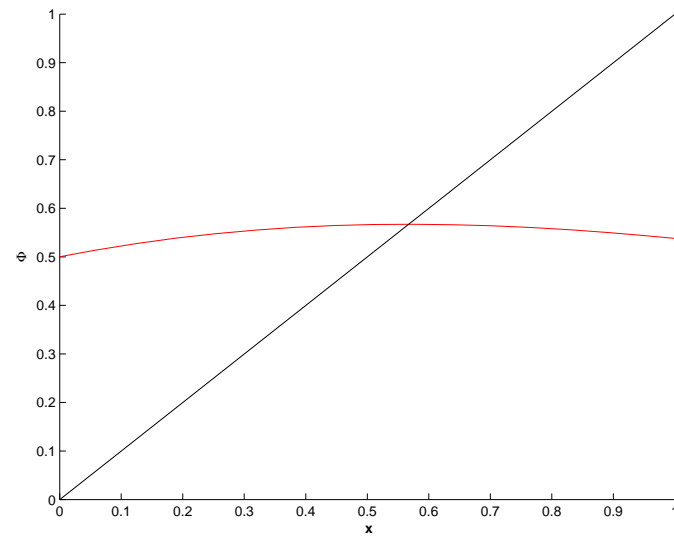
$$\Phi_2(x) = \frac{1+x}{1+e^x},$$

$$\Phi_3(x) = x + 1 - xe^x.$$

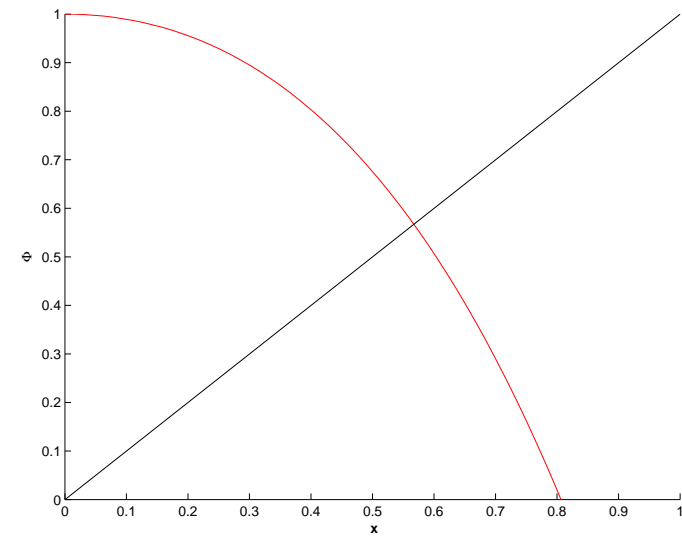




function Φ_1



function Φ_2



function Φ_3

| 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 |

| 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 | 0.0000000000000003 | 0.288178118764323 |
| 4 | 0.007078662470882 | 0.0000000000000000 | 0.723649245792953 |
| 5 | 0.004028858567431 | 0.0000000000000000 | 0.410183132337935 |
| 6 | 0.002280343429460 | 0.0000000000000000 | 1.186907542305364 |
| 7 | 0.001294757160282 | 0.0000000000000000 | 0.146569797006362 |
| 8 | 0.000733837662863 | 0.0000000000000000 | 0.310516641279937 |
| 9 | 0.000416343852458 | 0.0000000000000000 | 0.357777386500765 |
| 10 | 0.000236077474313 | 0.0000000000000000 | 0.974565695952037 |

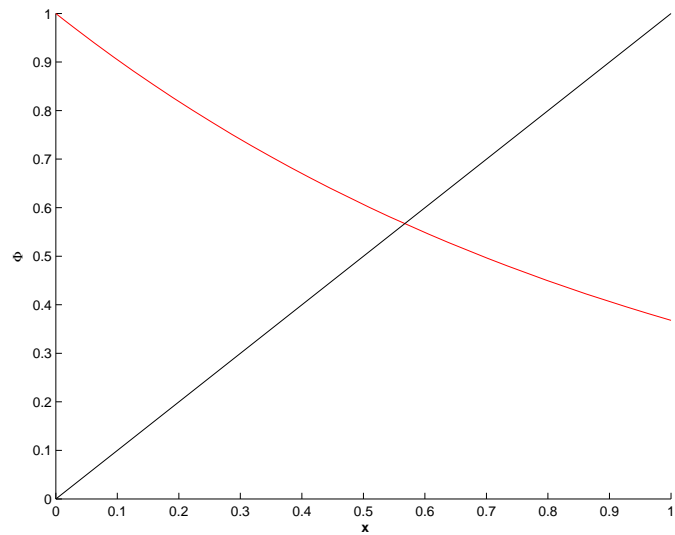
Observed: linear convergence of $x_1^{(k)}$, quadratic convergence of $x_2^{(k)}$,
no convergence (erratic behavior of $x_3^{(k)}$), $x_i^{(0)} = 0.5$.



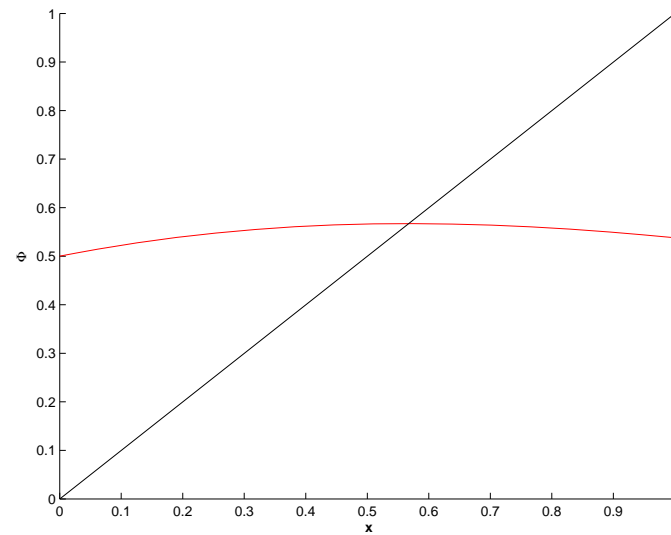
Question: can we explain/forecast the behaviour of the iteration?

4.2.2 Convergence of fixed point iterations

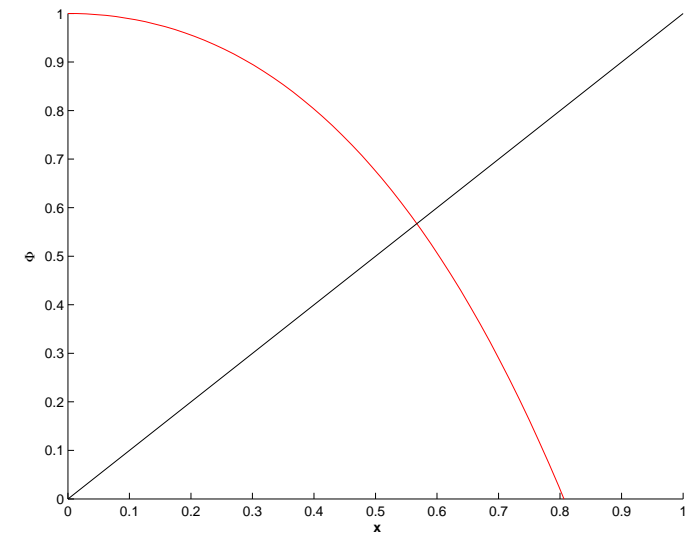
Ex. 4.2.3 revisited: vastly different behavior of different fixed point iterations for $n = 1$:



Φ_1 : linear convergence ?



Φ_2 : quadratic convergence ?

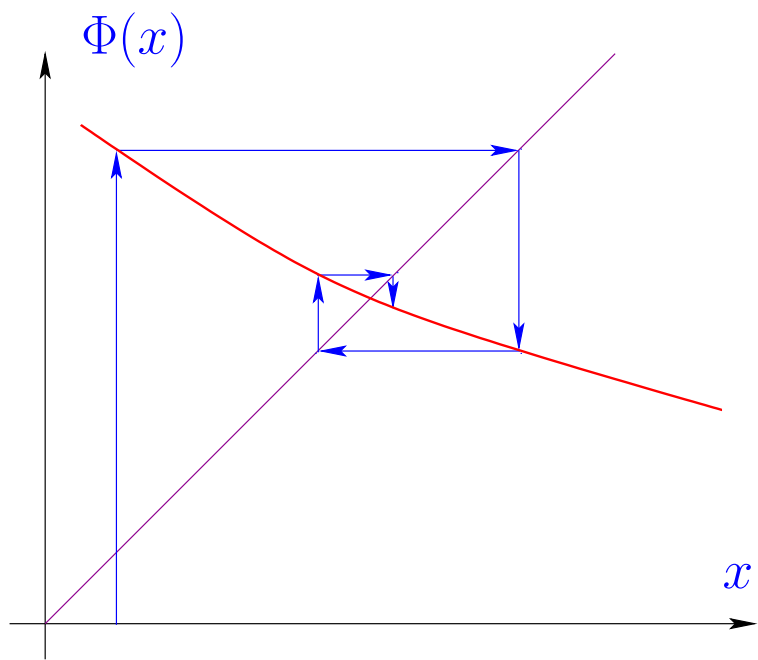


function Φ_3 : no convergence

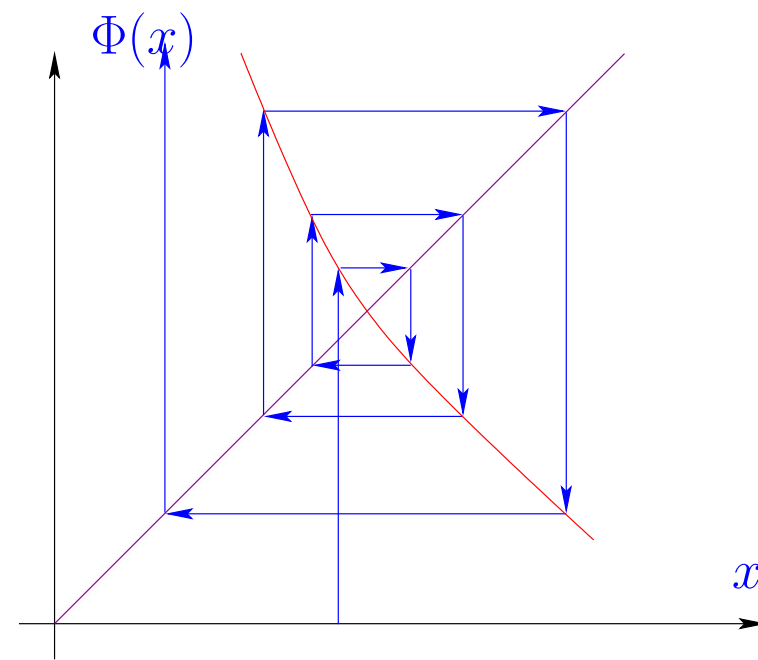
Example 4.2.4 (Fixed point iteration in 1D).

1D setting ($n = 1$): $\Phi : \mathbb{R} \mapsto \mathbb{R}$ continuously differentiable, $\Phi(x^*) = x^*$

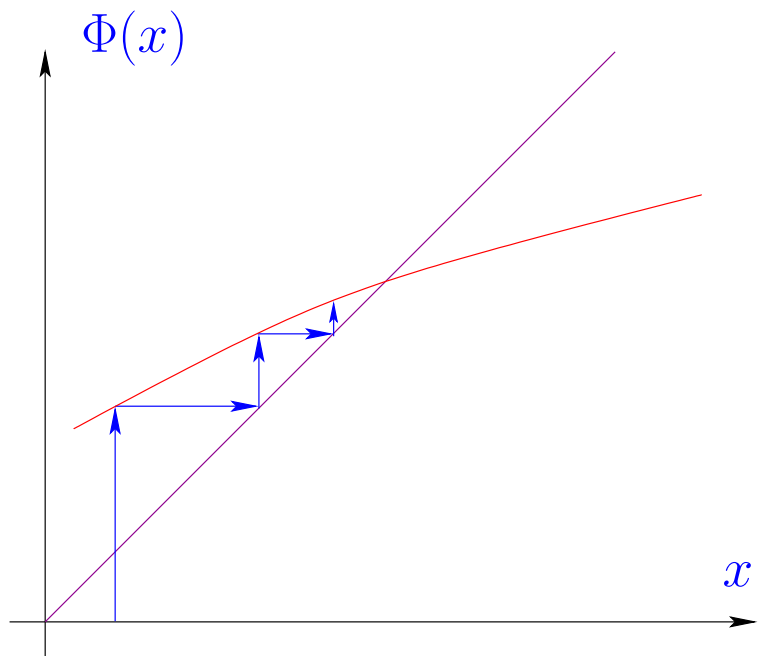
fixed point iteration: $x^{(k+1)} = \Phi(x^{(k)})$



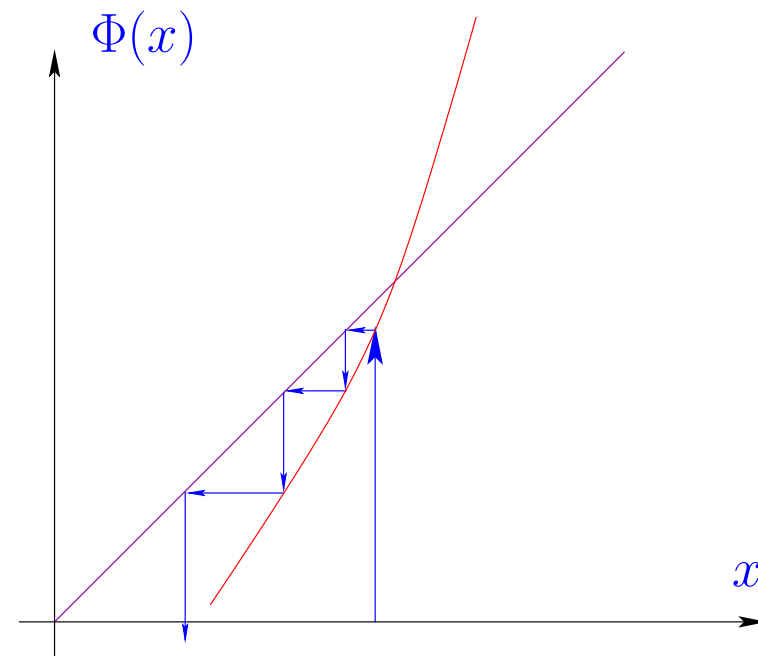
$-1 < \Phi'(x^*) \leq 0 \Rightarrow$ convergence



$\Phi'(x^*) < -1 \Rightarrow$ divergence



$0 \leq \Phi'(x^*) < 1 \Rightarrow$ convergence



$1 < \Phi'(x^*) \Rightarrow$ divergence

Lemma 4.2.8 (Sufficient condition for local linear convergence of fixed point iteration). → [35, Thm. 17.2], [13, Cor. 5.12]

If $\Phi : U \subset \mathbb{R}^n \mapsto \mathbb{R}^n$, $\Phi(\mathbf{x}^*) = \mathbf{x}^*$, Φ differentiable in \mathbf{x}^* , and $\|D\Phi(\mathbf{x}^*)\| < 1$, then the fixed point iteration (4.2.2) converges locally and at least linearly.

matrix norm, Def. 2.5.5 !

notation: $D\Phi(\mathbf{x}) \hat{=}$ **Jacobian** (ger.: Jacobi-Matrix) of Φ at $\mathbf{x} \in D$
→ [63, Sect. 7.6]

Remark 4.2.10 (Bound for asymptotic rate of linear convergence).

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 4.2.11 (Multidimensional fixed point iteration).

$$\begin{array}{l} \text{System of equations} \\ \left\{ \begin{array}{l} x_1 - c(\cos x_1 - \sin x_2) = 0 \\ (x_1 - x_2) - c \sin x_2 = 0 \end{array} \right. \end{array} \quad \text{in} \quad \Rightarrow \quad \begin{array}{l} \text{fixed point form:} \\ \left\{ \begin{array}{l} c(\cos x_1 - \sin x_2) = x_1 \\ c(\cos x_1 - 2 \sin x_2) = x_2 \end{array} \right. . \end{array}$$

Define: $\Phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = c \begin{pmatrix} \cos x_1 - \sin x_2 \\ \cos x_1 - 2 \sin x_2 \end{pmatrix} \Rightarrow D\Phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = -c \begin{pmatrix} \sin x_1 & \cos x_2 \\ \sin x_1 & 2 \cos x_2 \end{pmatrix} .$

Choose *appropriate* norm: $\|\cdot\| = \infty\text{-norm } \|\cdot\|_\infty$ (\rightarrow Example 2.5.6) ;

$$\text{if } c < \frac{1}{3} \Rightarrow \|D\Phi(\mathbf{x})\|_\infty < 1 \quad \forall \mathbf{x} \in \mathbb{R}^2 ,$$

➤ (at least) linear convergence of the fixed point iteration. ◇

What about higher order convergence (\rightarrow Def. 4.1.14, cf. Φ_2 in Ex. 4.2.3) ?

Refined convergence analysis for $n = 1$ (scalar case, $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$):

Theorem 4.2.12 (Taylor's formula). \rightarrow [63, Sect. 5.5]

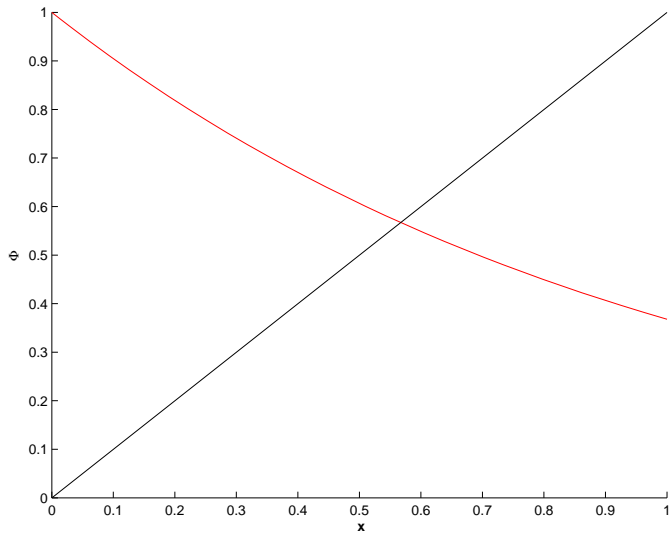
If $\Phi : U \subset \mathbb{R} \mapsto \mathbb{R}$, U interval, is $m + 1$ times continuously differentiable, $x \in U$

$$\Phi(y) - \Phi(x) = \sum_{k=1}^m \frac{1}{k!} \Phi^{(k)}(x)(y-x)^k + O(|y-x|^{m+1}) \quad \forall y \in U . \quad (4.2.13)$$

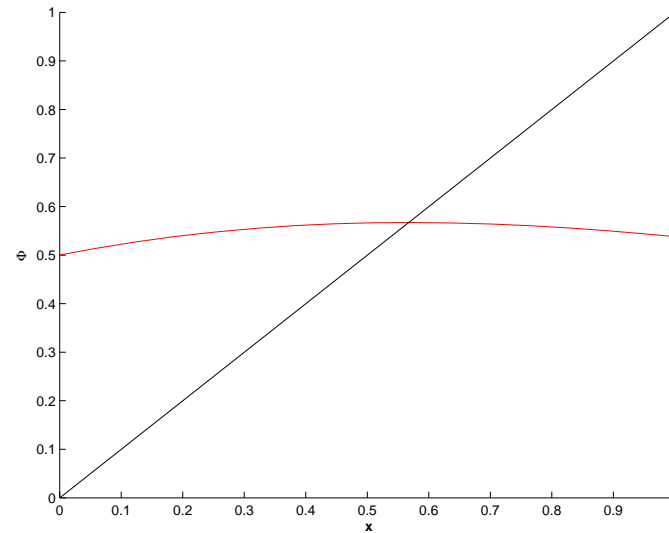
Lemma 4.2.15 (Higher order local convergence of fixed point iterations).

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 (4.2.2) converges locally to x^* with **order** $\geq m + 1$ (\rightarrow Def. 4.1.14).

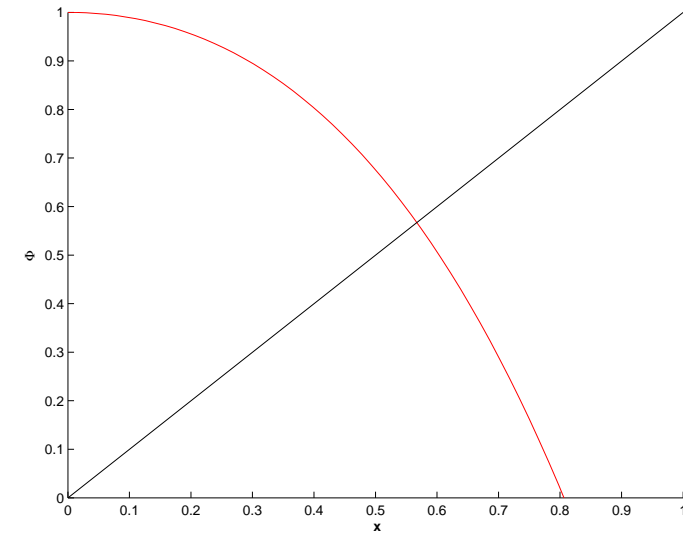
Example 4.2.3 continued:



function Φ_1



function Φ_2



function Φ_3

$$\Phi_2'(x) = \frac{1 - xe^x}{(1 + e^x)^2} = 0 \quad , \text{ if } \quad xe^x - 1 = 0 \quad \text{hence quadratic convergence ! .}$$

Example 4.2.3 continued: Since $x^*e^{x^*} - 1 = 0$

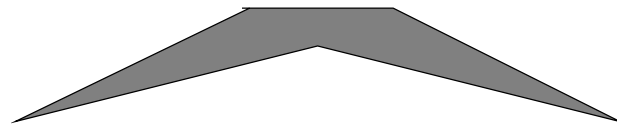
$$\Phi'_1(x) = -e^{-x} \Rightarrow \Phi'_1(x^*) = -x^* \approx -0.56 \quad \text{hence local linear convergence .}$$

$$\Phi'_3(x) = 1 - xe^x - e^x \Rightarrow \Phi'_3(x^*) = -\frac{1}{x^*} \approx -1.79 \quad \text{hence no convergence .}$$

Remark 4.2.16 (Termination criterion for contractive fixed point iteration).

Recap of Rem. 4.1.23:

$$\left\| \mathbf{x}^* - \mathbf{x}^{(k)} \right\| \leq \frac{L^{k-l}}{1-L} \left\| \mathbf{x}^{(l+1)} - \mathbf{x}^{(l)} \right\| . \quad (4.2.17)$$



Set $l = 0$ in (4.2.18)

a priori termination criterion

$$\left\| \mathbf{x}^* - \mathbf{x}^{(k)} \right\| \leq \frac{L^k}{1-L} \left\| \mathbf{x}^{(1)} - \mathbf{x}^{(0)} \right\| \quad (4.2.18)$$

Set $l = k - 1$ in (4.2.18)

a posteriori termination criterion

$$\left\| \mathbf{x}^* - \mathbf{x}^{(k)} \right\| \leq \frac{L}{1-L} \left\| \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right\| \quad (4.2.19)$$



4.3 Zero Finding

Now, focus on scalar case $n = 1$:

$F : I \subset \mathbb{R} \mapsto \mathbb{R}$ **continuous**, I interval

Sought:

$$x^* \in I: \quad F(x^*) = 0$$

4.3.1 Bisection [13, Sect. 5.5.1]

Idea: use ordering of real numbers & intermediate value theorem

Input: $a, b \in I$ such that $F(a)F(b) < 0$
(different signs !)

$\Rightarrow \exists x^* \in]\min\{a, b\}, \max\{a, b\}[$
 $F(x^*) = 0$.

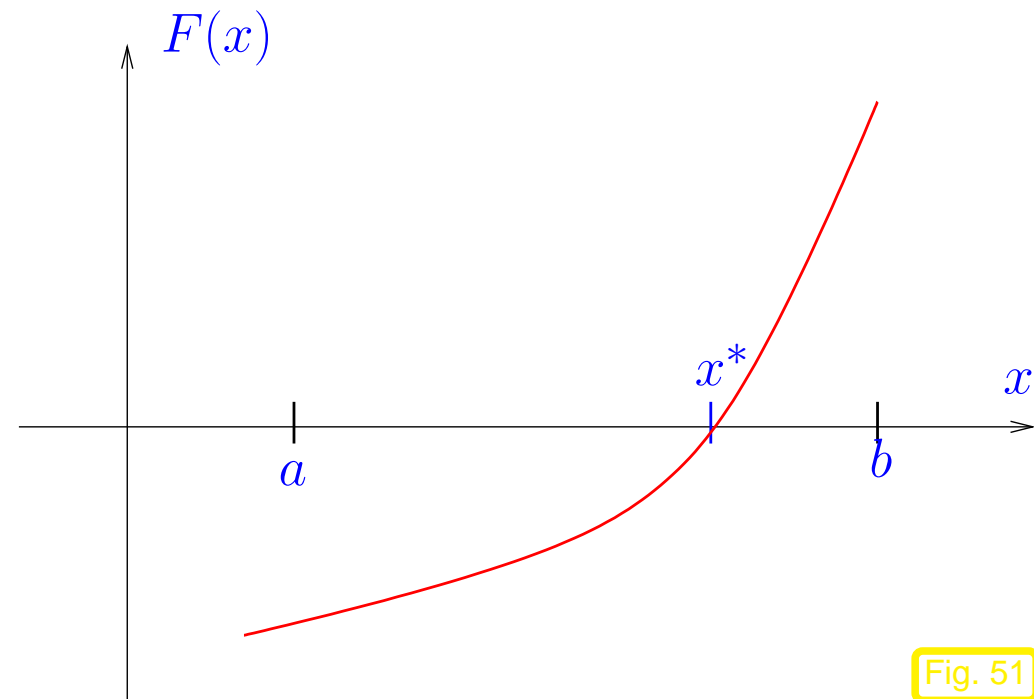


Fig. 51

Algorithm 4.3.1 (Bisection method).

MATLAB-CODE: bisection method

```
function x = bisect(F,a,b,tol)
% Searching zero by bisection
if (a>b), t=a; a=b; b=t; end;
fa = F(a); fb = F(b);
if (fa*fb>0)
    error('f(a), f(b) same sign'); end;
if (fa > 0), v=-1; else v = 1; end
x = 0.5*(b+a);
while((b-a > tol) & ((a<x) & (x<b)))
    if (v*F(x)>0), b=x; else a=x; end;
    x = 0.5*(a+b)
end
```

$tol \hat{=}$ absolute tolerance

Handle to MATLAB function providing F .

Avoid infinite loop, if $tol <$ resolution of M at zero x^* (“ M -based termination criterion”).



Advantages:

- “foolproof”
- requires only F evaluations



Drawbacks:

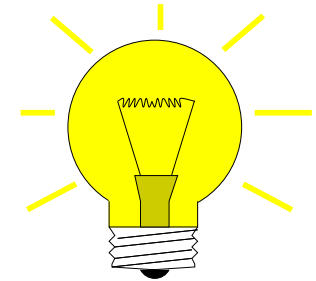
Merely “linear-type” convergence: $|x^{(k)} - x^*| \leq 2^{-k}|b - a|$

▶ $\log_2 \left(\frac{|b - a|}{tol} \right)$ steps necessary

4.3.2 Model function methods

$\hat{=}$ class of iterative methods for finding zeroes of F :

Idea: Given: approximate zeroes $x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}$



- ❶ replace F with **model function** \tilde{F}
(using function values/derivative values in $x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}$)
- ❷ $x^{(k+1)} :=$ zero of \tilde{F}
(has to be readily available \leftrightarrow analytic formula)

Distinguish (see (4.1.2)):

one-point methods : $x^{(k+1)} = \Phi_F(x^{(k)}), k \in \mathbb{N}$ (e.g., fixed point iteration \rightarrow Sect. 4.2)

multi-point methods : $x^{(k+1)} = \Phi_F(x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}), k \in \mathbb{N}, m = 2, 3, \dots$

4.3.2.1 Newton method in scalar case [35, Sect. 18.1], [13, Sect. 5.5.2]

Assume: $F : I \mapsto \mathbb{R}$ continuously differentiable

model function:= tangent at F in $x^{(k)}$:

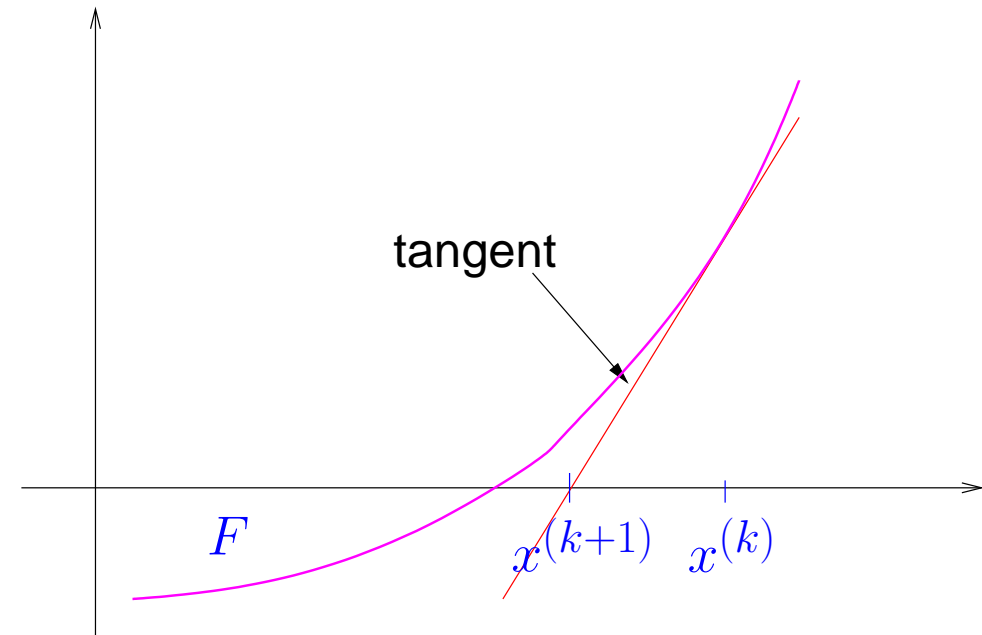
$$\tilde{F}(x) := F(x^{(k)}) + F'(x^{(k)})(x - x^{(k)})$$

take $x^{(k+1)} :=$ zero of tangent

We obtain **Newton iteration**

$$x^{(k+1)} := x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})}, \quad (4.3.3)$$

that requires $F'(x^{(k)}) \neq 0$.



Newton method locally quadratically convergent (\rightarrow Def. 4.1.14) to zero x^* , if $F'(x^*) \neq 0$

4.3.2.2 Special one-point methods

Example 4.3.5 (Halley's iteration). \rightarrow [35, Sect. 18.3]

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 (4.3.3) 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 (\rightarrow Sect. 4.3.3) \diamond

Example 4.3.6 (Adapted Newton method).

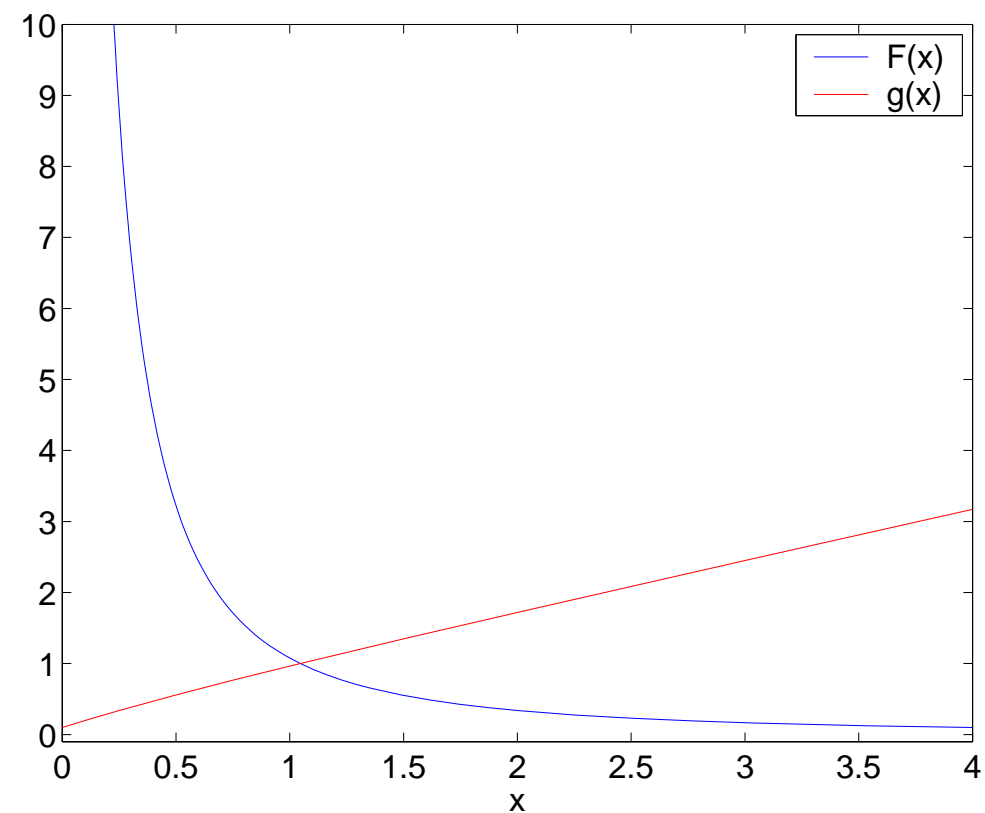
As in Ex. 4.3.5:

$$F(x) = \frac{1}{(x+1)^2} + \frac{1}{(x+0.1)^2} - 1, \quad x > 0 :$$

Observation:

$$F(x) + 1 \approx 2x^{-2} \text{ for } x \gg 1$$

and so $g(x) := \frac{1}{\sqrt{F(x) + 1}}$ “almost” linear for $x \gg 1$



Idea: instead of $F(x) \stackrel{!}{=} 0$ tackle $g(x) \stackrel{!}{=} 1$ with Newton's method (4.3.3).

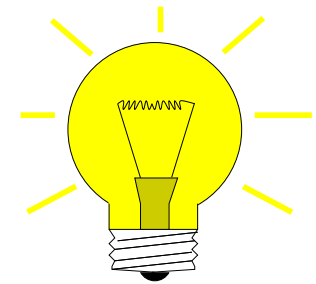
$$\begin{aligned} x^{(k+1)} &= x^{(k)} - \frac{g(x^{(k)}) - 1}{g'(x^{(k)})} = x^{(k)} + \left(\frac{1}{\sqrt{F(x^{(k)}) + 1}} - 1 \right) \frac{2(F(x^{(k)}) + 1)^{3/2}}{F'(x^{(k)})} \\ &= x^{(k)} + \frac{2(F(x^{(k)}) + 1)(1 - \sqrt{F(x^{(k)}) + 1})}{F'(x^{(k)})}. \end{aligned}$$

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 | 0.00000000000000 | 0.00000005485332 | -0.00000000000000 |



4.3.2.3 Multi-point methods



Idea:

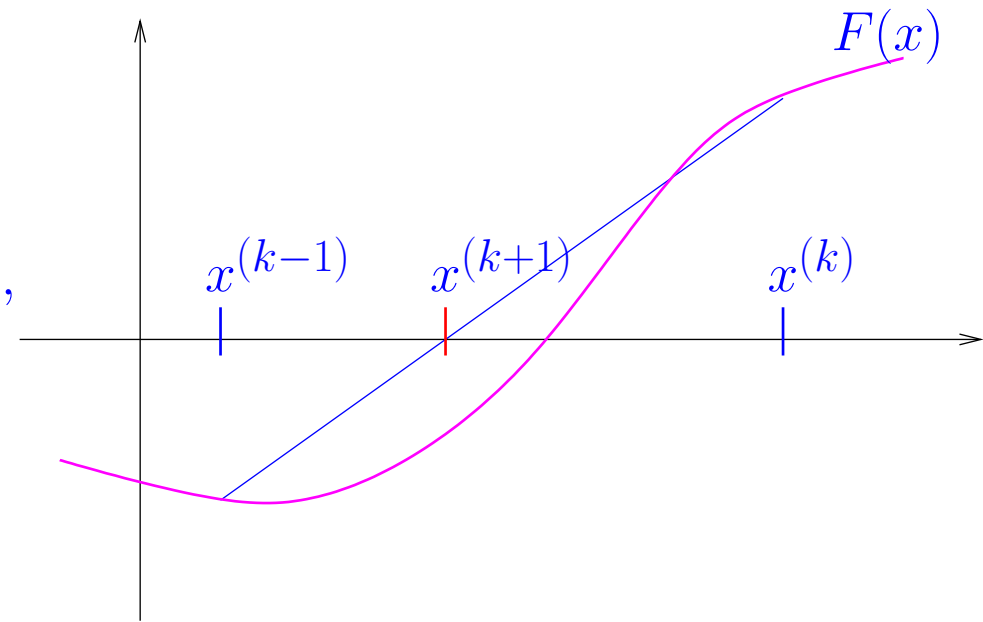
Replace F with **interpolating polynomial**
producing interpolatory model function methods

Simplest representative: **secant method** → [35,
Sect. 18.2], [13, Sect, 5.5.3]

$x^{(k+1)}$ = zero of secant

$$s(x) = F(x^{(k)}) + \frac{F(x^{(k)}) - F(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}(x - x^{(k)}), \quad (4.3.11)$$

►
$$x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})(x^{(k)} - x^{(k-1)})}{F(x^{(k)}) - F(x^{(k-1)})}. \quad (4.3.12)$$



Code 4.3.13: secant method

```

1 function x = secant(x0,x1,F,tol)
2 fo = F(x0);
3 for i=1:MAXIT
4     fn = F(x1);
5     s = fn*(x1-x0)/(fn-fo);
6     x0 = x1; x1 = x1-s;
7     if (abs(s) < tol), x = x1;
8         return; end
9     fo = fn;
9 end

```

secant method

(MATLAB implementation)

- Only one function evaluation per step
- **no derivatives required!**

Example 4.3.14 (secant method).

$$F(x) = xe^x - 1, \quad x^{(0)} = 0, \quad x^{(1)} = 5.$$

| 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 | -0.000000000000001 | -0.000000000000000 | |



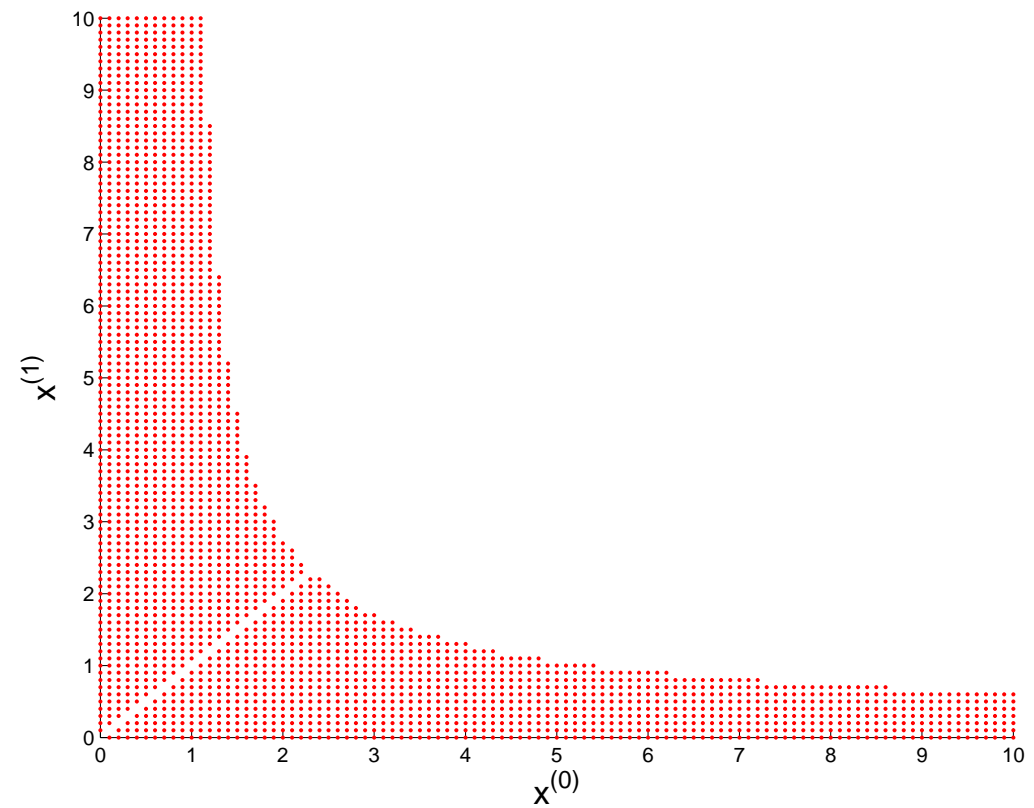
Example 4.3.18 (local convergence of secant method).

$$F(x) = \arctan(x)$$

• $\hat{=}$ secant method converges for a pair $(x^{(0)}, x^{(1)})$ of initial guesses.



= local convergence \rightarrow Def. 4.1.5



Another class of multi-point methods: *inverse interpolation*

Assume:

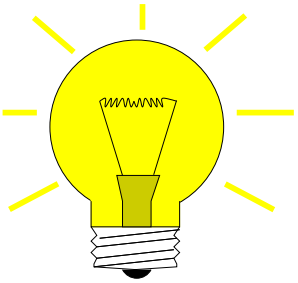
$F : I \subset \mathbb{R} \mapsto \mathbb{R}$ one-to-one

$$F(x^*) = 0 \Rightarrow F^{-1}(0) = x^* .$$

- Interpolate F^{-1} by polynomial p of degree d determined by

$$p(F(x^{(k-m)})) = x^{(k-m)} , \quad m = 0, \dots, d .$$

- New approximate zero $x^{(k+1)} := p(0)$



$$F(x^*) = 0 \Leftrightarrow F^{-1}(0) = x^*$$

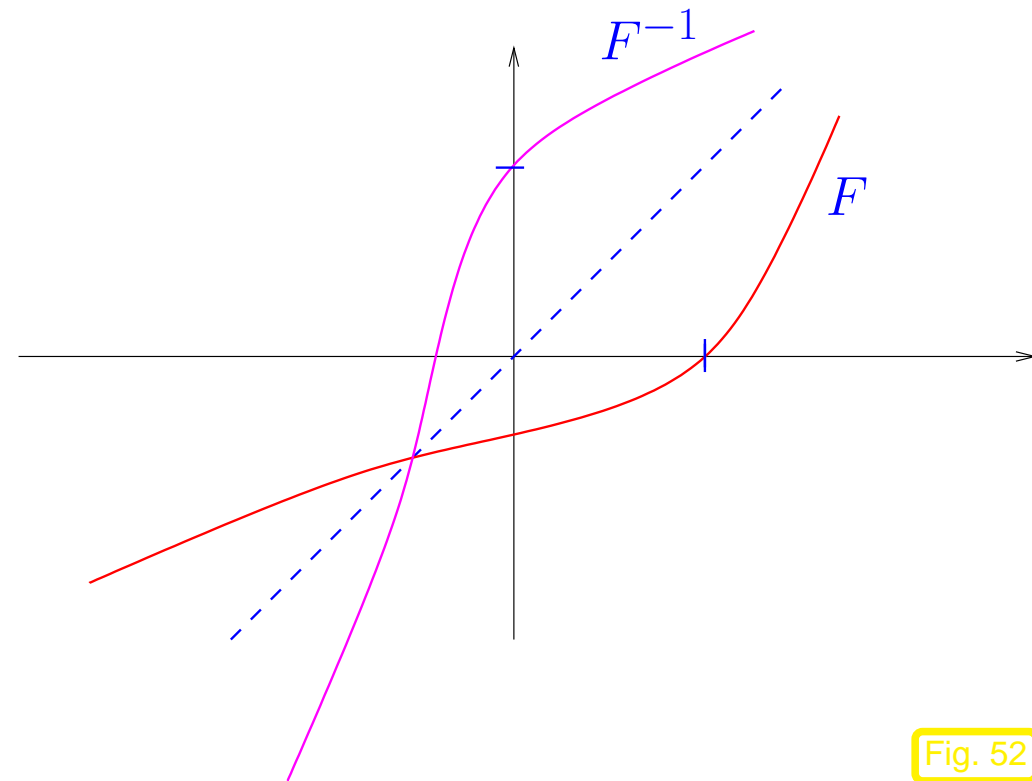


Fig. 52

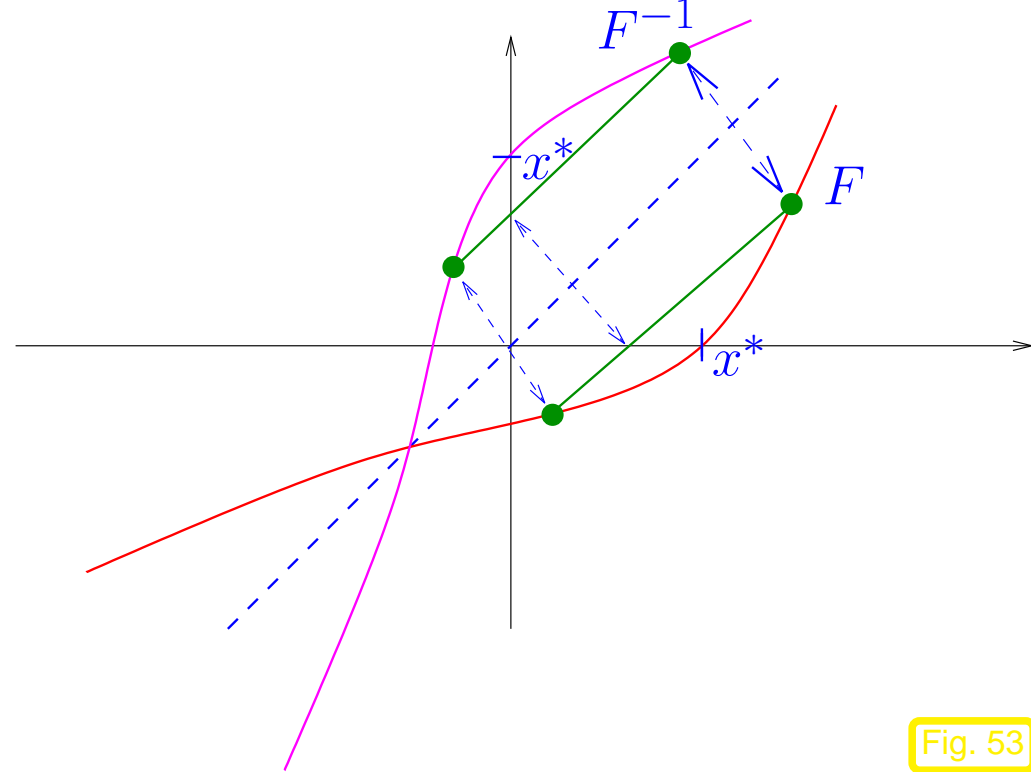


Fig. 53

Case $d = 1$ ➤ secant method

Case $d = 2$: quadratic inverse interpolation, see [45, Sect. 4.5]

MAPLE code: `p := x-> a*x^2+b*x+c;`
`solve({p(f0)=x0,p(f1)=x1,p(f2)=x2},{a,b,c});`
`assign(%); p(0);`

$$\blacktriangleright x^{(k+1)} = \frac{F_0^2(F_1x_2 - F_2x_1) + F_1^2(F_2x_0 - F_0x_2) + F_2^2(F_0x_1 - F_1x_0)}{F_0^2(F_1 - F_2) + F_1^2(F_2 - F_0) + F_2^2(F_0 - F_1)}.$$

$$(F_0 := F(x^{(k-2)}), F_1 := F(x^{(k-1)}), F_2 := F(x^{(k)}), x_0 := x^{(k-2)}, x_1 := x^{(k-1)}, x_2 := x^{(k)})$$

Example 4.3.19 (quadratic inverse interpolation). $F(x) = xe^x - 1$, $x^{(0)} = 0$, $x^{(1)} = 2.5$, $x^{(2)} = 5$.

| 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 | 0.000000000000003 | 0.000000000000001 | |



4.3.3 Note on Efficiency

Efficiency of an iterative method
(for solving $F(\mathbf{x}) = 0$)



computational effort to reach prescribed
number of significant digits in result.

Abstract:

$W \hat{=}$ computational effort per step

(e.g, $W \approx \frac{\#\{\text{evaluations of } F\}}{\text{step}} + n \cdot \frac{\#\{\text{evaluations of } F'\}}{\text{step}} + \dots$)

Crucial: number of steps $k = k(\rho)$ to achieve *relative reduction of error*

$$\|e^{(k)}\| \leq \rho \|e^{(0)}\|, \quad \rho > 0 \text{ prescribed ?} \tag{4.3.20}$$

Notice: $|\log \rho| \leftrightarrow$ No. of significant digits of $x^{(k)}$



Measure for efficiency:

$$\text{Efficiency} := \frac{\text{no. of digits gained}}{\text{total work required}} = \frac{|\log \rho|}{k(\rho) \cdot W}$$

(4.3.23)

► **asymptotic** efficiency w.r.t. $\rho \rightarrow 0 \rightarrow |\log \rho| \rightarrow \infty$):

$$\text{Efficiency}_{|\rho \rightarrow 0} = \begin{cases} \frac{\log C}{W} & , \text{ if } p = 1 , \\ \frac{\log p \cdot |\log \rho|}{W \log(|\log \rho|)} & , \text{ if } p > 1 . \end{cases} \tag{4.3.24}$$

Example 4.3.25 (Efficiency of iterative methods).

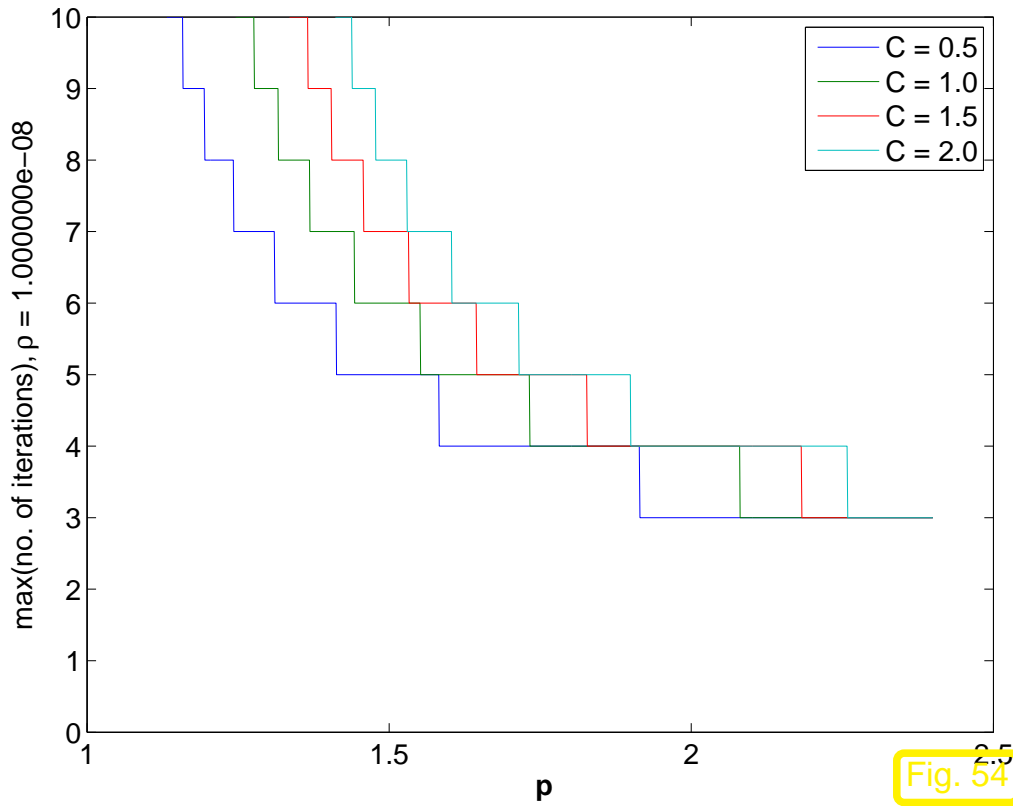


Fig. 54

Evaluation (4.3.22) for $\|e^{(0)}\| = 0.1, \rho = 10^{-8}$

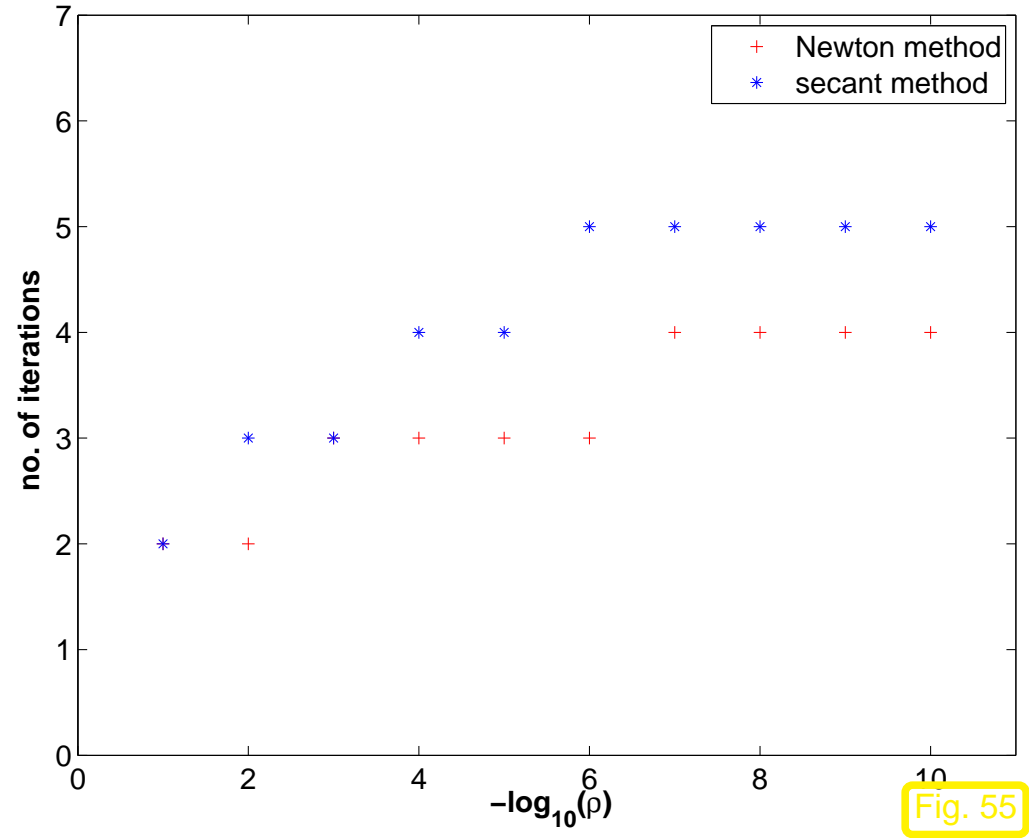


Fig. 55

Newton's method \leftrightarrow secant method, $C = 1$,
initial error $\|e^{(0)}\| = 0.1$

$$\begin{aligned} W_{\text{Newton}} &= 2W_{\text{secant}}, & p_{\text{Newton}} &= 2, & p_{\text{secant}} &= 1.62 \end{aligned} \quad \Rightarrow \quad \frac{\log p_{\text{Newton}}}{W_{\text{Newton}}} : \frac{\log p_{\text{secant}}}{W_{\text{secant}}} = 0.71 .$$

secant method is more efficient than Newton's method!



4.4 Newton's Method [35, Sect. 19], [13, Sect. 5.6]

Non-linear system of equations: for $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ find $\mathbf{x}^* \in D$: $F(\mathbf{x}^*) = 0$

Assume: $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ continuously differentiable

4.4.1 The Newton iteration

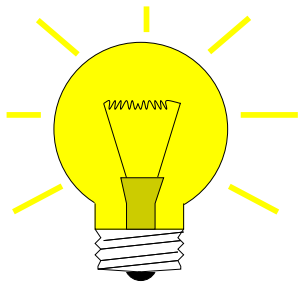
Idea (\rightarrow Sect. 4.3.2.1):

local linearization:

Given $\mathbf{x}^{(k)} \in D \succ \mathbf{x}^{(k+1)}$ as zero of affine linear model function

$$F(\mathbf{x}) \approx \tilde{F}(\mathbf{x}) := F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}),$$

$$DF(\mathbf{x}) \in \mathbb{R}^{n,n} = \text{Jacobian (ger.: Jacobi-Matrix)}, DF(\mathbf{x}) = \left(\frac{\partial F_j}{\partial x_k}(\mathbf{x}) \right)_{j,k=1}^n.$$



Newton iteration: (\leftrightarrow (4.3.3) for $n = 1$)

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - DF(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)}), \quad [\text{if } DF(\mathbf{x}^{(k)}) \text{ regular}] \quad (4.4.1)$$

Terminology: $-DF(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)}) = \text{Newton correction}$

MATLAB template for Newton
method:

Solve linear system:

$A \setminus b = A^{-1}b \rightarrow$ Chapter 2

F, DF : function handles

A posteriori termination criterion

MATLAB-CODE: Newton's method

```
function x = newton(x, F, DF, tr, ta)
for i=1:MAXIT
    s = DF(x) \ F(x);
    x = x-s;
    if ((norm(s)<tr*norm(x)) || (norm(s)<ta))
        return; end;
end
```



New aspect for $n \gg 1$ (compared to $n = 1$ -dimensional case, section. 4.3.2.1):

Computation of the Newton correction may be expensive!

Remark 4.4.2 (Affine invariance of Newton method).

An important property of the Newton iteration (4.4.1): **affine invariance** \rightarrow [15, Sect .1.2.2]

set $G(\mathbf{x}) := \mathbf{A}F(\mathbf{x})$ with regular $\mathbf{A} \in \mathbb{R}^{n,n}$ so that $F(\mathbf{x}^*) = 0 \Leftrightarrow G(\mathbf{x}^*) = 0$.

Affine invariance: Newton iteration for $G(\mathbf{x}) = 0$ is the same for all regular \mathbf{A} !

Code 4.4.7: simplified Newton method

```

1 function x = simpnewton(x,F,DF,rtol,atol)
2 % MATLAB template for simplified Newton method
3 [L,U] = lu(DF(x)); % one LU-decomposition
4 s = U\(L\F(x)); x = x-s;
5 % termination based on relative and absolute
   tolerance
6 ns = norm(s); nx = norm(x);
7 while ((ns > rtol*nx) && (ns > atol))
8     s = U\(L\F(x)); x = x-s;
9 end

```

Simplified Newton Method:

use the same $DF(\mathbf{x}^{(k)})$ for
all/several steps

➤ (usually) merely linear convergence instead of quadratic convergence



4.4.2 Convergence of Newton's method

Example 4.4.9 (Convergence of Newton's method in 2D).

$$F(\mathbf{x}) = \begin{pmatrix} x_1^2 - x_2^4 \\ x_1 - x_2^3 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2 \quad \text{with solution} \quad F \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 0.$$

| 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 | $(0.9999999999999998, 1.0000000000000000)^T$ | 2.11e-15 | |
| 6 | $(1, 1)^T$ | | |



4.4.3 Termination of Newton iteration

Practical a-posteriori termination criterion for Newton's method:

$$DF(\mathbf{x}^{(k-1)}) \approx DF(\mathbf{x}^{(k)}): \text{ quit as soon as } \left\| DF(\mathbf{x}^{(k-1)})^{-1} F(\mathbf{x}^{(k)}) \right\| < \tau_{\text{rel}} \left\| \mathbf{x}^{(k)} \right\|$$

affine invariant termination criterion

Terminology: $\Delta \bar{\mathbf{x}}^{(k)} := DF(\mathbf{x}^{(k-1)})^{-1} F(\mathbf{x}^{(k)}) \hat{=} \text{simplified Newton correction}$

Reuse of LU-factorization (\rightarrow Rem. 2.2.13) of $DF(\mathbf{x}^{(k-1)}) \blacktriangleright \Delta \bar{\mathbf{x}}^{(k)}$ available
with $O(n^2)$ operations

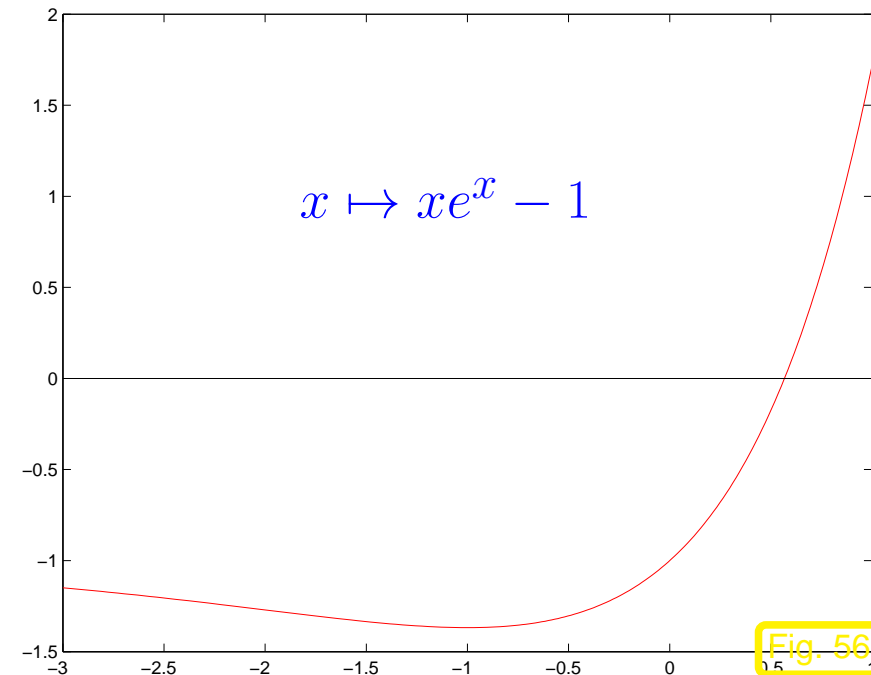
4.4.4 Damped Newton method [13, pp. 200]

Example 4.4.14 (Local convergence of Newton's method).

$$F(x) = xe^x - 1 \Rightarrow F'(-1) = 0$$

$$x^{(0)} < -1 \Rightarrow x^{(k)} \rightarrow -\infty$$

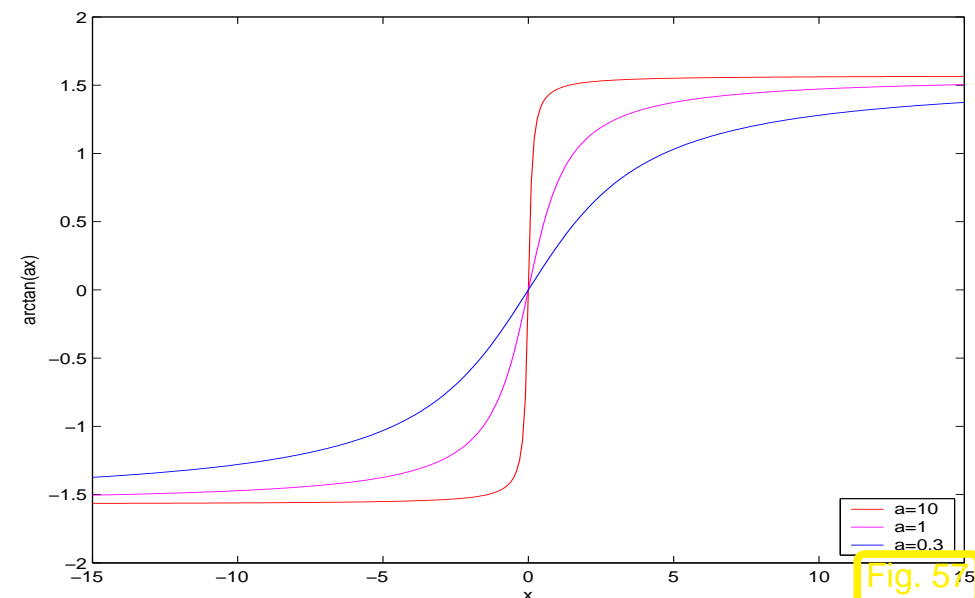
$$x^{(0)} > -1 \Rightarrow x^{(k)} \rightarrow x^*$$

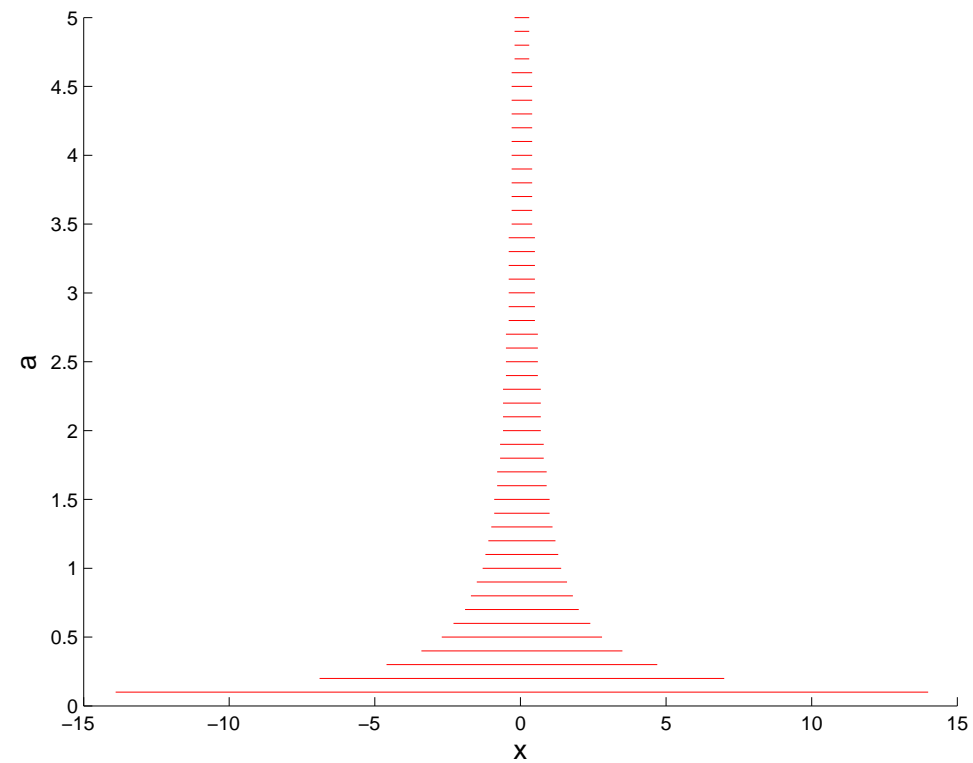
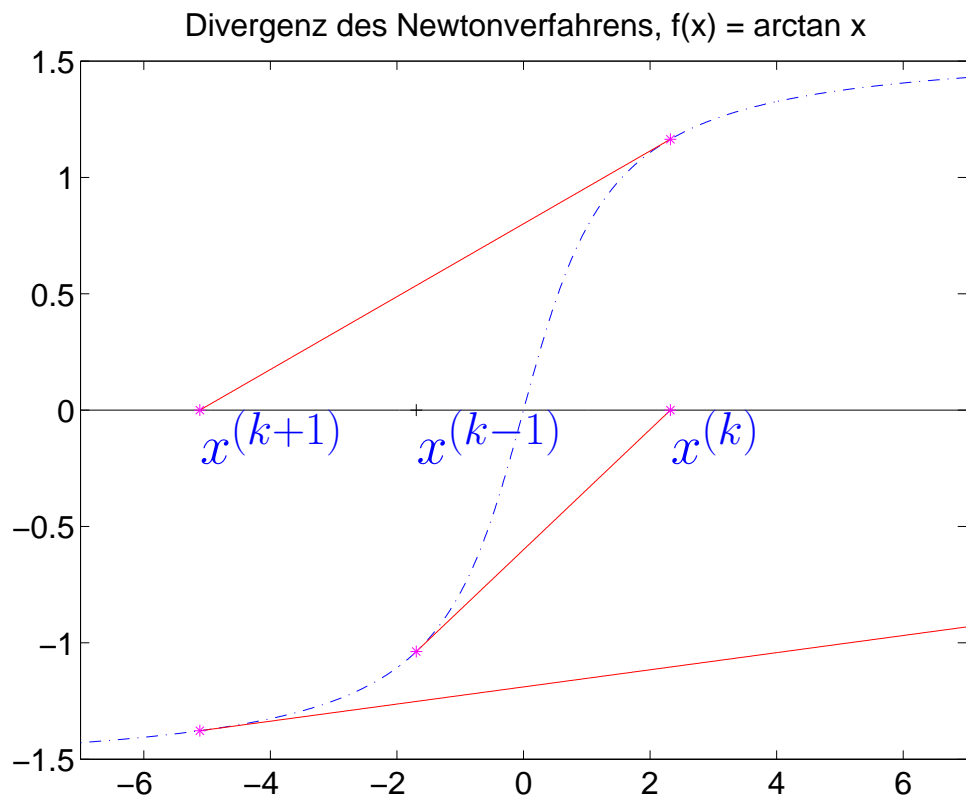


Example 4.4.15 (Region of convergence of Newton method).

$$F(x) = \arctan(ax), \quad a > 0, x \in \mathbb{R}$$

with zero $x^* = 0$.





red zone = $\{x^{(0)} \in \mathbb{R}, x^{(k)} \rightarrow 0\}$

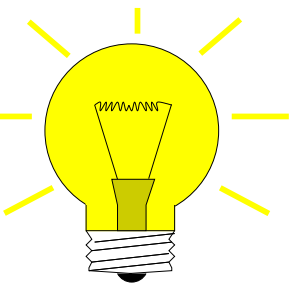
► we observe “overshooting” of Newton correction

Idea:

damping of Newton correction:

$$\text{With } \lambda^{(k)} > 0: \quad \mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \lambda^{(k)} DF(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)}) . \quad (4.4.16)$$

Terminology: $\lambda^{(k)}$ = damping factor



$$\text{“maximal” } 0 < \lambda^{(k)} \leq 1: \quad \left\| \Delta \bar{\mathbf{x}}(\lambda^{(k)}) \right\| \leq \left(1 - \frac{\lambda^{(k)}}{2}\right) \left\| \Delta \mathbf{x}^{(k)} \right\|_2 \quad (4.4.17)$$

where

$$\Delta \mathbf{x}^{(k)} := DF(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)}) \quad \rightarrow \text{current Newton correction ,}$$

$$\Delta \bar{\mathbf{x}}(\lambda^{(k)}) := DF(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)} + \lambda^{(k)} \Delta \mathbf{x}^{(k)}) \quad \rightarrow \text{tentative simplified Newton correction .}$$

Code 4.4.18: Damped Newton method

```

1 function [x,cvg] =
   dampnewton(x,F,DF,rtol,atol)
2 [L,U] = lu(DF(x)); s = U\(L\F(x));
3 xn = x-s; lambda = 1; cvg = 0;
4 f = F(xn); st = U\(L\f); stn = norm(st);
5 while ((stn>rtol*norm(xn)) && (stn > atol))
6   while (norm(st) > (1-lambda/2)*norm(s))
7     lambda = lambda/2;
8     if (lambda < LMIN), cvg = -1; return;
9     end
10    xn = x-lambda*s; f = F(xn);
11    st = U\(L\f);
12 end
13 x = xn; [L,U] = lu(DF(x)); s = U\(L\f);
14 lambda = min(2*lambda,1);
15 xn = x-lambda*s; f = F(xn); st =
   U\(L\f);
16 end
17 x = xn;

```

Reuse of LU-factorization, see
Rem. 2.2.13

a-posteriori termination
criterion (based on
simplified Newton correction,
cf. Sect. 4.4.3)

Natural monotonicity test
(4.4.17)

Reduce damping factor λ

Example 4.4.19 (Damped Newton method). (\rightarrow Ex. 4.4.15)

$$F(x) = \arctan(x),$$

- $x^{(0)} = 20$

- $q = \frac{1}{2}$

- LMIN = 0.001

Observation: asymptotic
quadratic convergence

| k | $\lambda^{(k)}$ | $x^{(k)}$ | $F(x^{(k)})$ |
|-----|-----------------|--------------------|--------------------|
| 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 | -0.000000000000001 | -0.000000000000001 |

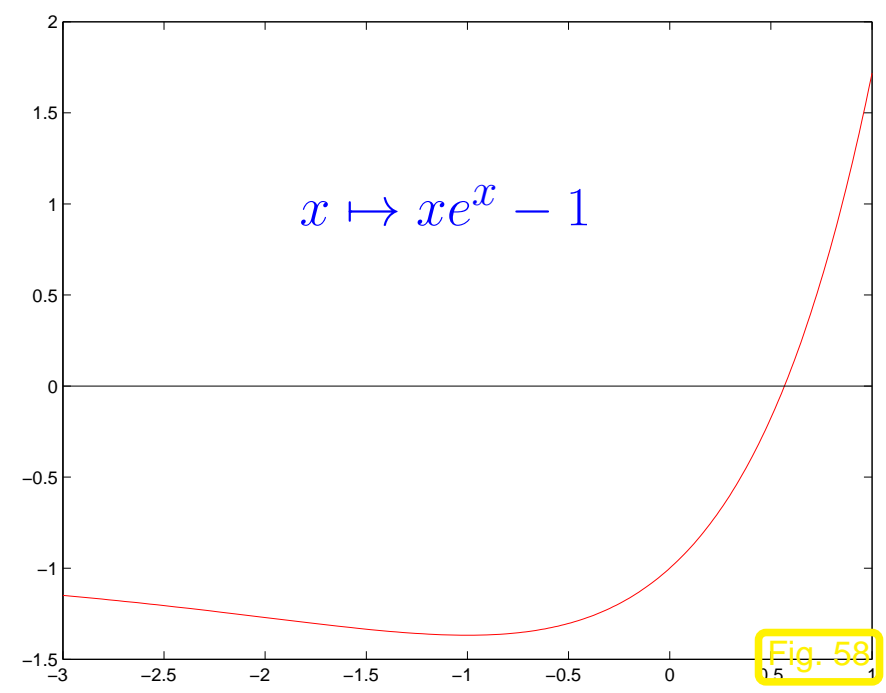


Example 4.4.20 (Failure of damped Newton method).

- As in Ex. 4.4.14:

$$F(x) = xe^x - 1,$$

- Initial guess for damped Newton method $x^{(0)} = -1.5$



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 |

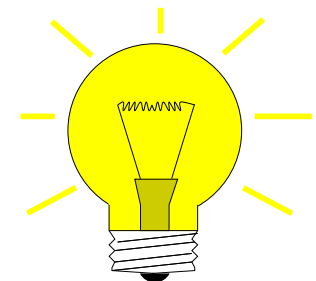
Bailed out because of $\lambda < \text{LMIN}$!



4.4.5 Quasi-Newton Method [51, Sect. 7.1.4]

What to do when $DF(\mathbf{x})$ is not available and numerical differentiation (see remark 4.4.8) is too expensive?

Idea: in one dimension ($n = 1$) apply the secant method (4.3.11) of section 4.3.2.3



$$F'(x^{(k)}) \approx \frac{F(x^{(k)}) - F(x^{(k-1)})}{x^{(k)} - x^{(k-1)}} \quad \text{"difference quotient"} \quad (4.4.21)$$

already computed ! → cheap



Generalisation for $n > 1$?

Idea: rewrite (4.4.21) as a **secant condition** for the approximation $\mathbf{J}_k \approx DF(\mathbf{x}^{(k)})$,
 $\mathbf{x}^{(k)} \hat{=}$ iterate:

$$\mathbf{J}_k(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) = F(\mathbf{x}^{(k)}) - F(\mathbf{x}^{(k-1)}) . \quad (4.4.22)$$

BUT: many matrices \mathbf{J}_k fulfill (4.4.22)

Hence: we need more conditions for $\mathbf{J}_k \in \mathbb{R}^{n,n}$

Idea: get \mathbf{J}_k by a **modification** of \mathbf{J}_{k-1}

Broyden conditions: $\mathbf{J}_k \mathbf{z} = \mathbf{J}_{k-1} \mathbf{z} \quad \forall \mathbf{z}: \mathbf{z} \perp (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) . \quad (4.4.23)$

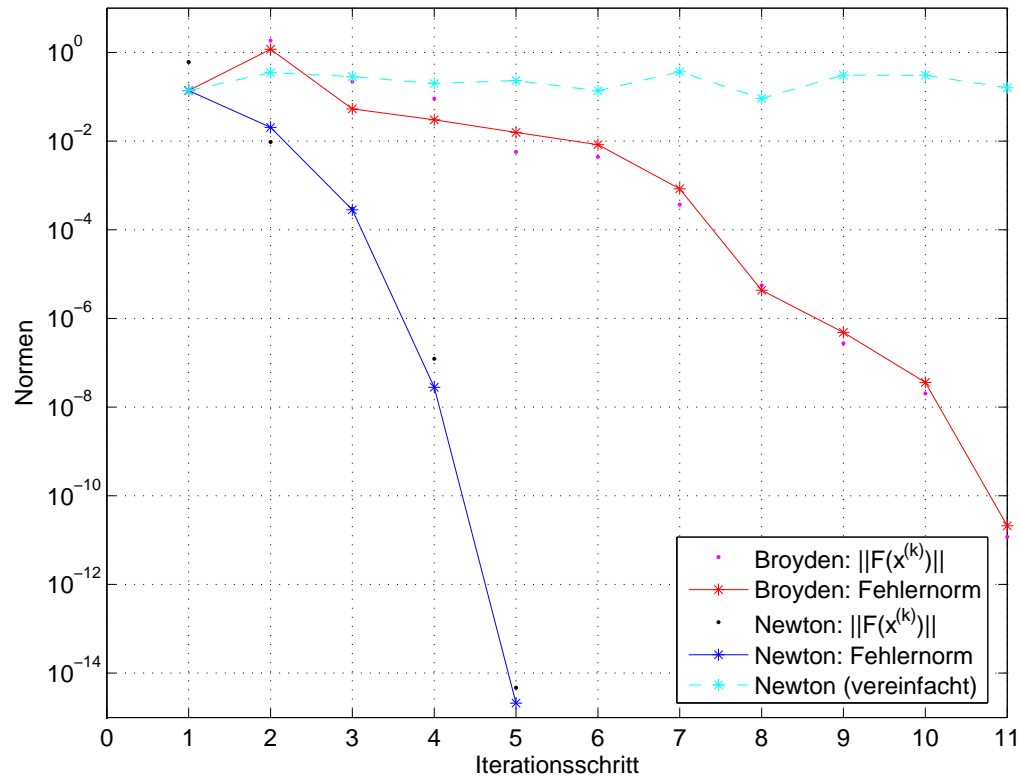
i.e.:
$$\mathbf{J}_k := \mathbf{J}_{k-1} + \frac{F(\mathbf{x}^{(k)})(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})^T}{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_2^2} \quad (4.4.24)$$

Broydens Quasi-Newton Method for solving $F(\mathbf{x}) = 0$:

$$\begin{aligned} \mathbf{x}^{(k+1)} &:= \mathbf{x}^{(k)} + \Delta \mathbf{x}^{(k)}, \quad \Delta \mathbf{x}^{(k)} := -\mathbf{J}_k^{-1} F(\mathbf{x}^{(k)}), \\ \mathbf{J}_{k+1} &:= \mathbf{J}_k + \frac{F(\mathbf{x}^{(k+1)})(\Delta \mathbf{x}^{(k)})^T}{\|\Delta \mathbf{x}^{(k)}\|_2^2} . \end{aligned} \quad (4.4.25)$$

Example 4.4.27 (Broydens Quasi-Newton Method: Convergence).

- In the non-linear system of the example 4.4.9, $n = 2$ take $\mathbf{x}^{(0)} = (0.7, 0.7)^T$ and $\mathbf{J}_0 = DF(\mathbf{x}^{(0)})$



convergence monitor

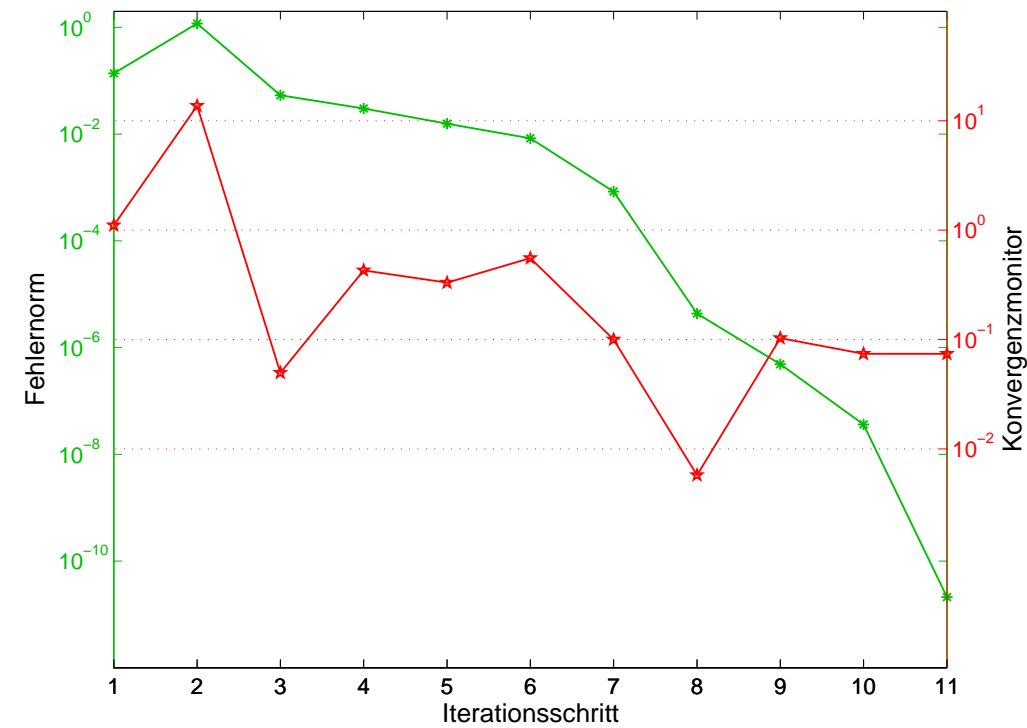
=

quantity that alerts to difficulties in the convergence of an iteration

Here:

$$\mu := \frac{\left\| \mathbf{J}_{k-1}^{-1} F(\mathbf{x}^{(k)}) \right\|}{\left\| \Delta \mathbf{x}^{(k-1)} \right\|}$$

Heuristics: no convergence whenever $\mu > 1$



```
function x = broyden(F,x,J,tol)
```

```
k = 1;
```

```
[L,U] = lu(J);
```

```
s = U\(L\F(x)); sn = dot(s,s);
```

```
dx = [s]; dxn = [sn];
```

```
x = x - s; f = F(x);
```

```
while (sqrt(sn) > tol), k=k+1
```

```
    w = U\(L\f);
```

```
    for l=2:k-1
```

```
        w = w+dx(:,l)*(dx(:,l-1)'*w)...  
            /dxn(l-1);
```

```
    end
```

```
    if (norm(w)>=sn)
```

```
        warning('Dubious step %d!',k);
```

```
    end
```

```
    z = s'*w; s = (1+z/(sn-z))*w; sn=s'*s;
```

```
    dx = [dx,s]; dxn = [dxn,sn];
```

```
    x = x - s; f = F(x);
```

```
end
```

unique LU-decomposition !

store $\Delta \mathbf{x}^{(k)}$, $\|\Delta \mathbf{x}^{(k)}\|_2^2$
(see (4.4.29))

solve two SLEs

Termination, see 4.4.2

construct $\mathbf{w} := \mathbf{J}_k^{-1} F(\mathbf{x}^{(k)})$
(\rightarrow recursion (4.4.29))

convergence monitor

$$\frac{\|\mathbf{J}_{k-1}^{-1} F(\mathbf{x}^{(k)})\|}{\|\Delta \mathbf{x}^{(k-1)}\|} < 1 ?$$

correction $\mathbf{s} = \mathbf{J}_k^{-1} F(\mathbf{x}^{(k)})$

- Computational cost : $O(N^2 \cdot n)$ operations with vectors, (Level I)
 N steps
- 1 LU-decomposition of J , $N \times$ solutions of SLEs, see section 2.2
 - N evaluations of F !

- Memory cost : LU-factors of J + auxiliary vectors $\in \mathbb{R}^n$
 N steps
- N vectors $\mathbf{x}^{(k)} \in \mathbb{R}^n$

Example 4.4.30 (Broyden method for a large non-linear system).

$$F(\mathbf{x}) = \begin{cases} \mathbb{R}^n \mapsto \mathbb{R}^n \\ \mathbf{x} \mapsto \text{diag}(\mathbf{x})\mathbf{A}\mathbf{x} - \mathbf{b} \end{cases},$$

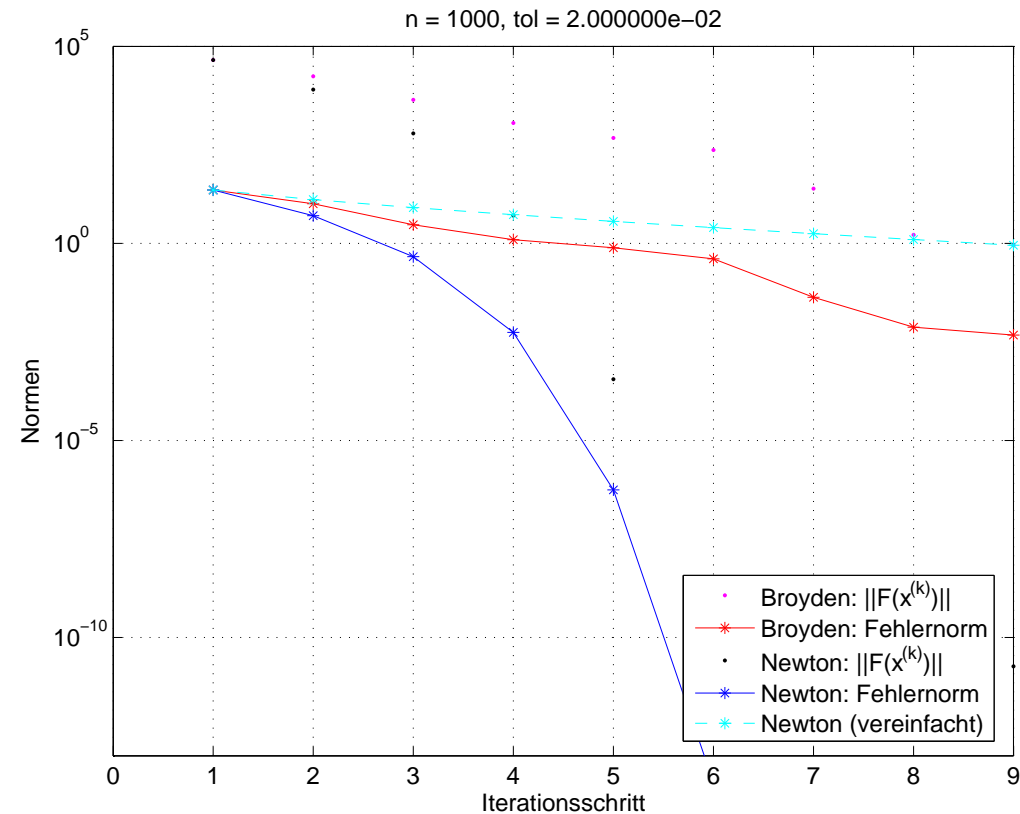
$$\mathbf{b} = (1, 2, \dots, n) \in \mathbb{R}^n,$$

$$\mathbf{A} = \mathbf{I} + \mathbf{a}\mathbf{a}^T \in \mathbb{R}^{n,n},$$

$$\mathbf{a} = \frac{1}{\sqrt{\mathbf{1} \cdot \mathbf{b} - \mathbf{1}}}(\mathbf{b} - \mathbf{1}).$$

The interpretation of the results resemble the example 4.4.27 \triangleright

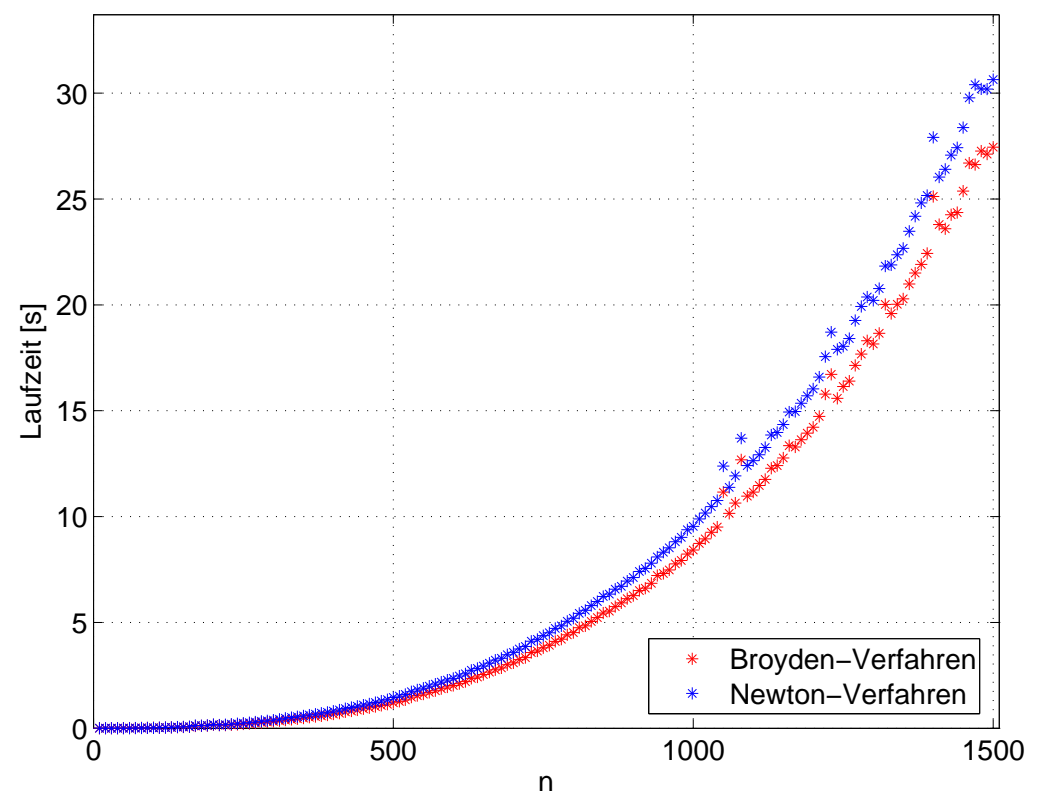
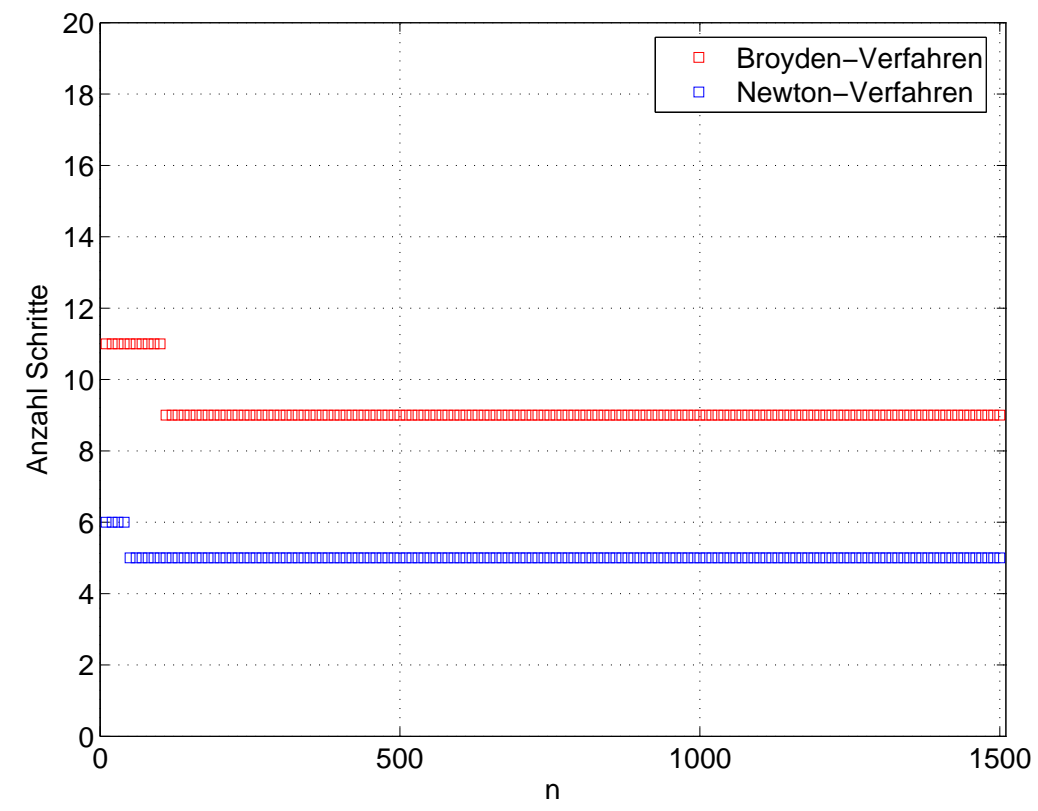
$$h = 2/n; \quad \mathbf{x}_0 = (2:h:4-h)';$$



Efficiency comparison:

Broyden method \longleftrightarrow Newton method:

(in case of dimension n use tolerance $\text{tol} = 2n \cdot 10^{-5}$, $h = 2/n$; $x_0 = (2:h:4-h)'$;)



5

Krylov Methods for Linear Systems of Equations

5.1 Descent Methods [51, 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. 2.7.9) system matrix \mathbf{A}

\rightarrow \mathbf{A} -inner product $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x}^\top \mathbf{A} \mathbf{y} \Rightarrow$ “ \mathbf{A} -geometry”

Definition 5.1.1 (Energy norm). \rightarrow [35, Def. 9.1]

A s.p.d. matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ induces an **energy norm**

$$\|\mathbf{x}\|_A := (\mathbf{x}^\top \mathbf{A} \mathbf{x})^{1/2}, \quad \mathbf{x} \in \mathbb{R}^n.$$

5.1.1 Quadratic minimization context

Lemma 5.1.3 (S.p.d. LSE and quadratic minimization problem). [13, (13.37)]

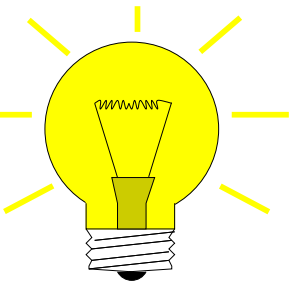
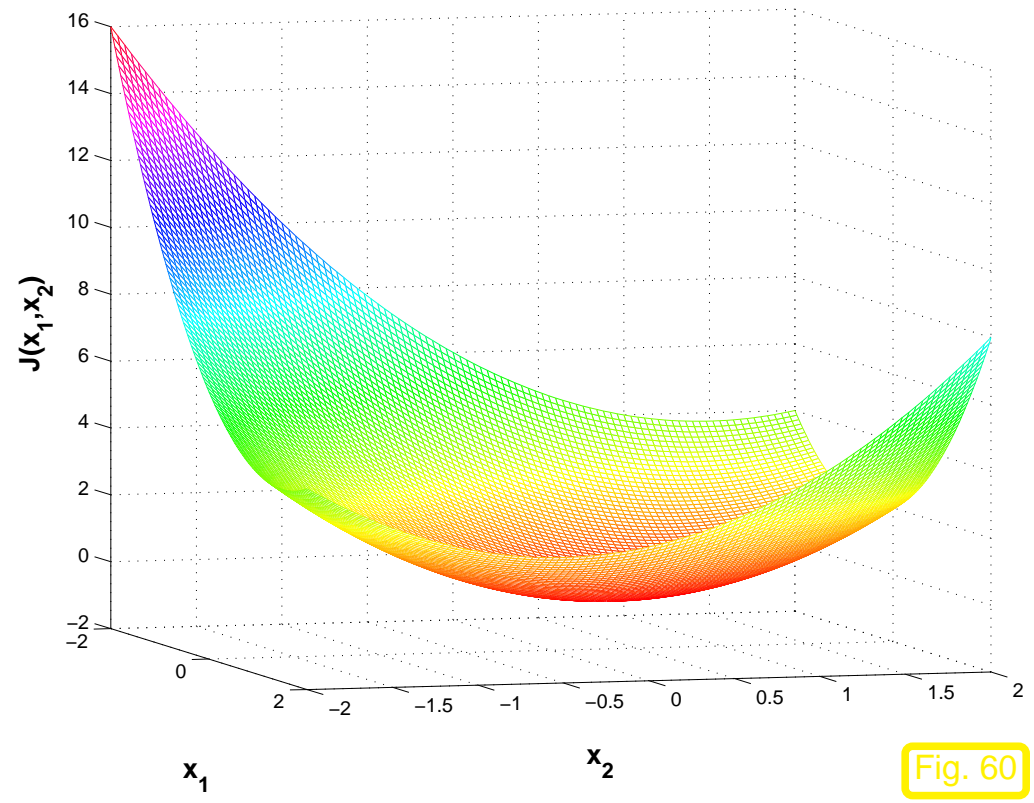
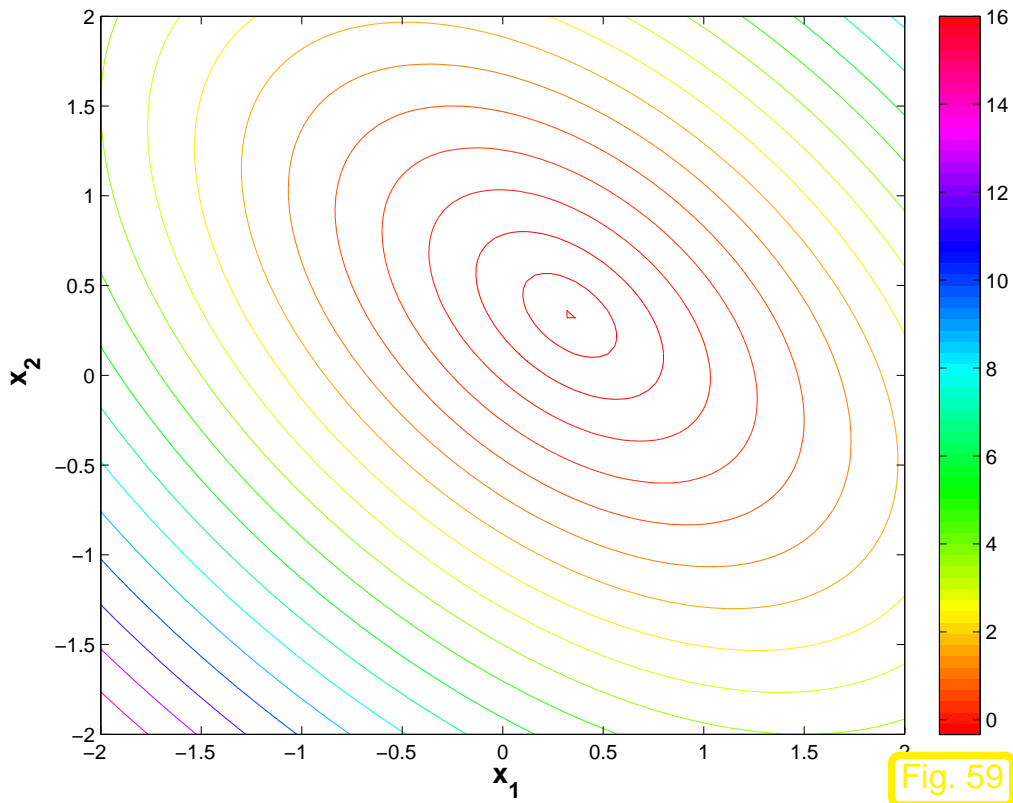
A LSE with $\mathbf{A} \in \mathbb{R}^{n,n}$ s.p.d. and $\mathbf{b} \in \mathbb{R}^n$ is equivalent to a minimization problem:

$$\mathbf{Ax} = \mathbf{b} \Leftrightarrow \mathbf{x} = \arg \min_{\mathbf{y} \in \mathbb{R}^n} J(\mathbf{y}), \quad J(\mathbf{y}) := \frac{1}{2} \mathbf{y}^\top \mathbf{A} \mathbf{y} - \mathbf{b}^\top \mathbf{y}. \quad (5.1.4)$$

A quadratic functional

Example 5.1.6 (Quadratic functional in 2D).

Plot of J from (5.1.4) for $\mathbf{A} = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$, $\mathbf{b} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$.



Algorithmic idea: (Lemma 5.1.3 \Rightarrow) Solve $\mathbf{Ax} = \mathbf{b}$ iteratively by successive solution of *simpler* minimization problems

5.1.2 Abstract steepest descent

Task: Given **continuously differentiable** $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}$,
 find **minimizer** $\mathbf{x}^* \in D$: $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in D} F(\mathbf{x})$

The most natural iteration:

Algorithm 5.1.7 (Steepest descent). (*ger.:* steilster Abstieg)

Initial guess $\mathbf{x}^{(0)} \in D, k = 0$

repeat

$$\mathbf{d}_k := -\operatorname{grad} F(\mathbf{x}^{(k)})$$

$$t^* := \operatorname{argmin}_{t \in \mathbb{R}} F(\mathbf{x}^{(k)} + t\mathbf{d}_k) \quad (\text{line search})$$

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + t^* \mathbf{d}_k$$

$$k := k + 1$$

until $\left(\begin{array}{l} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau_{\text{rel}} \|\mathbf{x}^{(k)}\| \text{ or} \\ \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau_{\text{abs}} \end{array} \right)$

- $\mathbf{d}_k \hat{=}$ *direction of steepest descent*
- linear search $\hat{=}$ 1D minimization: use Newton's method (\rightarrow Sect. 4.3.2.1) on derivative
- correction based a posteriori termination criterion, see Sect. 4.1.2 for a discussion.
($\tau \hat{=}$ prescribed tolerance)

The **gradient** (\rightarrow [63, Kapitel 7])

$$\mathbf{grad} F(\mathbf{x}) = \begin{pmatrix} \frac{\partial F}{\partial x_1} \\ \vdots \\ \frac{\partial F}{\partial x_n} \end{pmatrix} \in \mathbb{R}^n \quad (5.1.8)$$

provides the direction of **local** steepest ascent/descent of F

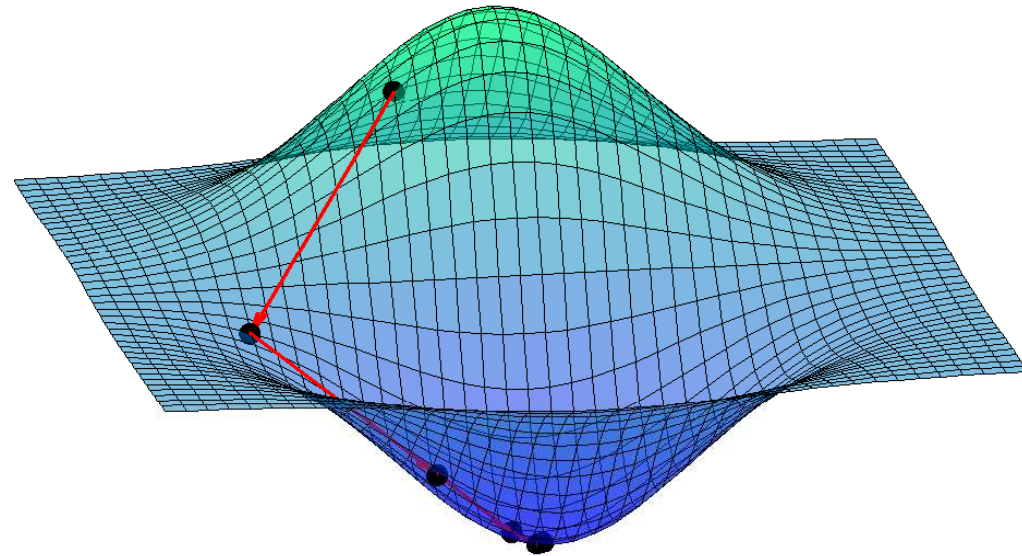


Fig. 61

5.1.3 Gradient method for s.p.d. linear system of equations

Steepest descent iteration = **gradient method** for LSE $\mathbf{Ax} = \mathbf{b}$, $\mathbf{A} \in \mathbb{R}^{n,n}$ s.p.d., $\mathbf{b} \in \mathbb{R}^n$:

Algorithm 5.1.11 (Gradient method for s.p.d. LSE).

Code 5.1.12: gradient method for $\mathbf{Ax} = \mathbf{b}$, \mathbf{A} s.p.d.Initial guess $\mathbf{x}^{(0)} \in \mathbb{R}^n$, $k = 0$ $\mathbf{r}_0 := \mathbf{b} - \mathbf{Ax}^{(0)}$ **repeat**

$$t^* := \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{A} \mathbf{r}_k}$$

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + t^* \mathbf{r}_k$$

$$\mathbf{r}_{k+1} := \mathbf{r}_k - t^* \mathbf{A} \mathbf{r}_k$$

$$k := k + 1$$

until $\left(\begin{array}{l} \left\| \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right\| \leq \tau_{\text{rel}} \left\| \mathbf{x}^{(k)} \right\| \text{ or} \\ \left\| \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right\| \leq \tau_{\text{abs}} \end{array} \right)$

```

1 function x =
   gradit(A,b,x,rtol,atol,maxit)
2 r = b-A*x; % residual → Def. 2.5.16
3 for k=1:maxit
4   p = A*r;
5   ts = (r'*r)/(r'*p); % cf. (5.1.10)
6   x = x + ts*r;
7   cn = (abs(ts)*norm(r)); % norm of
   correction
8   if ((cn < tol*norm(x)) ||
9       (cn < atol))
10      return; end
11   r = r - ts*p; %
12 end

```

One step of gradient method involves

- A single matrix \times vector product with \mathbf{A} ,
- 2 AXPY-operations (\rightarrow Sect. 1.4) on vectors of length n ,
- 2 dot products in \mathbb{R}^n .

Computational cost (per step) = cost(matrix \times vector) + $O(n)$

5.1.4 Convergence of the gradient method

Example 5.1.14 (Gradient method in 2D).

S.p.d. matrices $\in \mathbb{R}^{2,2}$:

$$\mathbf{A}_1 = \begin{pmatrix} 1.9412 & -0.2353 \\ -0.2353 & 1.0588 \end{pmatrix}, \quad \mathbf{A}_2 = \begin{pmatrix} 7.5353 & -1.8588 \\ -1.8588 & 0.5647 \end{pmatrix}$$

Eigenvalues: $\sigma(\mathbf{A}_1) = \{1, 2\}$, $\sigma(\mathbf{A}_2) = \{0.1, 8\}$

notation: **spectrum** of a matrix $\in \mathbb{K}^{n,n}$ $\sigma(\mathbf{M}) := \{\lambda \in \mathbb{C} : \lambda \text{ is eigenvalue of } \mathbf{M}\}$

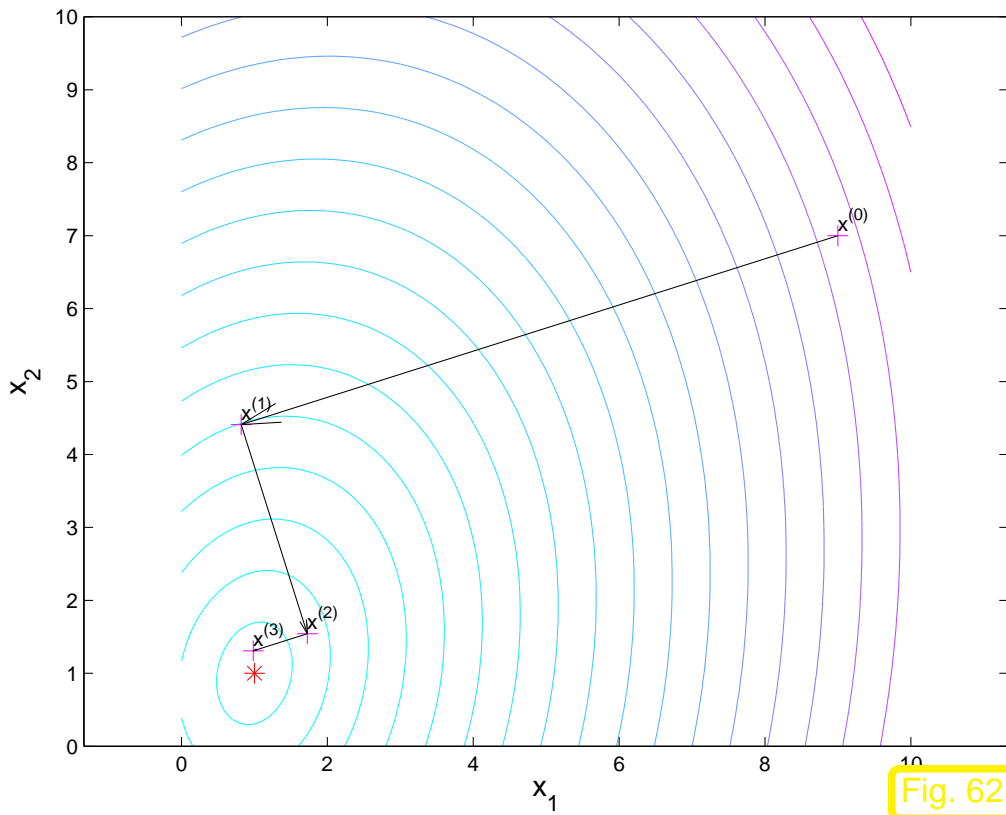


Fig. 62

iterates of Alg. 5.1.11 for A_1

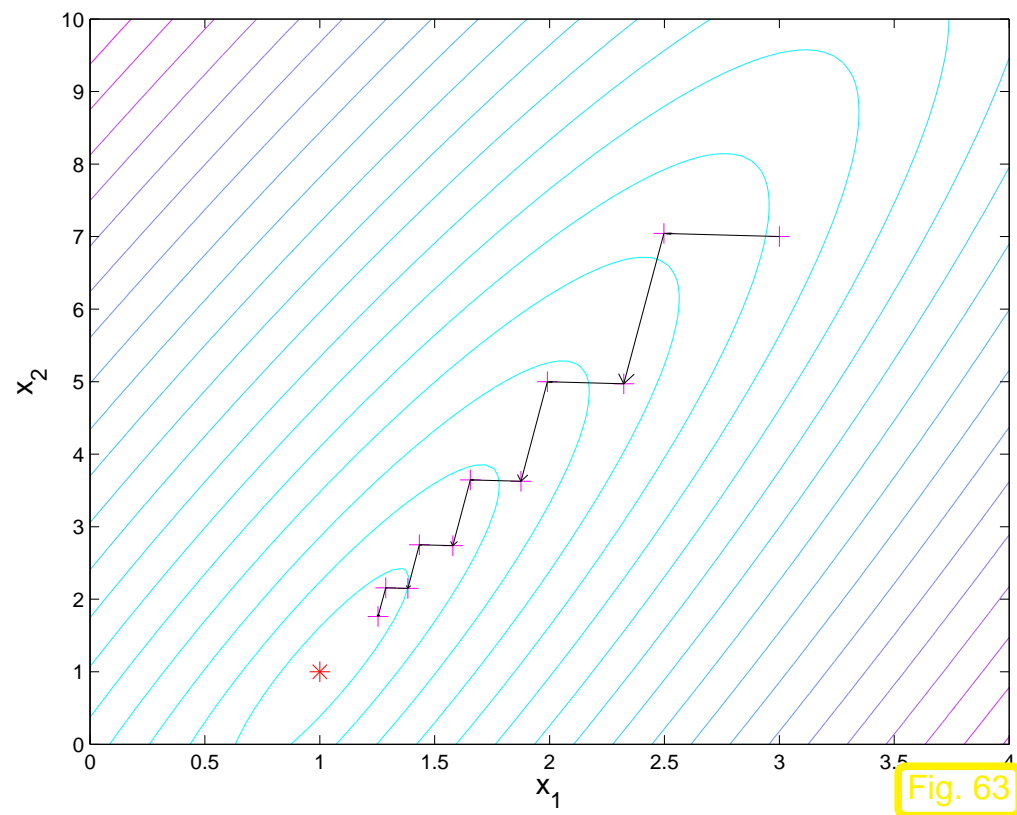


Fig. 63

iterates of Alg. 5.1.11 for A_2



Example 5.1.17 (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);

```

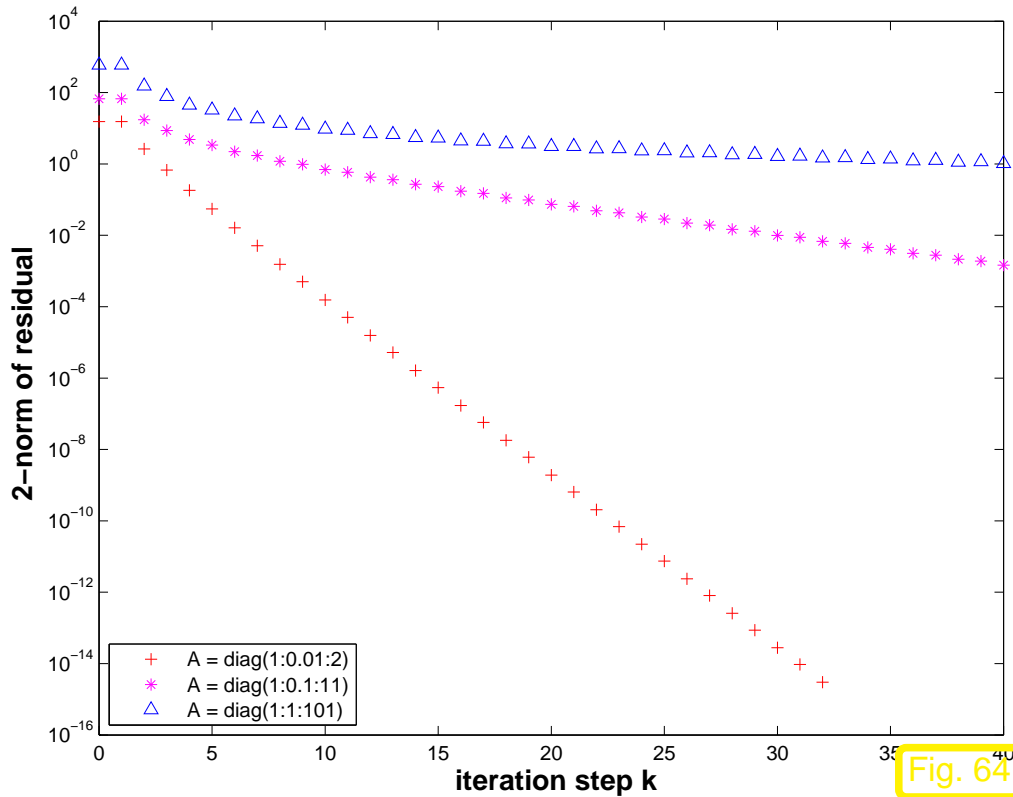



Fig. 64

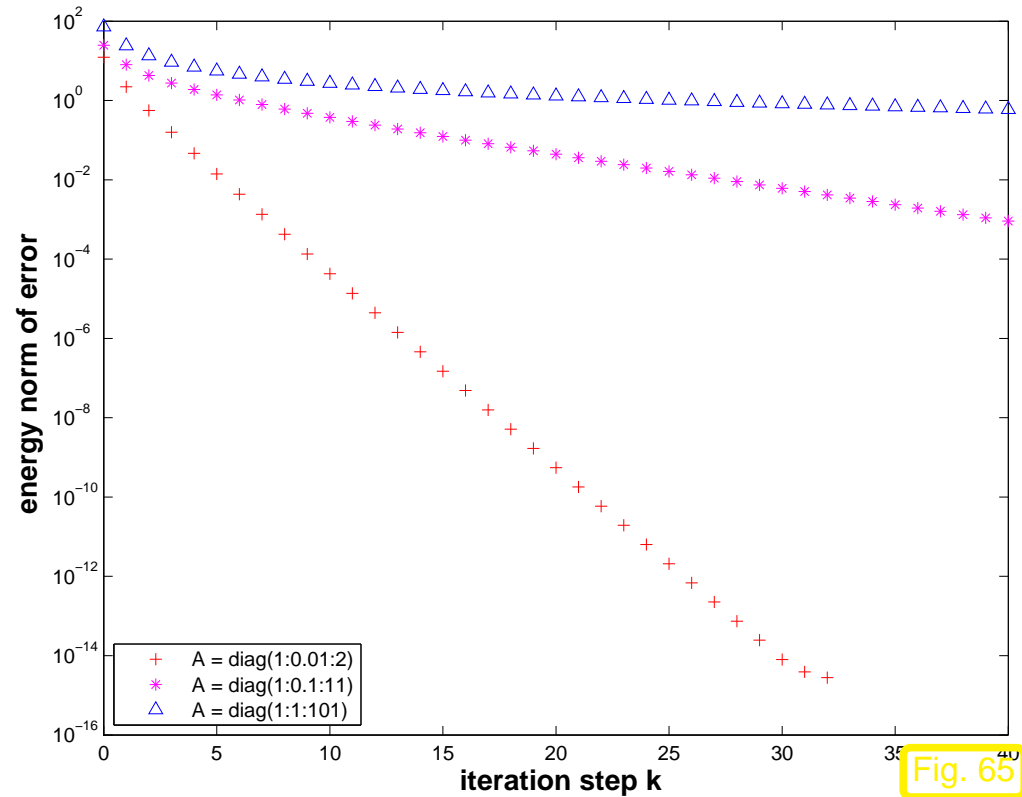


Fig. 65

Impact of distribution of diagonal entries (\leftrightarrow eigenvalues) of (diagonal matrix) A
 $(b = x^* = 0, x_0 = \cos((1:n)'))$

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

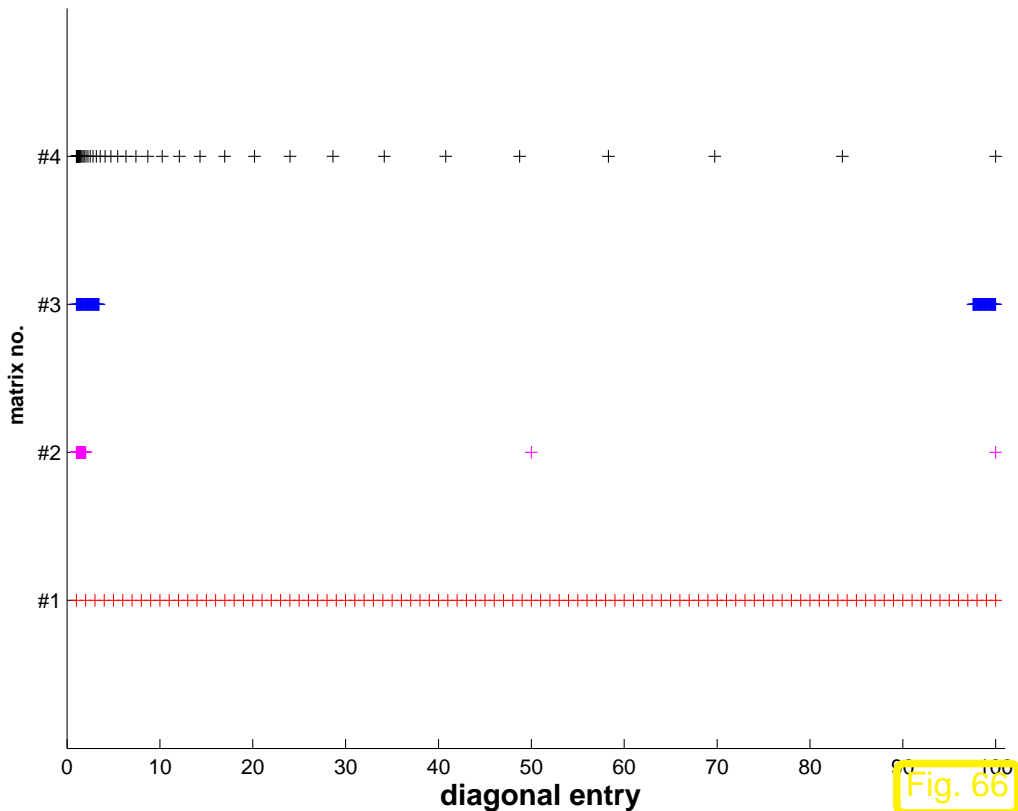
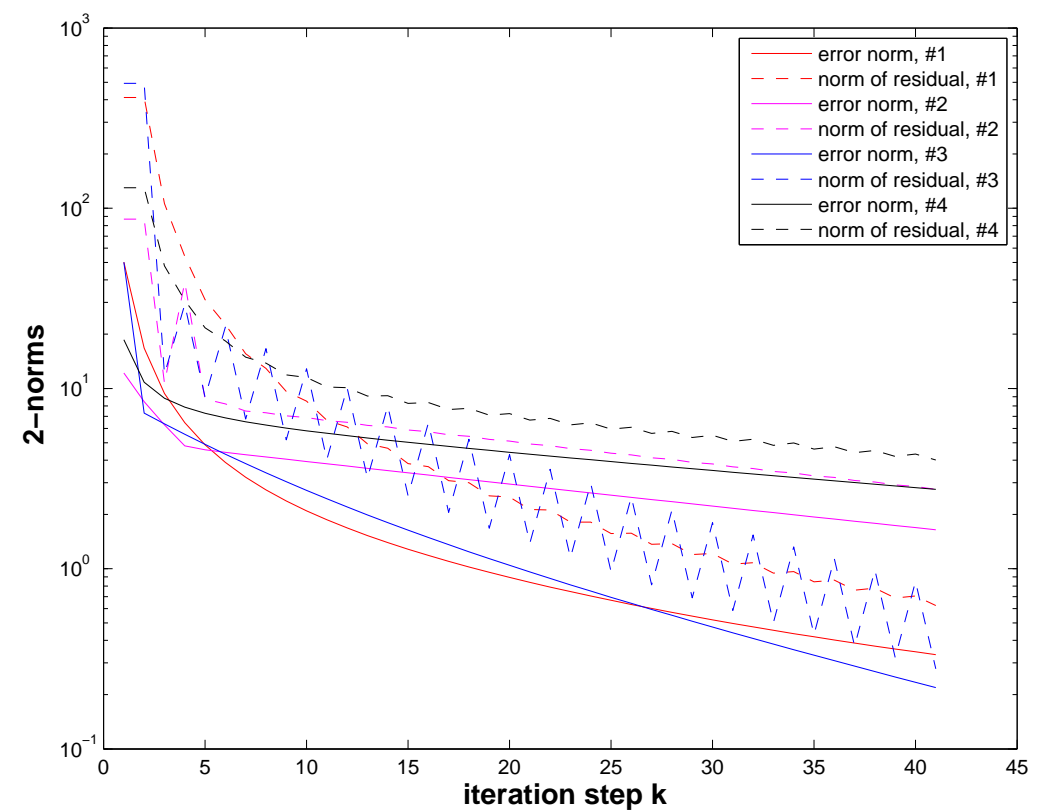


Fig. 66



Theorem 5.1.18 (Convergence of gradient method/steepest descent).

The iterates of the gradient method of Alg. 5.1.11 satisfy

$$\left\| \mathbf{x}^{(k+1)} - \mathbf{x}^* \right\|_A \leq L \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|_A, \quad L := \frac{\text{cond}_2(\mathbf{A}) - 1}{\text{cond}_2(\mathbf{A}) + 1},$$

that is, the iteration converges at least linearly (\rightarrow Def. 4.1.6) w.r.t. energy norm (\rightarrow Def. 5.1.1).

notation: $\text{cond}_2(\mathbf{A}) \hat{=}$ condition number (\rightarrow Def. 2.5.26) of \mathbf{A} induced by 2-norm

5.2 Conjugate gradient method (CG) [35, Ch. 9], [13, Sect. 13.4], [51, Sect. 4.3.4]

Liability of gradient method of Sect. 5.1.3:

NO MEMORY

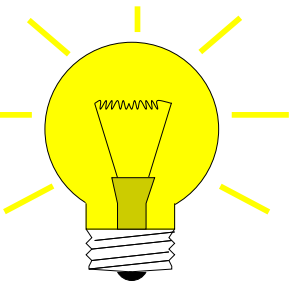
Idea:

Replace linear search with **subspace correction**

- Given:
- initial guess $\mathbf{x}^{(0)}$
 - **nested** subspaces $U_1 \subset U_2 \subset U_3 \subset \dots \subset U_n = \mathbb{R}^n$, $\dim U_k = k$

$$\mathbf{x}^{(k)} := \underset{\mathbf{x} \in U_k + \mathbf{x}^{(0)}}{\operatorname{argmin}} J(x), \quad (5.2.1)$$

quadratic functional from (5.1.4)





How to find suitable subspaces U_k ?

Idea: $U_{k+1} \leftarrow U_k +$ “local steepest descent direction”

given by $-\text{grad } J(\mathbf{x}^{(k)}) = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)} = \mathbf{r}_k$ (residual \rightarrow Def. 2.5.16)

$$U_{k+1} = \text{Span} \{U_k, \mathbf{r}_k\}, \quad \mathbf{x}^{(k)} \text{ from (5.2.1)}. \quad (5.2.2)$$

(5.2.1) and (5.2.2) define the **conjugate gradient method** (CG) for the iterative solution of

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

Definition 5.2.6 (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

$$\mathcal{K}_l(\mathbf{A}, \mathbf{z}) := \text{Span} \left\{ \mathbf{z}, \mathbf{A}\mathbf{z}, \dots, \mathbf{A}^{l-1}\mathbf{z} \right\} .$$

Lemma 5.2.7. The subspaces $U_k \subset \mathbb{R}^n$, $k \geq 1$, defined by (5.2.1) and (5.2.2) satisfy

$$U_k = \text{Span} \left\{ \mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \dots, \mathbf{A}^{k-1}\mathbf{r}_0 \right\} = \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0) ,$$

where $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$ is the initial residual.

5.2.2 Implementation of CG

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 $\mathbf{Ax} = \mathbf{b}$

$$\mathbf{r}_0 := \mathbf{b} - \mathbf{Ax}^{(0)};$$

$$\text{for } j = 1 \text{ to } l \text{ do } \left\{ \mathbf{x}^{(j)} := \mathbf{x}^{(j-1)} + \frac{\mathbf{p}_j^\top \mathbf{r}_0}{\mathbf{p}_j^\top \mathbf{A} \mathbf{p}_j} \mathbf{p}_j \right\} \quad (5.2.9)$$

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 \blacktriangleright short recursions

Algorithm 5.2.17 (CG method for solving $\mathbf{Ax} = \mathbf{b}$, \mathbf{A} s.p.d.). \rightarrow [13, Alg. 13.27]

Input : initial guess $\mathbf{x}^{(0)} \in \mathbb{R}^n$

Output : approximate solution $\mathbf{x}^{(l)} \in \mathbb{R}^n$

$$\mathbf{p}_1 = \mathbf{r}_0 := \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)};$$

for $j = 1$ to l do {

$$\mathbf{x}^{(j)} := \mathbf{x}^{(j-1)} + \frac{\mathbf{p}_j^T \mathbf{r}_{j-1}}{\mathbf{p}_j^T \mathbf{A}\mathbf{p}_j} \mathbf{p}_j;$$

$$\mathbf{r}_j = \mathbf{r}_{j-1} - \frac{\mathbf{p}_j^T \mathbf{r}_{j-1}}{\mathbf{p}_j^T \mathbf{A}\mathbf{p}_j} \mathbf{A}\mathbf{p}_j;$$

$$\mathbf{p}_{j+1} = \mathbf{r}_j - \frac{(\mathbf{A}\mathbf{p}_j)^T \mathbf{r}_j}{\mathbf{p}_j^T \mathbf{A}\mathbf{p}_j} \mathbf{p}_j;$$

}

Input: initial guess $\mathbf{x} \hat{=} \mathbf{x}^{(0)} \in \mathbb{R}^n$

tolerance $\tau > 0$

Output: approximate solution $\mathbf{x} \hat{=} \mathbf{x}^{(l)}$

$$\mathbf{p} := \mathbf{r}_0 := \mathbf{r} := \mathbf{b} - \mathbf{A}\mathbf{x};$$

for $j = 1$ to l_{\max} do {

$$\beta := \mathbf{r}^T \mathbf{r};$$

$$\mathbf{h} := \mathbf{A}\mathbf{p};$$

$$\alpha := \frac{\beta}{\mathbf{p}^T \mathbf{h}};$$

$$\mathbf{x} := \mathbf{x} + \alpha \mathbf{p};$$

$$\mathbf{r} := \mathbf{r} - \alpha \mathbf{h};$$

if $\|\mathbf{r}\| \leq \tau \|\mathbf{r}_0\|$ then stop;

$$\beta := \frac{\mathbf{r}^T \mathbf{r}}{\beta};$$

$$\mathbf{p} := \mathbf{r} + \beta \mathbf{p};$$

}

1 matrix \times vector product, 3 dot products, 3 **AXPY-operations** per step:

If **A** sparse, $\text{nnz}(\mathbf{A}) \sim n$ \blacktriangleright computational effort $O(n)$ per step

Code 5.2.18: basic CG iteration for solving $\mathbf{Ax} = \mathbf{b}$, Alg. 5.2.17

```

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 Rem. 5.2.19
12    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

5.2.3 Convergence of CG

Example 5.2.21 (Impact of roundoff errors on CG). \rightarrow [51, Rem. 4.3]

Numerical experiment: $A = \text{hilb}(20)$,
 $\mathbf{x}^{(0)} = \mathbf{0}$, $\mathbf{b} = (1, \dots, 1)^\top$

Hilbert-Matrix: extremely ill-conditioned

residual norms during CG iteration

$$\mathbf{R} = [\mathbf{r}_0, \dots, \mathbf{r}^{(10)}]$$

▷:

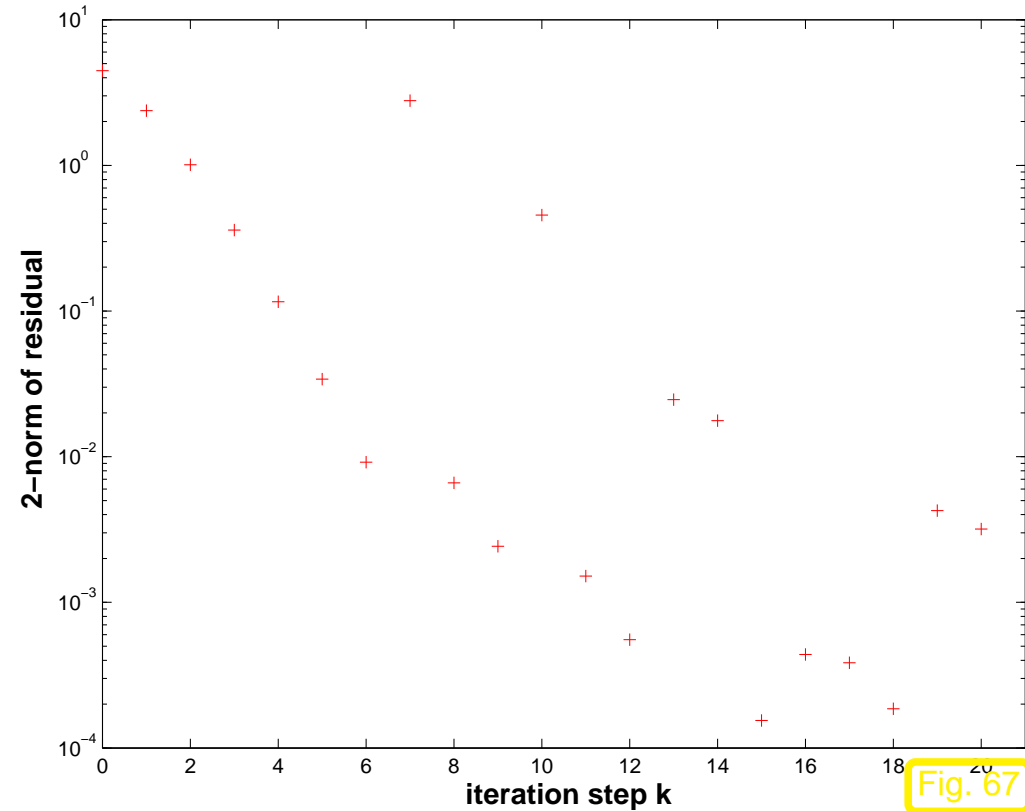


Fig. 67

$$\mathbf{R}^\top \mathbf{R} =$$

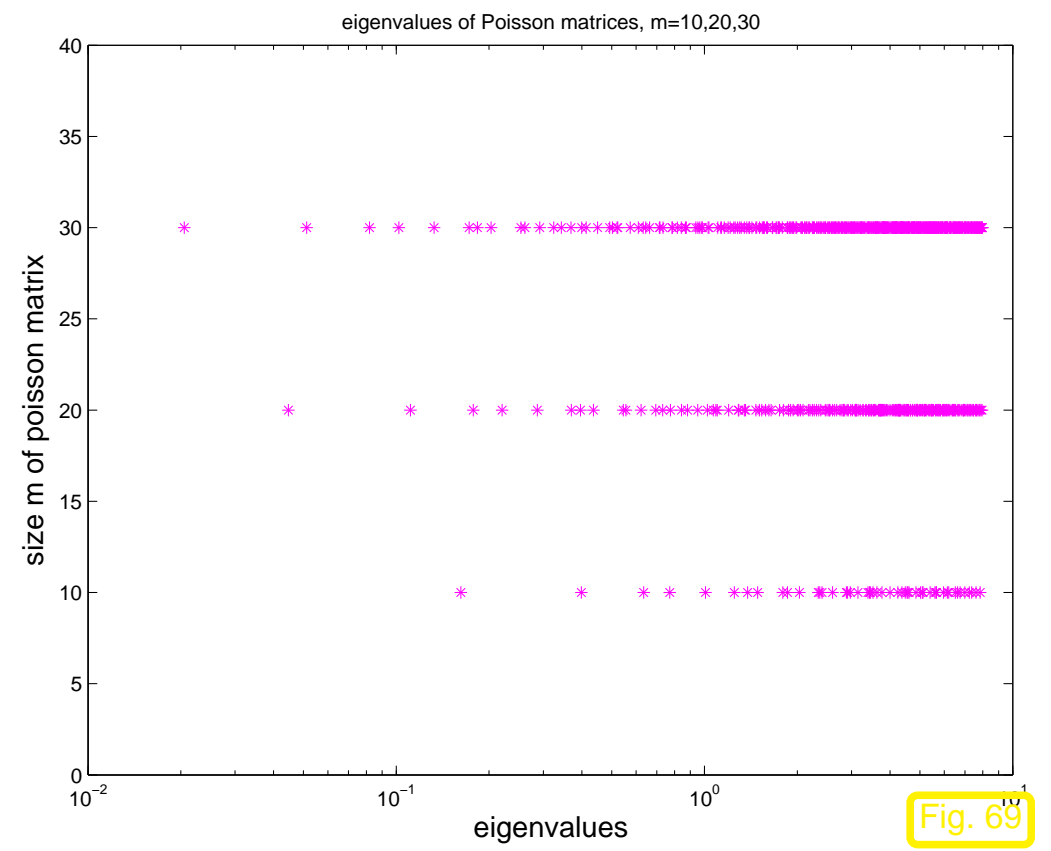
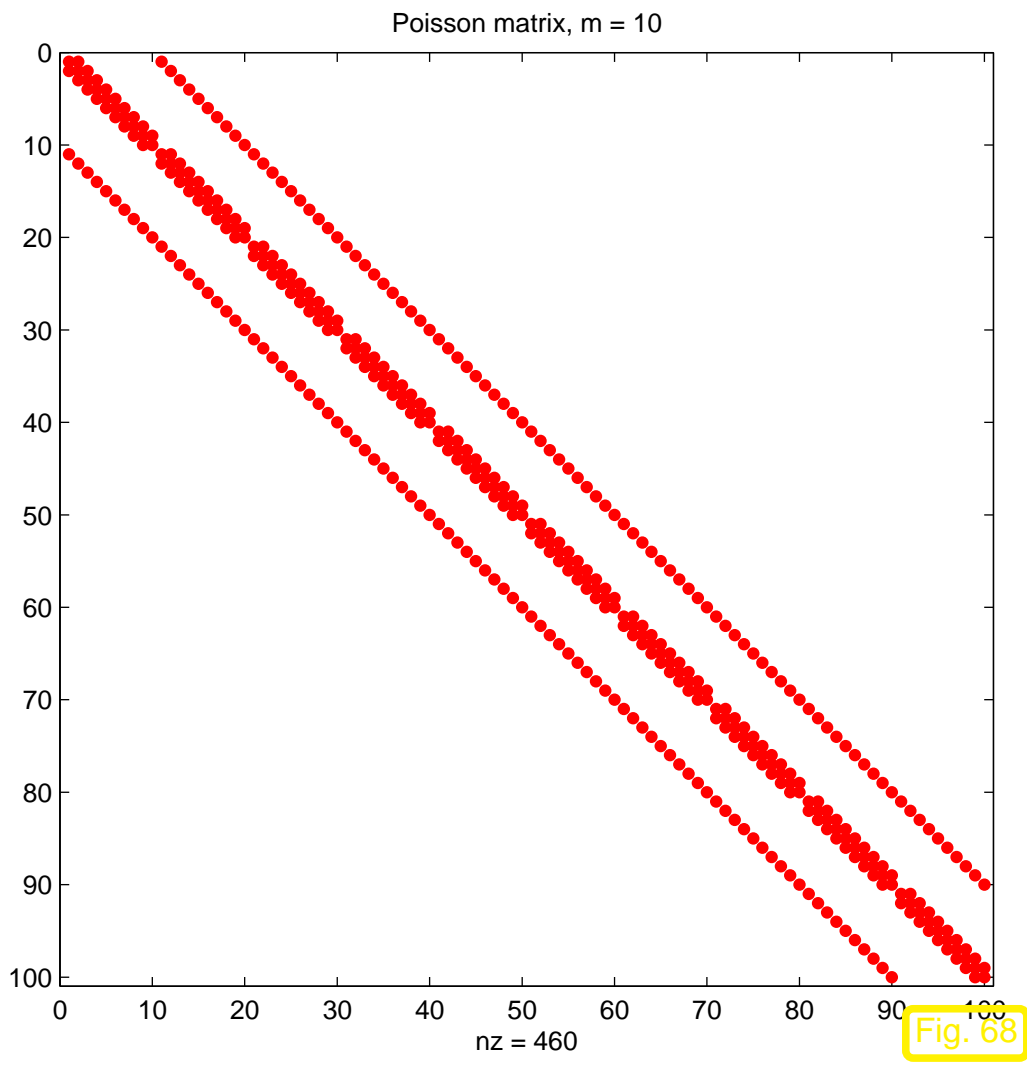
| | | | | | | | | | |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 1.000000 | -0.000000 | 0.000000 | -0.000000 | 0.000000 | -0.000000 | 0.016019 | -0.795816 | -0.430569 | 0.348133 |
| -0.000000 | 1.000000 | -0.000000 | 0.000000 | -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 |

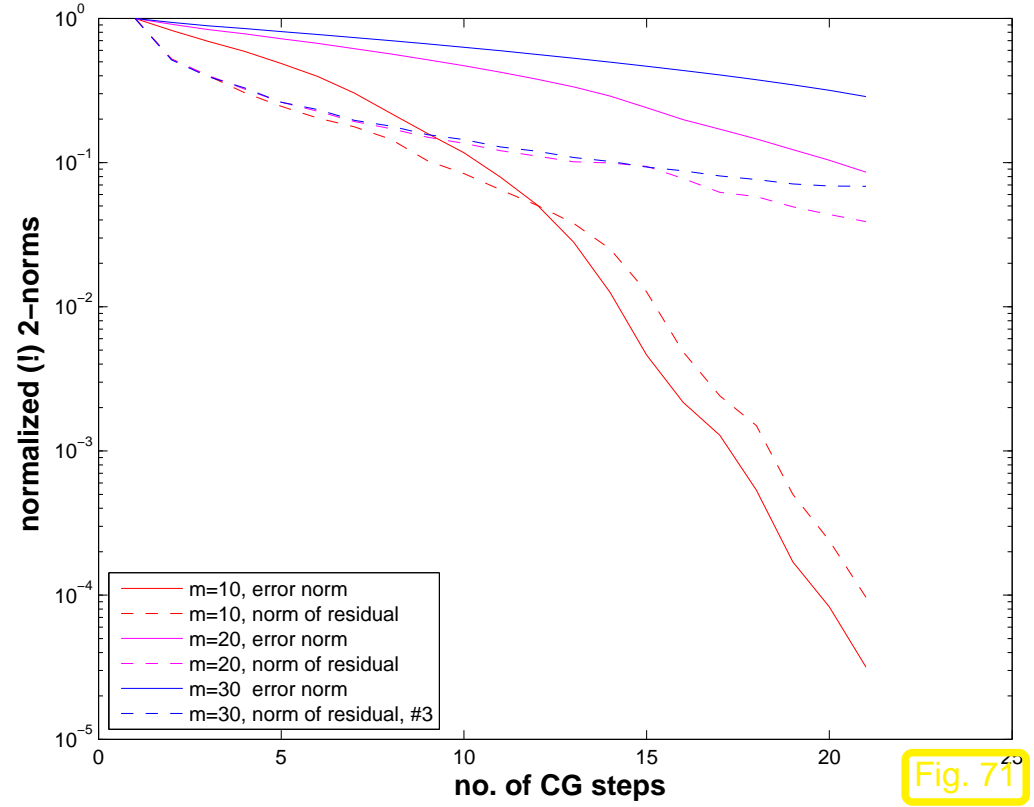
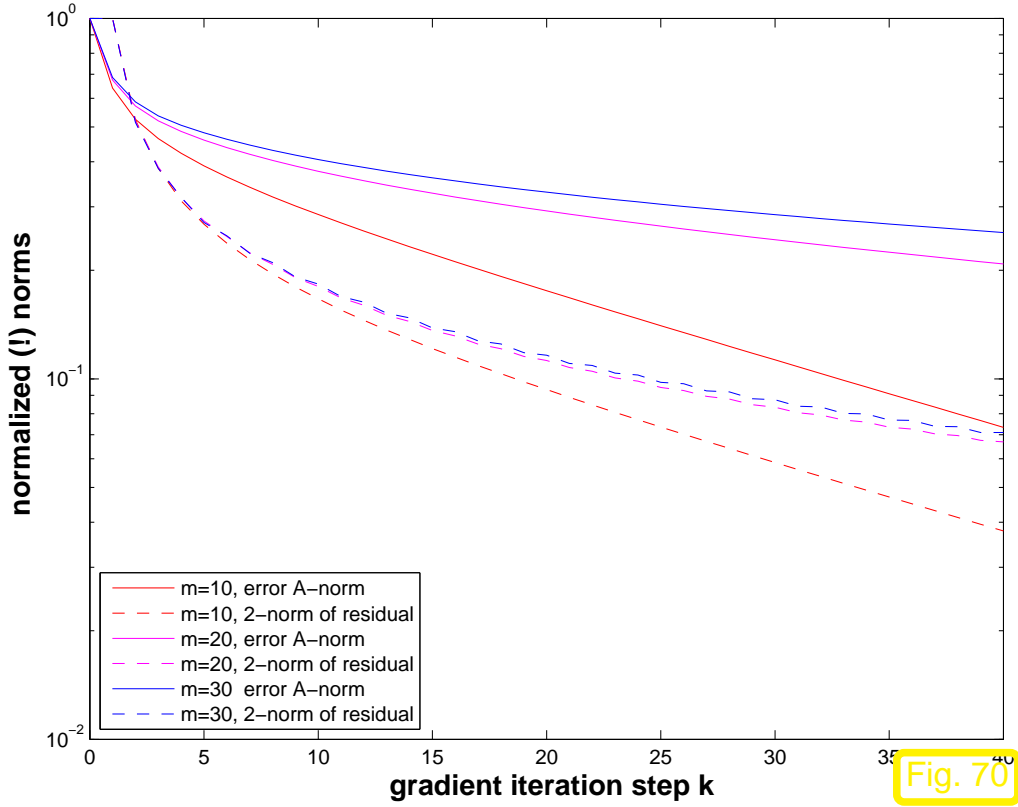


Example 5.2.22 (Convergence of CG as iterative solver).

CG (Code 5.2.17) & gradient method (Code 5.1.11) for LSE with sparse s.p.d. "Poisson matrix"

```
A = gallery('poisson',m); x0 = (1:n)'; b = zeros(n,1);
```





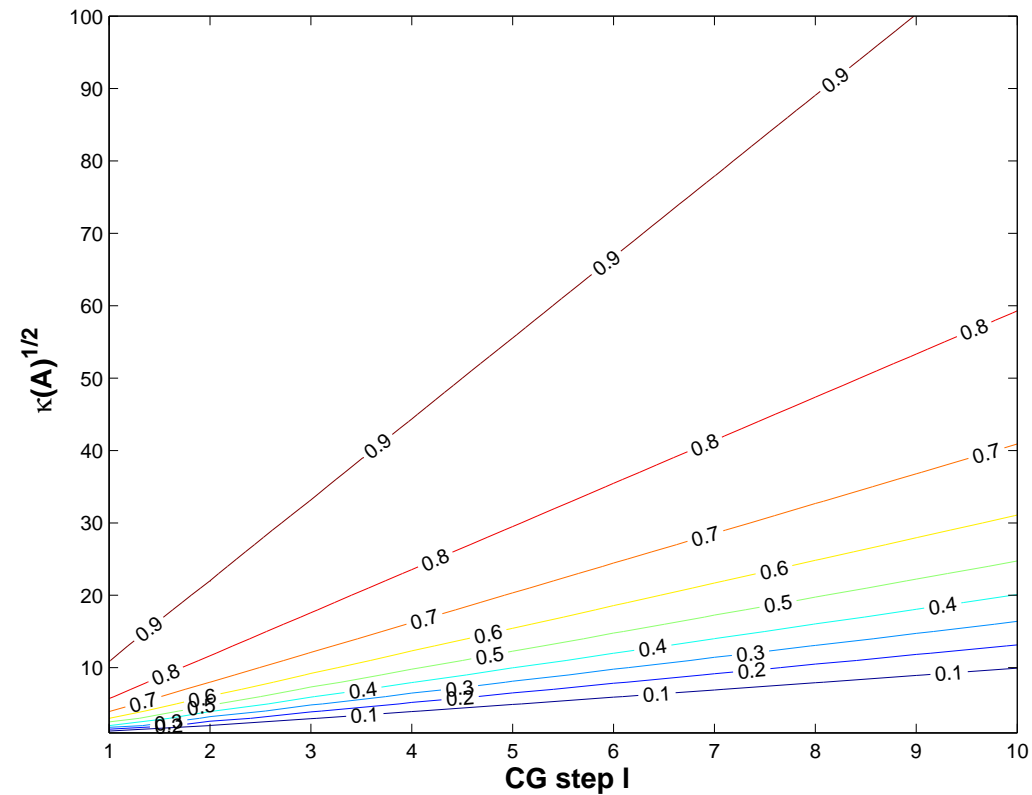
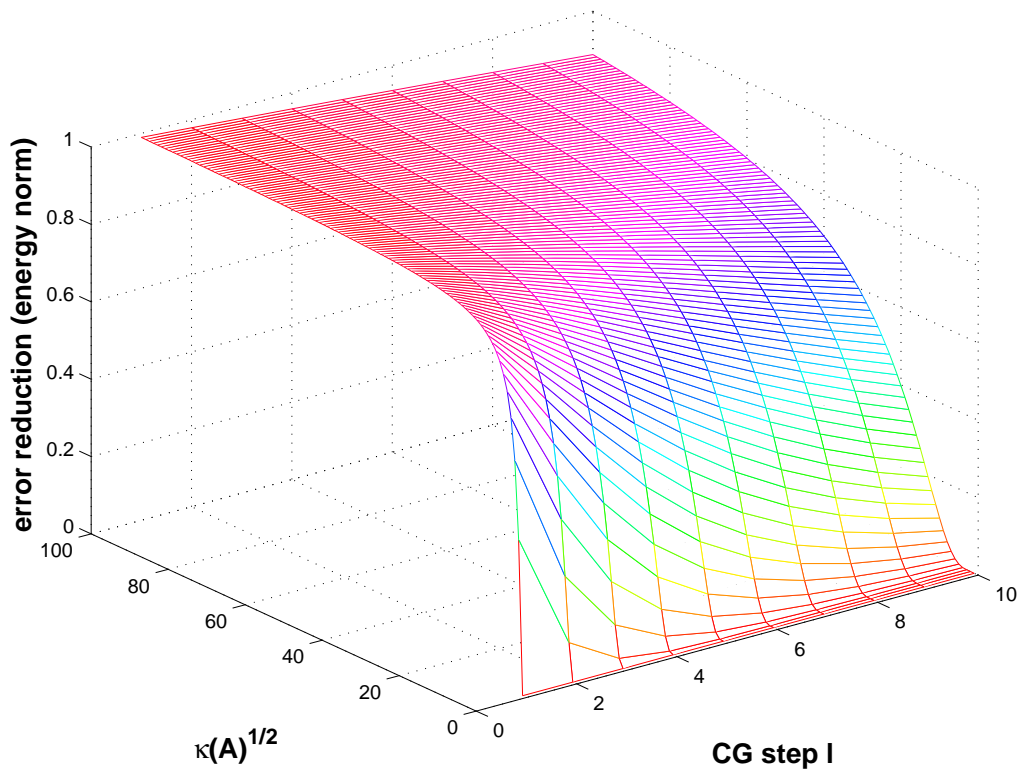
Theorem 5.2.25 (Convergence of CG method).

The iterates of the CG method for solving $\mathbf{Ax} = \mathbf{b}$ (see Code 5.2.17) 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: $\kappa(\mathbf{A}) = \text{spectral condition number of } \mathbf{A}$, $\kappa(\mathbf{A}) = \text{cond}_2(\mathbf{A})$)

Plots of bounds for error reduction (in energy norm) during CG iteration from Thm. 5.2.25:



Example 5.2.27 (Convergence rates for CG method).

Code 5.2.28: 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);

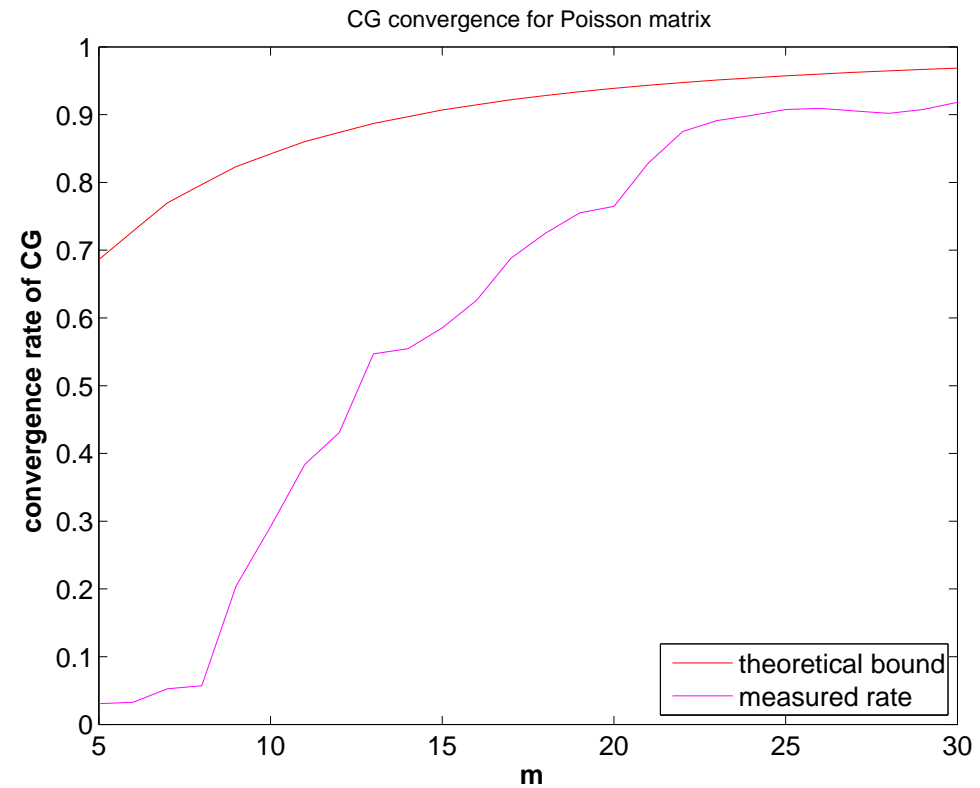
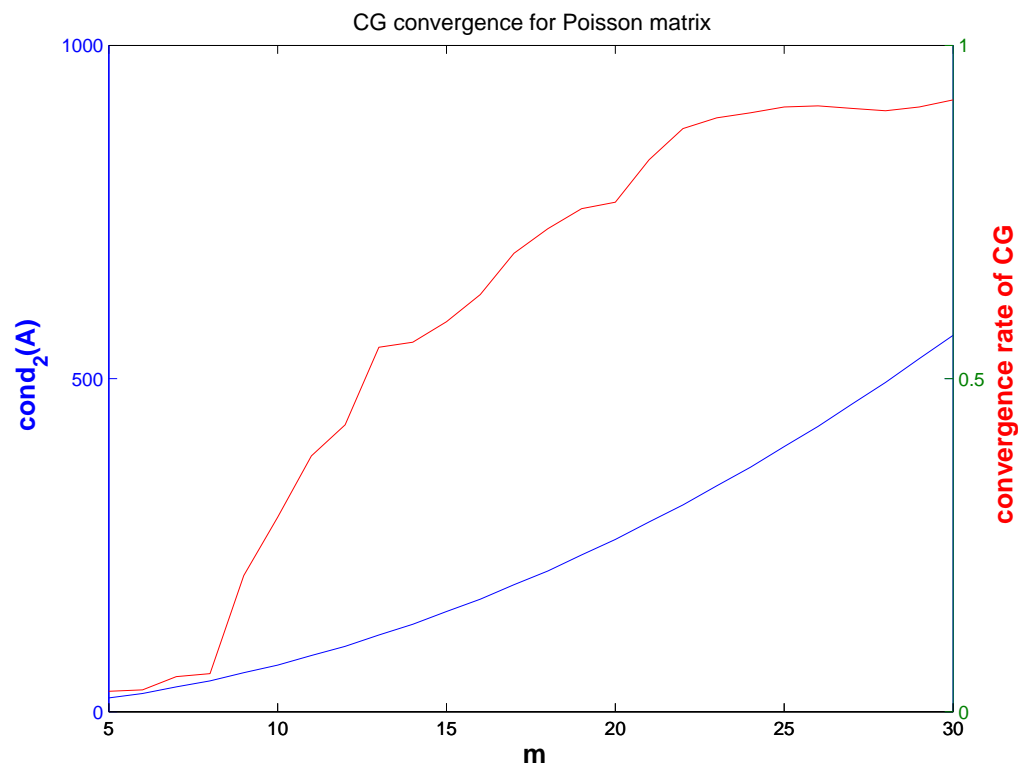
```

Measurement

of

rate of (linear) convergence:

$$\text{rate} \approx \sqrt[10]{\frac{\|\mathbf{r}_{30}\|_2}{\|\mathbf{r}_{20}\|_2}}. \quad (5.2.29)$$



Example 5.2.30 (CG convergence and spectrum). → Ex. 5.1.17

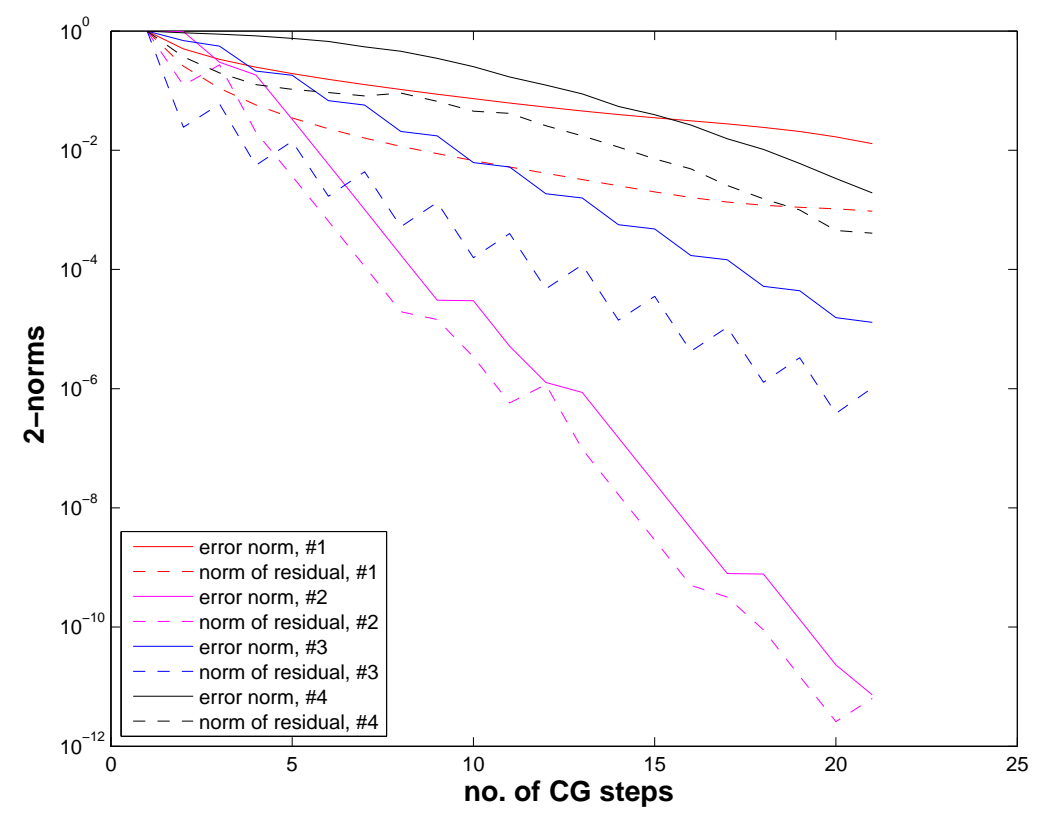
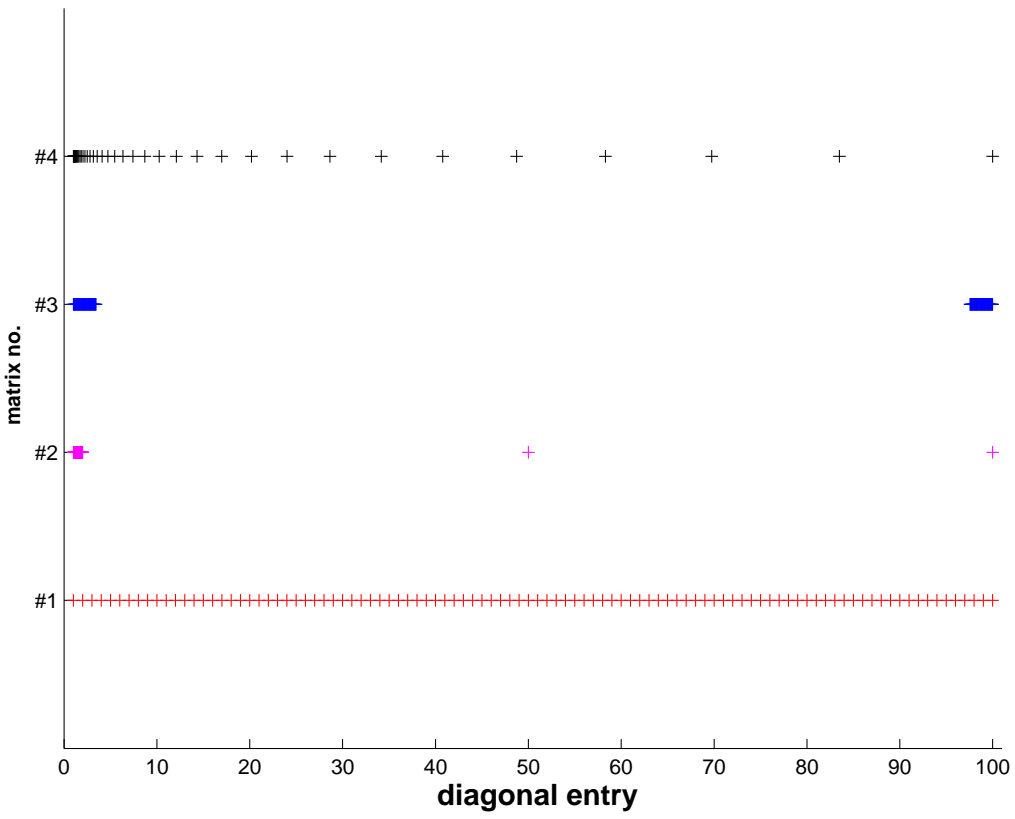
Test matrix #1: $A = \text{diag}(d)$; $d = (1:100)$;

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

$x_0 = \cos((1:n)')$; $b = \text{zeros}(n,1)$;



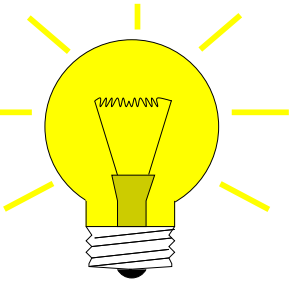
Idea:

Preconditioning

Apply CG method to **transformed linear system**

$$\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}} \quad , \quad \tilde{\mathbf{A}} := \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2} \quad , \quad \tilde{\mathbf{x}} := \mathbf{B}^{1/2}\mathbf{x} \quad , \quad \tilde{\mathbf{b}} := \mathbf{B}^{-1/2}\mathbf{b} \quad , \quad (5.3.1)$$

with “small” $\kappa(\tilde{\mathbf{A}})$, $\mathbf{B} = \mathbf{B}^T \in \mathbb{R}^{N,N}$ s.p.d. $\hat{=}$ **preconditioner**.



Notion 5.3.3 (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×vector multiplication $\mathbf{A}\mathbf{x}$, $\mathbf{x} \in \mathbb{R}^n$.

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 $\mathbf{x}^{(l)} \in \mathbb{R}^n$

$$\mathbf{B}^{1/2}\tilde{\mathbf{r}}_0 := \mathbf{B}^{1/2}\tilde{\mathbf{b}} - \mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{x}}^{(0)};$$

$$\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_1 := \mathbf{B}^{-1}(\mathbf{B}^{1/2}\tilde{\mathbf{r}}_0);$$

for $j = 1$ to l do {

$$\alpha := \frac{(\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \mathbf{B}^{1/2}\tilde{\mathbf{r}}_{j-1}}{(\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j}$$

$$\mathbf{B}^{-1/2}\tilde{\mathbf{x}}^{(j)} := \mathbf{B}^{-1/2}\tilde{\mathbf{x}}^{(j-1)} + \alpha \mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j;$$

$$\mathbf{B}^{1/2}\tilde{\mathbf{r}}_j = \mathbf{B}^{1/2}\tilde{\mathbf{r}}_{j-1} - \alpha \mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j;$$

$$\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_{j+1} = \mathbf{B}^{-1}(\mathbf{B}^{-1/2}\tilde{\mathbf{r}}_j) - \frac{(\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \mathbf{A}\mathbf{B}^{-1}(\mathbf{B}^{1/2}\tilde{\mathbf{r}}_j)}{(\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j)^T \mathbf{A}\mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j} \mathbf{B}^{-1/2}\tilde{\mathbf{p}}_j;$$

}

Algorithm 5.3.5 (Preconditioned CG method (PCG)).

[13, Alg. 13.32], [35, Alg. 10.1]]

Input: initial guess $\mathbf{x} \in \mathbb{R}^n \hat{=} \mathbf{x}^{(0)} \in \mathbb{R}^n$, tolerance $\tau > 0$

Output: approximate solution $\mathbf{x} \hat{=} \mathbf{x}^{(l)}$

$\mathbf{p} := \mathbf{r} := \mathbf{b} - \mathbf{A}\mathbf{x}$; $\mathbf{p} := \mathbf{B}^{-1}\mathbf{r}$; $\mathbf{q} := \mathbf{p}$; $\tau_0 := \mathbf{p}^\top \mathbf{r}$;

for $l = 1$ **to** l_{\max} **do** {

$\beta := \mathbf{r}^\top \mathbf{q}$; $\mathbf{h} := \mathbf{A}\mathbf{p}$; $\alpha := \frac{\beta}{\mathbf{p}^\top \mathbf{h}}$;

$\mathbf{x} := \mathbf{x} + \alpha\mathbf{p}$;

$\mathbf{r} := \mathbf{r} - \alpha\mathbf{h}$;

(5.3.6)

$\mathbf{q} := \mathbf{B}^{-1}\mathbf{r}$; $\beta := \frac{\mathbf{r}^\top \mathbf{q}}{\beta}$;

if $|\mathbf{q}^\top \mathbf{r}| \leq \tau \cdot \tau_0$ **then stop**;

$\mathbf{p} := \mathbf{q} + \beta\mathbf{p}$;

}

Computational effort per step: 1 evaluation $\mathbf{A} \times$ vector,
1 evaluation $\mathbf{B}^{-1} \times$ vector,
3 dot products, 3 **AXPY-operations**

Example 5.3.11 (Tridiagonal preconditioning).

Code 5.3.12: LSE for Ex. 5.3.11

```
1 A =  
    spdiags(repmat([1/n,-1,2+2/n,-1,1/n],n,1),[-n/2,-1,0,1,n/2],n,n);  
2 b = ones(n,1); x0 = ones(n,1); tol = 1.0E-4; maxit = 1000;  
3 evalA = @(x) A*x;  
4  
5 % no preconditioning, see Code 5.3.6  
6 invB = @(x) x; [x,rn] = pcgbase(evalA,b,tol,maxit,invB,x0);  
7  
8 % tridiagonal preconditioning, see Code 5.3.6  
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); %
```

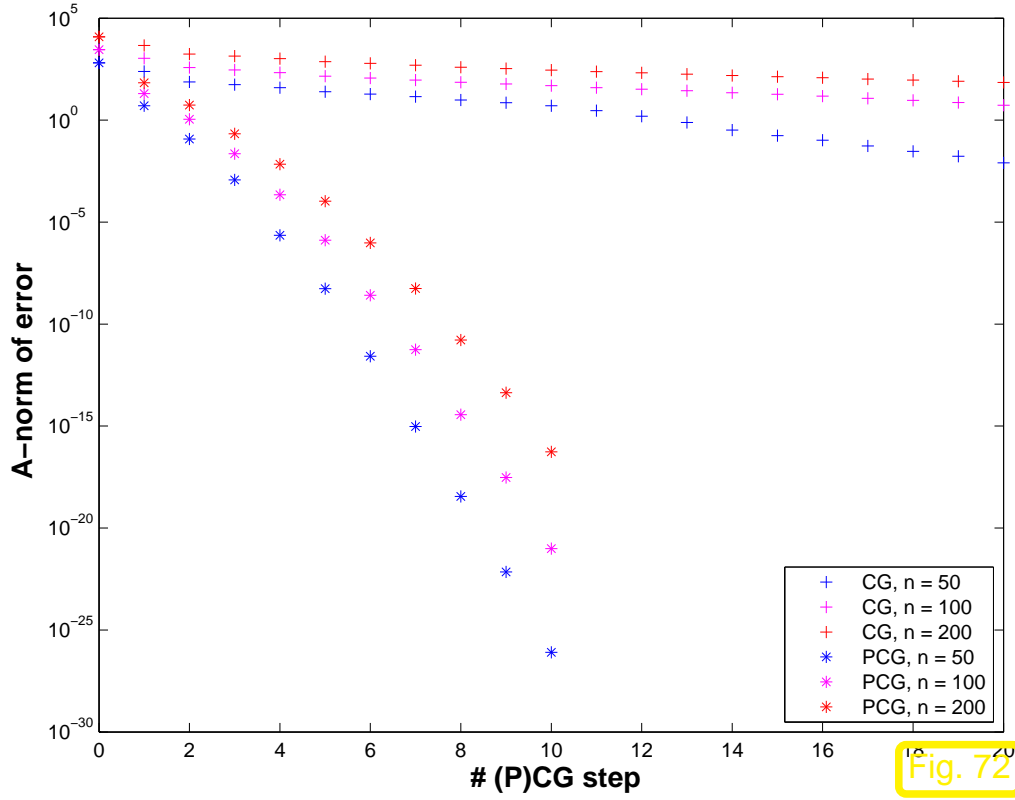


Fig. 72

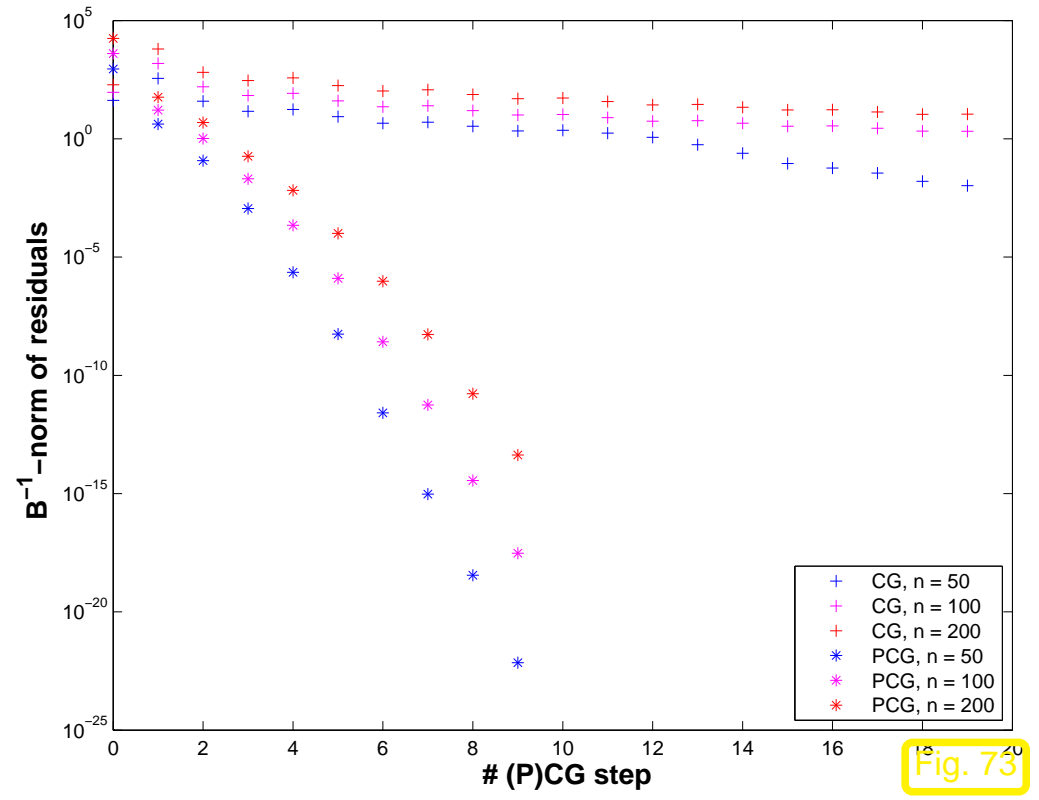


Fig. 73

| 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 |

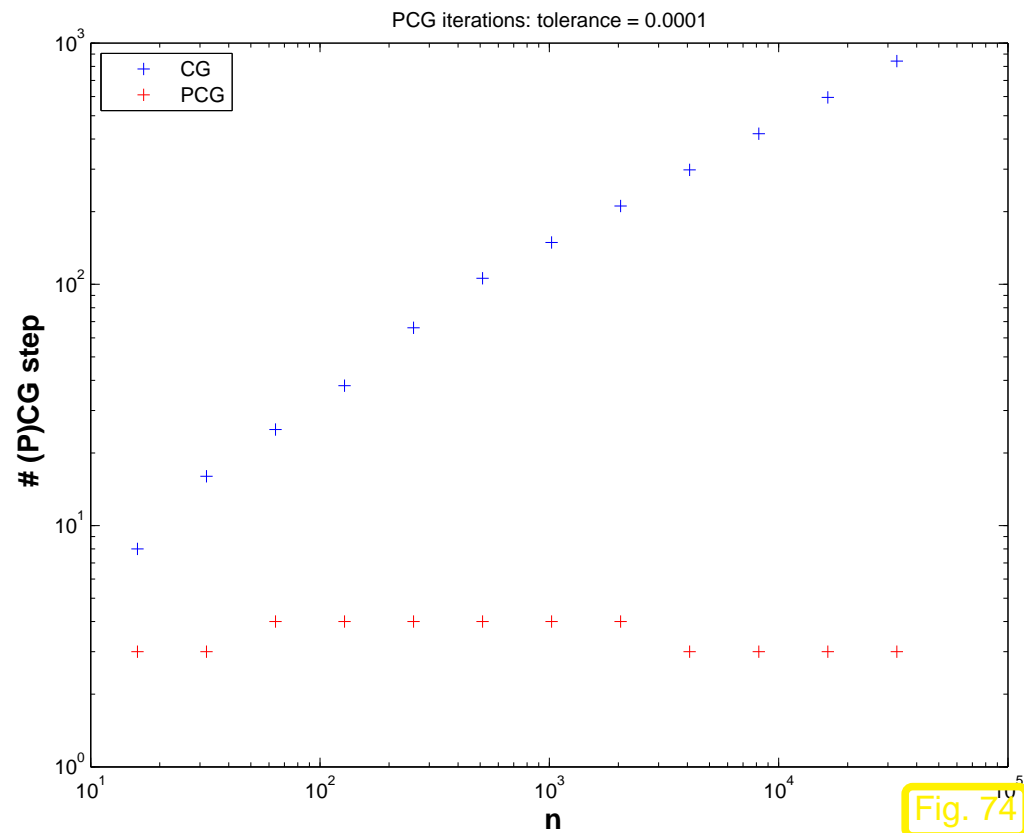


Fig. 74



Remark 5.3.13 (Termination of PCG).

available during PCG iteration (5.3.6)

$$\frac{1}{\kappa(\mathbf{B}^{-1}\mathbf{A})} \frac{\|\mathbf{e}^{(l)}\|_A^2}{\|\mathbf{e}^{(0)}\|_A^2} \leq \frac{(\mathbf{B}^{-1}\mathbf{r}_l)^\top \mathbf{r}_l}{(\mathbf{B}^{-1}\mathbf{r}_0)^\top \mathbf{r}_0} \leq \kappa(\mathbf{B}^{-1}\mathbf{A}) \frac{\|\mathbf{e}^{(l)}\|_A^2}{\|\mathbf{e}^{(0)}\|_A^2} \tag{5.3.14}$$

MATLAB-function: `[x,flag,relr,it,rv] = pcg(A,b,tol,maxit,B,[],x0);`
(A , B may be handles to functions providing Ax and $B^{-1}x$, resp.)

Remark 5.3.15 (Termination criterion in MATLAB-`pcg`). → [51, Sect. 4.6]

Implementation (skeleton) of MATLAB built-in `pcg`:


```
function x = pcg(Afun,b,tol,maxit,Binvfun,x0)
x = x0; r = b - feval(Afun,x); rho = 1;
for i = 1 : maxit
    y = feval(Binvfun,r);
    rho1 = rho; rho = r' * y;
    if (i == 1)
        p = y;
    else
        beta = rho / rho1;
        p = y + beta * p;
    end
    q = feval(Afun,p);
    alpha = rho / (p' * q);
    x = x + alpha * p;
    if (norm(b - evalf(Afun,x)) <= tol*b*norm(b)), return; end
    r = r - alpha * q;
end
```

Dubious termination criterion !

5.4.1 Minimal residual methods

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 $\mathbf{A} \in \mathbb{R}^{n,n}$:

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

5.4.2 Iterations with short recursions [51, Sect. 4.5]

Iterative methods for *general* regular system matrix \mathbf{A} :

MATLAB-function: • `[x,flag,r,it,rv] = bicgstab(A,b,tol,maxit,B,[],x0)`
 • `[...] = bicgstab(Afun,b,tol,maxit,Binvfun,[],x0);`

Computational costs : $2 \mathbf{A} \times \text{vector}$, $2 \mathbf{B}^{-1} \times \text{vector}$, 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 : $2 \mathbf{A} \times \text{vector}$, $2 \mathbf{B}^{-1} \times \text{vector}$, 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 5.4.3 (Failure of Krylov iterative solvers).

$$\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 & \cdots & & \cdots & 0 \\ 0 & 0 & 1 & 0 & & & \vdots \\ \vdots & \cdots & \cdots & \cdots & & & \\ & & & & \cdots & \vdots & \\ \vdots & & & \cdots & \cdots & 0 & \\ 0 & & & & 0 & 1 & \\ 1 & 0 & \cdots & & \cdots & 0 & \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \quad \blacktriangleright \quad \mathbf{x} = \mathbf{e}_1.$$

$$\mathbf{x}^{(0)} = 0 \quad \blacktriangleright \quad \mathbf{r}_0 = \mathbf{e}_n \quad \blacktriangleright \quad \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0) = \text{Span} \{ \mathbf{e}_n, \mathbf{e}_{n-1}, \dots, \mathbf{e}_{n-l+1} \}$$

$$\blacktriangleright \quad \min \{ \|\mathbf{y} - \mathbf{x}\|_2 : \mathbf{y} \in \mathcal{K}_l(\mathbf{A}, \mathbf{r}_0) \} = \begin{cases} 1 & , \text{ if } l \leq n, \\ 0 & , \text{ for } l = n. \end{cases}$$

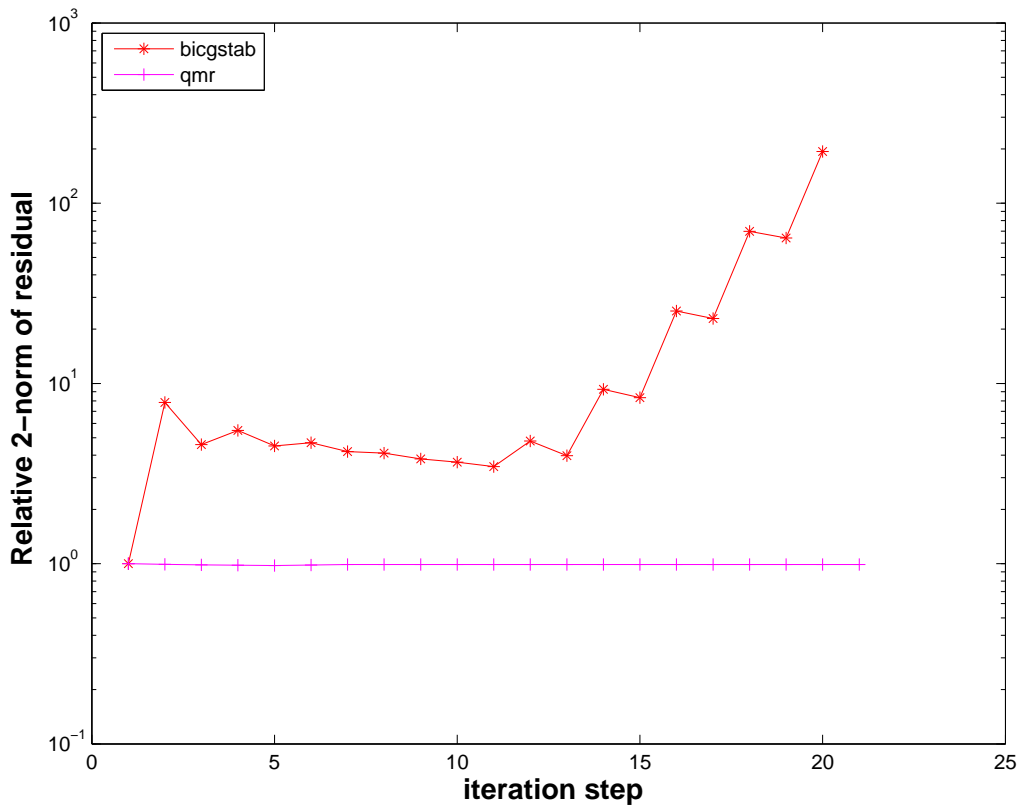


TRY & PRAY

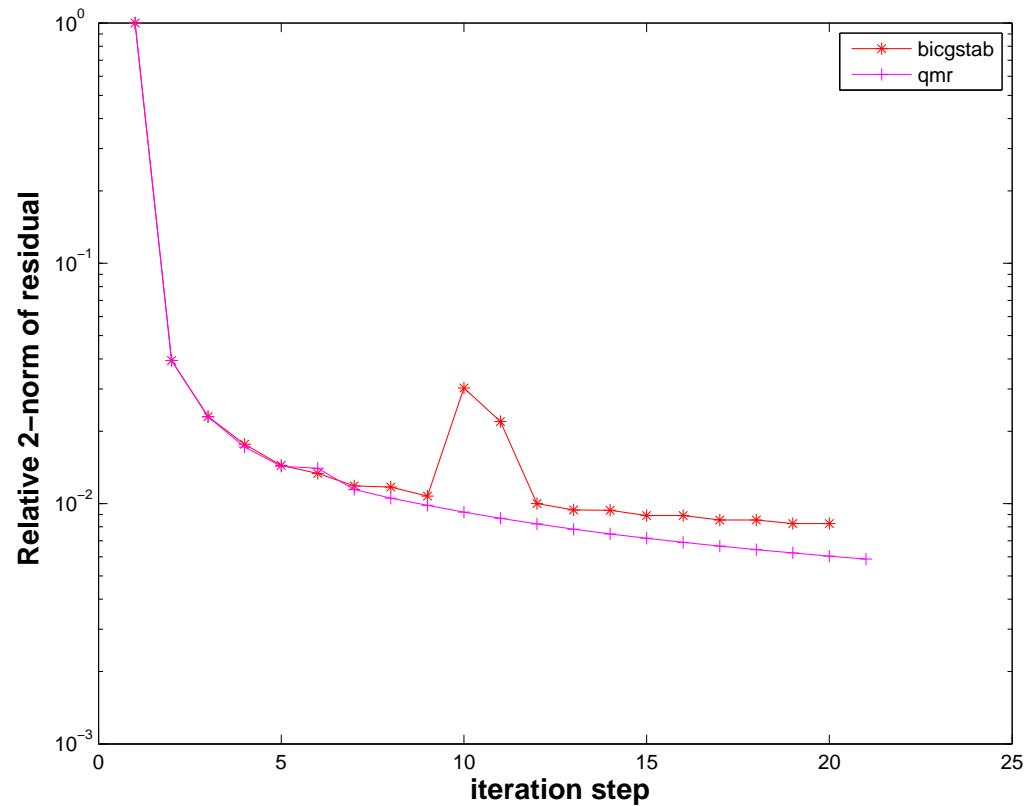
Example 5.4.4 (Convergence of Krylov subspace methods for non-symmetric system matrix).

```
A = gallery('tridiag',-0.5*ones(n-1,1),2*ones(n,1),-1.5*ones(n-1,1));
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**



6

Eigenvalues

Example 6.0.1 (Resonances of linear electric circuits).

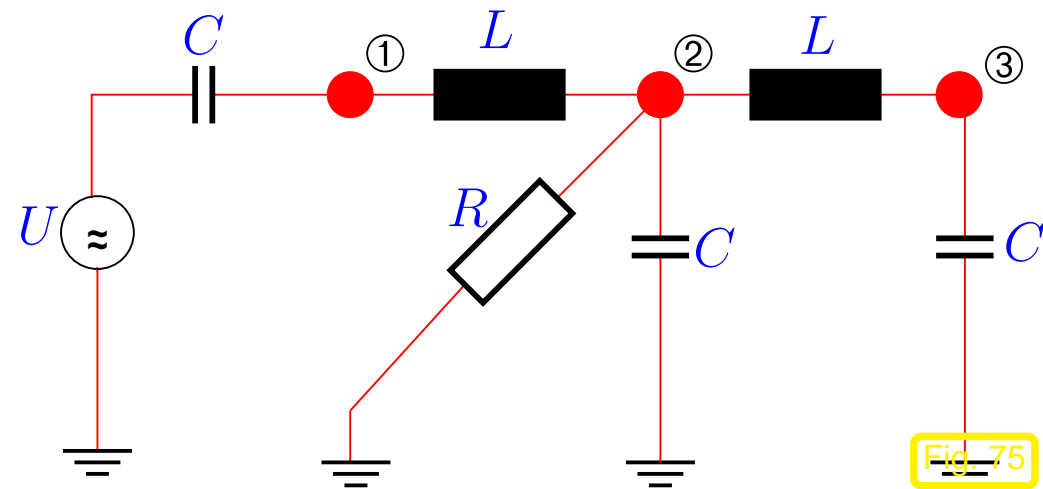


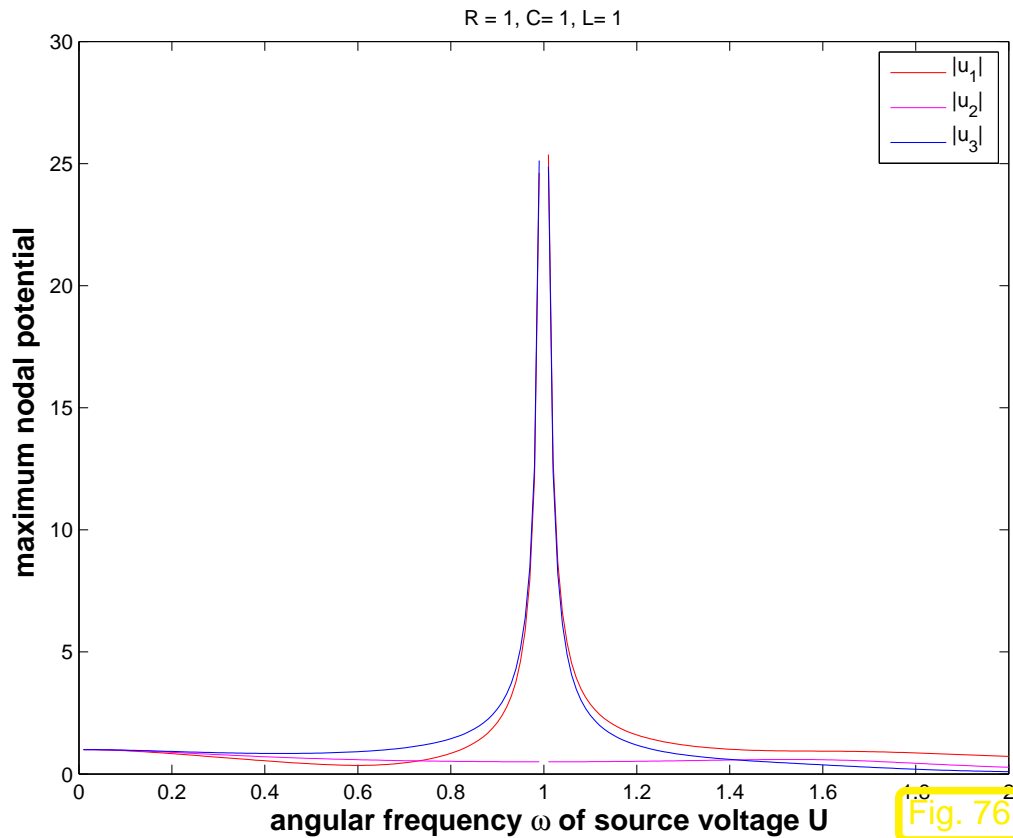
Fig. 75

➤ system matrix from nodal analysis at angular frequency $\omega > 0$:

$$\mathbf{A} = \begin{pmatrix} \omega C + \frac{1}{\omega L} & -\frac{1}{\omega L} & 0 \\ -\frac{1}{\omega L} & \omega C + \frac{1}{R} + \frac{2}{\omega L} & -\frac{1}{\omega L} \\ 0 & -\frac{1}{\omega L} & \omega C + \frac{1}{\omega L} \end{pmatrix}$$

$$= \begin{pmatrix} 0 & 0 & 0 \\ 0 & \frac{1}{R} & 0 \\ 0 & 0 & 0 \end{pmatrix} + \omega \begin{pmatrix} C & 0 & 0 \\ 0 & C & 0 \\ 0 & 0 & C \end{pmatrix} + 1/\omega \begin{pmatrix} \frac{1}{L} & -\frac{1}{L} & 0 \\ -\frac{1}{L} & \frac{2}{L} & -\frac{1}{L} \\ 0 & -\frac{1}{L} & \frac{1}{L} \end{pmatrix}.$$

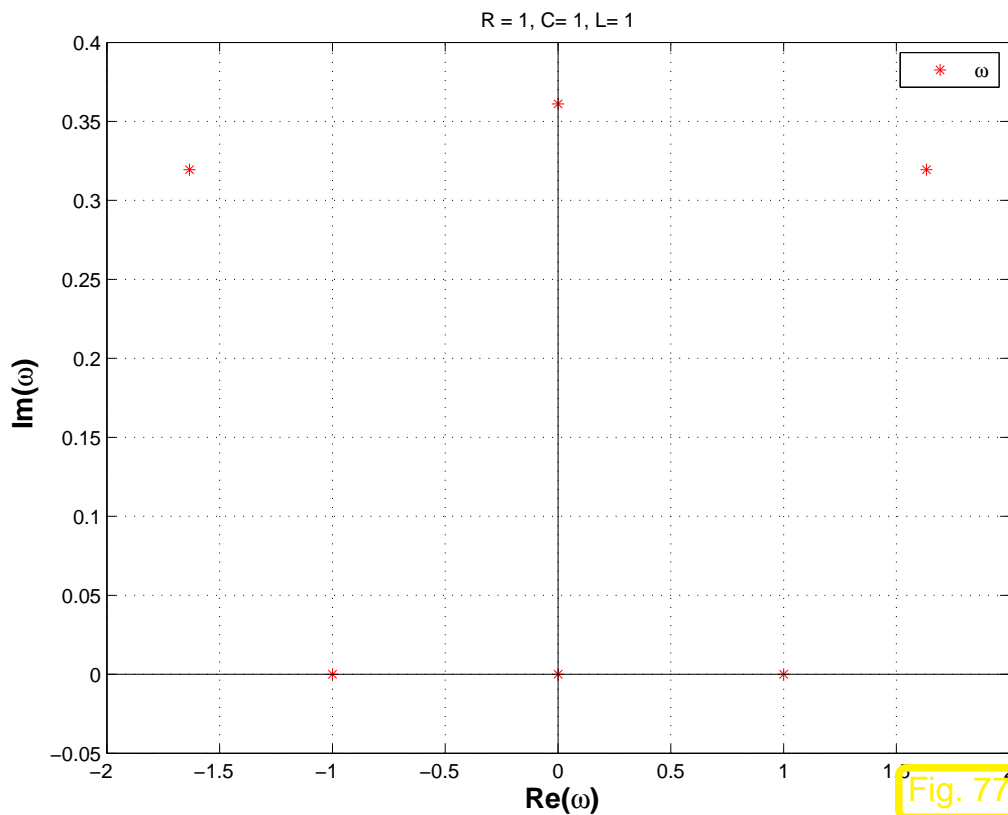
$$\mathbf{A}(\omega) := \mathbf{W} + i\omega\mathbf{C} - i\omega^{-1}\mathbf{S} \quad , \quad \mathbf{W}, \mathbf{C}, \mathbf{S} \in \mathbb{R}^{n,n} \text{ symmetric .} \quad (6.0.2)$$



◁ plot of $|u_i(U)|$, $i = 1, 2, 3$ for $R = L = C = 1$
(scaled model)

Blow-up of some nodal potentials for certain ω !

resonant frequencies = $\omega \in \{\omega \in \mathbb{R} : \mathbf{A}(\omega) \text{ singular}\}$



◁ resonant frequencies for circuit from Fig. 6.0.2
(including decaying modes with $\text{Im}(\omega) > 0$)



Example 6.0.7 (Analytic solution of homogeneous linear ordinary differential equations). → [63, Remark 5.6.1], [27, Sect. 10.1], [48, Sect. 8.1], [13, Ex. 7.3]

Autonomous homogeneous linear ordinary differential equation (ODE):

$$\dot{\mathbf{y}} = \mathbf{A}\mathbf{y} \quad , \quad \mathbf{A} \in \mathbb{C}^{n,n} . \quad (6.0.8)$$

$$\mathbf{A} = \mathbf{S} \underbrace{\begin{pmatrix} \lambda_1 & & \\ & \cdots & \\ & & \lambda_n \end{pmatrix}}{=: \mathbf{D}} \mathbf{S}^{-1}, \quad \mathbf{S} \in \mathbb{C}^{n,n} \text{ regular} \implies \left(\dot{\mathbf{y}} = \mathbf{A}\mathbf{y} \quad \begin{array}{c} \mathbf{z} = \mathbf{S}^{-1}\mathbf{y} \\ \longleftrightarrow \end{array} \quad \dot{\mathbf{z}} = \mathbf{D}\mathbf{z} \right).$$



6.1 Theory of eigenvalue problems [48, Ch. 7], [27, Ch. 9], [51, Sect. 1.7]

Eigenvalue

- problems:*
- ❶ Given $\mathbf{A} \in \mathbb{K}^{n,n}$ find **all eigenvalues** (= spectrum of \mathbf{A}).
 - (EVPs) ❷ Given $\mathbf{A} \in \mathbb{K}^{n,n}$ find $\sigma(\mathbf{A})$ plus **all eigenvectors**.
 - ❸ 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 $\mathbf{x} \neq 0$, $\lambda \in \mathbb{C}$

$$\mathbf{Ax} = \lambda \mathbf{Bx} \Leftrightarrow \mathbf{B}^{-1} \mathbf{Ax} = \lambda \mathbf{x} . \quad (6.1.10)$$

$\mathbf{x} \hat{=}$ generalized eigenvector, $\lambda \hat{=}$ generalized eigenvalue

6.2 “Direct” Eigensolvers

Purpose: solution of eigenvalue problems ❶, ❷ for **dense** matrices “up to machine precision”

MATLAB-function:

`eig`

`d = eig(A)` : computes spectrum $\sigma(\mathbf{A}) = \{d_1, \dots, d_n\}$ of $\mathbf{A} \in \mathbb{C}^{n,n}$

`[V,D] = eig(A)` : computes $\mathbf{V} \in \mathbb{C}^{n,n}$, *diagonal* $\mathbf{D} \in \mathbb{C}^{n,n}$ such that $\mathbf{AV} = \mathbf{VD}$

Similar functionality for generalized EVP $\mathbf{Ax} = \lambda \mathbf{Bx}$, $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n,n}$

$d = \text{eig}(A,B)$: computes all generalized eigenvalues
 $[V,D] = \text{eig}(A,B)$: computes $V \in \mathbb{C}^{n,n}$, diagonal $D \in \mathbb{C}^{n,n}$ such that $AV = BVD$

Remark 6.2.4 (Computational effort for eigenvalue computations).

Computational effort (#elementary operations) for $\text{eig}()$:

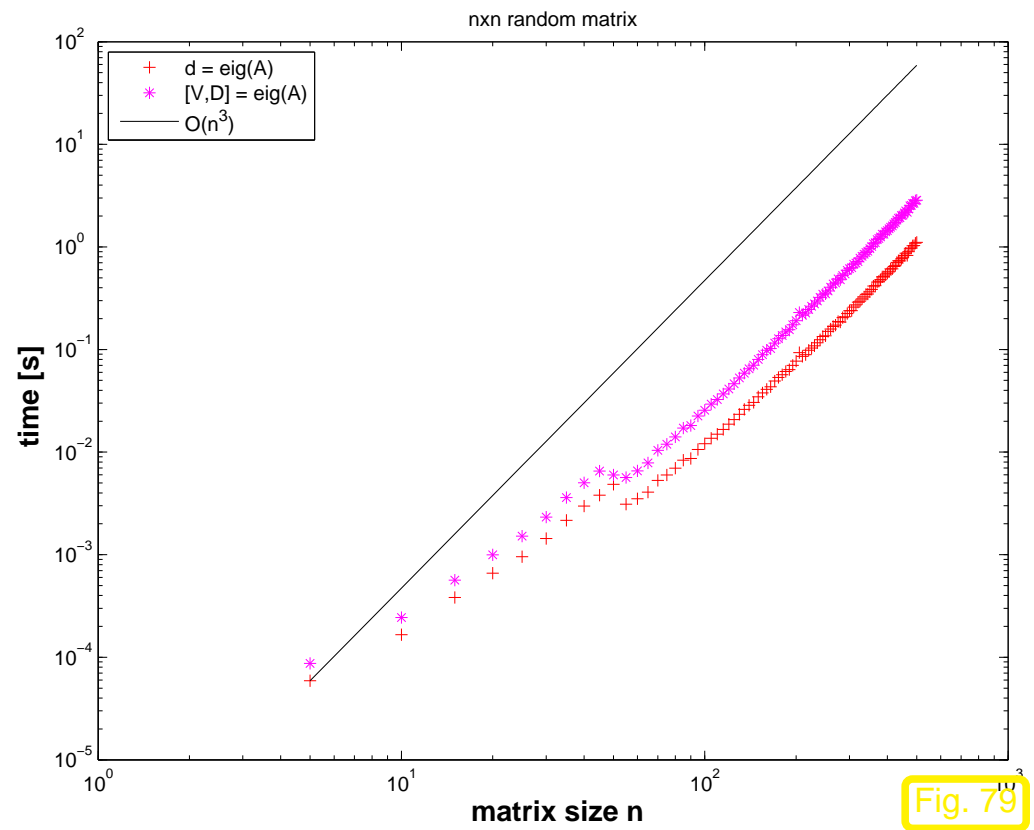
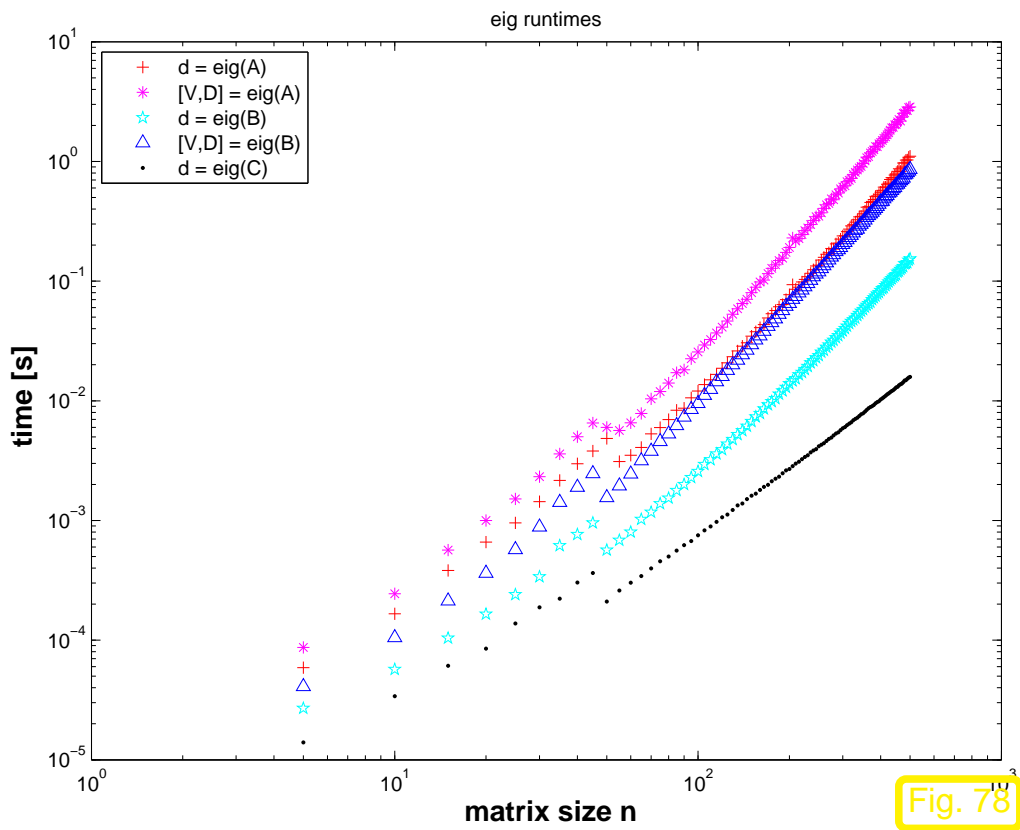
| | | |
|---|--------------------------------|-------------|
| eigenvalues & eigenvectors of $A \in \mathbb{K}^{n,n}$ | $\sim 25n^3 + O(n^2)$ | } $O(n^3)!$ |
| only eigenvalues of $A \in \mathbb{K}^{n,n}$ | $\sim 10n^3 + O(n^2)$ | |
| eigenvalues and eigenvectors $A = A^H \in \mathbb{K}^{n,n}$ | $\sim 9n^3 + O(n^2)$ | |
| only eigenvalues of $A = A^H \in \mathbb{K}^{n,n}$ | $\sim \frac{4}{3}n^3 + O(n^2)$ | |
| only eigenvalues of tridiagonal $A = A^H \in \mathbb{K}^{n,n}$ | $\sim 30n^2 + O(n)$ | |

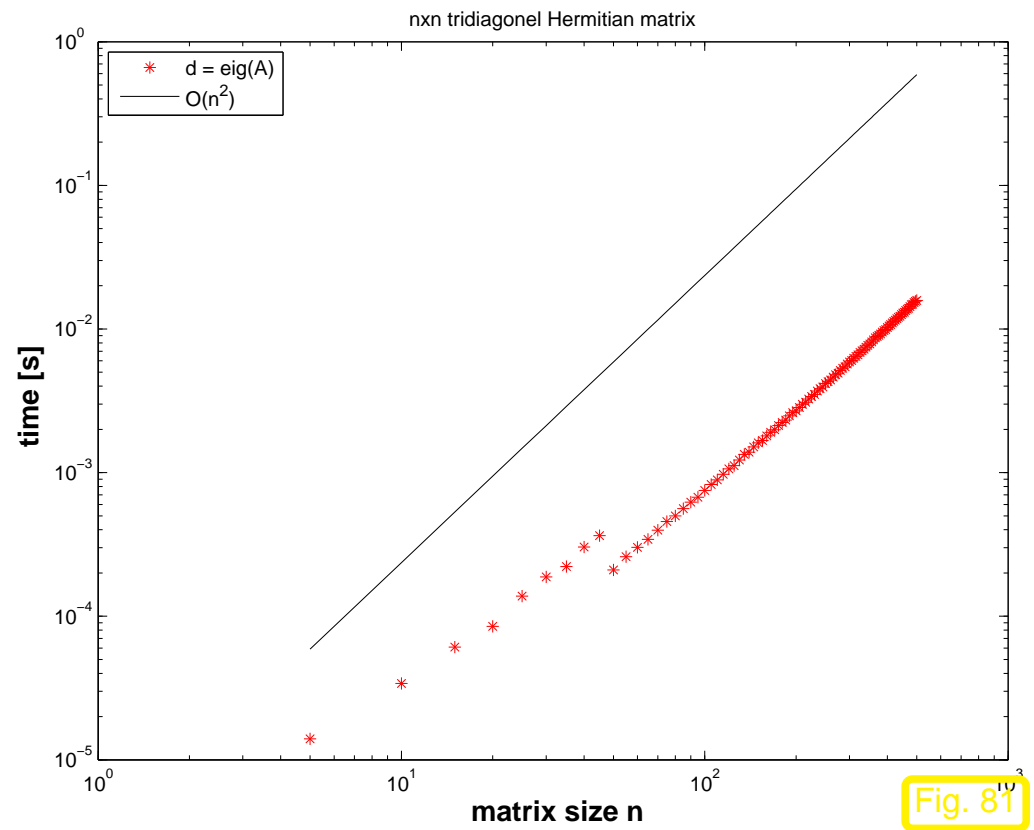
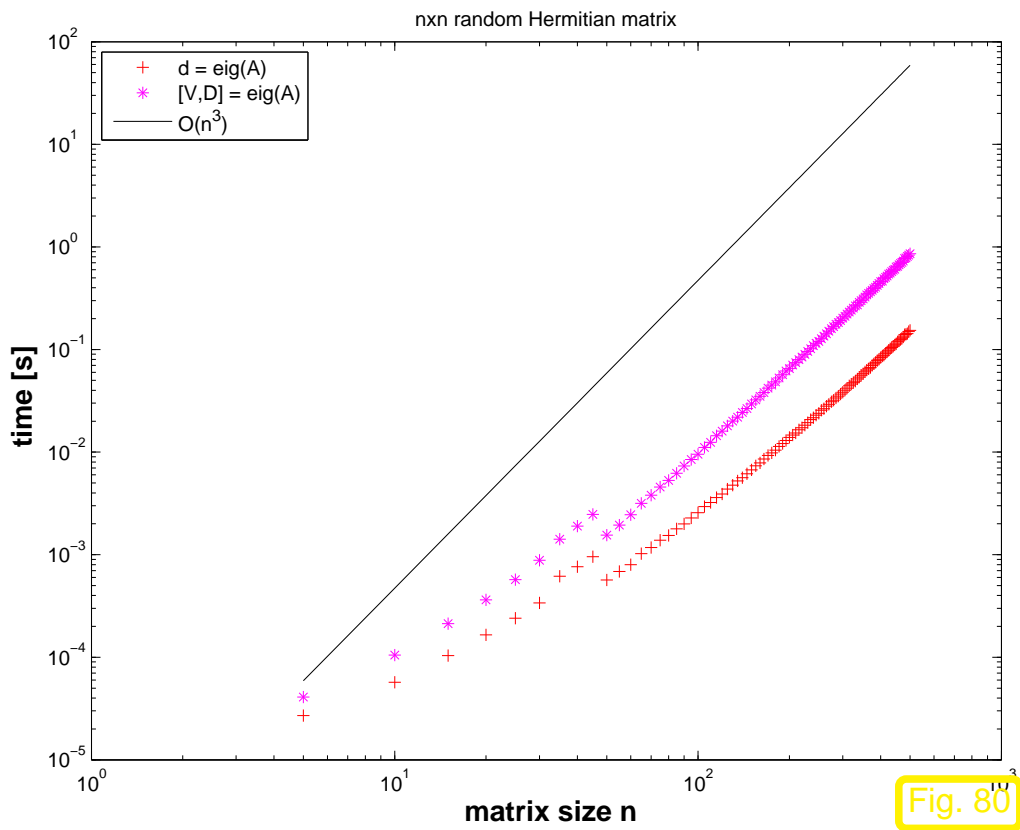


Example 6.2.5 (Runtimes of eig).

```
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. 2.7.9) matrix
- **C** s.p.d. *tridiagonal* matrix





6.3 Power Methods

6.3.1 Direct power method [13, Sect. 7.5], [51, Sect. 5.3.1], [51, Sect. 5.3]

Example 6.3.1 ((Simplified) Page rank algorithm). → [42]

Model: **Random surfer** visits a web page, stays there for fixed time Δt , and then

- ❶ either follows each of ℓ links on a page with probability $1/\ell$.
- ❷ or resumes surfing at a randomly (with equal probability) selected page

Option ❷ is chosen with probability d , $0 \leq d \leq 1$, option ❶ with probability $1 - d$.

Question: Fraction of time spent by random surfer on i -th page (= **page rank** $x_i \in [0, 1]$)

Method: Stochastic simulation

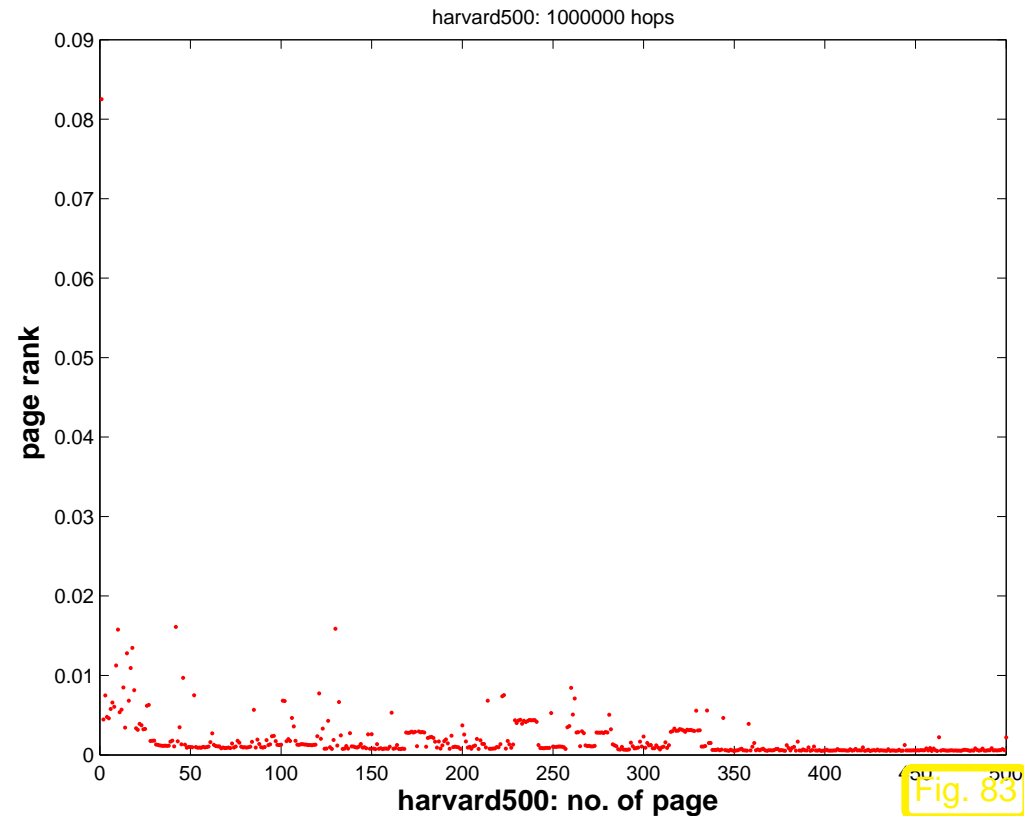
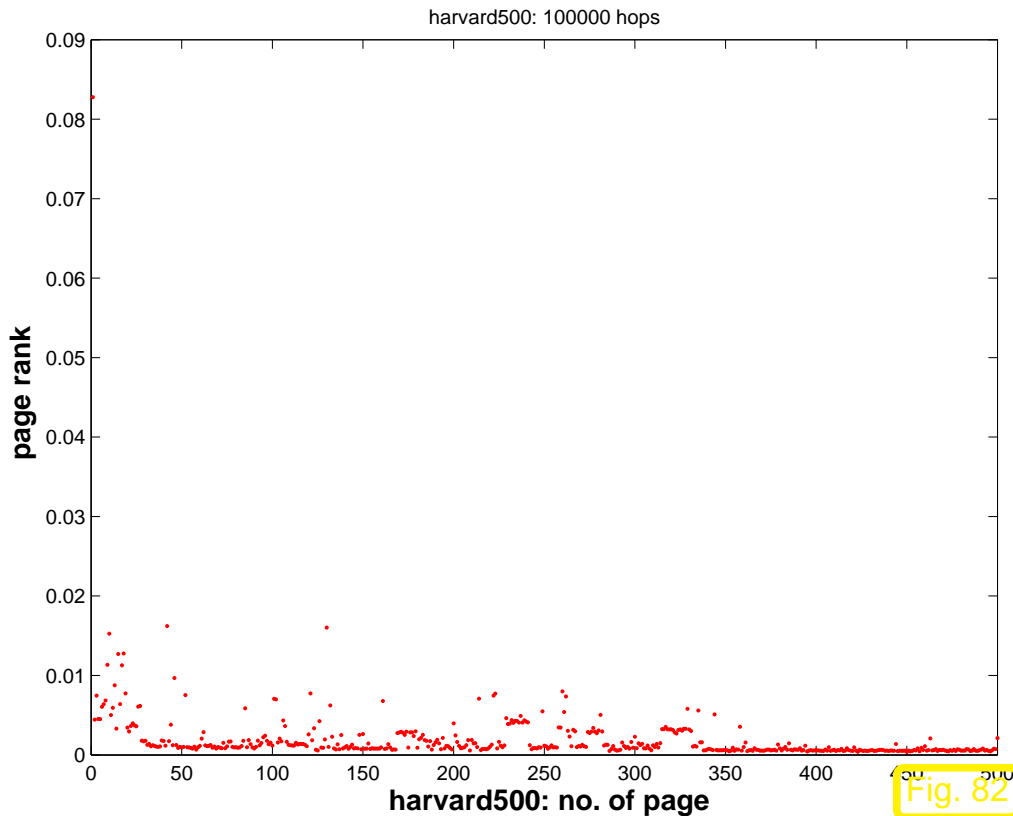


Code 6.3.2: stochastic page rank simulation

```

1 function prstochsim(Nhops)
2 % Load web graph data stored in N x 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. '); axis([0 N+1 0 0.1]);
19    xlabel('\bf harvard500: no. of page','fontsize',14);
20    ylabel('\bf page rank','fontsize',14);
21    title(sprintf('\bf page rank, harvard500: %d
22    hops}',n),'fontsize',14);
23    drawnow;
end
```

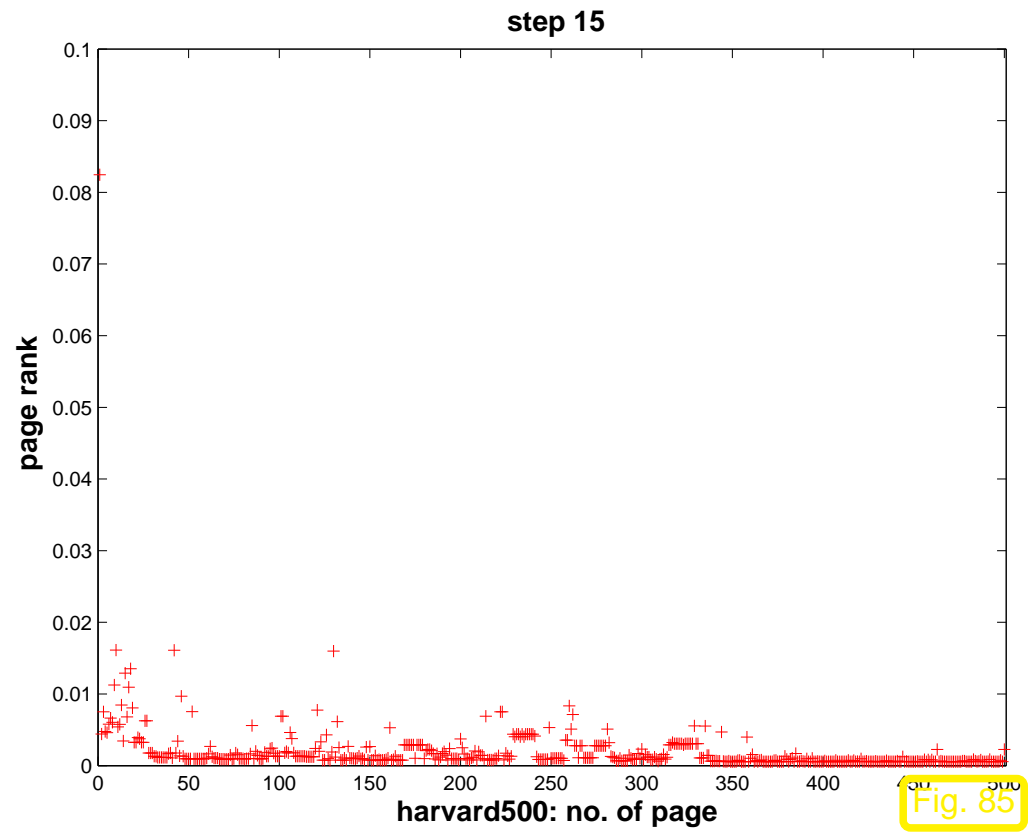
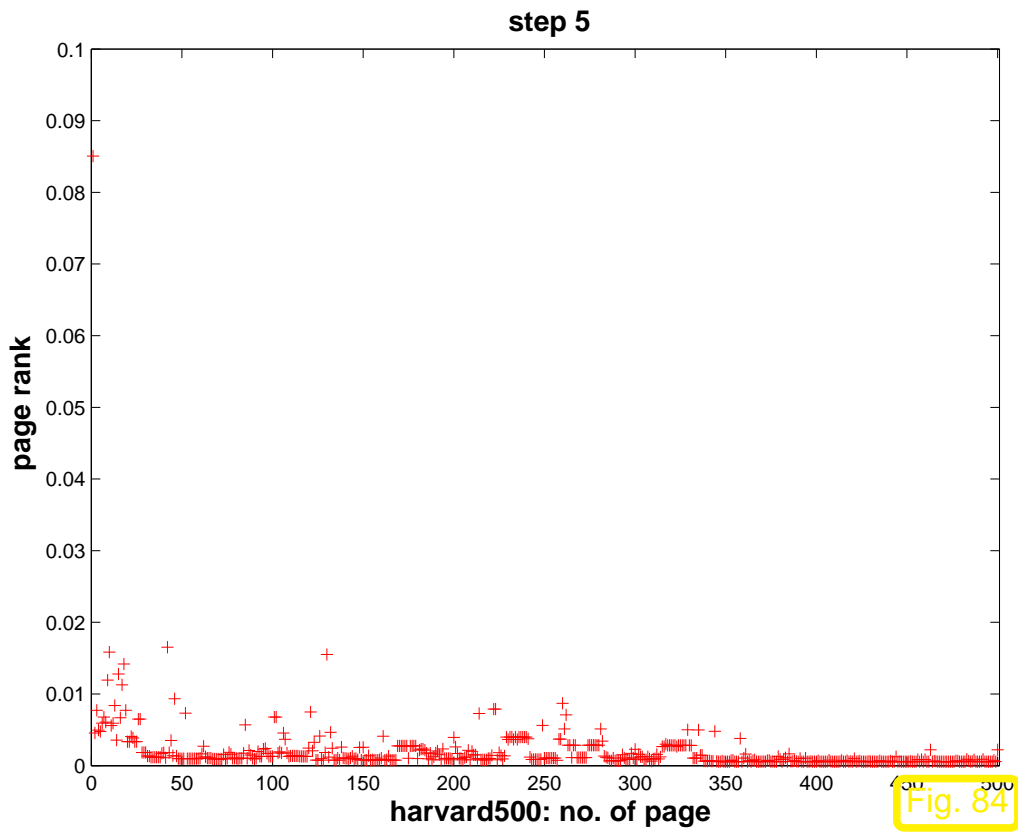


- Numbering of pages $1, \dots, N$, $\ell_i \hat{=}$ number of links from page i
- $N \times N$ -matrix of transition probabilities page $j \rightarrow$ page i : $\mathbf{A} = (a_{ij})_{i,j=1}^N \in \mathbb{R}^{N,N}$

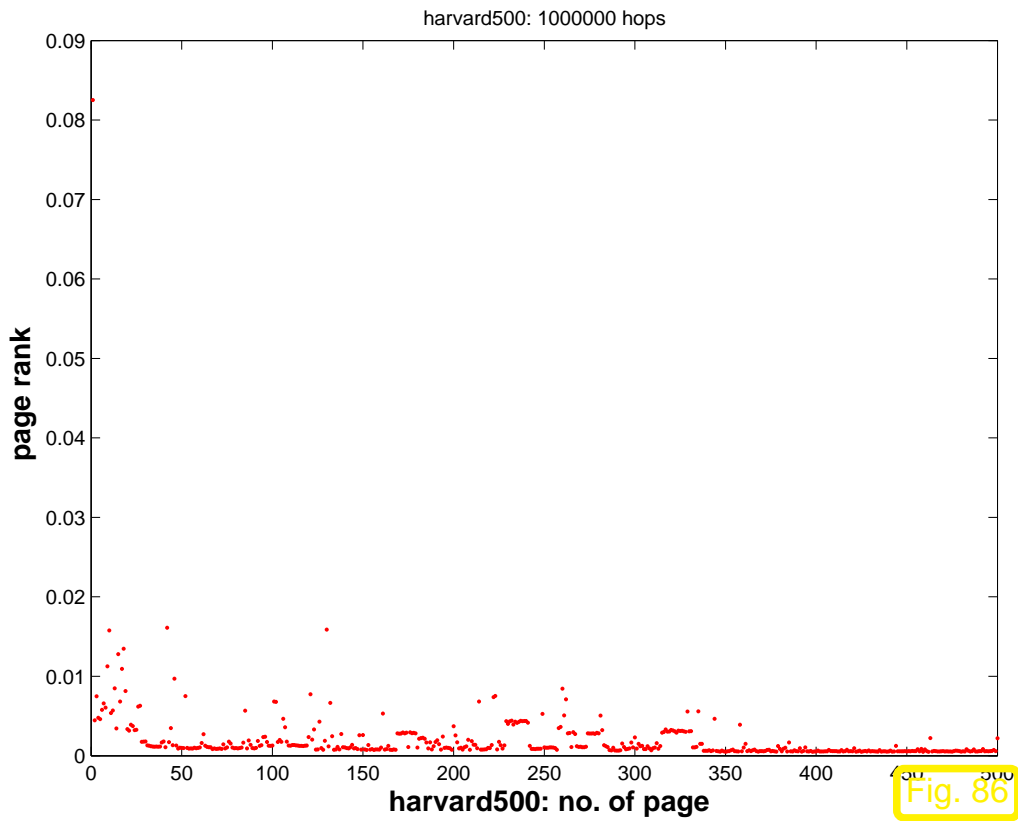
$a_{ij} \in [0, 1] \hat{=}$ probability to jump from page j to page i .

$$\Rightarrow \sum_{i=1}^N a_{ij} = 1 . \tag{6.3.3}$$

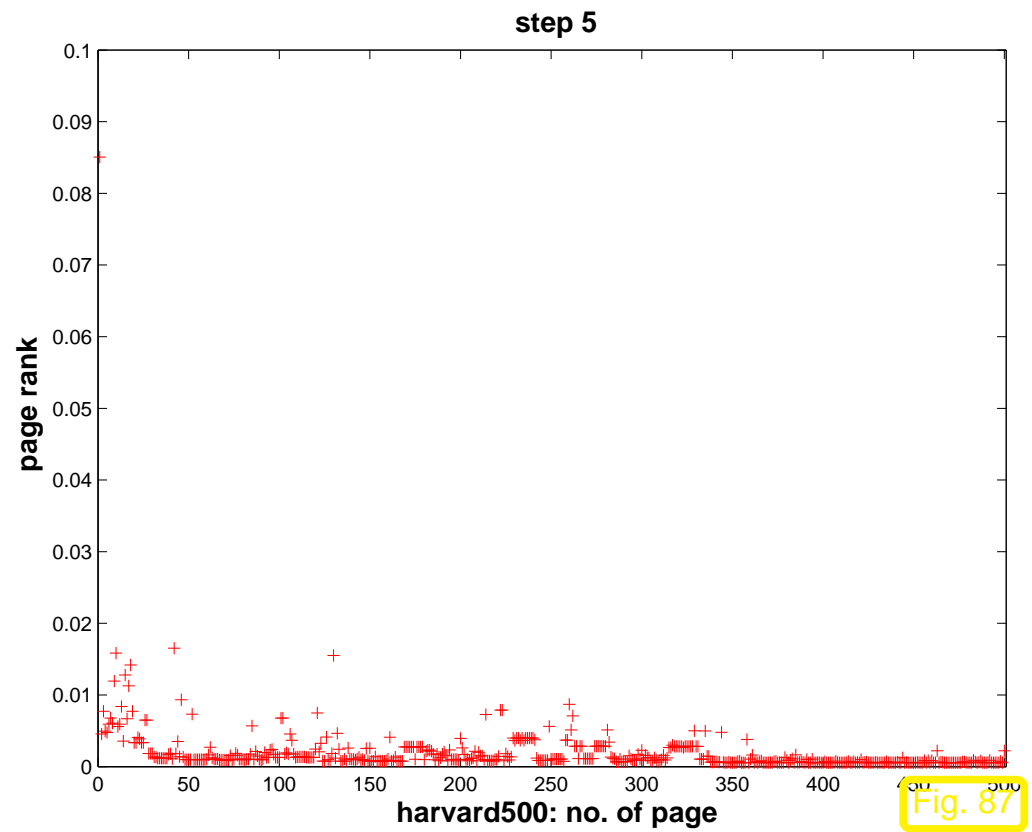

```
1 function prpowitsim(d,Nsteps)
2 % MATLAB way of specifying Default arguments
3 if (nargin < 2), Nsteps = 5; end
4 if (nargin < 1), d = 0.15; end
5 % load connectivity matrix and build transition matrix
6 load harvard500.mat; A = prbuildA(G,d);
7 N = size(A,1); x = ones(N,1)/N;
8
9 figure('position',[0 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:



Single surfer stochastic simulation



Power method, Code 6.3.6

Code 6.3.9: computing page rank vector \mathbf{r} via `eig`

```

1 function prevp
2 load harvard500.mat; d = 0.15;
3 [V,D] = eig(prbuildA(G,d));
4
5 figure; r = V(:,1); N = length(r);
6 plot(1:N,r/sum(r),'m. '); axis([0 N+1 0 0.1]);
7 xlabel('{\bf harvard500: no. of
   page}','fontsize',14);
8 ylabel('{\bf entry of r-vector}','fontsize',14);
9 title('harvard 500: Perron-Frobenius vector');
10 print -depsc2 '../PICTURES/prevp.eps';

```

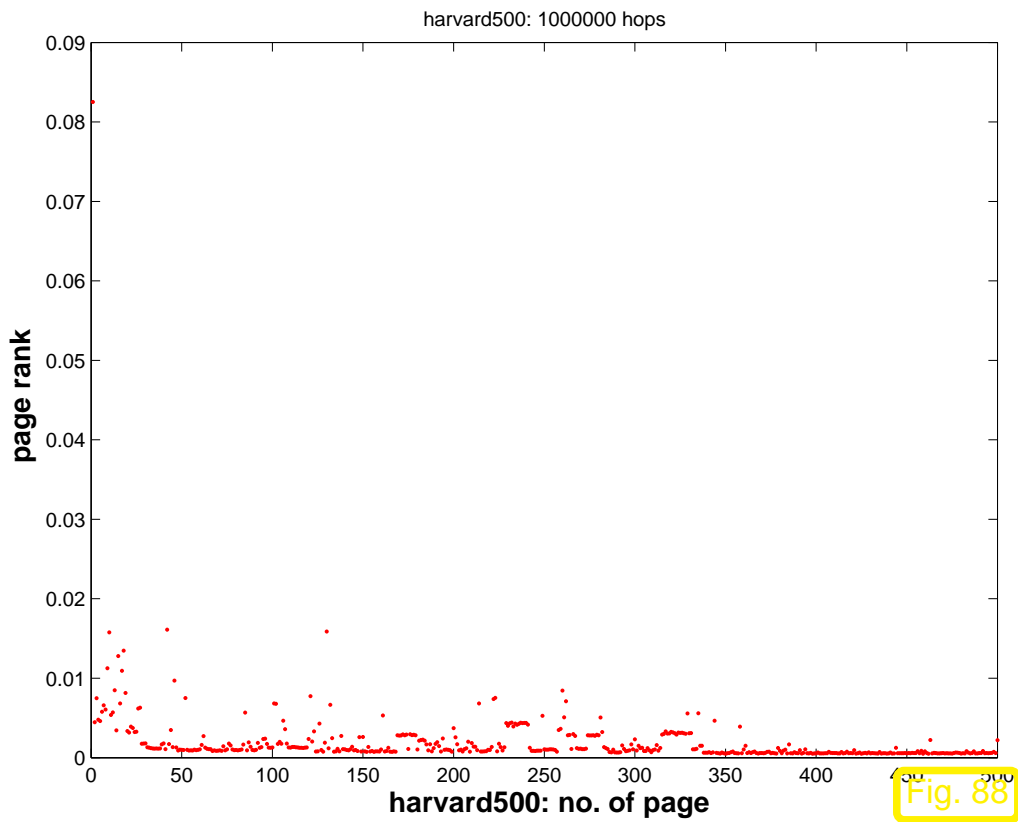
Plot of entries of
unique vector $\mathbf{r} \in \mathbb{R}^N$ with

$$0 \leq (\mathbf{r})_i \leq 1,$$

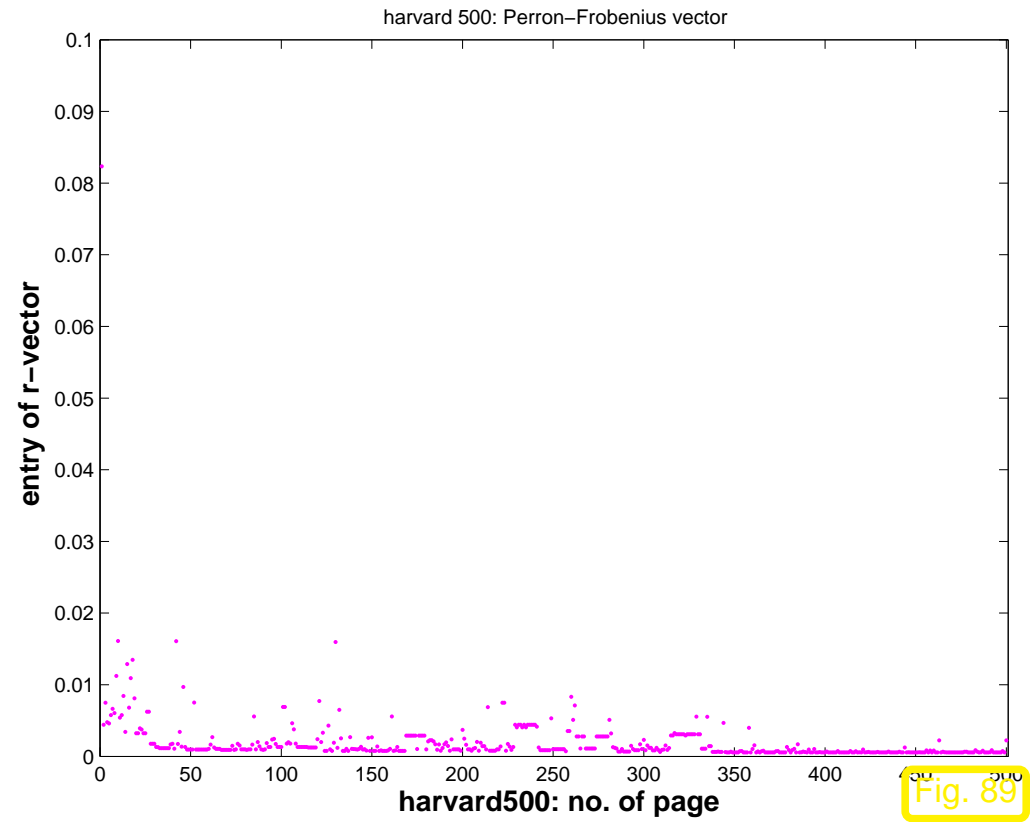
$$\|\mathbf{r}\|_1 = 1,$$

$$\mathbf{A}\mathbf{r} = \mathbf{r}.$$

Inefficient implementation!



stochastic simulation



eigenvector computation

$\mathbf{A} \hat{=}$ page rank transition probability matrix, see Code 6.3.4, $d = 0.15$, harvard500 example.

Errors:

$$\left\| \mathbf{A}^k \mathbf{x}_0 - \mathbf{r} \right\|_1,$$

with $\mathbf{x}_0 = \mathbf{1}/N$, $N = 500$.

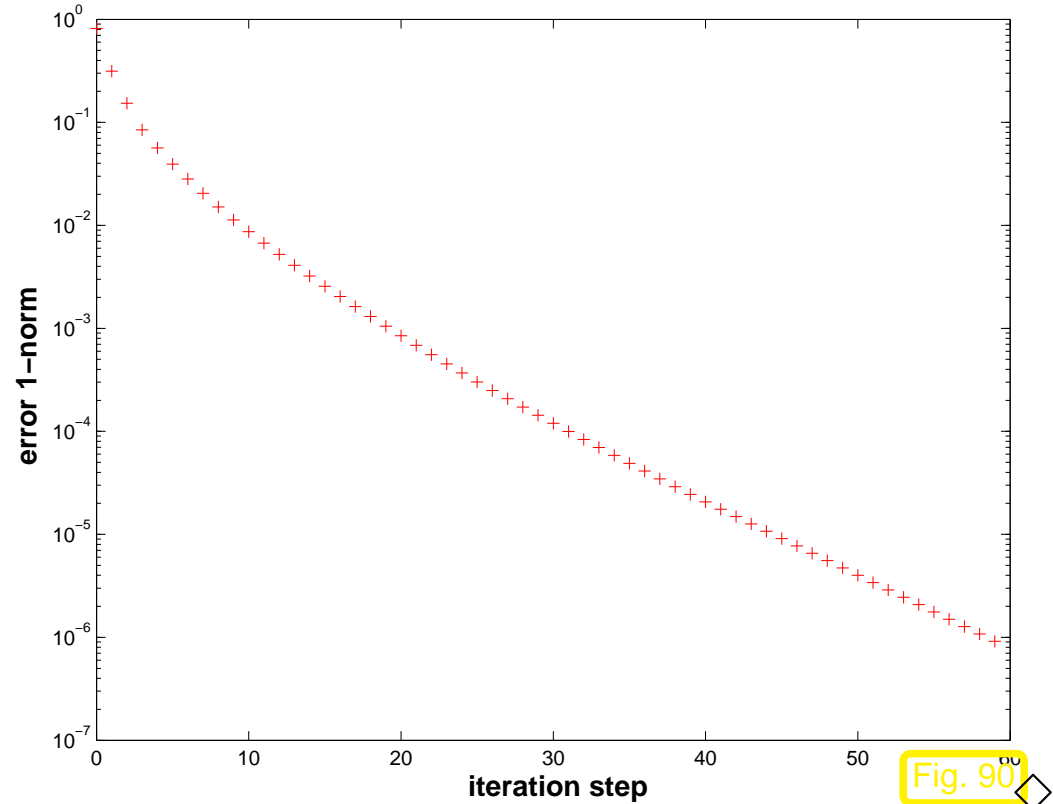
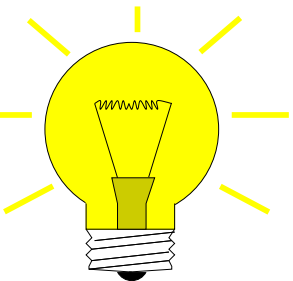


Fig. 90

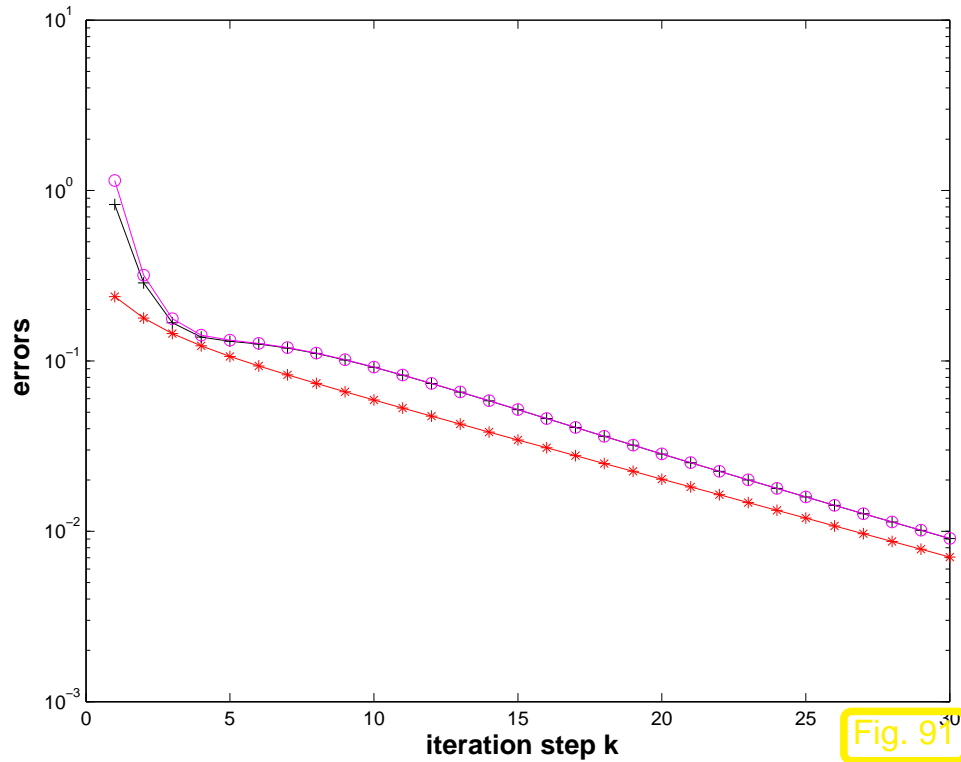
Task: given $\mathbf{A} \in \mathbb{K}^{n,n}$, find **largest** (in modulus) eigenvalue of \mathbf{A} and (an) associated eigenvector.

Idea: (suggested by page rank computation, Code 6.3.6)

Iteration: $\mathbf{z}^{(k+1)} = \mathbf{A}\mathbf{z}^{(k)}, \mathbf{z}^{(0)}$ arbitrary



Example 6.3.10 (Power iteration). → Ex. 6.3.1



```
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: $(\mathbf{S})_{:,10} \hat{=}$ eigenvector for eigenvalue 10)

$\mathbf{z}^{(0)}$ = random vector



Suggests **direct power method** (*ger.:* Potenzmethode): iterative method (→ Sect. 4.1)

initial guess: $\mathbf{z}^{(0)}$ “arbitrary”,

next iterate: $\mathbf{w} := \mathbf{A}\mathbf{z}^{(k-1)}$, $\mathbf{z}^{(k)} := \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, $k = 1, 2, \dots$. (6.3.11)

Computational effort: $1 \times \text{matrix} \times \text{vector}$ per step \rightarrow inexpensive for sparse matrices

When (6.3.11) has converged, two common ways to recover λ_{\max} \rightarrow [13, Alg. 7.20]

$$\textcircled{1} \quad \mathbf{Az}^{(k)} \approx \lambda_{\max} \mathbf{z}^{(k)} \quad \rightarrow \quad |\lambda_n| \approx \frac{\|\mathbf{Az}^{(k)}\|}{\|\mathbf{z}^{(k)}\|} \quad (\text{modulus only!})$$

$$\textcircled{2} \quad \lambda_{\max} \approx \underset{\theta \in \mathbb{R}}{\operatorname{argmin}} \left\| \mathbf{Az}^{(k)} - \theta \mathbf{z}^{(k)} \right\|_2^2 \quad \rightarrow \quad \lambda_{\max} \approx \frac{(\mathbf{z}^{(k)})^H \mathbf{Az}^{(k)}}{\|\mathbf{z}^{(k)}\|_2^2}.$$

Definition 6.3.15. For $\mathbf{A} \in \mathbb{K}^{n,n}$, $\mathbf{u} \in \mathbb{K}^n$ the *Rayleigh quotient* is defined by

$$\rho_{\mathbf{A}}(\mathbf{u}) := \frac{\mathbf{u}^H \mathbf{A} \mathbf{u}}{\mathbf{u}^H \mathbf{u}}.$$

Example 6.3.17 (Direct power method). \rightarrow Ex. 6.3.17 cnt'd

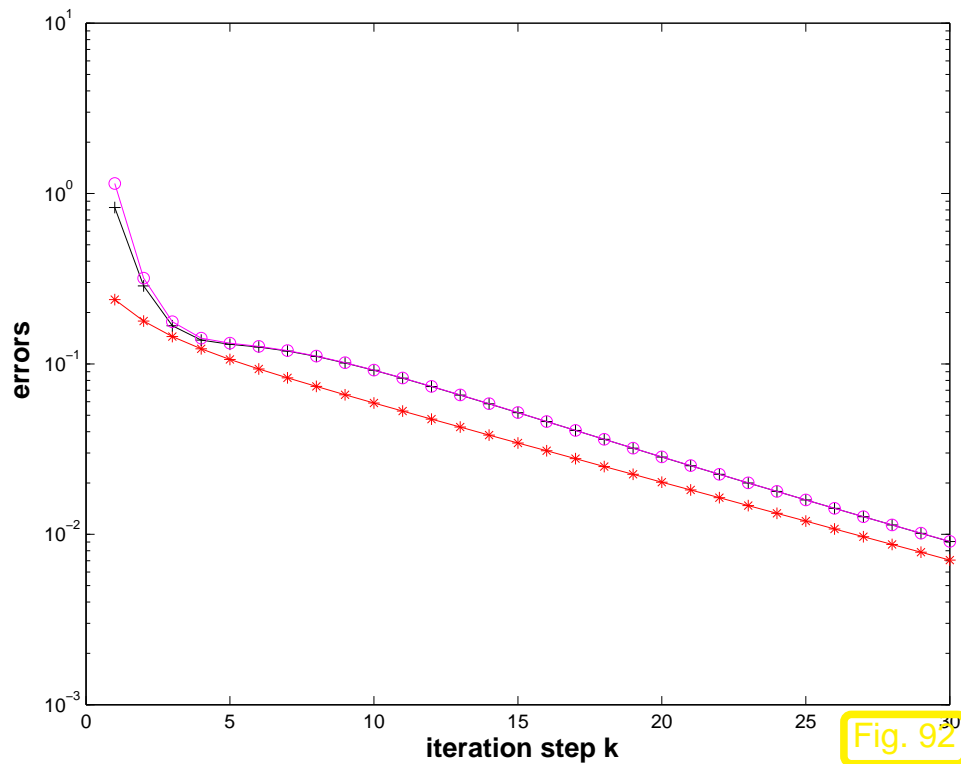


Fig. 92

```
n = length(d);
S = triu(diag(n:-1:1,0)+...
ones(n,n));
A = S*diag(d,0)*inv(S);
```

```
d = (1:10)';
```

○ : error $|\lambda_n - \rho_{\mathbf{A}}(\mathbf{z}^{(k)})|$

△ : error norm $\|\mathbf{z}^{(k)} - \mathbf{s}_{\cdot,n}\|$

+ : $\left| \lambda_n - \frac{\|\mathbf{A}\mathbf{z}^{(k-1)}\|_2}{\|\mathbf{z}^{(k-1)}\|_2} \right|$

$\mathbf{z}^{(0)}$ = random vector

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 $d = 1 - 2.^{-(1:0.5:5)'}$; $\Rightarrow |\lambda_{n-1}| : |\lambda_n| = 0.9866$

$$\rho_{EV}^{(k)} := \frac{\|\mathbf{z}^{(k)} - \mathbf{s}_{\cdot,n}\|}{\|\mathbf{z}^{(k-1)} - \mathbf{s}_{\cdot,n}\|},$$

$$\rho_{EW}^{(k)} := \frac{|\rho_{\mathbf{A}}(\mathbf{z}^{(k)}) - \lambda_n|}{|\rho_{\mathbf{A}}(\mathbf{z}^{(k-1)}) - \lambda_n|}.$$

| | ① | | ② | | ③ | |
|-----|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| k | $\rho_{EV}^{(k)}$ | $\rho_{EW}^{(k)}$ | $\rho_{EV}^{(k)}$ | $\rho_{EW}^{(k)}$ | $\rho_{EV}^{(k)}$ | $\rho_{EW}^{(k)}$ |
| 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 |



Theorem 6.3.19 (Convergence of direct power method). \rightarrow [13, 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 λ_n normalized according to $\|\mathbf{y}\|_2 = \|\mathbf{v}\|_2 = 1$. Then there is convergence

$$\left\| \mathbf{A} \mathbf{z}^{(k)} \right\|_2 \rightarrow \lambda_n, \quad \mathbf{z}^{(k)} \rightarrow \pm \mathbf{v} \quad \text{linearly with rate } \frac{|\lambda_{n-1}|}{|\lambda_n|},$$

where $\mathbf{z}^{(k)}$ are the iterates of the direct power iteration and $\mathbf{y}^H \mathbf{z}^{(0)} \neq 0$ is assumed.

6.3.2 Inverse Iteration [13, Sect. 7.6], [51, Sect. 5.3.2]

Example 6.3.22 (Image segmentation).

Given: gray-scale image: intensity matrix $\mathbf{P} \in \{0, \dots, 255\}^{m,n}$, $m, n \in \mathbb{N}$
 $((\mathbf{P})_{ij} \leftrightarrow \text{pixel}, 0 \hat{=} \text{black}, 255 \hat{=} \text{white})$

Code 6.3.23: loading and displaying an image

Loading and displaying im-
ages in MATLAB

```

1 M = imread('eth.pbm');
2 [m,n] = size(M);
3 fprintf('%dx%d grey scale pixel image\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

Next: Statement of (rigorously defined) problem, *cf.* Sect. 2.5.2:

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:

$$\mathbf{W} \in \mathbb{R}^{N,N}, \quad N := mn, \quad (6.3.24)$$

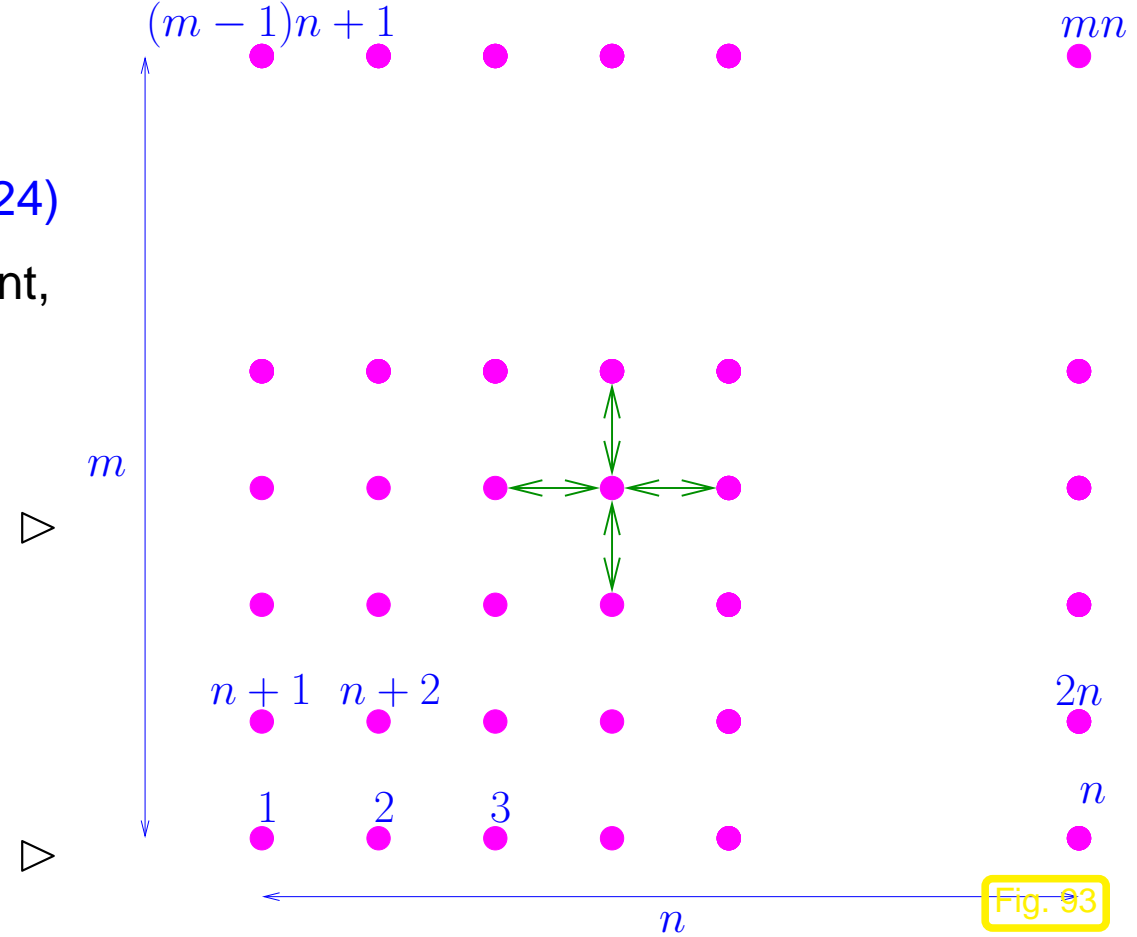
$$(\mathbf{W})_{ij} = \begin{cases} 0 & , \text{ if pixels } i, j \text{ not adjacent,} \\ 0 & , \text{ if } i = j, \\ \sigma(p_i, p_j) & , \text{ if pixels } i, j \text{ adjacent.} \end{cases}$$

$\leftrightarrow \hat{=}$ 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



Definition 6.3.25 (Normalized cut). (\rightarrow [58, Sect. 2])

For $\mathcal{X} \subset \mathcal{V}$ define the **normalized cut** as

$$\text{Ncut}(\mathcal{X}) := \frac{\text{cut}(\mathcal{X})}{\text{weight}(\mathcal{X})} + \frac{\text{cut}(\mathcal{X})}{\text{weight}(\mathcal{V} \setminus \mathcal{X})},$$

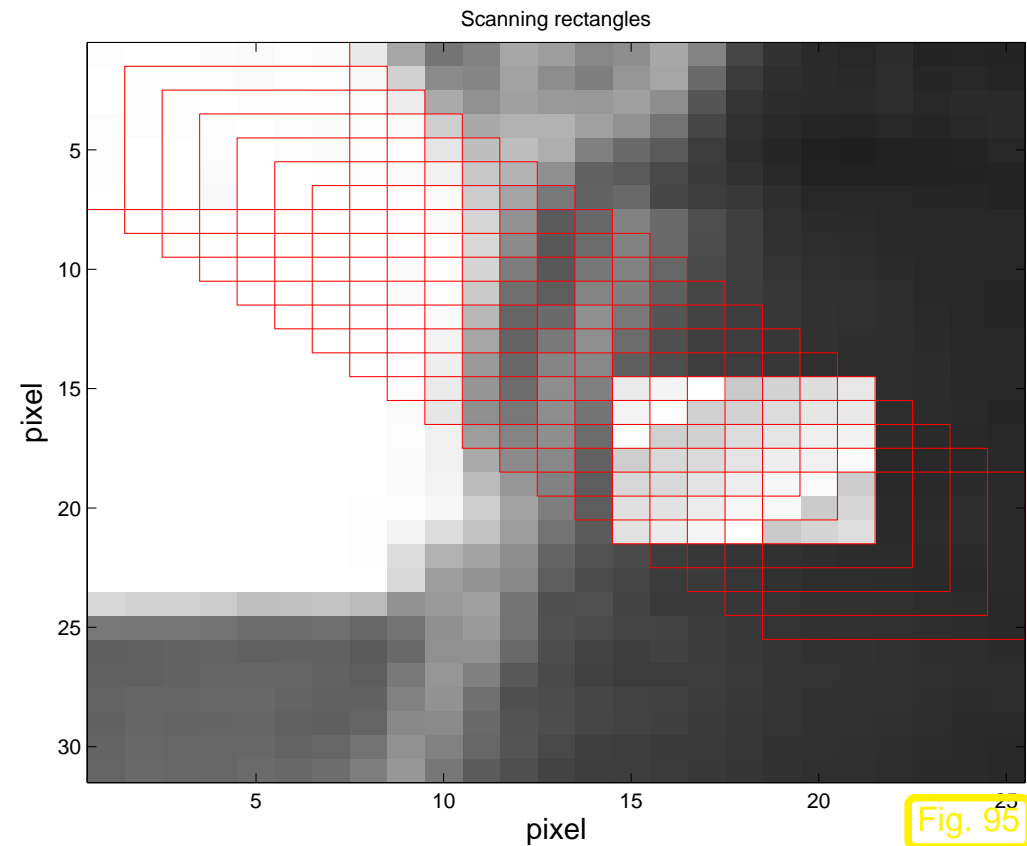
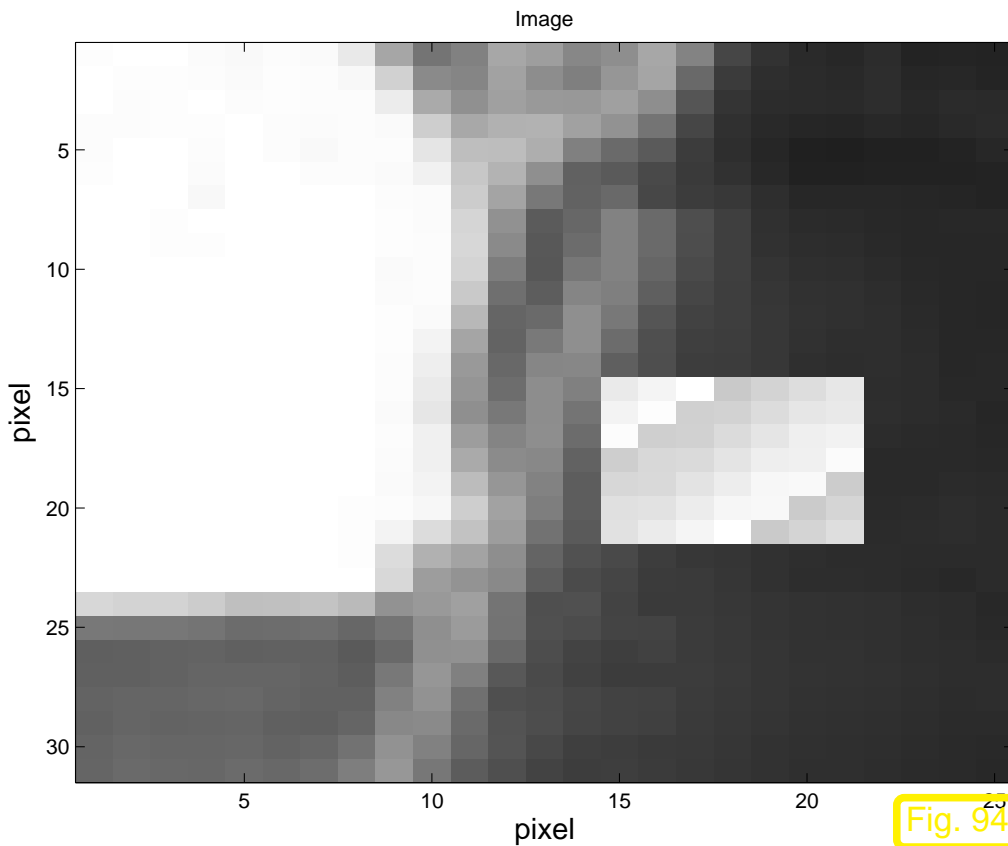
with

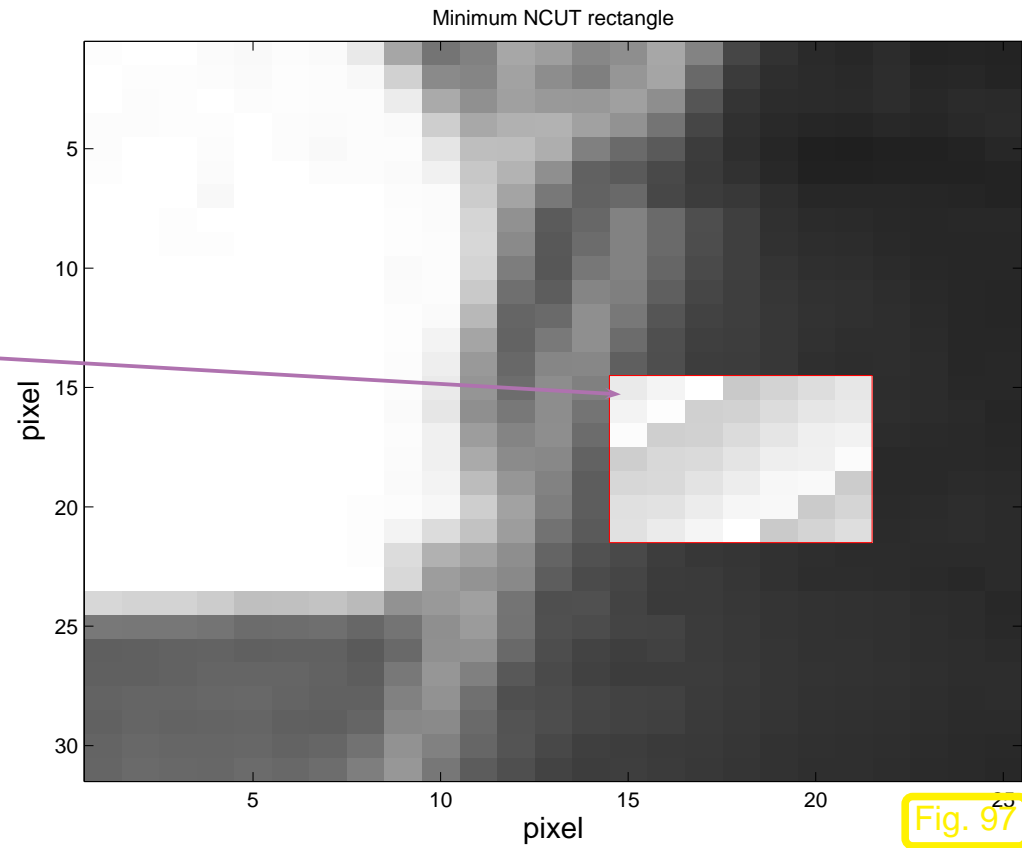
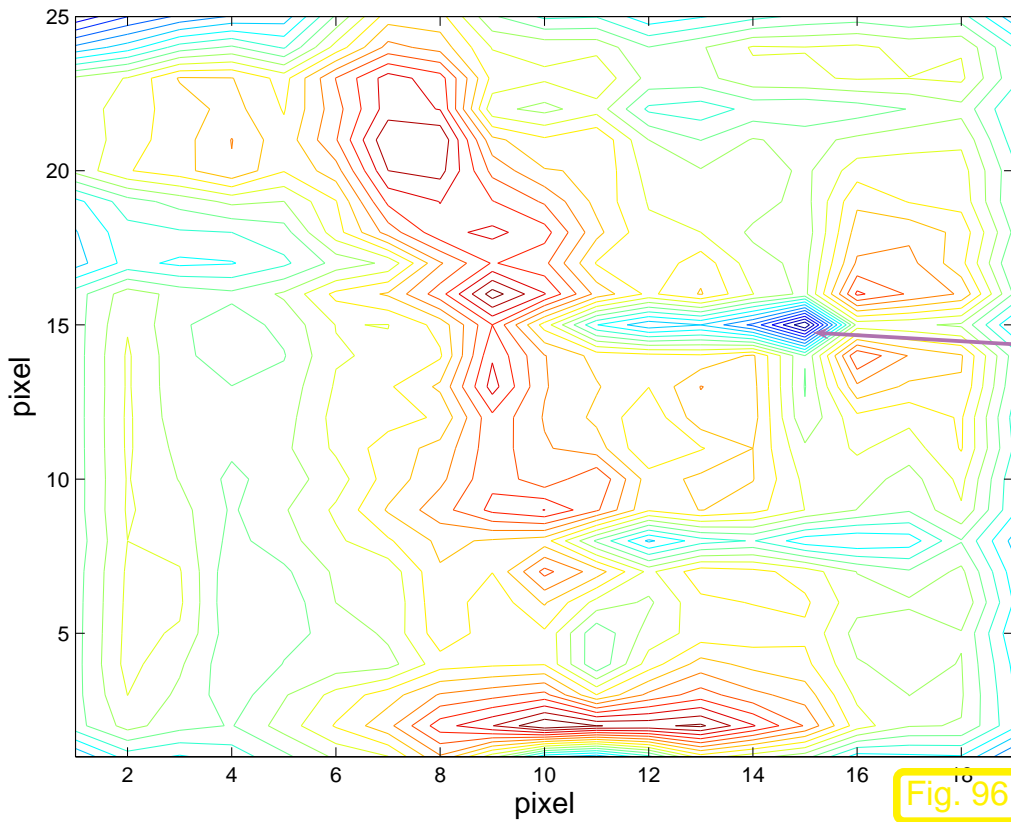
$$\text{cut}(\mathcal{X}) := \sum_{i \in \mathcal{X}, j \notin \mathcal{X}} w_{ij}, \quad \text{weight}(\mathcal{X}) := \sum_{i \in \mathcal{X}, j \in \mathcal{X}} w_{ij}.$$

$$\text{find } \mathcal{X}^* \subset \mathcal{V}: \mathcal{X}^* = \underset{\mathcal{X} \subset \mathcal{V}}{\operatorname{argmin}} \operatorname{Ncut}(\mathcal{X}) . \quad (6.3.26)$$



NP-hard combinatorial optimization problem !





Equivalent reformulation:

indicator function: $z : \{1, \dots, N\} \mapsto \{-1, 1\}$, $z_i := z(i) = \begin{cases} 1 & , \text{if } i \in \mathcal{X} , \\ -1 & , \text{if } i \notin \mathcal{X} . \end{cases}$ (6.3.27)

▶
$$\text{Ncut}(\mathcal{X}) = \frac{\sum_{z_i > 0, z_j < 0} -w_{ij} z_i z_j}{\sum_{z_i > 0} d_i} + \frac{\sum_{z_i > 0, z_j < 0} -w_{ij} z_i z_j}{\sum_{z_i < 0} d_i},$$
 (6.3.28)

$$d_i = \sum_{j \in \mathcal{V}} w_{ij} = \text{weight}(\{i\}) . \tag{6.3.29}$$

Sparse matrices:

$$\mathbf{D} := \text{diag}(d_1, \dots, d_N) \in \mathbb{R}^{N,N} , \quad \mathbf{A} := \mathbf{D} - \mathbf{W} = \mathbf{A}^\top . \tag{6.3.30}$$

Lemma 6.3.33 (Ncut and Rayleigh quotient). (→ [58, Sect. 2])

With $\mathbf{z} \in \{-1, 1\}^N$ according to (6.3.27) 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 $\rho_{\mathbf{A}, \mathbf{D}}(\mathbf{y})$

► segmentation problem (6.3.26) $\Leftrightarrow \underset{\mathbf{y} \in \{2, -2\beta\}^N, \mathbf{1}^\top \mathbf{D} \mathbf{y} = 0}{\text{argmin}} \rho_{\mathbf{A}, \mathbf{D}}(\mathbf{y}) . \tag{6.3.35}$

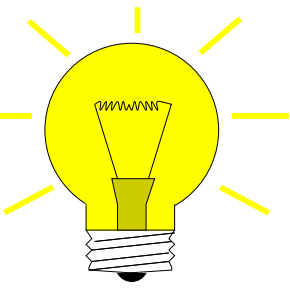
still NP-hard

Idea:

Relaxation

Discrete optimization problem \rightarrow continuous optimization problem

$$(6.3.35) \rightarrow \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^N, \mathbf{y} \neq 0, \mathbf{1}^\top \mathbf{D} \mathbf{y} = 0} \rho_{\mathbf{A}, \mathbf{D}}(\mathbf{y}). \quad (6.3.36)$$



Theorem 6.3.37 (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

$$\operatorname{Eig}_{\mathbf{A}}(\lambda_1) = \operatorname{argmin}_{\mathbf{y} \in \mathbb{C}^{n,n} \setminus \{0\}} \rho_{\mathbf{A}}(\mathbf{y}) \quad \text{and} \quad \operatorname{Eig}_{\mathbf{A}}(\lambda_m) = \operatorname{argmax}_{\mathbf{y} \in \mathbb{C}^{n,n} \setminus \{0\}} \rho_{\mathbf{A}}(\mathbf{y}).$$

Algorithm 6.3.49 (Binary grayscale image segmentation (outline)).

① Given similarity function σ compute (sparse!) matrices $\mathbf{W}, \mathbf{D}, \mathbf{A} \in \mathbb{R}^{N,N}$, see (6.3.24), (6.3.30).

② Compute \mathbf{y}^* , $\|\mathbf{y}^*\|_2 = 1$, as eigenvector belonging to the *smallest eigenvalue* of $\hat{\mathbf{A}} := \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} + 2(\mathbf{D}^{1/2} \mathbf{1})(\mathbf{D}^{1/2} \mathbf{1})^\top$.

③ Set $\mathbf{x}^* = \mathbf{D}^{-1/2}\mathbf{y}^*$ and define the image segment as pixel set

$$\mathcal{X} := \{i \in \{1, \dots, N\} : x_i^* > \frac{1}{N} \sum_{i=1}^N x_i^*\} . \quad (6.3.50)$$

↑
mean value of entries of \mathbf{x}^*

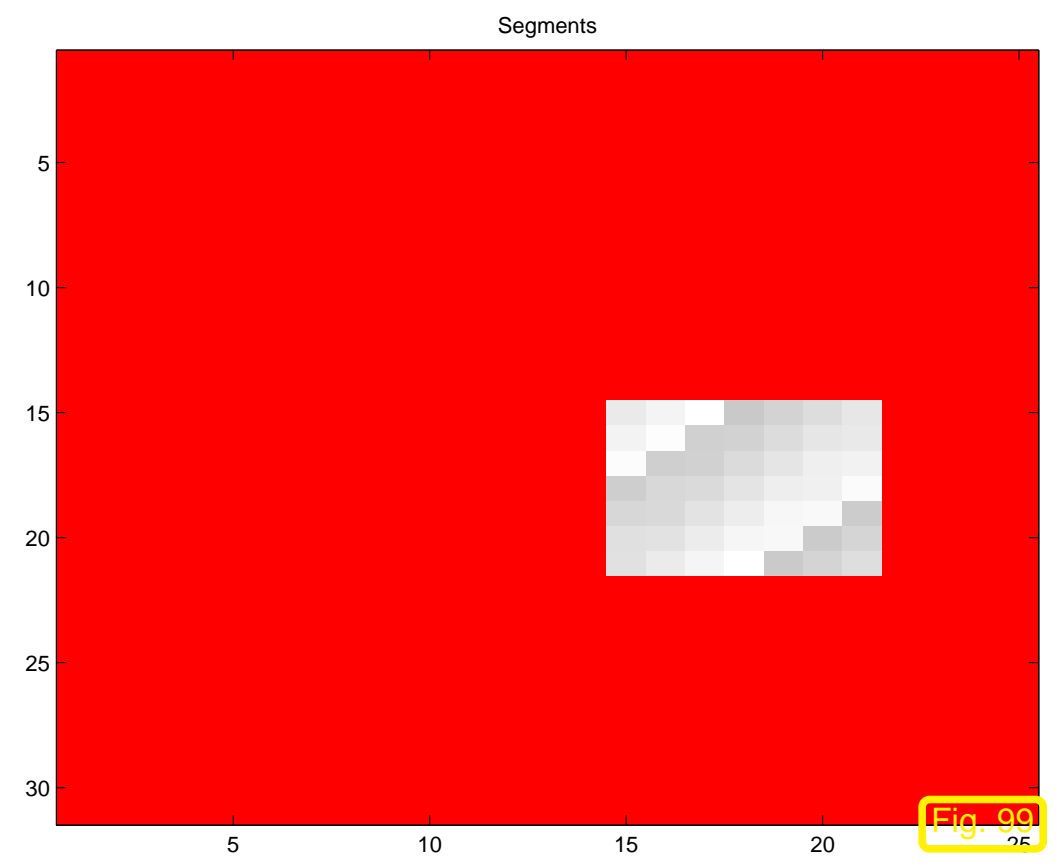
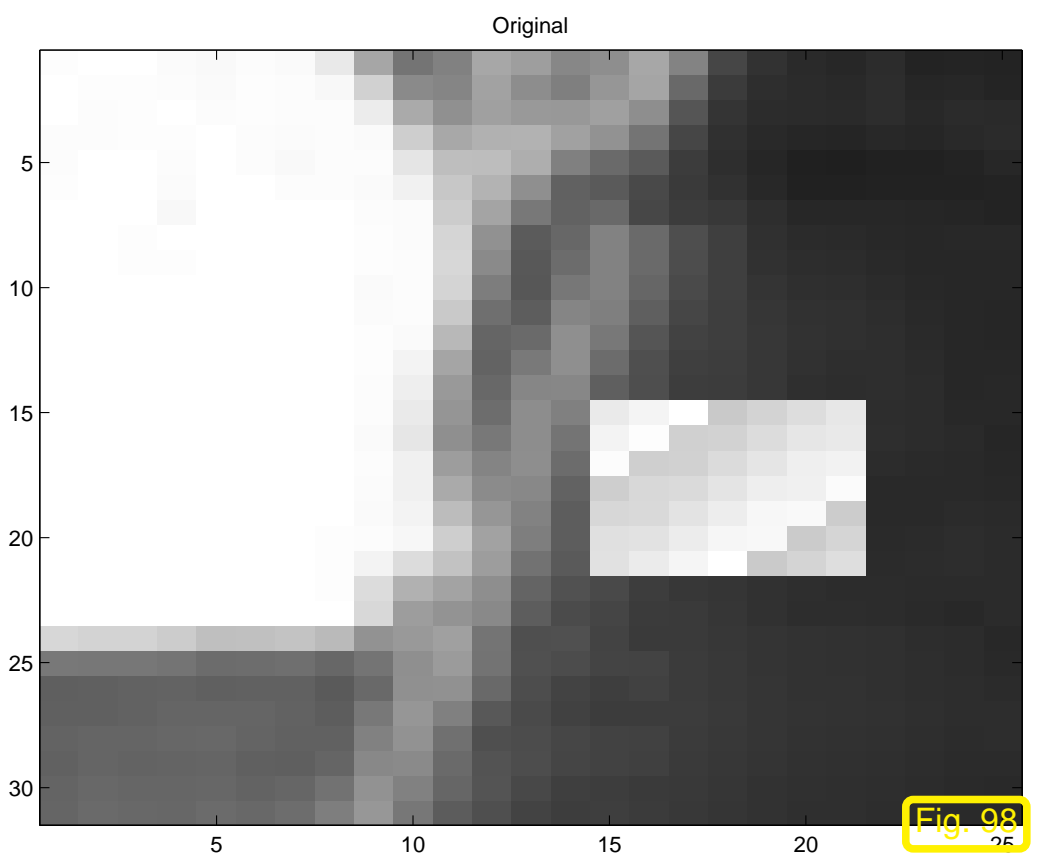
Code 6.3.51: 1st stage of segmentation of grayscale image

```

1 % Read image and build matrices, see Code 6.3.31 and (6.3.30)
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  $\mathbf{x}^*$  belonging to 2nd smallest generalized
12 % eigenvalue of  $\mathbf{A}$  and  $\mathbf{D}$ 
13 x = W(:,1); x = Dm*v;
14 % Extract segmented image
15 xs = reshape(x,m,n); Xidx = find(xs > (sum(sum(xs))/(n*m)));

```

1st-stage of segmentation of 31×25 grayscale pixel image (root .pbm, red pixels $\hat{=}$ \mathcal{X} , $\sigma(x, y) = \exp(-(x-y/10)^2)$)



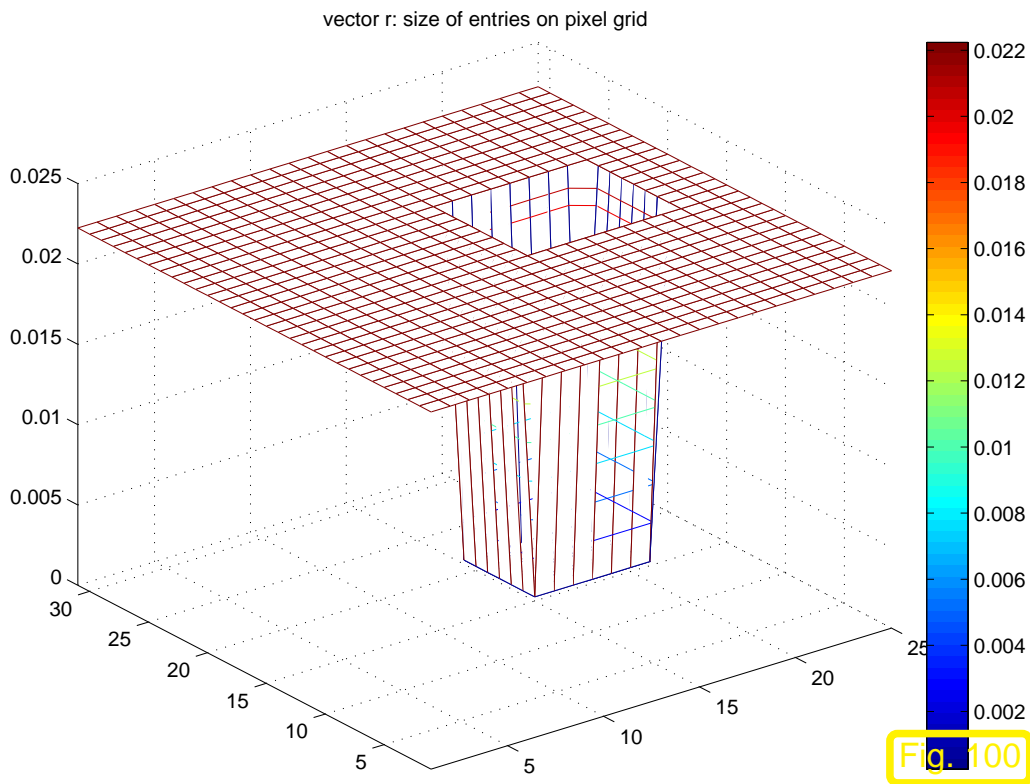


Image from Fig. 111:

◁ eigenvector \mathbf{x}^* plotted on pixel grid

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:

$$\text{Smallest (in modulus) EV of } \mathbf{A} = \left(\text{Largest (in modulus) EV of } \mathbf{A}^{-1} \right)^{-1}$$

Code 6.3.52: 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. 2.2.13
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:  $y = \mathbf{A}^{-1}\mathbf{x}$ ,
8     lmin = 1/norm(y); % new approxmation of  $\lambda_{\min}(\mathbf{A})$ 
9     y = y*lmin; % normalization  $y := \frac{y}{\|y\|_2}$ 
10 end
    
```

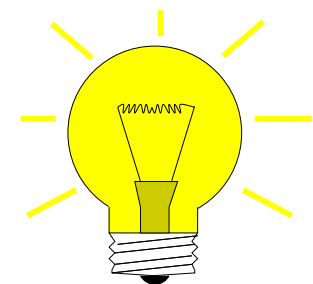
Remark 6.3.53 (Shifted inverse iteration).



Idea:

A posteriori adaptation of shift

Use $\alpha := \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})$ in k -th step of inverse iteration.



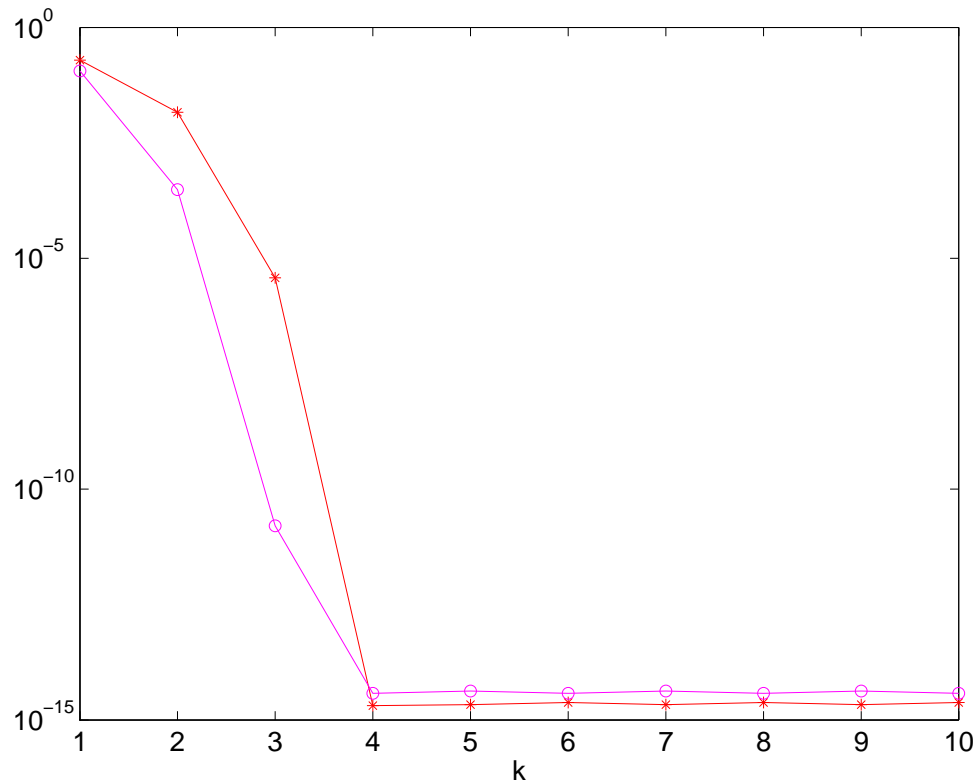
Code 6.3.57: Rayleigh quotient iteration (for *normal* $\mathbf{A} \in \mathbb{R}^{n,n}$)

```

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); %  $\mathbf{z}^{(0)}$ 
4 for i=1:maxit
5     z = (A-alpha*speye(n))\z; %  $\mathbf{z}^{(k+1)} = (\mathbf{A} - \rho_{\mathbf{A}}(\mathbf{z}^{(k)})\mathbf{I})^{-1}\mathbf{x}^{(k)}$ 
6     z = z/norm(z); lmin=dot(A*z,z); % Computation of  $\rho_{\mathbf{A}}(\mathbf{z}^{(k+1)})$ 
7     if (abs(alpha-lmin) < tol*lmin) % Desired relative accuracy reached ?
8         break; end;
9     alpha = lmin;
10 end

```

Example 6.3.58 (Rayleigh quotient iteration).



```
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';
```

- : $|\lambda_{\min} - \rho_{\mathbf{A}}(\mathbf{z}^{(k)})|$
- * : $\|\mathbf{z}^{(k)} - \mathbf{x}_j\|, \lambda_{\min} = \lambda_j, \mathbf{x}_j \in \text{Eig}_{\mathbf{A}}(\lambda_j),$
- : $\|\mathbf{x}_j\|_2 = 1$

| 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 | 0.00000000001783 | 0.00000411928759 |
| 4 | 0.00000000000000 | 0.00000000000000 |
| 5 | 0.00000000000000 | 0.00000000000000 |

Theorem 6.3.59. \rightarrow [35, Thm. 25.4]
If $\mathbf{A} = \mathbf{A}^H$, then $\rho_{\mathbf{A}}(\mathbf{z}^{(k)})$ converges locally of order 3 (\rightarrow Def. 4.1.14) to the smallest eigenvalue (in modulus), when $\mathbf{z}^{(k)}$ are generated by the Rayleigh quotient iteration (6.3.56).



6.3.3 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 (\rightarrow Code 6.3.51) and Rayleigh quotient iteration (6.3.56).

?

What if direct solution of $\mathbf{Ax} = \mathbf{b}$ not feasible ?

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** $\mathbf{B}^{-1} \approx \mathbf{A}^{-1}$

\triangleright $\mathbf{B} \hat{=} \mathbf{Preconditioner}$ for \mathbf{A} , see Def. 5.3.3

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



[**Preconditioned inverse iteration** (PINVIT) for s.p.d. \mathbf{A}]

$\mathbf{z}^{(0)}$ arbitrary,

$$\mathbf{w} = \mathbf{z}^{(k-1)} - \mathbf{B}^{-1}(\mathbf{A}\mathbf{z}^{(k-1)} - \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})\mathbf{z}^{(k-1)}), \quad k = 1, 2, \dots \quad (6.3.60)$$

$$\mathbf{z}^{(k)} = \frac{\mathbf{w}}{\|\mathbf{w}\|_2},$$

Code 6.3.61: preconditioned inverse iteration (6.3.60)

```

1 function [lmin,z,res] =
   pinvit(evalA,n,invB,tol,maxit)
2 % invB  $\hat{=}$  handle to function implementing preconditioner  $B^{-1}$ 
3 z = (1:n)'; z = z/norm(z); % initial guess
4 res = []; rho = 0;
5 for i=1:maxit
6   v = evalA(z); rhon = dot(v,z); % Rayleigh quotient
7   r = v - rhon*z; % residual
8   z = z - invB(r); % iteration according to (6.3.60)
9   z = z/norm(z); % normalization
0   res = [res; rhon]; % tracking iteration
1   if (abs(rho-rhon) < tol*abs(rhon)), break;
2   else rho = rhon; end
3 end
4 lmin = dot(evalA(z),z); res = [res; lmin],

```

Computational
effort:

1 matrix \times vector

1 evaluation of
preconditioner

A few

AXPY-operations

Example 6.3.62 (Convergence of PINVIT).

S.p.d. matrix $A \in \mathbb{R}^{n,n}$, tridiagonal preconditioner, see Ex. 5.3.11

```

1 A = spdiags(repmat([1/n,-1,2*(1+1/n),-1,1/n],n,1),
   [-n/2,-1,0,1,n/2],n,n);
2 evalA = @(x) A*x;
3 % inverse iteration

```

```

4 invB = @(x) A\x;
5 % tridiagonal preconditioning
6 B = spdiags(spdiags(A, [-1,0,1]), [-1,0,1], n, n); invB = @(x) B\x;
    
```

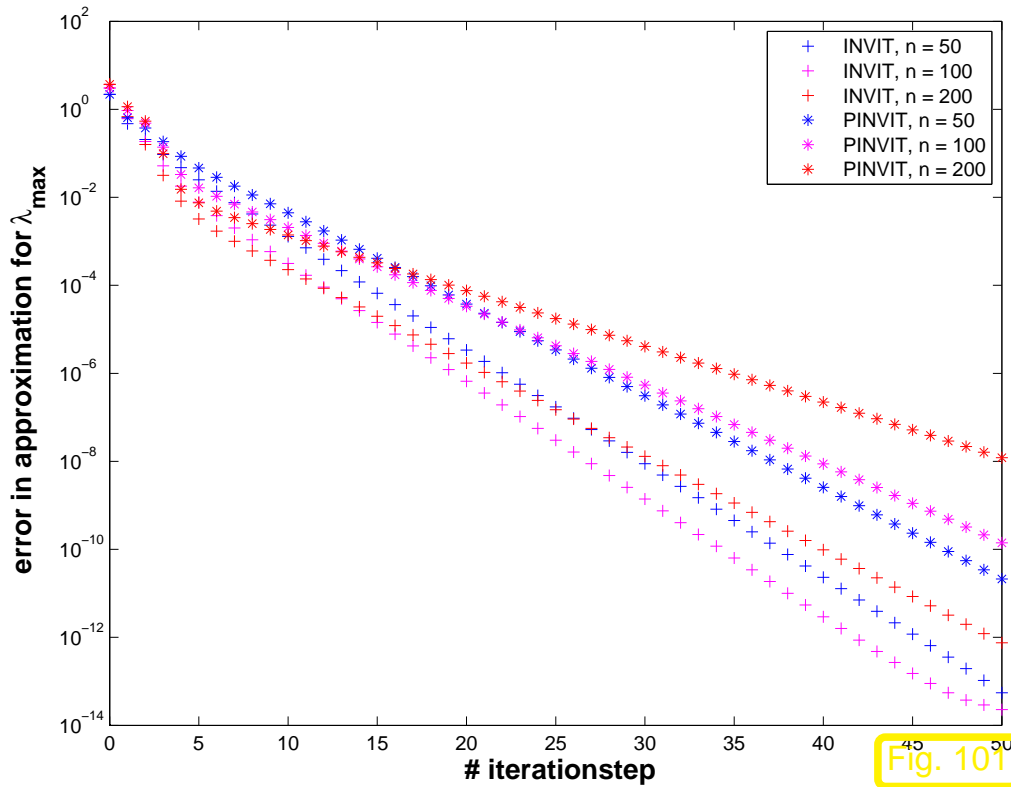


Fig. 101

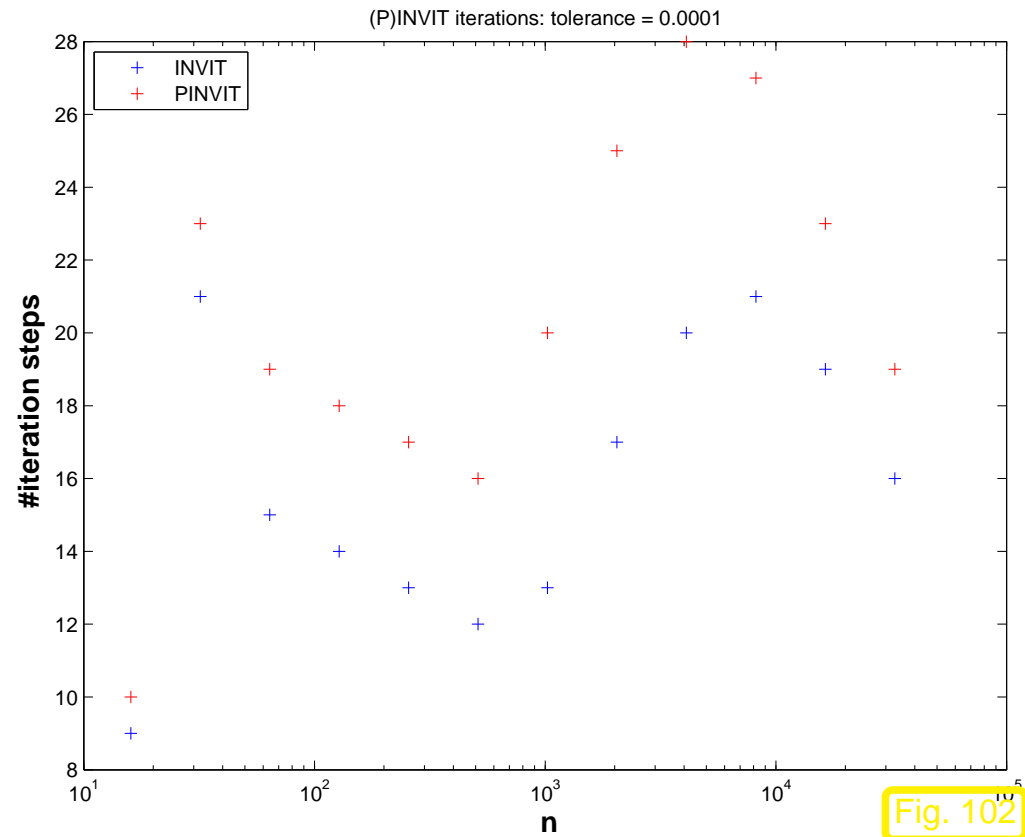


Fig. 102

Theory [47, 46]:

- linear convergence of (6.3.60)
- fast convergence, if spectral condition number $\kappa(\mathbf{B}^{-1}\mathbf{A})$ small

6.3.4 Subspace iterations

Example 6.3.68 (Vibrations of a truss structure).

cf. [35, Sect. 3], MATLAB's `truss` demo

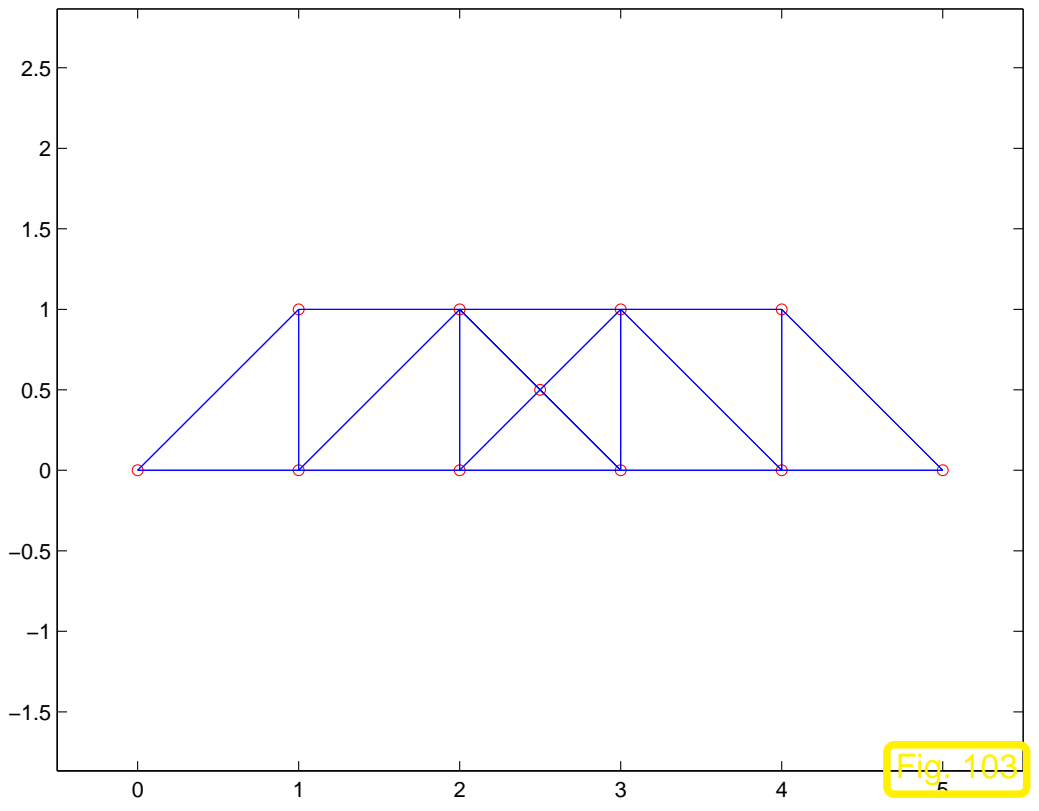


Fig. 103

◁ a “bridge” truss

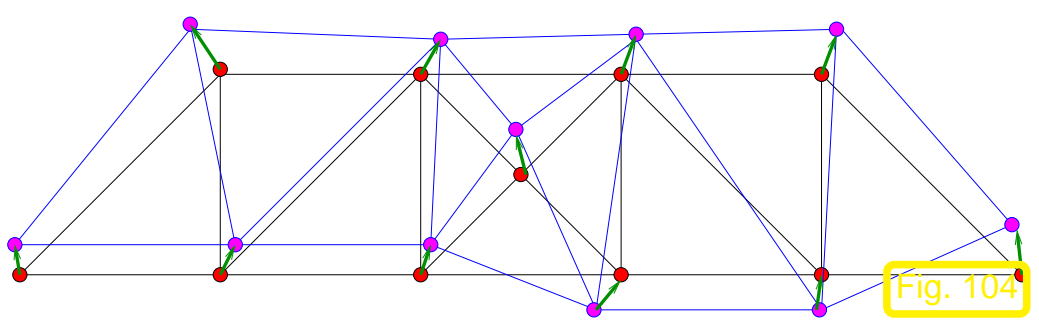


Fig. 104

◁ deformed truss:

● $\hat{=}$ point masses at positions \mathbf{p}^i

→ $\hat{=}$ displacement vectors \mathbf{u}^i

● $\hat{=}$ shifted masses at $\mathbf{p}^i + \mathbf{u}^i$

(6.3.78)



$$\mathbf{M} \frac{d\mathbf{u}}{dt^2}(t) + \mathbf{A}\mathbf{u}(t) = \mathbf{f}(t)$$

(6.3.79)

truss resonant frequencies

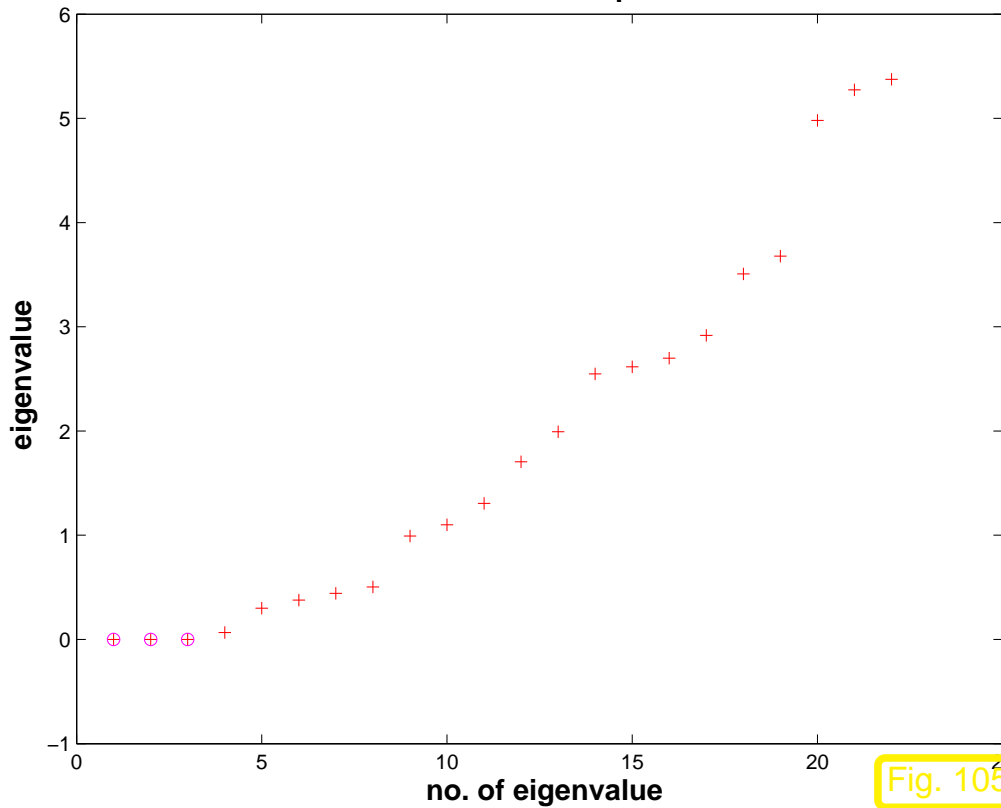
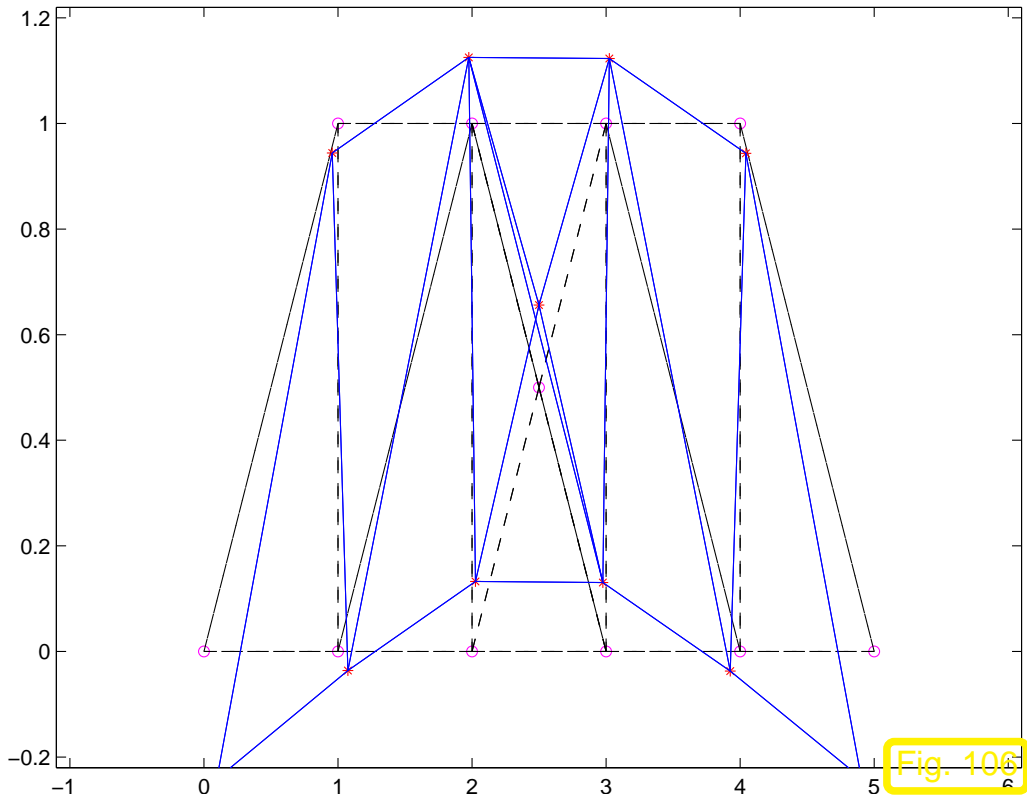


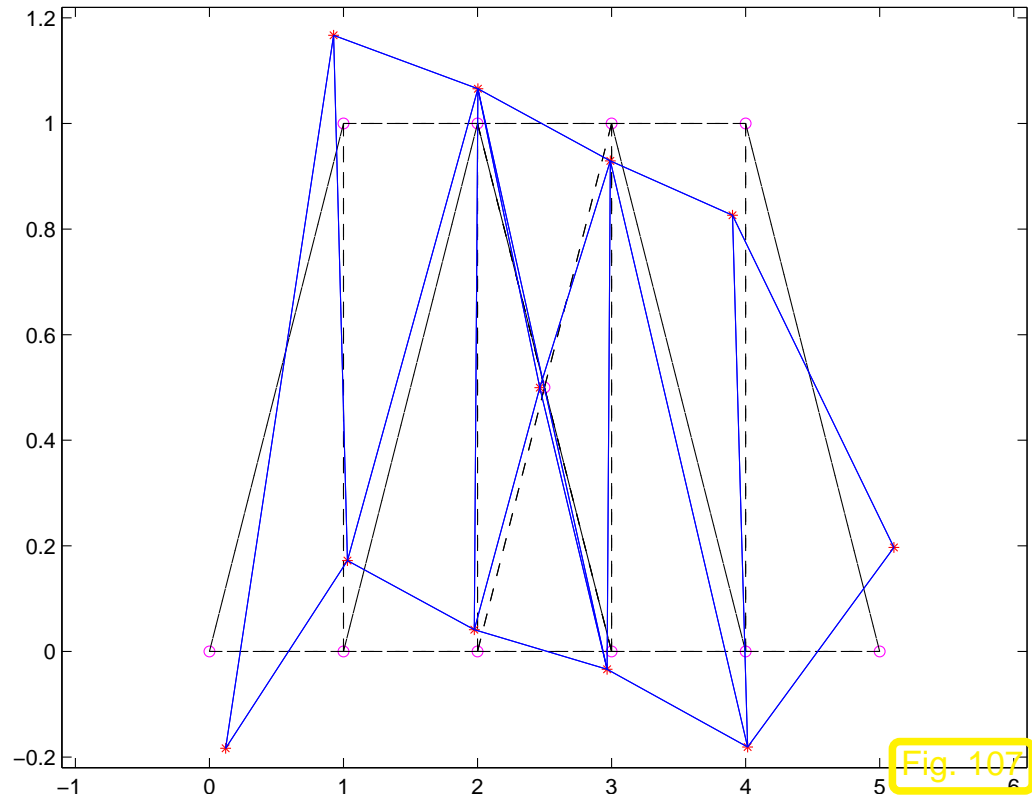
Fig. 105

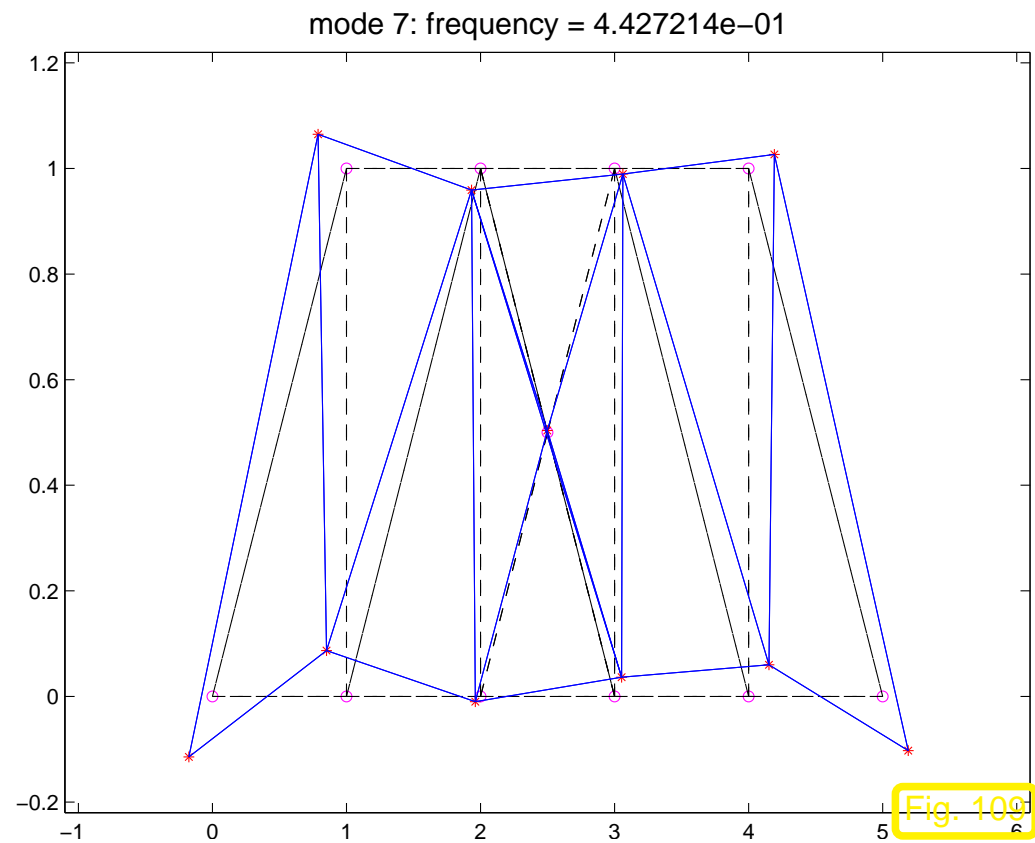
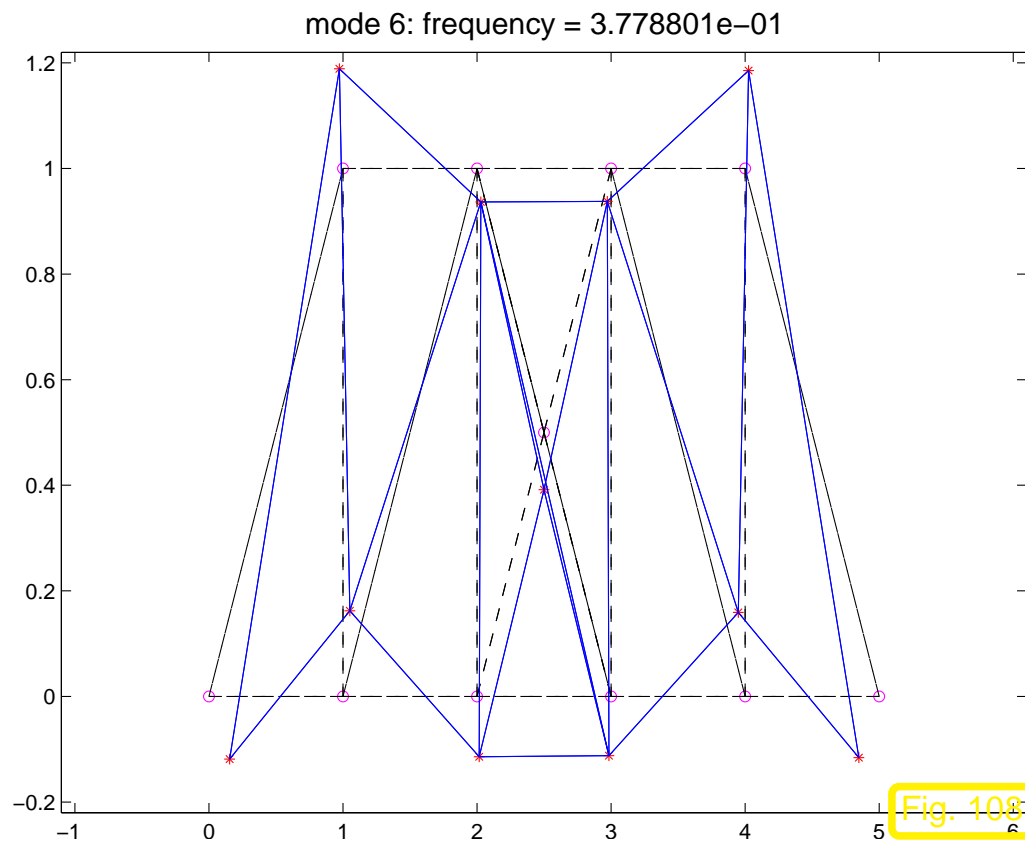
◁ resonant frequencies of bridge truss from Fig. 116.

mode 4: frequency = $6.606390e-02$



mode 5: frequency = $3.004450e-01$





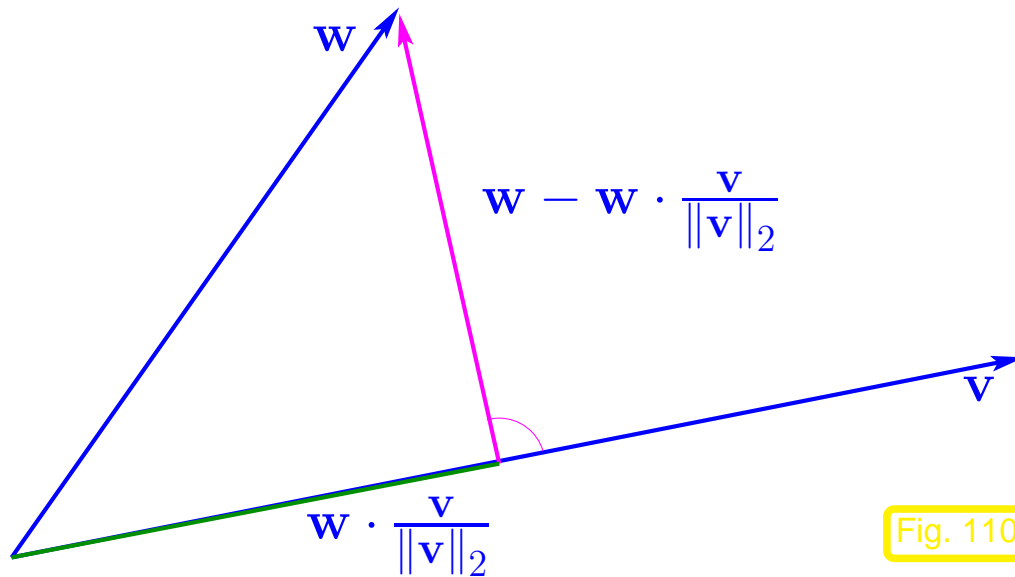
Task: Compute m , $m \ll n$, of the smallest/largest (in modulus) eigenvalues of $\mathbf{A} = \mathbf{A}^H \in \mathbb{C}^{n,n}$ and associated eigenvectors.

6.3.4.1 Orthogonalization

```

1 function [v,w] = sspowitstep(A,v,w)
2 v = A*v; w = A*w; % "power iteration", cf. (6.3.11)
3 % orthogonalization, cf. Gram-Schmidt orthogonalization (5.2.11)
4 v = v/norm(v); w = w - dot(v,w)*v; w = w/norm(w); % now w ⊥ v

```



◁ Orthogonalization of two vectors
(see Line 4 of Code 6.3.81)

Fig. 110

Example 6.3.83 (Subspace power iteration with orthogonal projection).

☞ construction of matrix $\mathbf{A} = \mathbf{A}^\top$ as in Ex. 6.3.58

$\sigma(\mathbf{A}) = \{1, 2, \dots, 10\}$:

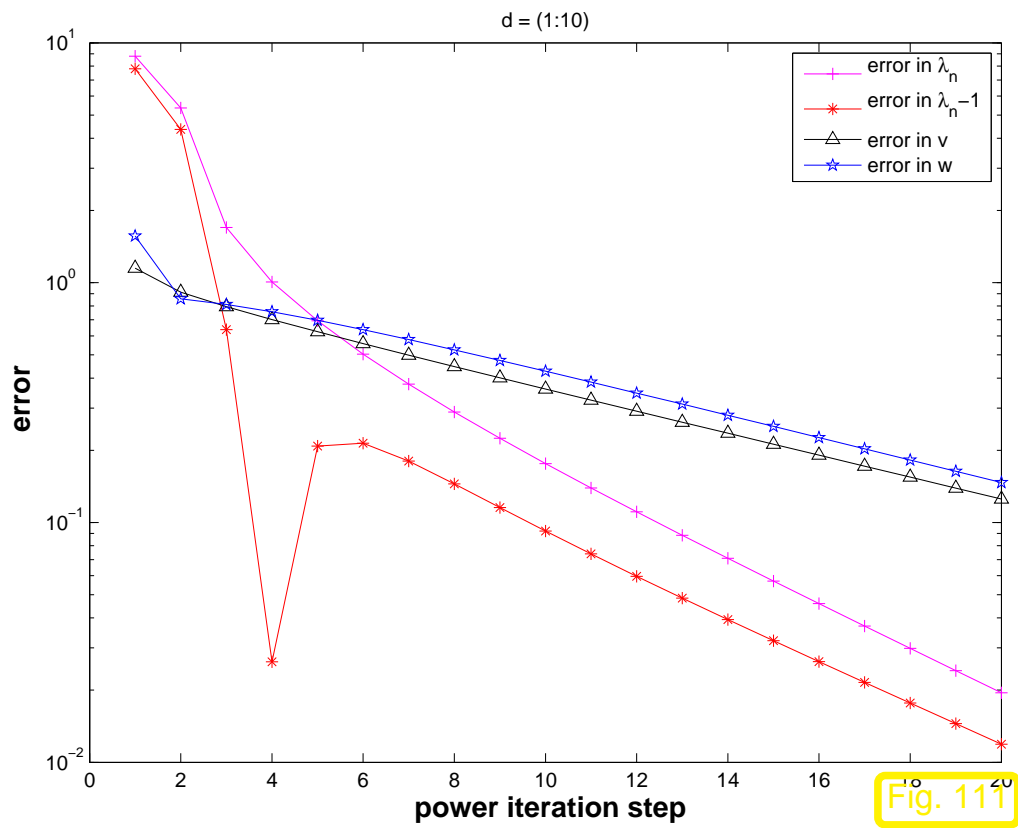


Fig. 111

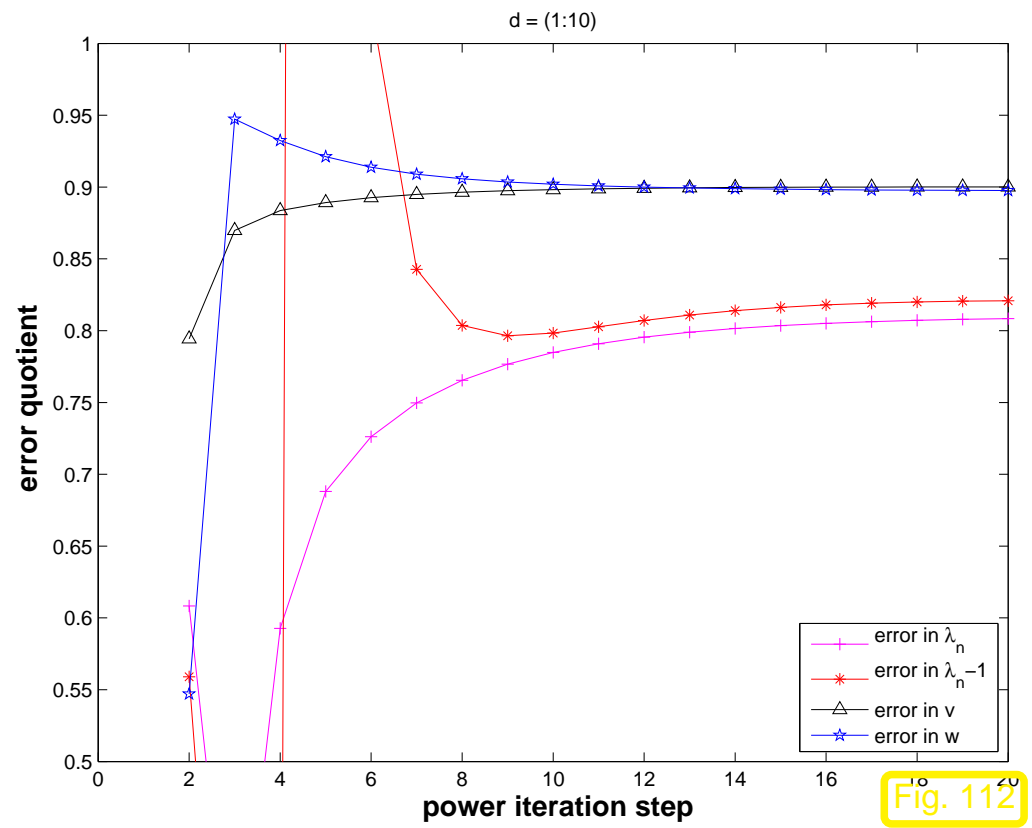


Fig. 112

$\sigma(\mathbf{A}) = \{0.5, 1, \dots, 4, 9.5, 10\}$:

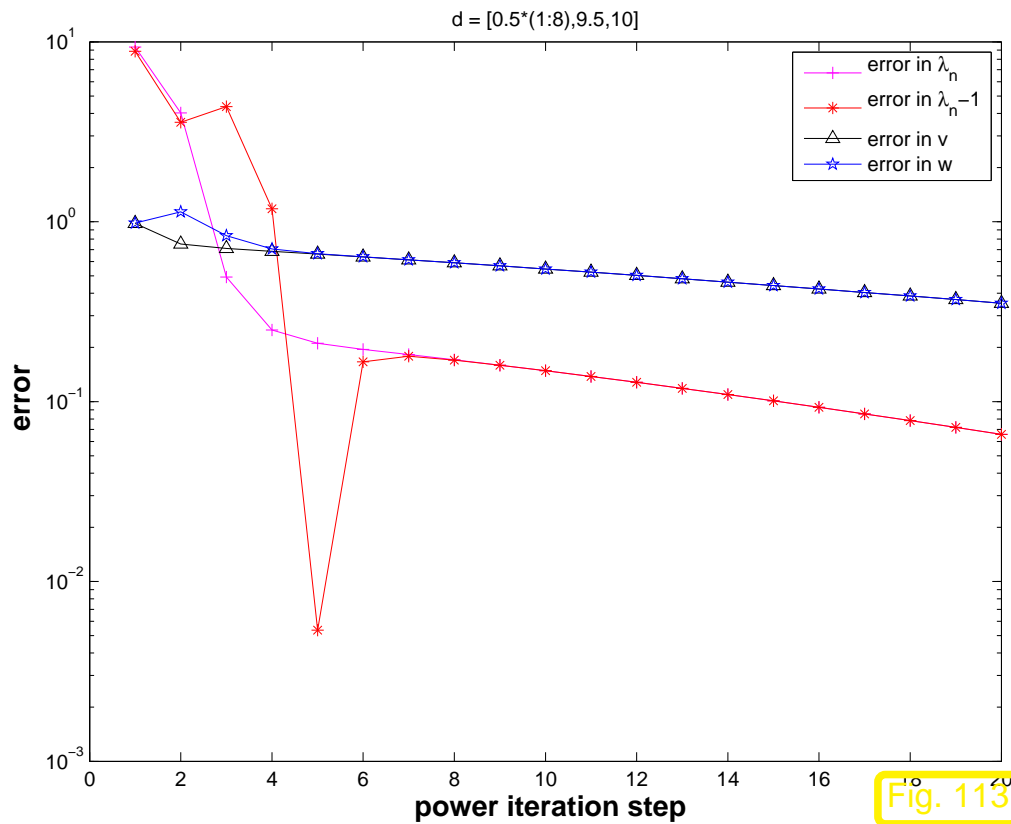


Fig. 113

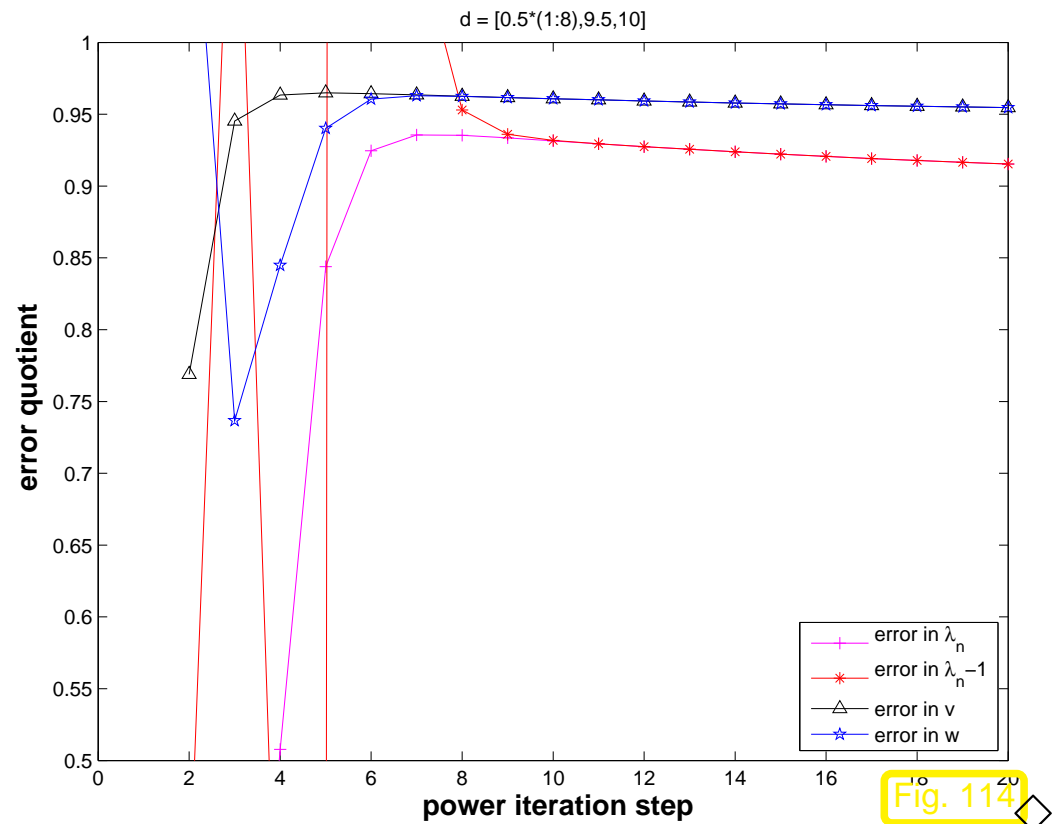


Fig. 114

Code 6.3.89: **Gram-Schmidt orthonormalization** (do not use, unstable algorithm!)

```

1 function Q = gso(V)
2 % Gram-Schmidt orthonormalization of the columns of  $V \in \mathbb{R}^{n,m}$ , see
3 % (6.3.87). The vectors  $q_1, \dots, q_m$  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

```

$[Q,R] = \text{qr}(V,0)$ ← (Asymptotic computational cost: $O(m^2n)$)

dummy return value (for our purposes) dummy argument

Code 6.3.91: 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 (6.3.87)

```

6.3.4.2 Ritz projection

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 $\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$.

Code 6.3.98: one step of subspace power iteration with Ritz projection, matrix version

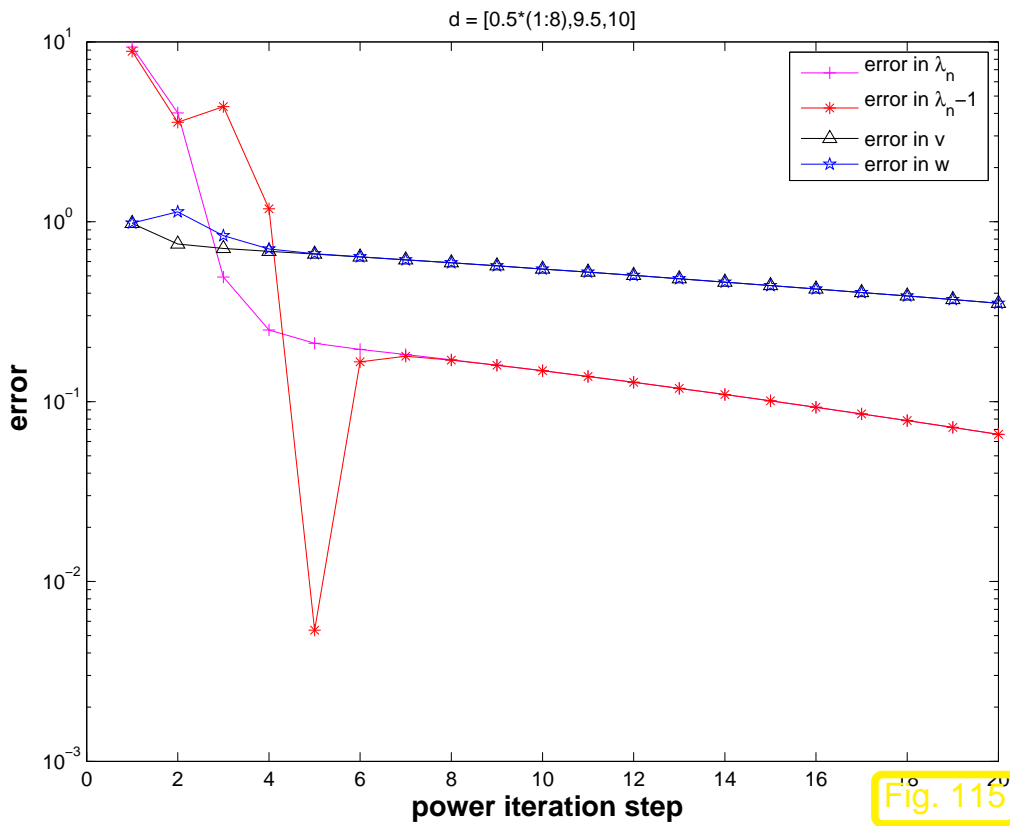
NumCSE,
autumn
2010

```
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 Sect. 6.3.4.1
4 [U,D] = eig(Q'*A*Q); % Solve Ritz projected  $m \times m$  eigenvalue problem
5 V = Q*U;          % recover approximate eigenvectors
6 ev = diag(D);     % approximate eigenvalues
```

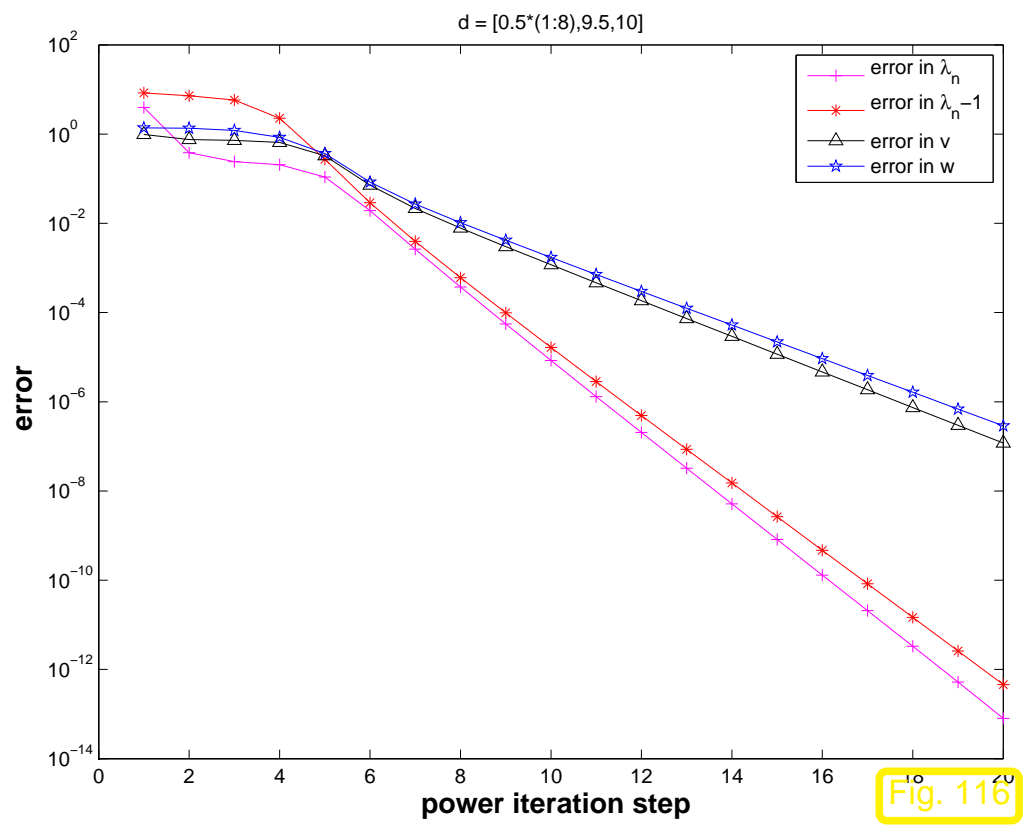
Example 6.3.99 (Power iteration with Ritz projection).

Matrix as in Ex. 6.3.83, $\sigma(\mathbf{A}) = \{0.5, 1, \dots, 4, 9.5, 10\}$:

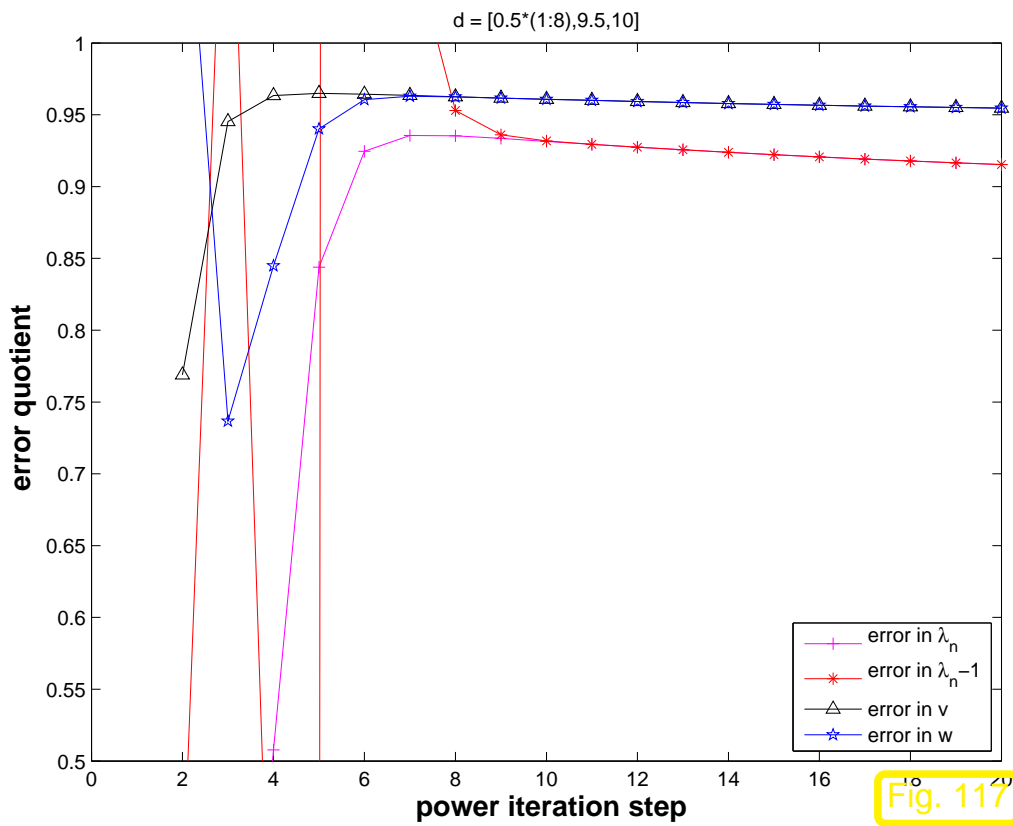
R. Hiptmair
rev 38355,
November
10, 2011



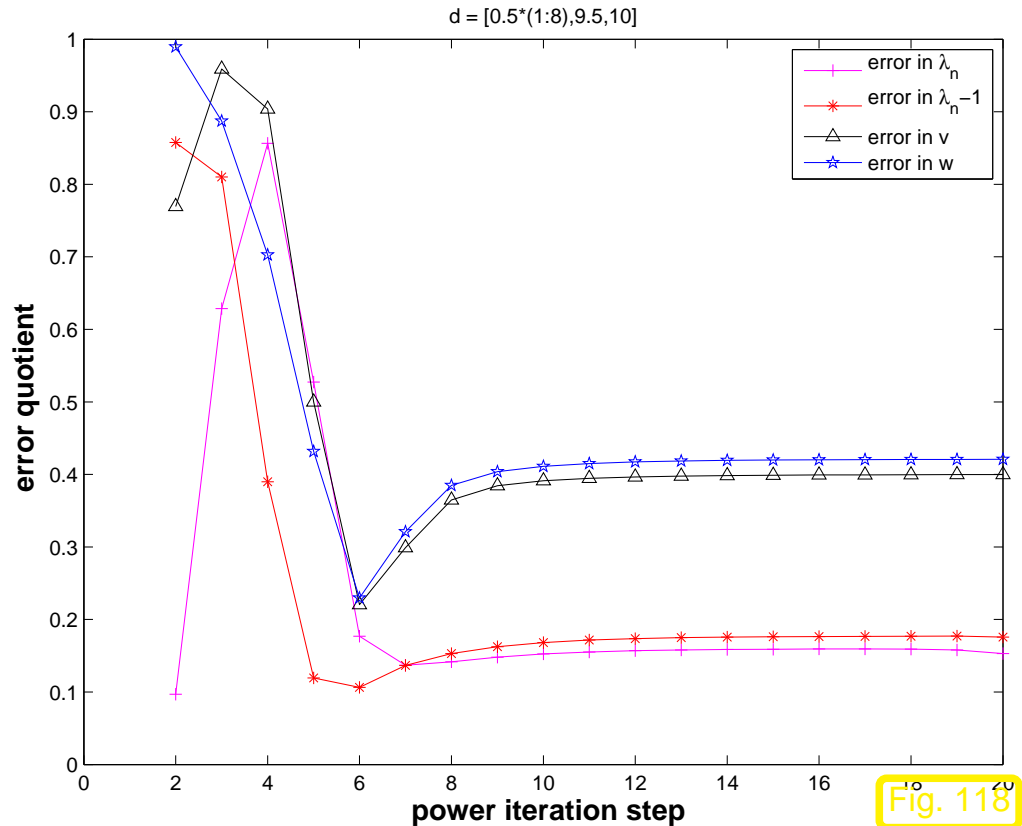
simple orthonormalization, Ex. 6.3.83



with Ritz projection, Code 6.3.97



simple orthonormalization, Ex. 6.3.83



with Ritz projection, Code 6.3.97



Example 6.3.102 (Convergence of subspace variant of direct power method).

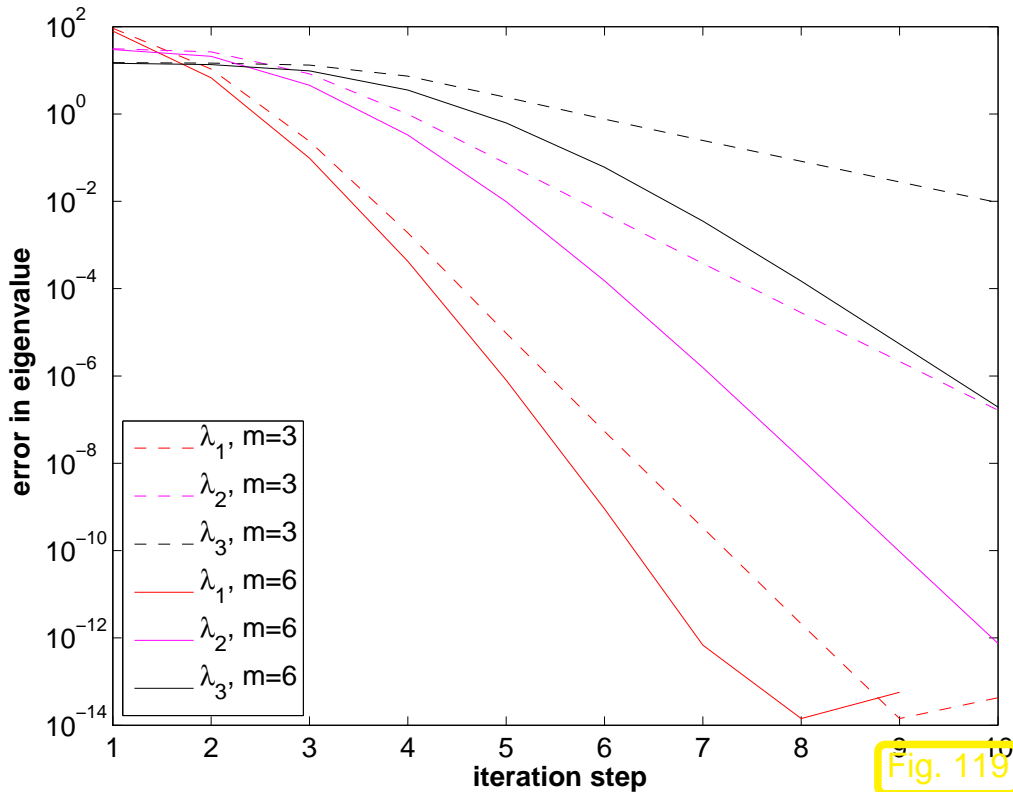


Fig. 119

S.p.d. test matrix: $a_{ij} := \min\{\frac{i}{j}, \frac{j}{i}\}$

`n=200; A = gallery('lehmer', n);`

“Initial eigenvector guesses”:

`V = eye(n, m);`

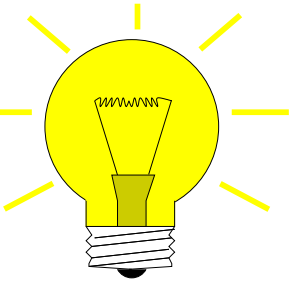
◇ R. Hiptmair
rev 38355,
November
7, 2011

Remark 6.3.103 (Subspace power methods).

Analogous to Alg. 6.3.101: construction of subspace variants of inverse iteration (\rightarrow Code 6.3.51), PINVIT (6.3.60), and Rayleigh quotient iteration (6.3.56).



6.4 Krylov Subspace Methods [35, Sect. 30]



Idea: Better $\mathbf{z}^{(k)}$ from Ritz projection onto $V := \text{Span} \{ \mathbf{z}^{(0)}, \dots, \mathbf{z}^{(k)} \}$
 (= space spanned by previous iterates)

For direct power method (6.3.11):

$$\mathbf{z}^{(k)} \parallel \mathbf{A}^k \mathbf{z}^{(0)}$$

$$V = \text{Span} \{ \mathbf{z}^{(0)}, \mathbf{A}\mathbf{z}^{(0)}, \dots, \mathbf{A}^{(k)}\mathbf{z}^{(0)} \} = \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{z}^{(0)}) \text{ a Krylov space, } \rightarrow \text{Def. 5.2.6 . (6.4.2)}$$

Code 6.4.3: Ritz projections onto Krylov space (6.4.2)

```

1 function [V,D] = kryleig(A,m)
2 % Ritz projection onto Krylov subspace. An
   orthonormal basis of  $\mathcal{K}_m(\mathbf{A}, \mathbf{1})$  is assembled into
   the columns of  $\mathbf{V}$ .
3 n = size(A,1); V = (1:n)'; V = V/norm(V);
4 for l=1:m-1
5     V = [V,A*V(:,end)]; [Q,R] = qr(V,0);
6     [W,D] = eig(Q'*A*Q); V = Q*W;
7 end
```

◁ direct power method with Ritz projection onto Krylov space from (6.4.2), cf. Alg 6.3.101.

Note: implementation for demonstration purposes only (inefficient for sparse matrix **A**!)

Example 6.4.4 (Ritz projections onto Krylov space).

```

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,  $\sigma(\mathbf{A}) = \{1,2,3,\dots,100\}$ 

```

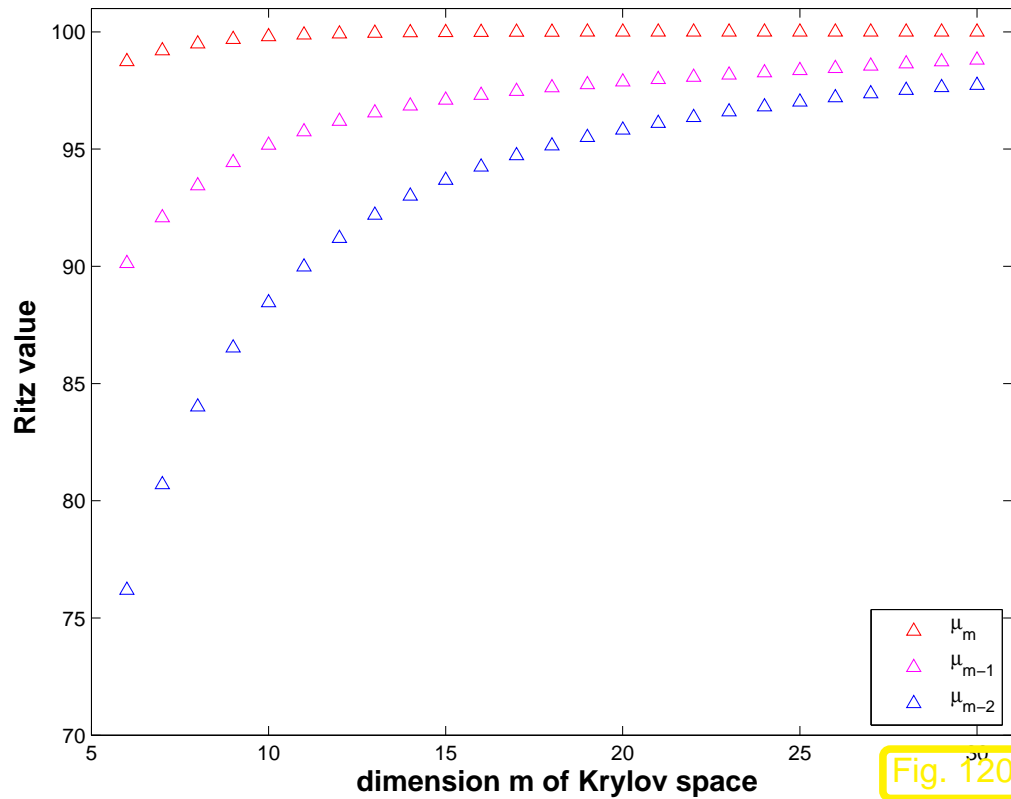


Fig. 120

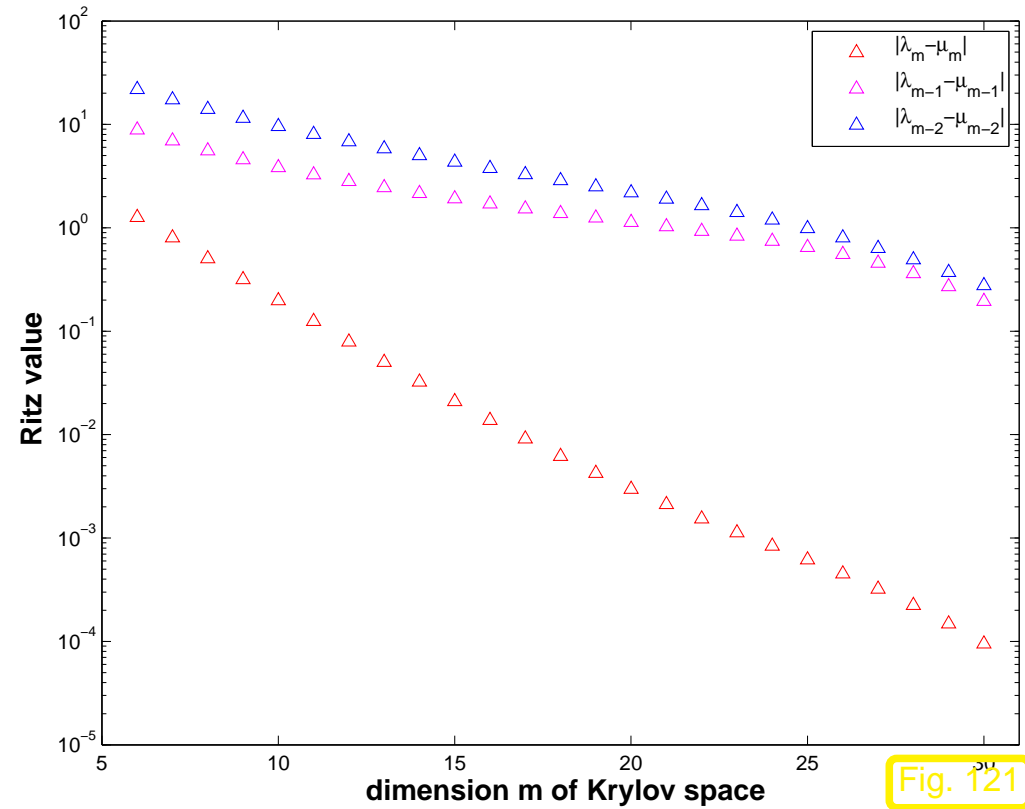
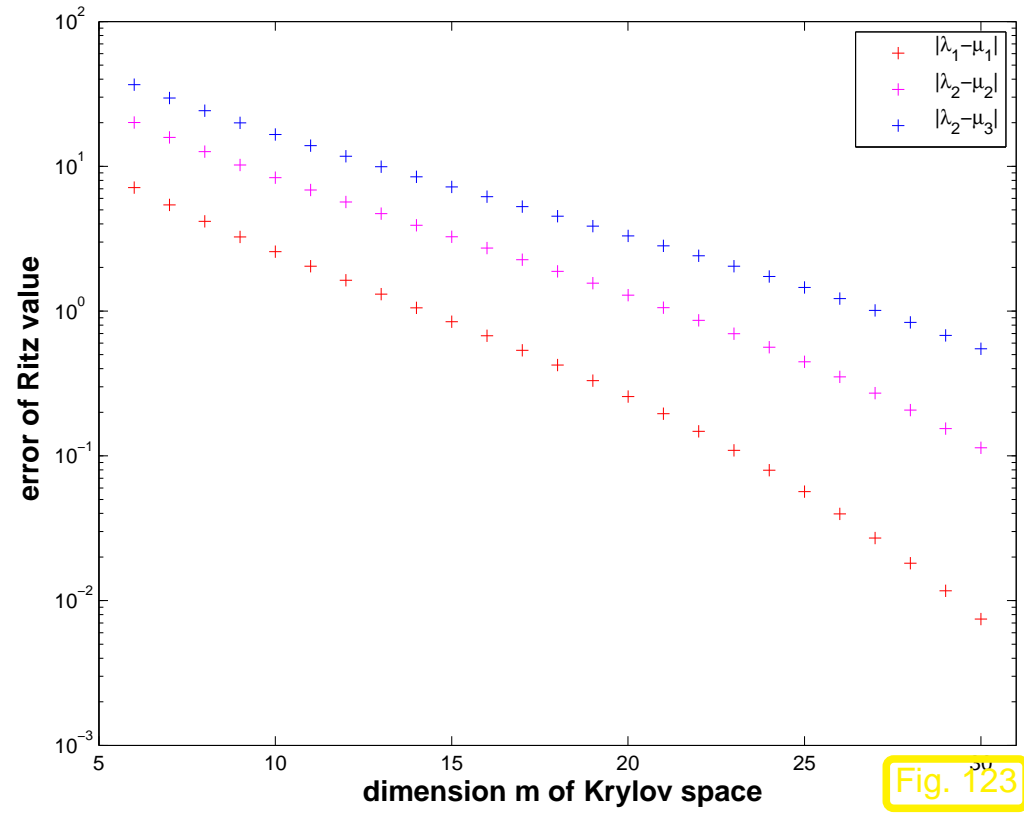
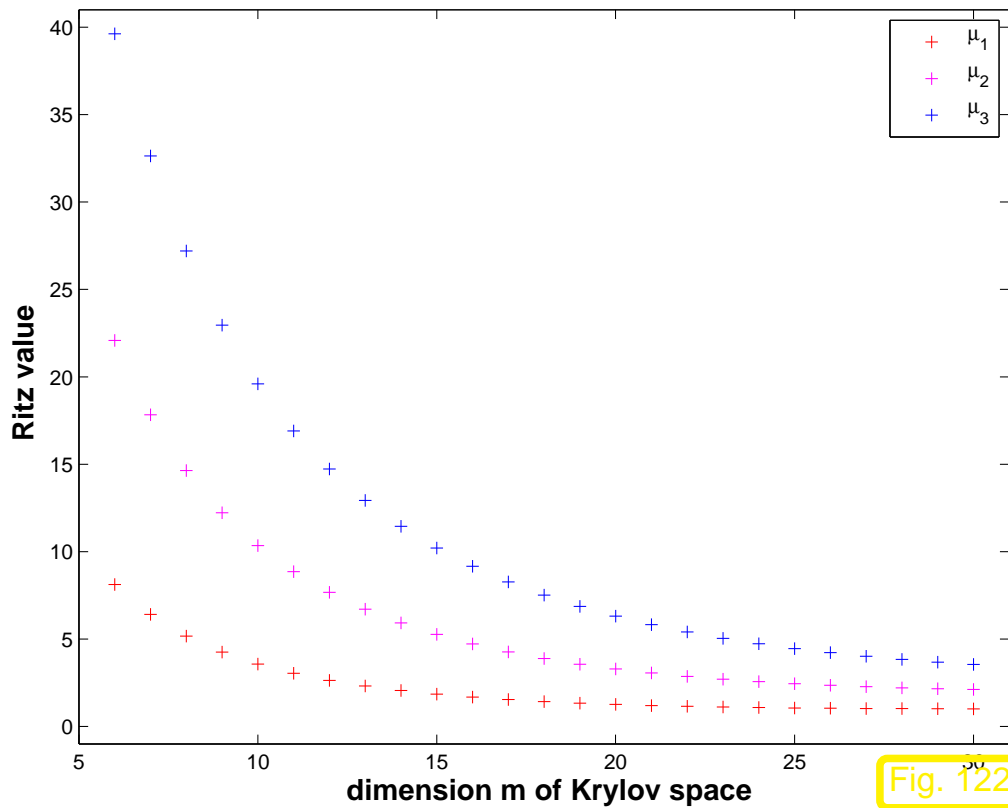


Fig. 121



Recall from Sect. 5.2.2 , Lemma 5.2.12:

Residuals $\mathbf{r}_0, \dots, \mathbf{r}_{m-1}$ generated in CG iteration, Alg. 5.2.17 applied to $\mathbf{Ax} = \mathbf{z}$ with $\mathbf{x}^{(0)} = \mathbf{0}$, provide *orthogonal basis* for $\mathcal{K}_m(\mathbf{A}, \mathbf{z})$ (, if $\mathbf{r}_k \neq \mathbf{0}$).

► Inexpensive Ritz projection of $\mathbf{Ax} = \lambda \mathbf{x}$ onto $\mathcal{K}_m(\mathbf{A}, \mathbf{z})$:

orthogonal matrix

$$\mathbf{V}_m^T \mathbf{A} \mathbf{V}_m \mathbf{x} = \lambda \mathbf{x}, \quad \mathbf{V}_m := \left(\frac{\mathbf{r}_0}{\|\mathbf{r}_0\|}, \dots, \frac{\mathbf{r}_{m-1}}{\|\mathbf{r}_{m-1}\|} \right) \in \mathbb{R}^{n,m}. \quad (6.4.8)$$

recall: residuals generated by *short recursions*, see Alg. 5.2.17

Lemma 6.4.9 (Tridiagonal Ritz projection from CG residuals).

$\mathbf{V}_m^T \mathbf{A} \mathbf{V}_m$ is a tridiagonal matrix.

$$\mathbf{V}_l^H \mathbf{A} \mathbf{V}_l = \begin{pmatrix} \alpha_1 & \beta_1 & & & & \\ \beta_1 & \alpha_2 & \beta_2 & & & \\ & \beta_2 & \alpha_3 & \ddots & & \\ & & \ddots & \ddots & & \\ & & & & \ddots & \\ & & & \ddots & \ddots & \beta_{k-1} \\ & & & & \beta_{k-1} & \alpha_k \end{pmatrix} =: \mathbf{T}_l \in \mathbb{K}^{k,k} \quad \text{[tridiagonal matrix]} \quad (6.4.10)$$

Algorithm for computing V_l and T_l :

Lanczos process

Computational effort/step:

- 1 × $A \times$ vector
- 2 dot products
- 2 AXPY-operations
- 1 division

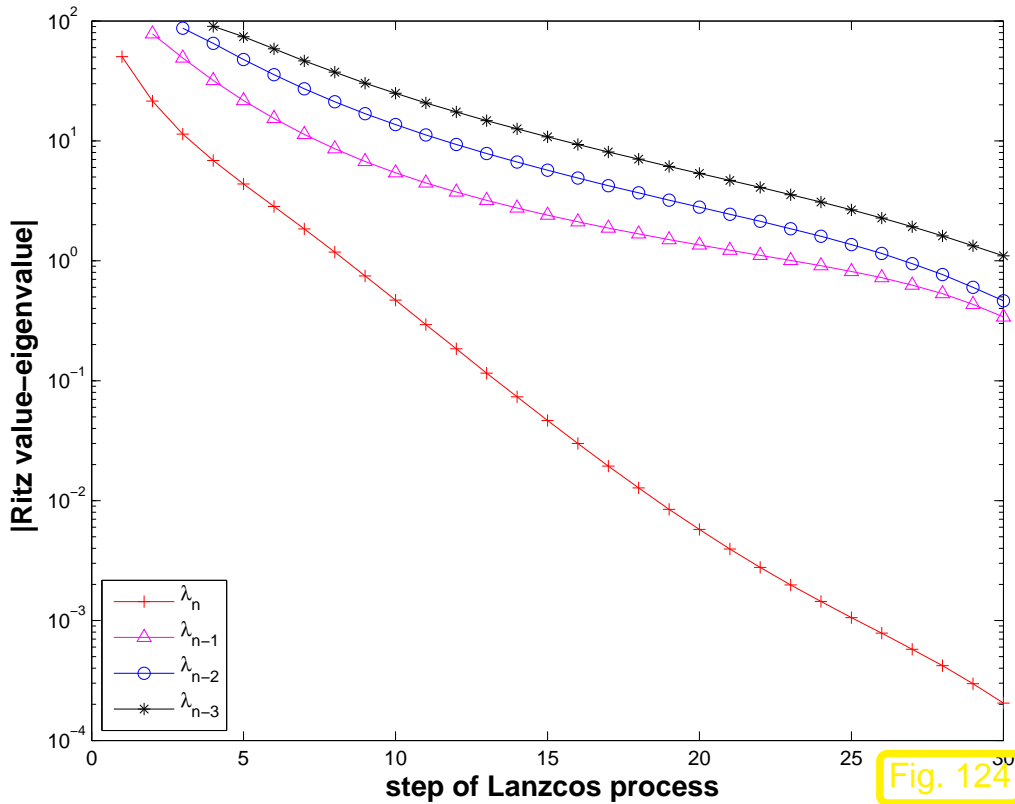
```

1 function [V,alph,bet] = lanczos(A,k,z0)
2 % Note: this implementation of the Lanczos process
   also records the orthonormal CG residuals in the
   columns of the matrix V, which is not needed when
   only eigenvalue approximations are desired.
3 V = z0/norm(z0);
4 % Vectors storing entries of tridiagonal matrix (6.4.10)
5 alph=zeros(k,1); bet = zeros(k,1);
6 for j=1:k
7     q = A*V(:,j); alph(j) = dot(q,V(:,j));
8     w = q - alph(j)*V(:,j);
9     if (j > 1), w = w - bet(j-1)*V(:,j-1); end
10    bet(j) = norm(w); V = [V,w/bet(j)];
11 end
12 bet = bet(1:end-1);

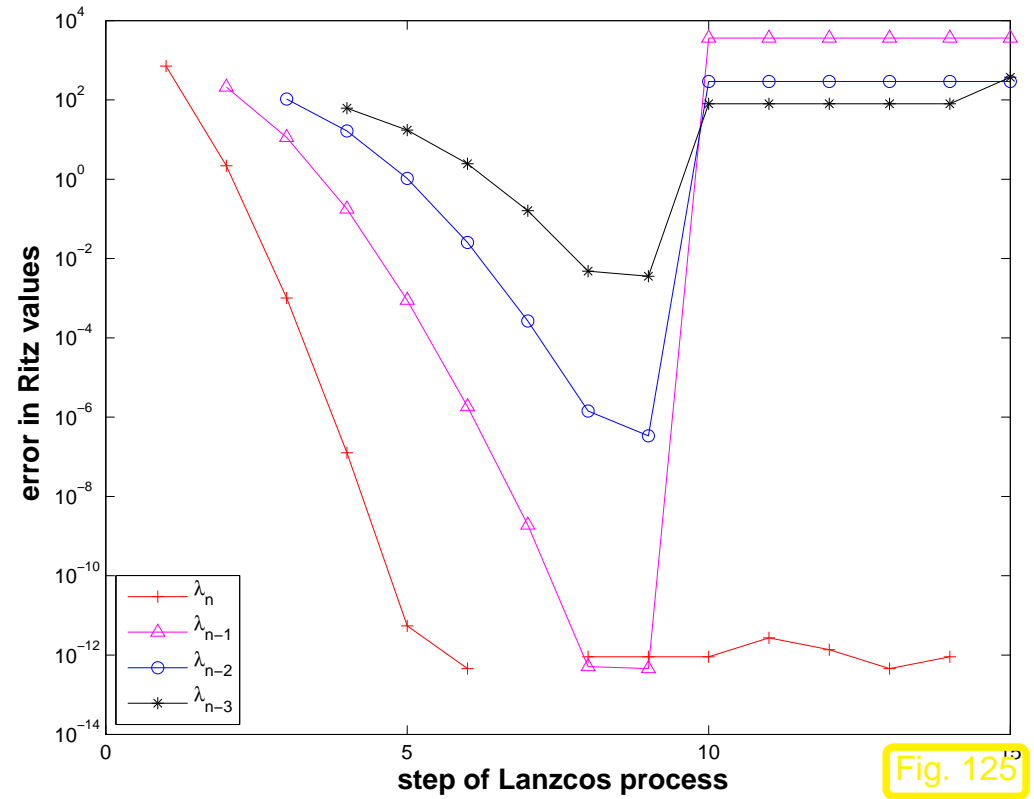
```

Example 6.4.12 (Lanczos process for eigenvalue computation).

A from Ex. 6.4.4



A = gallery('minij',100);



Example 6.4.13 (Impact of roundoff on Lanczos process).

$$A \in \mathbb{R}^{10,10}, \quad a_{ij} = \min\{i, j\}.$$

```
A = gallery('minij',10);
```

Computed by `[V,alpha,beta] = lanczos(A,n,ones(n,1));`, see Code 6.4.10:

$$\mathbf{T} = \begin{pmatrix} 38.500000 & 14.813845 & & & & & & & & & \\ 14.813845 & 9.642857 & 2.062955 & & & & & & & & \\ & 2.062955 & 2.720779 & 0.776284 & & & & & & & \\ & & 0.776284 & 1.336364 & 0.385013 & & & & & & \\ & & & 0.385013 & 0.826316 & 0.215431 & & & & & \\ & & & & 0.215431 & 0.582380 & 0.126781 & & & & \\ & & & & & 0.126781 & 0.446860 & 0.074650 & & & \\ & & & & & & 0.074650 & 0.363803 & 0.043121 & & \\ & & & & & & & 0.043121 & 3.820888 & 11.991094 & \\ & & & & & & & & 11.991094 & 41.254286 & \end{pmatrix}$$

$$\sigma(\mathbf{A}) = \{0.255680, 0.273787, 0.307979, 0.366209, 0.465233, 0.643104, 1.000000, 1.873023, 5.048917, 44.766069\}$$

$$\sigma(\mathbf{T}) = \{0.263867, 0.303001, 0.365376, 0.465199, 0.643104, 1.000000, 1.873023, 5.048917, 44.765976, 44.766069\}$$

▶ Uncanny cluster of computed eigenvalues of \mathbf{T} (“ghost eigenvalues”, [23, Sect. 9.2.5])

NumCSE,
autumn
2010

$$\mathbf{V}^H \mathbf{V} = \begin{pmatrix} 1.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000251 & 0.258801 & 0.883711 \\ 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000106 & 0.109470 & 0.373799 \\ 0.000000 & -0.000000 & 1.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000005 & 0.005373 & 0.018347 \\ 0.000000 & 0.000000 & 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000096 & 0.000328 \\ 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000001 & 0.000003 \\ 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000000 \\ 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 \\ 0.000251 & 0.000106 & 0.000005 & 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 \\ 0.258801 & 0.109470 & 0.005373 & 0.000096 & 0.000001 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & 0.000000 \\ 0.883711 & 0.373799 & 0.018347 & 0.000328 & 0.000003 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 1.000000 \end{pmatrix}$$

| 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 |

R. Hiptmair
rev 38355,
November
7, 2011

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*.

MATLAB-functions:

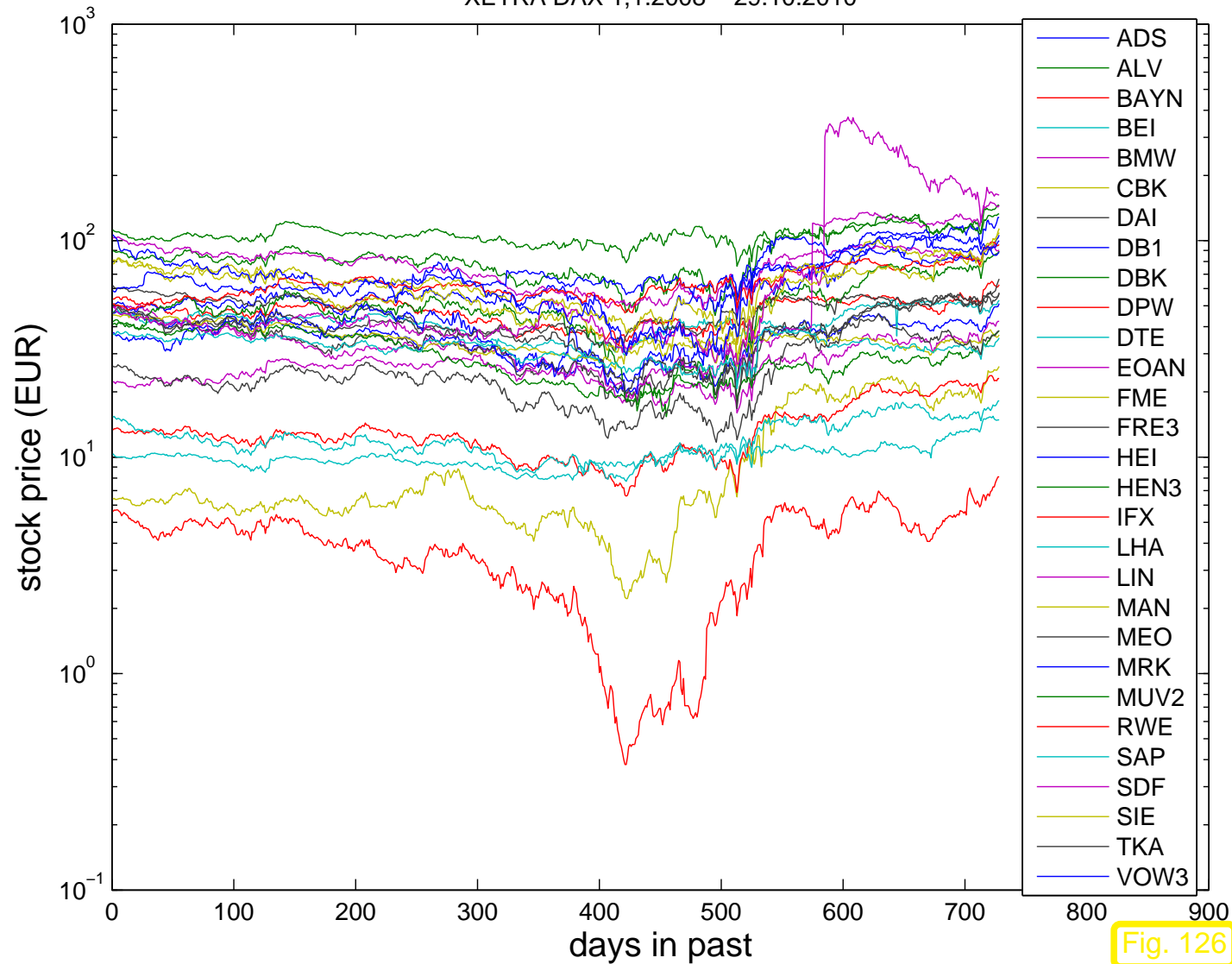
`d = eigs(A,k,sigma)` : k largest/smallest eigenvalues of \mathbf{A}

`d = eigs(A,B,k,sigma)`: k largest/smallest eigenvalues for generalized EVP $\mathbf{Ax} = \lambda\mathbf{Bx}$, \mathbf{B} s.p.d.

`d = eigs(Afun,n,k)` : \mathbf{Afun} = handle to function providing matrix \times vector for $\mathbf{A}/\mathbf{A}^{-1}/\mathbf{A} - \alpha\mathbf{I}/(\mathbf{A} - \alpha\mathbf{B})^{-1}$. (Use flags to tell `eigs` about special properties of matrix behind \mathbf{Afun} .)

6.5 Singular Value Decomposition

Example 6.5.1 (Trend analysis).



◁ (end of day) stock prizes

Are there underlying governing trends ?



Example 6.5.2 (Classification from measured data).

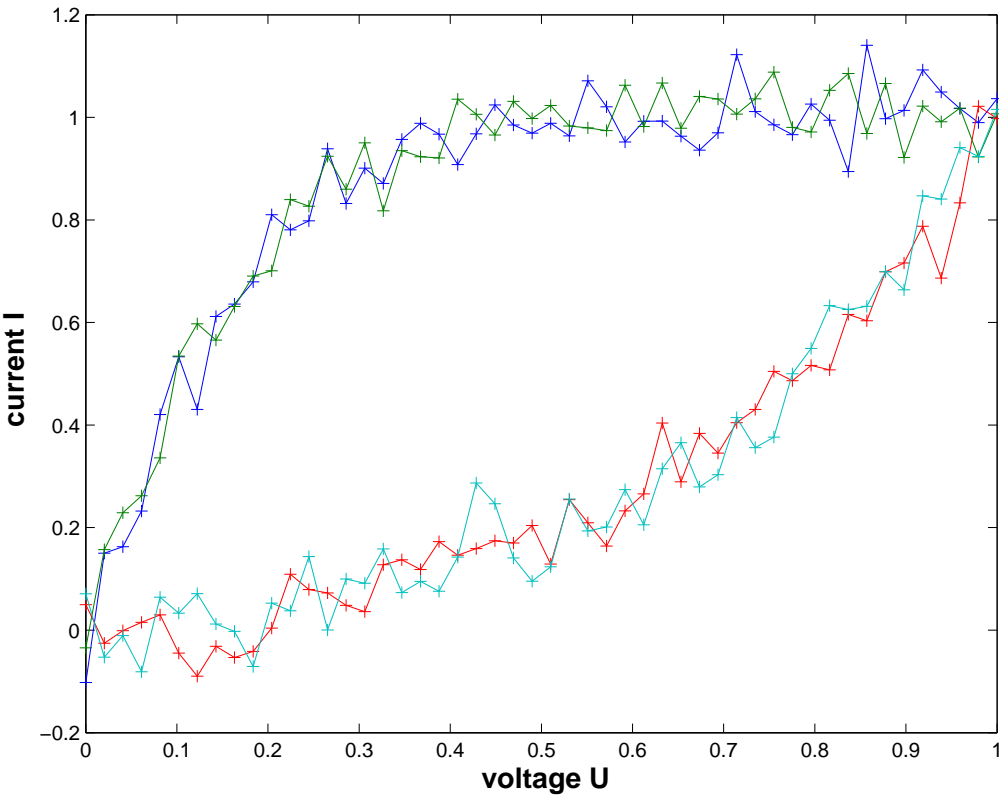


Given: measured U - I characteristics of diodes

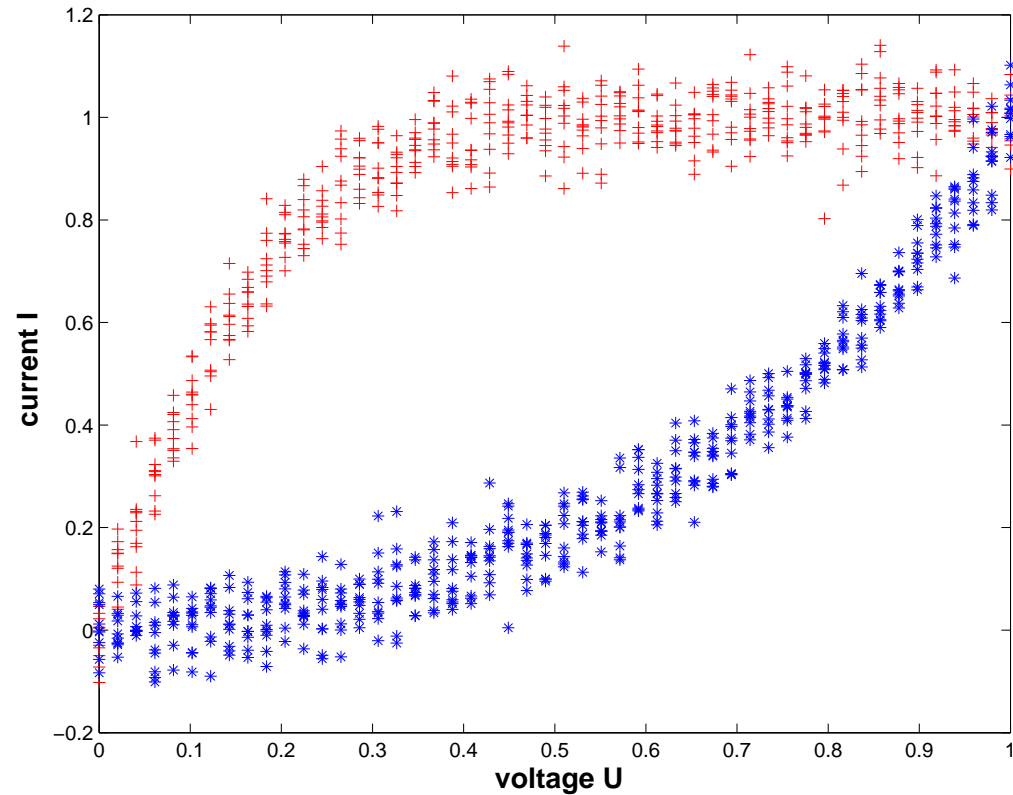
Find out

- how many types,
- the U - I characteristic of each type.

measured U-I characteristics for some diodes



measured U-I characteristics for all diodes



Theorem 6.5.6 (singular value decomposition). $\rightarrow [48, \text{Thm. 9.6}], [27, \text{Thm. 11.1}]$

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 $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_p) \in \mathbb{R}^{m,n}$, $p := \min\{m, n\}$, $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ such that

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H.$$

(*): $\mathbf{\Sigma}$ (generalized) diagonal matrix $\Leftrightarrow (\mathbf{\Sigma})_{i,j} = 0$, if $i \neq j$, $1 \leq i \leq m$, $1 \leq j \leq n$.

$$\begin{pmatrix} \boxed{A} \end{pmatrix} = \begin{pmatrix} \boxed{U} \end{pmatrix} \begin{pmatrix} \boxed{\Sigma} \end{pmatrix} \begin{pmatrix} \boxed{V^H} \end{pmatrix}$$

$$\begin{pmatrix} \boxed{A} \end{pmatrix} = \begin{pmatrix} \boxed{U} \end{pmatrix} \begin{pmatrix} \boxed{\Sigma} \end{pmatrix} \begin{pmatrix} \boxed{V^H} \end{pmatrix}$$

Definition 6.5.8 (Singular value decomposition (SVD)).

The decomposition $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H$ of Thm. 6.5.6 is called *singular value decomposition (SVD)* of \mathbf{A} . The diagonal entries σ_i of $\mathbf{\Sigma}$ are the *singular values* of \mathbf{A} .

MATLAB-functions (for algorithms see [23, Sect. 8.3]):

- $\mathbf{s} = \text{svd}(\mathbf{A})$: computes singular values of matrix \mathbf{A}
- $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{A})$: computes singular value decomposition according to Thm. 6.5.6
- $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{A}, 0)$: “economical” singular value decomposition for $m > n$: : $\mathbf{U} \in \mathbb{K}^{m,n}$, $\mathbf{\Sigma} \in \mathbb{R}^{n,n}$, $\mathbf{V} \in \mathbb{K}^{n,n}$
- $\mathbf{s} = \text{svds}(\mathbf{A}, k)$: k largest singular values (important for sparse $\mathbf{A} \rightarrow$ Def. 2.6.1)
- $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svds}(\mathbf{A}, k)$: partial singular value decomposition: $\mathbf{U} \in \mathbb{K}^{m,k}$, $\mathbf{V} \in \mathbb{K}^{n,k}$, $\mathbf{\Sigma} \in \mathbb{R}^{k,k}$ diagonal with k largest singular values of \mathbf{A} .

Explanation: “economical” singular value decomposition:

$$\begin{pmatrix} \text{A} \end{pmatrix} = \begin{pmatrix} \text{U} \end{pmatrix} \begin{pmatrix} \Sigma \end{pmatrix} \begin{pmatrix} \text{V}^H \end{pmatrix}$$

(MATLAB) algorithm for computing SVD is (numerically) stable \rightarrow Def. 2.5.11

Complexity:

$$2mn^2 + 2n^3 + O(n^2) + O(mn) \quad \text{for } \mathbf{s} = \text{svd}(\mathbf{A}),$$

$$4m^2n + 22n^3 + O(mn) + O(n^2) \quad \text{for } [\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{A}),$$

$$O(mn^2) + O(n^3) \quad \text{for } [\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{A}, 0), m \gg n.$$

- Application of SVD: computation of rank (\rightarrow Def. 2.0.2), kernel and range of a matrix

Illustration:

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U} \end{pmatrix} \begin{pmatrix} \Sigma_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{V}^H \end{pmatrix} \quad (6.5.15)$$

columns = ONB of $\text{Im}(\mathbf{A})$

rows = ONB of $\text{Ker}(\mathbf{A})$

Code 6.5.16: Computing an ONB of the kernel of a matrix

```

1 function V = kerncomp(A,tol)
2 % computes an orthonormal basis of Ker(A) using Lemma 6.5.14
3 % kernel selection with relative tolerance tol
4 if (nargin < 2), tol = eps; end
5 [U,S,V] = svd(A); % Singular value decomposition
6 s = diag(S); % Extract vector of singular values
7 % find singular values of relative (w.r.t.  $\sigma_1$ ) size  $\leq$  tol
8 r = min(find(s/s(1) <= tol)); % "Numerical rank" +1
9 V = V(:,r:end); % rightmost columns of V

```

$$\begin{pmatrix} \\ \\ \\ \\ \end{pmatrix} = \sigma_1 \begin{pmatrix} \\ \mathbf{u}_1 \\ \\ \\ \end{pmatrix} \left(\begin{array}{|c|} \hline \\ \hline \mathbf{v}_1^\top \\ \hline \\ \hline \\ \hline \\ \hline \end{array} \right) + \sigma_2 \begin{pmatrix} \\ \\ \mathbf{u}_2 \\ \\ \end{pmatrix} \left(\begin{array}{|c|} \hline \\ \hline \\ \hline \mathbf{v}_2^\top \\ \hline \\ \hline \\ \hline \\ \hline \end{array} \right) + \dots$$

$$\begin{pmatrix} \\ \\ \\ \\ \end{pmatrix} = \sigma_1 \begin{pmatrix} \\ \mathbf{u}_1 \\ \\ \\ \end{pmatrix} \left(\begin{array}{|c|} \hline \\ \hline \mathbf{v}_1^\top \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \end{array} \right) + \sigma_2 \begin{pmatrix} \\ \\ \mathbf{u}_2 \\ \\ \end{pmatrix} \left(\begin{array}{|c|} \hline \\ \hline \\ \hline \mathbf{v}_2^\top \\ \hline \\ \hline \\ \hline \\ \hline \end{array} \right) + \dots$$

Example 6.5.18 (Principal component analysis for data classification).

→ Ex. 6.5.2 cnt'd

measured U-I characteristics for some diodes

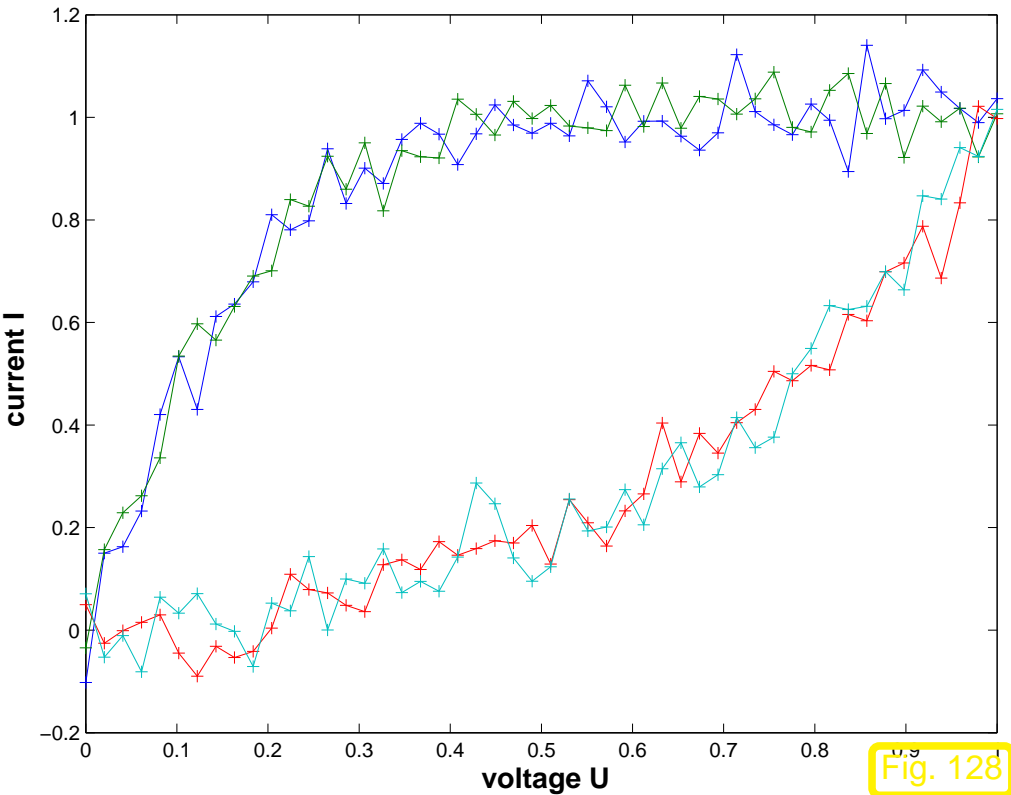


Fig. 128

measured U-I characteristics for all diodes

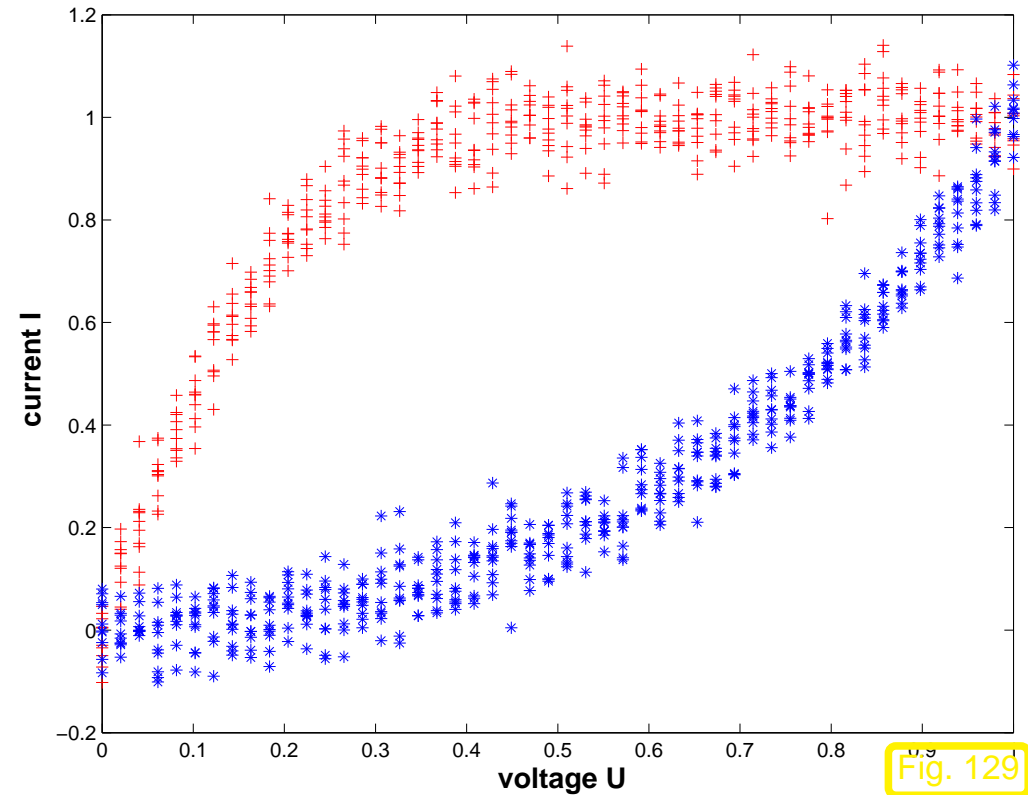


Fig. 129

Code 6.5.19: Generation of synthetic perturbed $U-I$ characteristics

```

1 % Generate synthetic measurement curves with random multiplicative and additive
2 % perturbations supposed to reflect measurement errors and manufacturing
  tolerances
3
4 % Voltage-current characteristics of both kinds of diodes
5 i1 = @(u) (2./(1+exp(-10*u)) - 1);
6 i2 = @(u) ((exp(3*u)-1)/(exp(3)-1));
7 % Simulated measurements
8 m = 50; % number of measurements for different input voltages
9 n = 10; % no. of diodes of each kind
0 na = 0.05; % level of additive noise (normally distributed)

```

```

1 nm = 0.02; % level of multiplicative noise (normally distributed)
2
3 uvals = (0:1/(m-1):1);
4 D1 = (1+nm*randn(n,m)).*(i1 repmat(uvals,n,1))+na*randn(n,m);
5 D2 = (1+nm*randn(n,m)).*(i2 repmat(uvals,n,1))+na*randn(n,m);
6 A = ([D1;D2])'; A = A(1:size(A,1),randperm(1:size(A,2)));

```

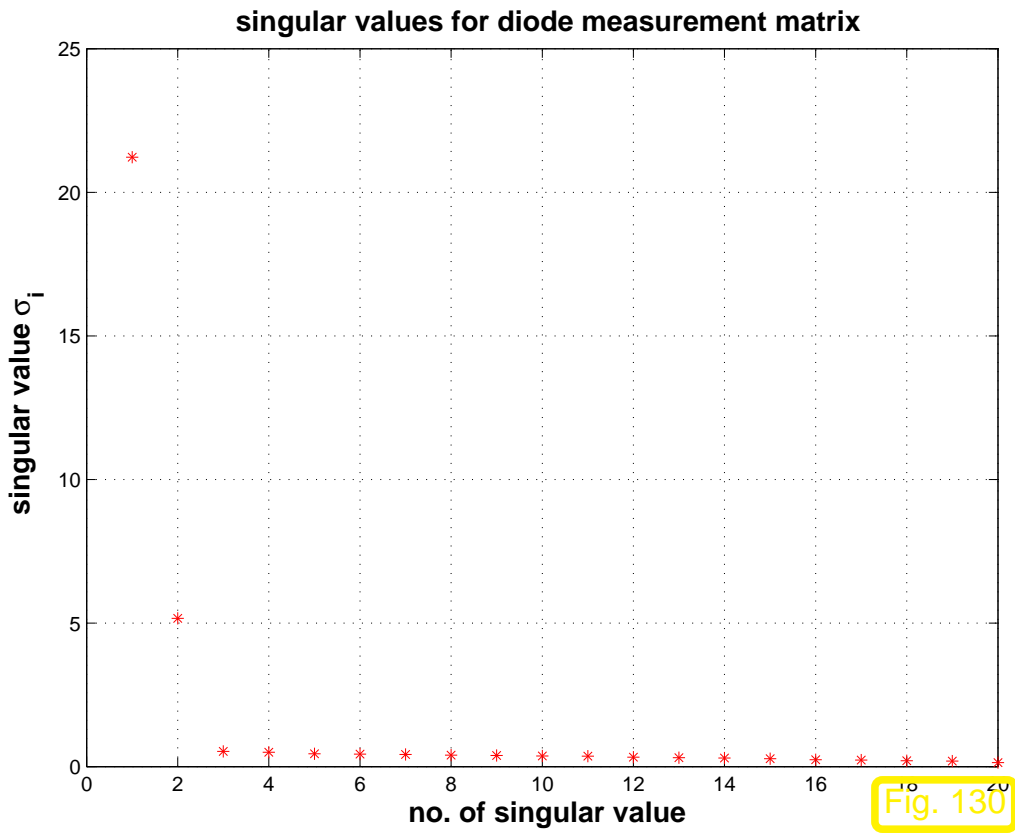


Fig. 130

← distribution of singular values of matrix

two dominant singular values !



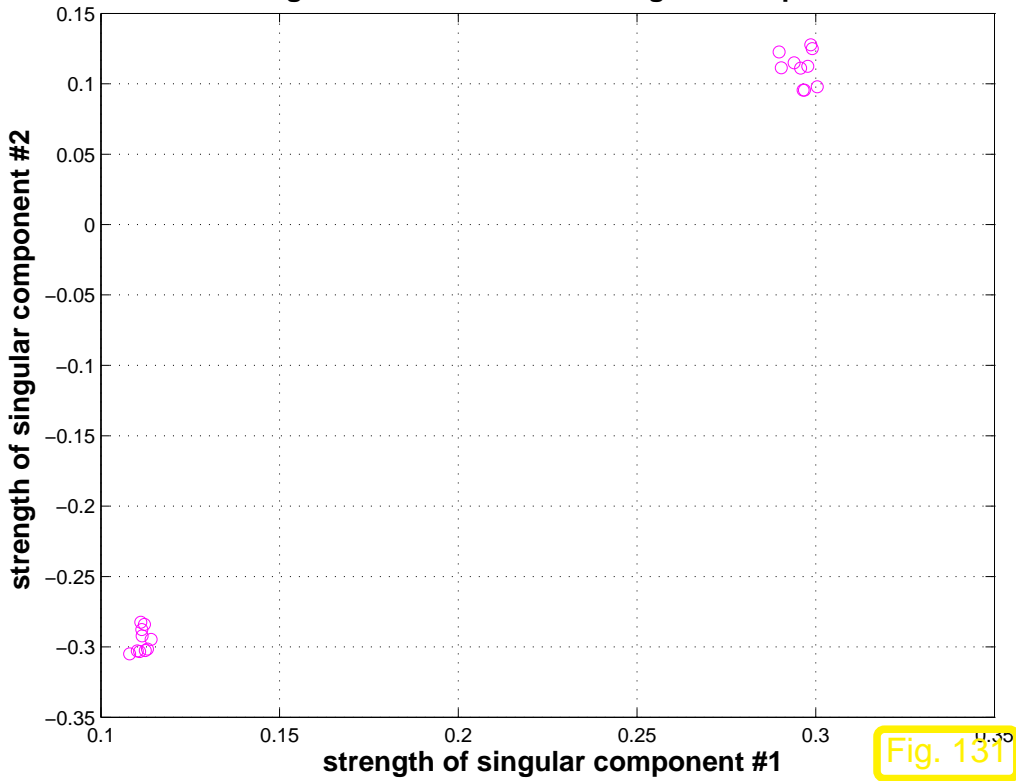
measurements display linear correlation with **two**
principal components



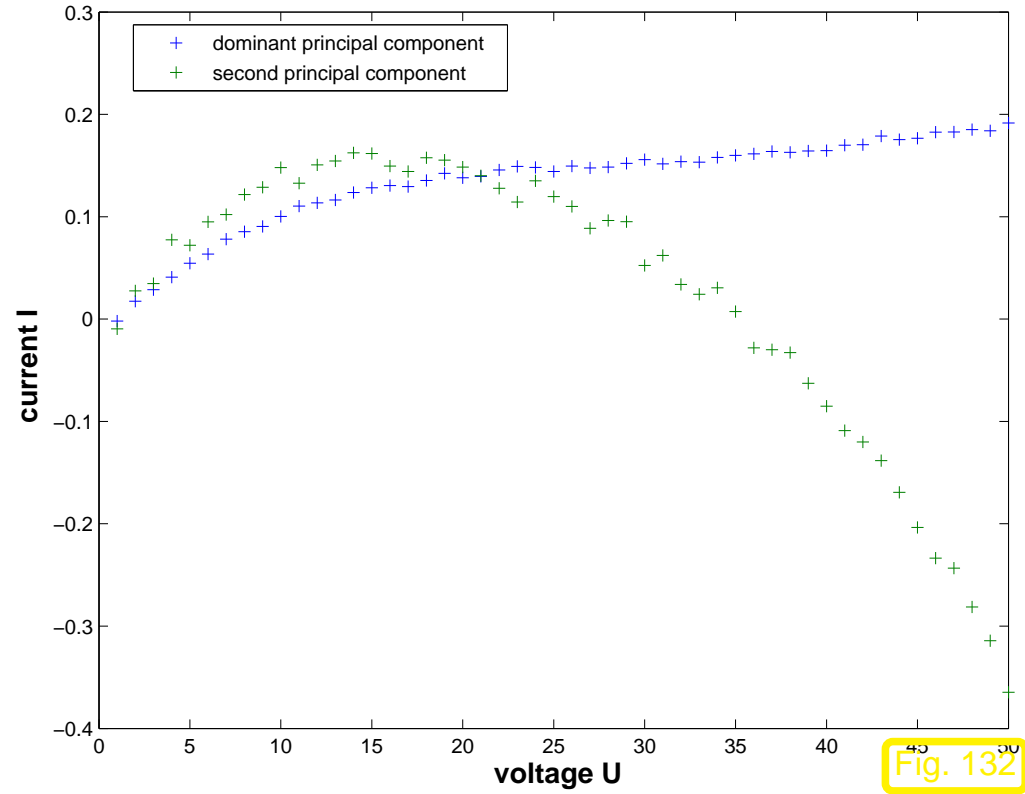
two types of diodes in batch



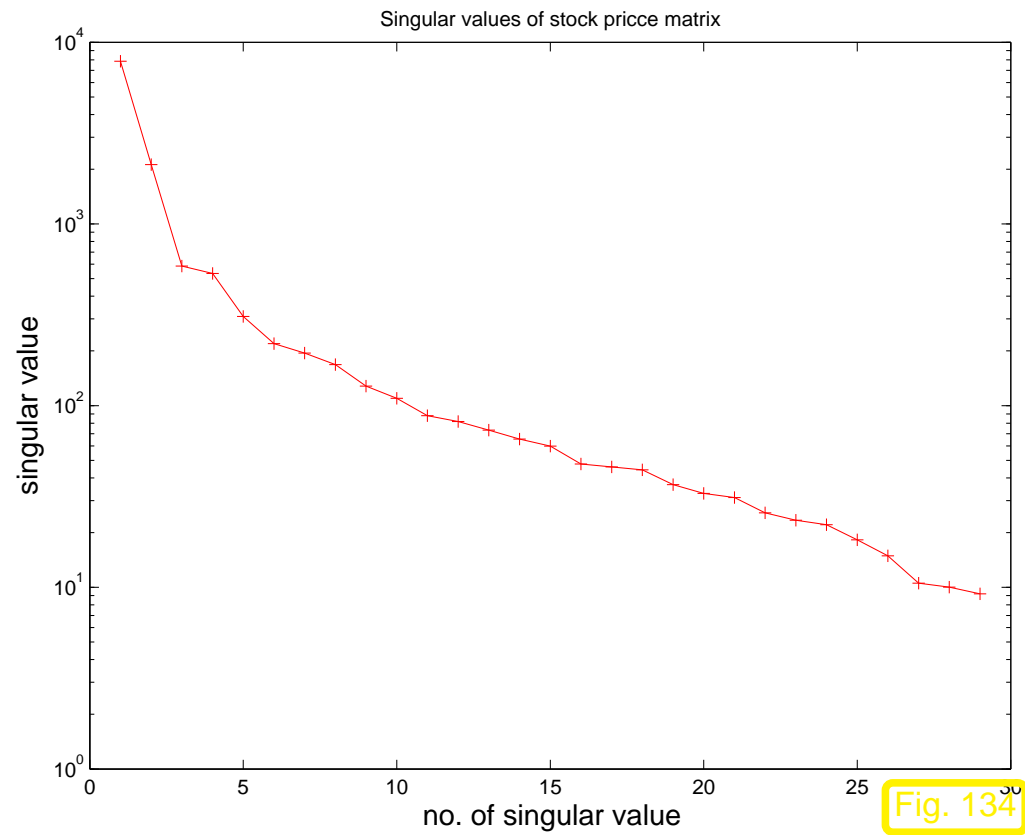
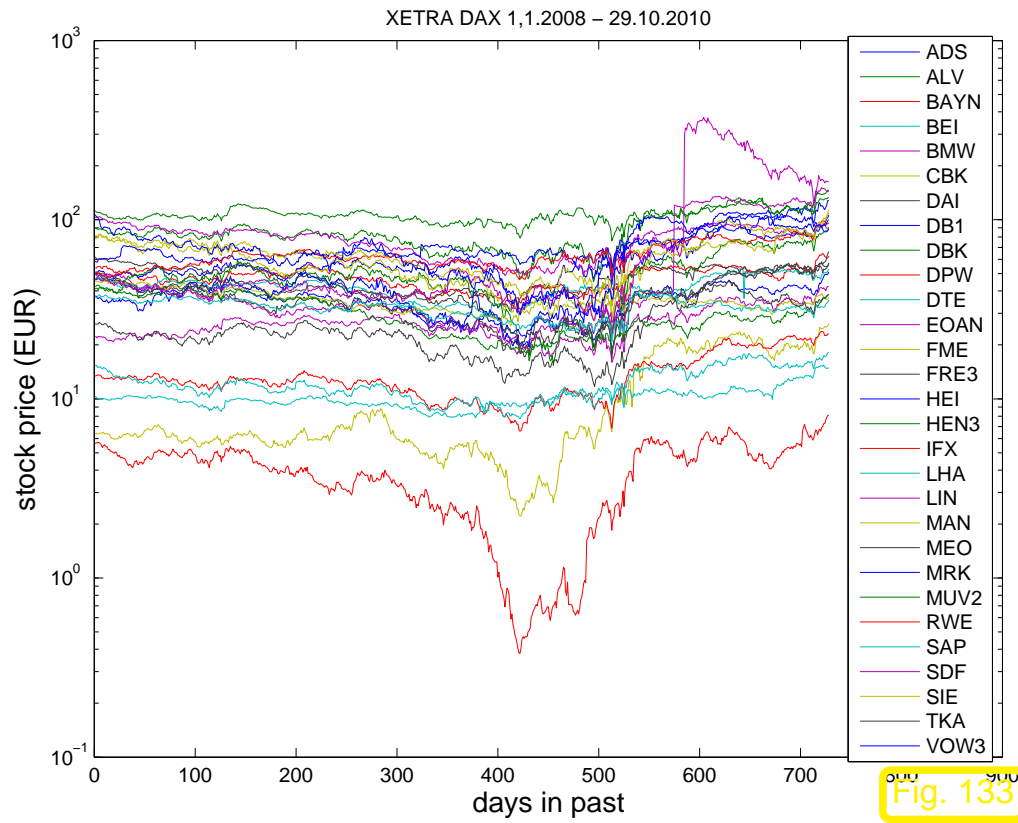
strengths of contributions of singular components

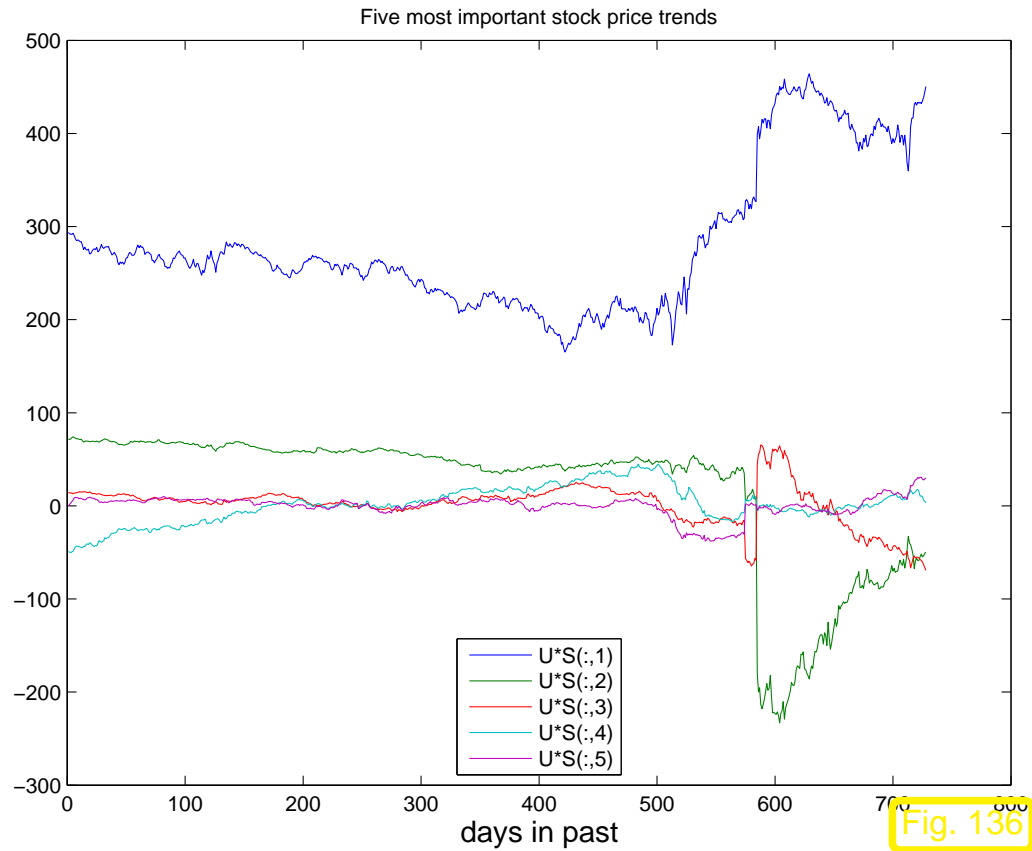
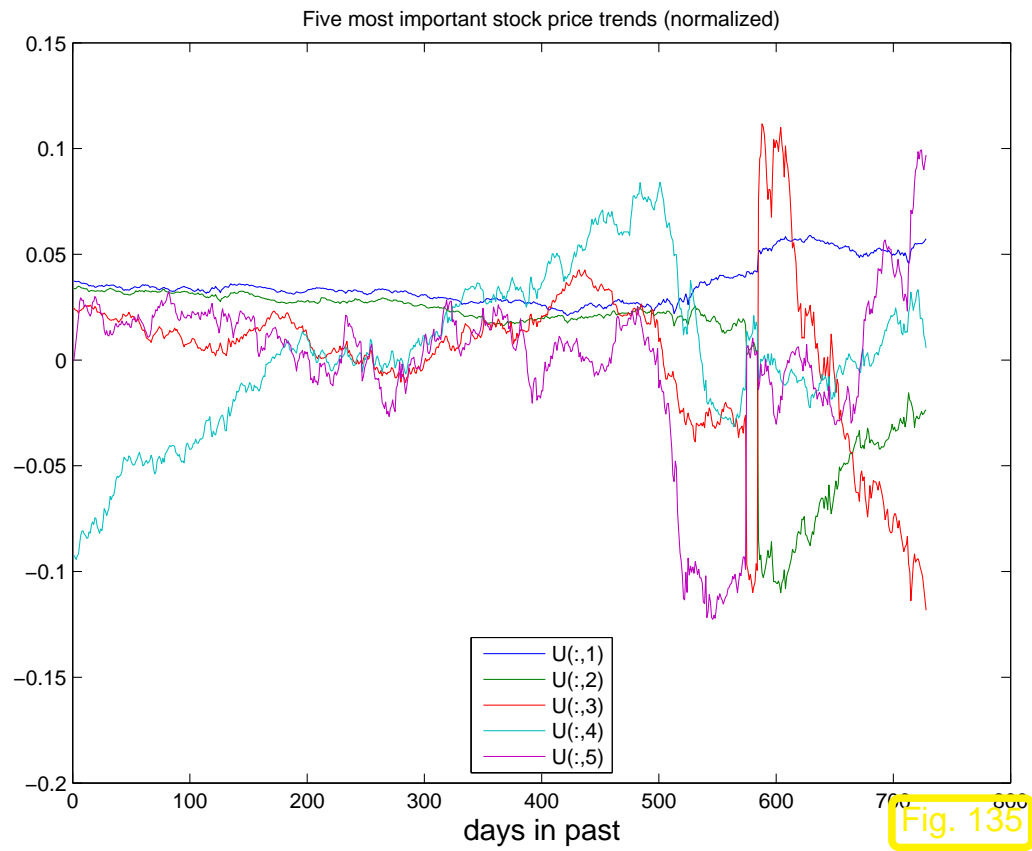


principal components (trend vectors) for diode measurements



Example 6.5.29 (PCA of stock prices). → Ex. 6.5.1







Trends in BMW stock, 1.1.2008 – 29.10.2010

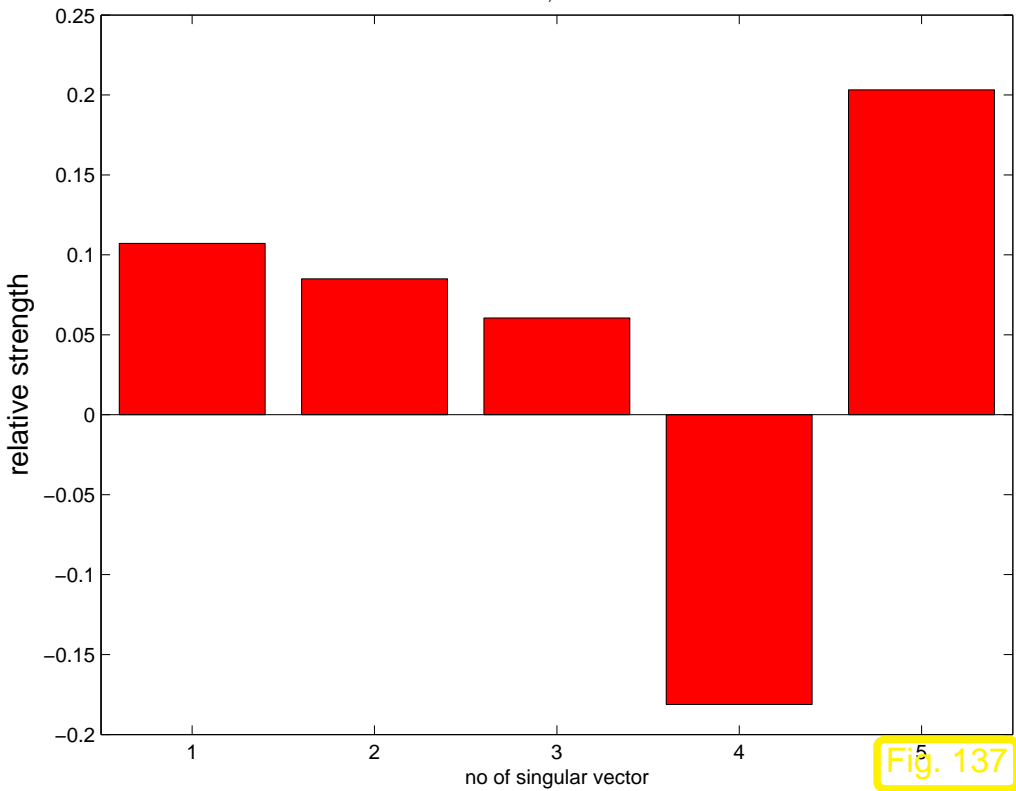


Fig. 137

Trends in Daimler stock, 1.1.2008 – 29.10.2010

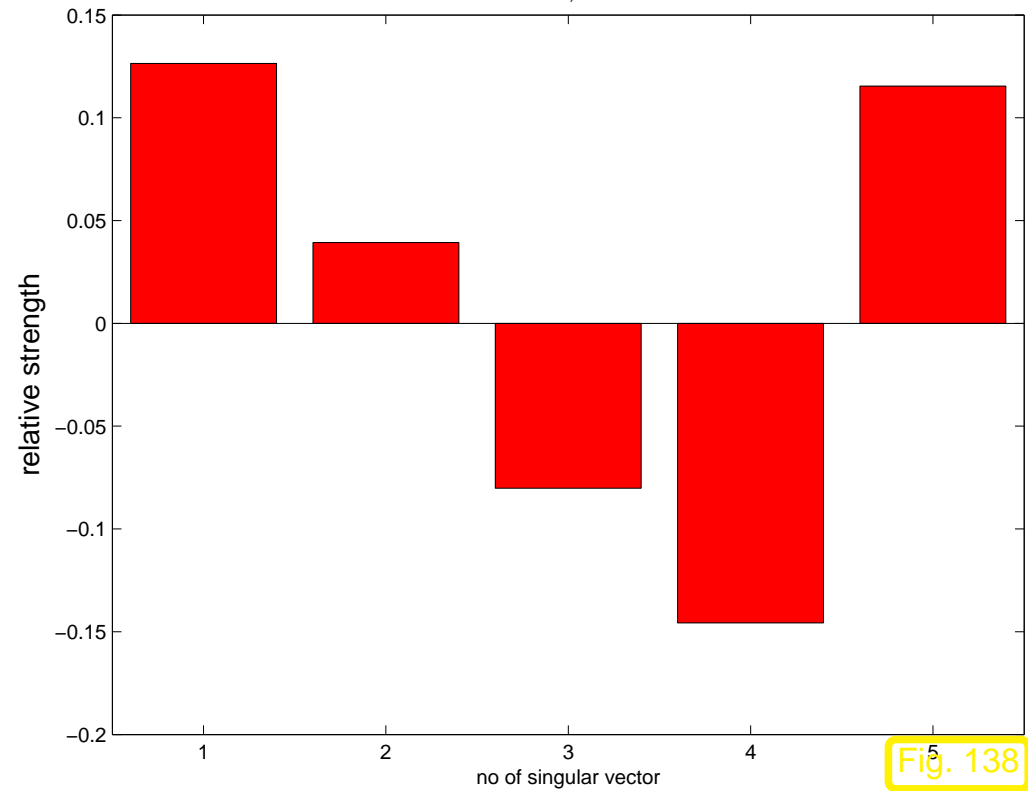


Fig. 138

- Application of SVD: extrema of quadratic forms on the unit sphere
- Application of SVD: *best low rank approximation*

Theorem 6.5.36 (best low rank approximation). \rightarrow [27, Thm. 11.6]

Let $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H$ be the SVD of $\mathbf{A} \in \mathbb{K}^{m,n}$ (\rightarrow Thm. 6.5.6). For $1 \leq k \leq \text{rank}(\mathbf{A})$ set $\mathbf{U}_k := [\mathbf{u}_{\cdot,1}, \dots, \mathbf{u}_{\cdot,k}] \in \mathbb{K}^{m,k}$, $\mathbf{V}_k := [\mathbf{v}_{\cdot,1}, \dots, \mathbf{v}_{\cdot,k}] \in \mathbb{K}^{n,k}$, $\mathbf{\Sigma}_k := \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbb{K}^{k,k}$. Then, for $\|\cdot\| = \|\cdot\|_F$ and $\|\cdot\| = \|\cdot\|_2$, holds true

$$\left\| \mathbf{A} - \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^H \right\| \leq \|\mathbf{A} - \mathbf{F}\| \quad \forall \mathbf{F} \in \mathcal{R}_k(m, n).$$

Example 6.5.39 (Image compression).

Image composed of $m \times n$ pixels (greyscale, BMP format) \leftrightarrow matrix $\mathbf{A} \in \mathbb{R}^{m,n}$, $a_{ij} \in \{0, \dots, 255\}$, see Ex. 6.3.22

► Thm. 6.5.37 \triangleright best low rank approximation of image: $\tilde{\mathbf{A}} = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}^T$

Code 6.5.40: SVD based image compression

```

1 P = double(imread('eth.pbm'));
2 [m,n] = size(P); [U,S,V] = svd(P); s = diag(S);
3 k = 40; S(k+1:end,k+1:end) = 0; PC = U*S*V';
4
5 figure('position',[0 0 1600 1200]); col = [0:1/215:1]'*[1,1,1];
6 subplot(2,2,1); image(P); title('original image'); colormap(col);

```

```

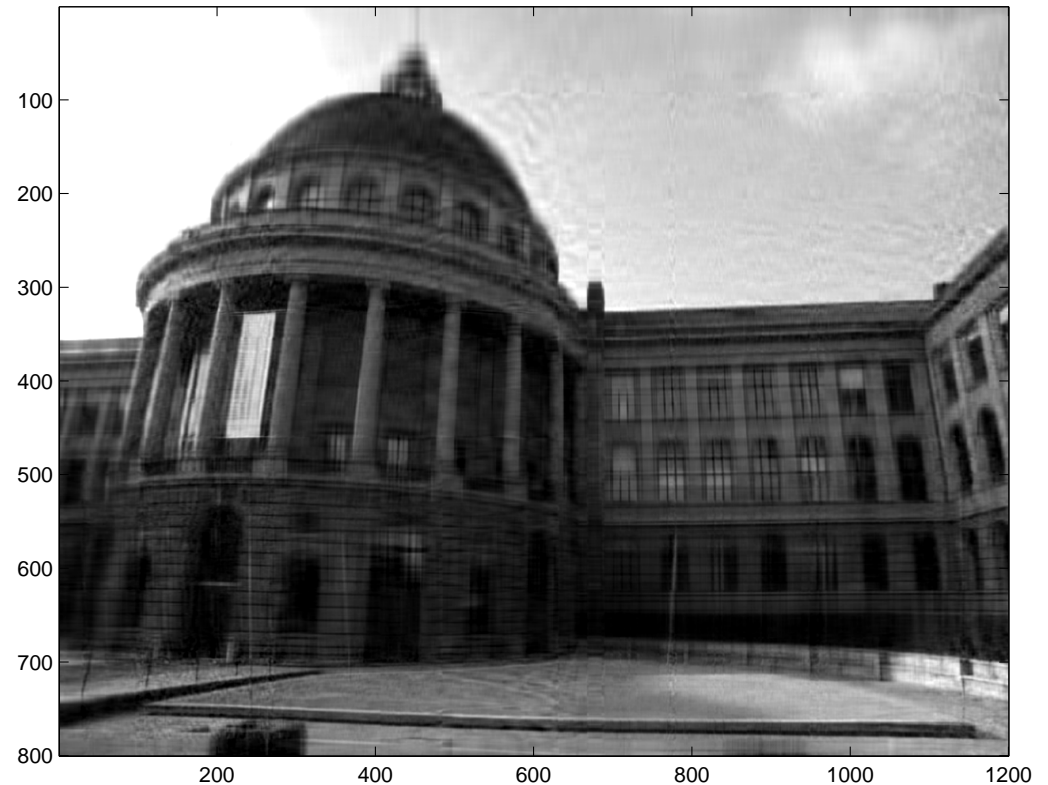
7 subplot(2,2,2); image(PC); title('compressed (40 S.V.)');
  colormap(col);
8 subplot(2,2,3); image(abs(P-PC)); title('difference');
  colormap(col);
9 subplot(2,2,4); cla; semilogy(s); hold on; plot(k,s(k),'ro');

```

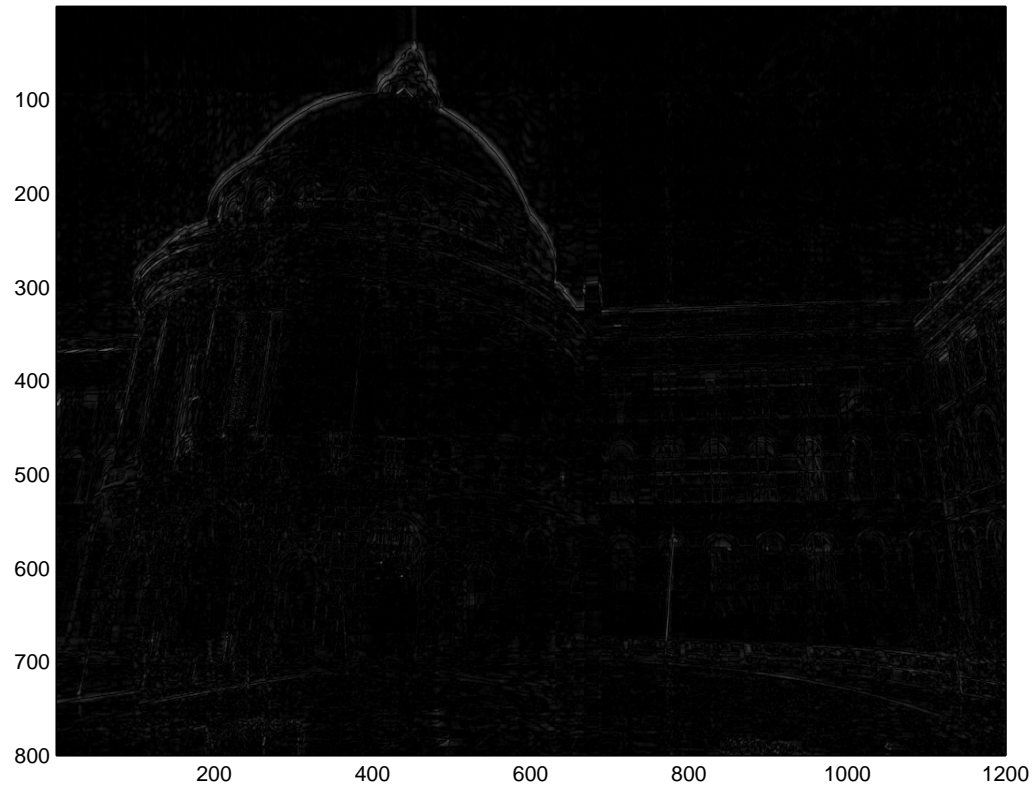
View of ETH Zurich main building



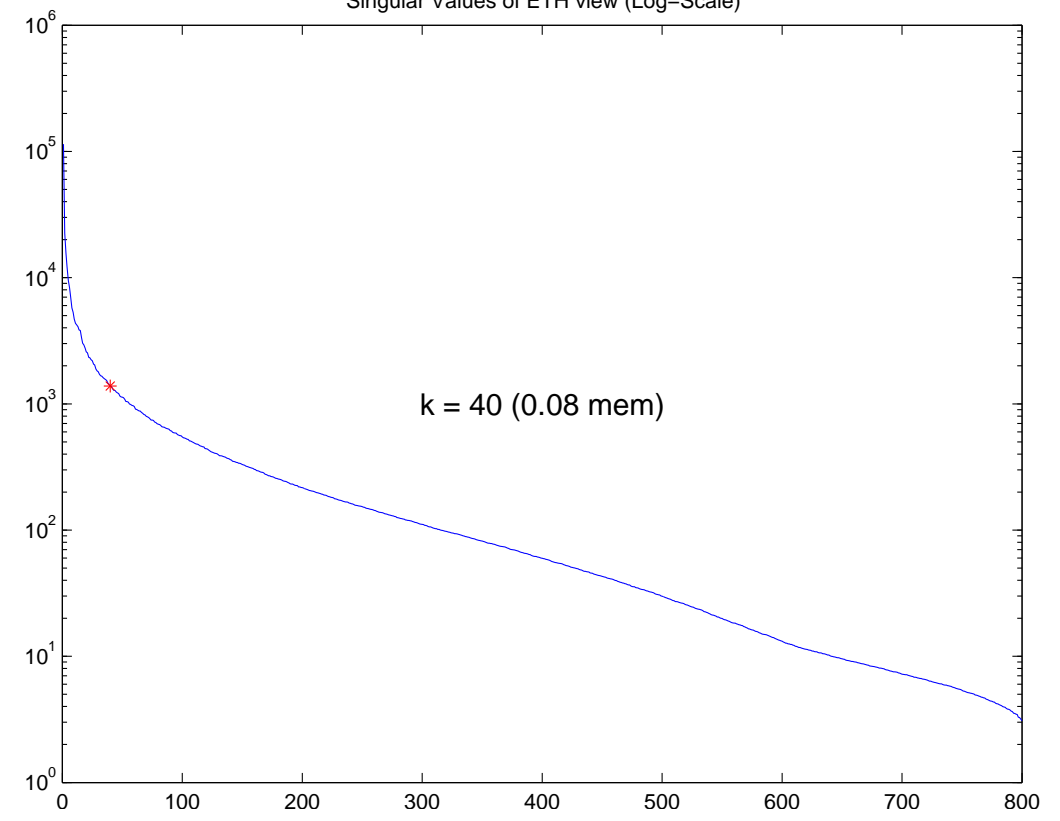
Compressed image, 40 singular values used



Difference image: |original – approximated|



Singular Values of ETH view (Log-Scale)



7

Least Squares

Example 7.0.1 (Least squares data fitting).

Non-linear least squares fitting problem: [13, Sect. 6.1]

- Given:
- data points (t_i, y_i) , $i = 1, \dots, m$
 - (symbolic formula) for parameterized function

$$f(x_1, \dots, x_n; \cdot) : D \subset \mathbb{R} \mapsto \mathbb{R}, \quad n < m$$

Sought: parameter values $x_1^*, \dots, x_n^* \in \mathbb{R}$ such that

$$(x_1^*, \dots, x_n^*) = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m |f(x_1, \dots, x_n; t_i) - y_i|^2. \quad (7.0.2)$$

Example 7.0.3 (Linear data fitting). \rightarrow [13, Sect. 4.1] \rightarrow [13, Ex. 4.2 & Ex. 4.3]

Linear least squares fitting problem:

- Given:
- data points (t_i, y_i) , $i = 1, \dots, m$
 - basis functions $b_j : I \mapsto \mathbb{K}$, $j = 1, \dots, n$, $n < m$

Sought: coefficients $x_j^* \in \mathbb{R}$, $j = 1, \dots, n$, such that

$$(x_1^*, \dots, x_n^*) = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m \left| \sum_{j=1}^n x_j b_j(t_i) - y_i \right|^2 . \quad (7.0.5)$$



Example 7.0.7 (Polybomial interpolation vs. polynomial fitting).

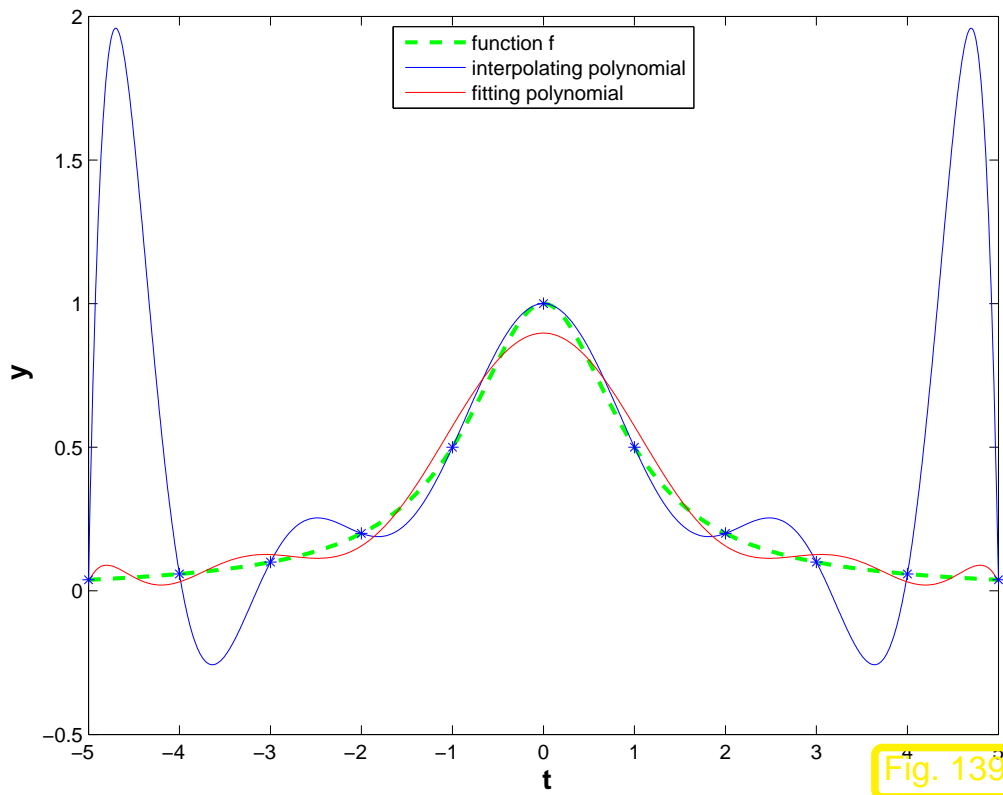


Fig. 139

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. 3.5.1,
- fitting to data points $(t_j, f(t_j))$, $j = 0, \dots, 3d$, $t_j = -5 + j/3$.



(Full rank linear) **least squares problem**: [13, 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 $\|\mathbf{Ax} - \mathbf{b}\|_2 = \inf\{\|\mathbf{Ay} - \mathbf{b}\|_2 : \mathbf{y} \in \mathbb{R}^n\}$

(7.0.14)



$$\mathbf{x} = \underset{\mathbf{y} \in \mathbb{R}^n}{\text{argmin}} \|\mathbf{Ay} - \mathbf{b}\|_2$$

(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: $\mathbf{x} \in \mathbb{R}^n$ such that

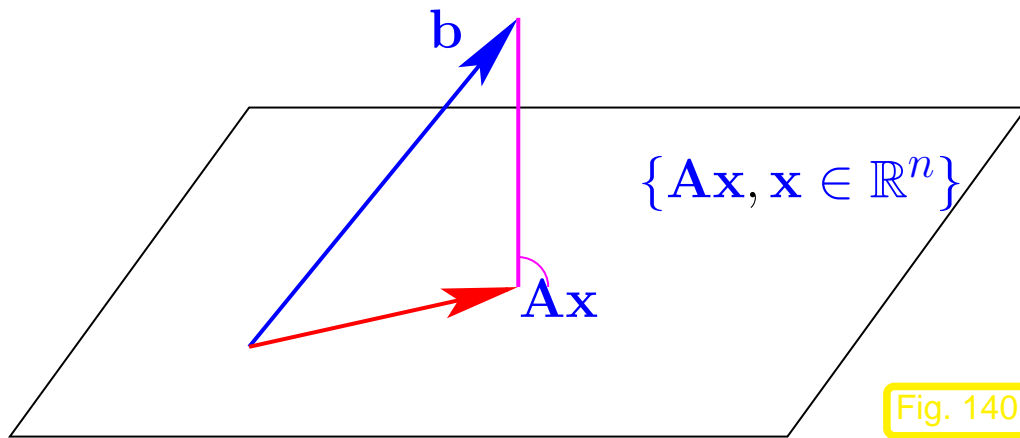
$$(i) \quad \|\mathbf{Ax} - \mathbf{b}\|_2 = \inf\{\|\mathbf{Ay} - \mathbf{b}\|_2 : \mathbf{y} \in \mathbb{R}^n\}, \quad (7.0.16)$$

(ii) $\|\mathbf{x}\|_2$ is minimal under the condition (i).

MATLAB “black-box” solver for general linear least squares problems (7.0.16)

$\mathbf{x} = \mathbf{A} \backslash \mathbf{b}$ (“backslash”) solves (7.0.16) for $\mathbf{A} \in \mathbb{K}^{m,n}$, $m \neq n$.

7.1 Normal Equations [13, Sect. 4.2], [35, Ch. 11]



Geometric interpretation of linear least squares problem (7.0.16):

$\hat{\mathbf{x}} \triangleq$ orthogonal projection of \mathbf{b} on the subspace $\text{Im}(\mathbf{A}) := \text{Span} \{(\mathbf{A})_{:,1}, \dots, (\mathbf{A})_{:,n}\}$.

Example 7.1.4 (Instability of normal equations). \rightarrow [13, Ex. 4.12]

Caution: loss of information in the computation of $\mathbf{A}^H \mathbf{A}$, e.g.



$$\mathbf{A} = \begin{pmatrix} 1 & 1 \\ \delta & 0 \\ 0 & \delta \end{pmatrix} \Rightarrow \mathbf{A}^H \mathbf{A} = \begin{pmatrix} 1 + \delta^2 & 1 \\ 1 & 1 + \delta^2 \end{pmatrix}$$

```
1 >> A = [1 1; ...
2         sqrt(eps) 0; ...
3         0 sqrt(eps)];
4 >> rank(A)
5         ans = 2
6 >> rank(A' * A)
7         ans = 1
```



7.2 Orthogonal Transformation Methods [13, Sect. 4.4.2]

QR-decomposition: $\mathbf{A} = \mathbf{Q}\mathbf{R}$, $\mathbf{Q} \in \mathbb{K}^{m,m}$ unitary, $\mathbf{R} \in \mathbb{K}^{m,n}$ (regular) upper triangular matrix.

$$\|\mathbf{Ax} - \mathbf{b}\|_2 = \left\| \mathbf{Q}(\mathbf{Rx} - \mathbf{Q}^H\mathbf{b}) \right\|_2 = \left\| \mathbf{Rx} - \tilde{\mathbf{b}} \right\|_2, \quad \tilde{\mathbf{b}} := \mathbf{Q}^H\mathbf{b}.$$

$$\|\mathbf{Ax} - \mathbf{b}\|_2 \rightarrow \min \Leftrightarrow \left\| \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} - \begin{pmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_m \end{pmatrix} \right\|_2 \rightarrow \min.$$

$$\mathbf{x} = \left(\begin{array}{c} \text{yellow bars} \\ \mathbf{R} \end{array} \right)^{-1} \begin{pmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_n \end{pmatrix}, \quad \text{residuum } \mathbf{r} = \mathbf{Q} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \tilde{b}_{n+1} \\ \vdots \\ \tilde{b}_m \end{pmatrix}.$$

Code 7.2.1: QR-based solver for full rank linear least squares problem (7.0.14)

```

1 function [x,res] = qrllsqsolve(A,b)
2 % Solution of linear least squares problem (7.0.14) by means of QR-decomposition
3 % Note:  $\mathbf{A} \in \mathbb{R}^{m,n}$  with  $m > n$ ,  $\text{rank}(\mathbf{A}) = n$  is assumed
4 [m,n] = size(A);
5 R = triu(qr([A,b],0)), % economical QR-decomposition of extended matrix
6 x = R(1:n,1:n)\R(1:n,n+1); %  $\hat{\mathbf{x}} = (\mathbf{R})_{1:n,1:n}^{-1}(\mathbf{Q}^T \mathbf{b})_{1:n}$ 
7 res = R(n+1,n+1); % =  $\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|_2$  (why ?)

```

Alternative: Solving linear least squares problem (7.0.16) by SVD (\rightarrow Def. 6.5.8)

Most general setting: $\mathbf{A} \in \mathbb{K}^{m,n}$, $\text{rank}(\mathbf{A}) = r \leq \min\{m, n\}$:

SVD: $\mathbf{A} = [\mathbf{U}_1 \quad \mathbf{U}_2] \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{pmatrix}$

(7.2.2)

R. Hiptmair
rev 38355,
November
18, 2011

$$\|\mathbf{Ax} - \mathbf{b}\|_2 = \left\| [\mathbf{U}_1 \quad \mathbf{U}_2] \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{pmatrix} \mathbf{x} - \mathbf{b} \right\|_2 = \left\| \begin{pmatrix} \Sigma_r \mathbf{V}_1^H \mathbf{x} \\ 0 \end{pmatrix} - \begin{pmatrix} \mathbf{U}_1^H \mathbf{b} \\ \mathbf{U}_2^H \mathbf{b} \end{pmatrix} \right\|_2 \quad (7.2.3)$$

► solution $\mathbf{x} = \mathbf{V}_1 \Sigma_r^{-1} \mathbf{U}_1^H \mathbf{b}$, $\|\mathbf{r}\|_2 = \left\| \mathbf{U}_2^H \mathbf{b} \right\|_2$. (7.2.6)

Code 7.2.7: Solving LSQ problem via SVD

Practical implementation:

“numerical rank” test:

$$r = \max\{i: \sigma_i/\sigma_1 > \text{tol}\}$$

```

1 function y = lsqsvd(A,b)
2 [U,S,V] = svd(A,0);
3 sv = diag(S);
4 r = max(find(sv/sv(1) > eps));
5 y = V(:,1:r)*(diag(1./sv(1:r))*...
6           (U(:,1:r)'*b));

```

7.3 Total Least Squares

Total least squares problem:

Given: $\mathbf{A} \in \mathbb{R}^{m,n}$, $m \geq n$, $\text{rank}(\mathbf{A}) = n$, $\mathbf{b} \in \mathbb{R}^m$,find: $\hat{\mathbf{A}} \in \mathbb{R}^{m,n}$, $\hat{\mathbf{b}} \in \mathbb{R}^m$ with

$$\left\| \underbrace{\begin{bmatrix} \mathbf{A} & \mathbf{b} \end{bmatrix}}_{=:\mathbf{C}} - \underbrace{\begin{bmatrix} \hat{\mathbf{A}} & \hat{\mathbf{b}} \end{bmatrix}}_{=:\hat{\mathbf{C}}} \right\|_F \rightarrow \min, \quad \hat{\mathbf{b}} \in \text{Im}(\hat{\mathbf{A}}).$$

(7.3.1)

$$\hat{\mathbf{b}} \in \text{Im}(\hat{\mathbf{A}}) \Rightarrow \text{rank}(\hat{\mathbf{C}}) = n \quad \blacktriangleright \quad (7.3.1) \Rightarrow \min_{\text{rank}(\hat{\mathbf{C}})=n} \|\mathbf{C} - \hat{\mathbf{C}}\|_F .$$

Code 7.3.5: Total least squares via SVD

```

1 function x = lsqtotal(A,b);
2 % computes only solution x of fitted consistent LSE
3 [m,n]=size(A);
4 [U, Sigma, V] = svd([A,b]); % see (7.3.2)
5 s = V(n+1,n+1);
6 if s == 0, error('No solution'); end
7 x = -V(1:n,n+1)/s; % see (7.3.4)

```

7.4 Constrained Least Squares

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: $\mathbf{x} \in \mathbb{R}^n$ with $\|\mathbf{Ax} - \mathbf{b}\|_2 \rightarrow \min$, $\mathbf{Cx} = \mathbf{d}$. (7.4.1)

Linear constraint

7.4.1 Solution via normal equations

Idea: coupling the constraint using the **Lagrange multiplier** $\mathbf{m} \in \mathbb{R}^p$

$$\mathbf{x} = \underset{\mathbf{x} \in \mathbb{R}^n}{\text{argmin}} \max_{\mathbf{m} \in \mathbb{R}^p} L(\mathbf{x}, \mathbf{m}), \quad (7.4.2)$$

$$L(\mathbf{x}, \mathbf{m}) := \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 + \mathbf{m}^H (\mathbf{Cx} - \mathbf{d}). \quad (7.4.3)$$

Solution of min-max problem:

saddle point

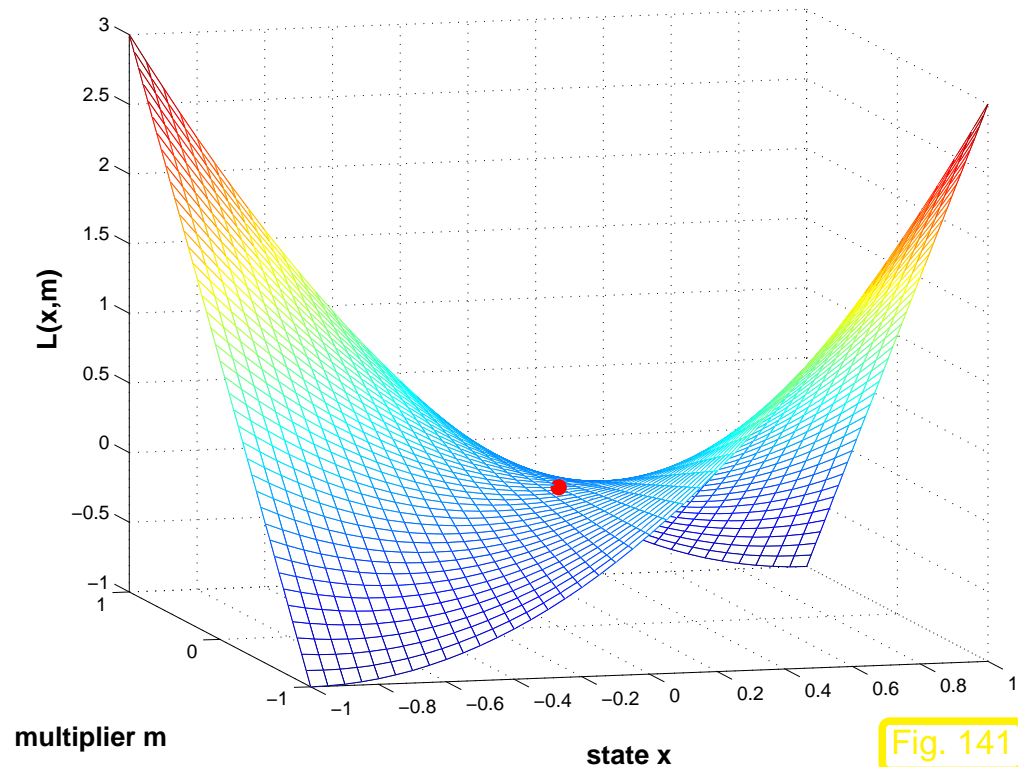


Fig. 141

Necessary (and sufficient) condition for the minimizer (→ Section 7.1)

$$\frac{\partial L}{\partial \mathbf{x}}(\mathbf{x}, \mathbf{m}) = \mathbf{A}^H(\mathbf{A}\mathbf{x} - \mathbf{b}) + \mathbf{C}^H\mathbf{m} \stackrel{!}{=} 0, \quad \frac{\partial L}{\partial \mathbf{m}}(\mathbf{x}, \mathbf{m}) = \mathbf{C}\mathbf{x} - \mathbf{d} \stackrel{!}{=} 0. \quad (7.4.4)$$

$$\begin{pmatrix} \mathbf{A}^H\mathbf{A} & \mathbf{C}^H \\ \mathbf{C} & 0 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{m} \end{pmatrix} = \begin{pmatrix} \mathbf{A}^H\mathbf{b} \\ \mathbf{d} \end{pmatrix} \quad \text{Augmented normal equations} \\ \text{(matrix saddle point problem)} \quad (7.4.5)$$

7.4.2 Solution via SVD

① Compute orthonormal basis of $\text{Ker}(\mathbf{C})$ using SVD (\rightarrow Lemma 6.5.14, (7.2.2)):

$$\mathbf{C} = \mathbf{U} \begin{bmatrix} \Sigma & 0 \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix}, \quad \mathbf{U} \in \mathbb{R}^{p,p}, \Sigma \in \mathbb{R}^{p,p}, \mathbf{V}_1 \in \mathbb{R}^{n,p}, \mathbf{V}_2 \in \mathbb{R}^{n,n-p}$$

$$\blacktriangleright \quad \text{Ker}(\mathbf{C}) = \text{Im}(\mathbf{V}_2).$$

and the **particular solution** of the constraint equation

$$\mathbf{x}_0 := \mathbf{V}_1 \Sigma^{-1} \mathbf{U}^H \mathbf{d}.$$

Representation of the solution \mathbf{x} of (7.4.1): $\mathbf{x} = \mathbf{x}_0 + \mathbf{V}_2 \mathbf{y}, \quad \mathbf{y} \in \mathbb{R}^{n-p}.$

② Insert this representation in (7.4.1) \blacktriangleright standard linear least squares

$$\|\mathbf{A}(\mathbf{x}_0 + \mathbf{V}_2 \mathbf{y}) - \mathbf{b}\|_2 \rightarrow \min \quad \Leftrightarrow \quad \|\mathbf{A} \mathbf{V}_2 \mathbf{y} - (\mathbf{b} - \mathbf{A} \mathbf{x}_0)\| \rightarrow \min.$$

7.5 Non-linear Least Squares [13, Ch. 6]

Non-linear least squares problem

Given: $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^m$, $m, n \in \mathbb{N}$, $m > n$.

Find: $\mathbf{x}^* \in D$: $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in D} \Phi(\mathbf{x})$, $\Phi(\mathbf{x}) := \frac{1}{2} \|F(\mathbf{x})\|_2^2$. (7.5.3)

Terminology: $D \hat{=}$ parameter space, $x_1, \dots, x_n \hat{=}$ parameter.

We require “independence for each parameter”: \rightarrow Rem. 7.0.18

\exists neighbourhood $\mathcal{U}(\mathbf{x}^*)$ such that $DF(\mathbf{x})$ has full rank $n \quad \forall \mathbf{x} \in \mathcal{U}(\mathbf{x}^*)$. (7.5.4)

(It means: the columns of the Jacobi matrix $DF(\mathbf{x})$ are linearly independent.)

7.5.1 (Damped) Newton method

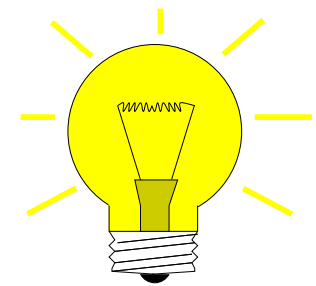
$\Phi(\mathbf{x}^*) = \min \Rightarrow \operatorname{grad} \Phi(\mathbf{x}) = 0$, $\operatorname{grad} \Phi(\mathbf{x}) := \left(\frac{\partial \Phi}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial \Phi}{\partial x_n}(\mathbf{x}) \right)^T \in \mathbb{R}^n$.

Simple idea: use Newton's method (\rightarrow Sect. 4.4) to determine a zero of $\text{grad } \Phi : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$.

► For Newton iterate $\mathbf{x}^{(k)}$: Newton correction $\mathbf{s} \in \mathbb{R}^n$ from LSE

$$\underbrace{\left(DF(\mathbf{x}^{(k)})^T DF(\mathbf{x}^{(k)}) + \sum_{j=1}^m F_j(\mathbf{x}^{(k)}) D^2 F_j(\mathbf{x}^{(k)}) \right)}_{=H\Phi(\mathbf{x}^{(k)})} \mathbf{s} = - \underbrace{DF(\mathbf{x}^{(k)})^T F(\mathbf{x}^{(k)})}_{=\text{grad } \Phi(\mathbf{x}^{(k)})}. \quad (7.5.6)$$

7.5.2 Gauss-Newton method



Idea: **local linearization** of F : $F(\mathbf{x}) \approx F(\mathbf{y}) + DF(\mathbf{y})(\mathbf{x} - \mathbf{y})$

➤ sequence of *linear* least squares problems

Initial guess $\mathbf{x}^{(0)} \in D$

$$\mathbf{x}^{(k+1)} := \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left\| F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}) \right\|_2. \quad (7.5.9)$$

linear least squares problem

MATLAB-`\` used to solve linear least squares problem in each step:

for $\mathbf{A} \in \mathbb{R}^{m,n}$

$$\mathbf{x} = \mathbf{A} \backslash \mathbf{b}$$



\mathbf{x} minimizer of $\|\mathbf{Ax} - \mathbf{b}\|_2$

with minimal 2-norm

Example 7.5.11 (Non-linear data fitting (II)). → Ex. 7.5.1

Code 7.5.10: template for Gauss-Newton method

```

1 function x = gn(x,F,J,tol)
2 s = J(x)\F(x); %
3 x = x-s;
4 while (norm(s) > tol*norm(x)) %
5     s = J(x)\F(x); %
6     x = x-s;
7 end
    
```

Non-linear data fitting problem (7.5.2) for $f(t) = x_1 + x_2 \exp(-x_3 t)$.

$$F(\mathbf{x}) = \begin{pmatrix} x_1 + x_2 \exp(-x_3 t_1) - y_1 \\ \vdots \\ x_1 + x_2 \exp(-x_3 t_m) - y_m \end{pmatrix} : \mathbb{R}^3 \mapsto \mathbb{R}^m, \quad DF(\mathbf{x}) = \begin{pmatrix} 1 & e^{-x_3 t_1} & -x_2 t_1 e^{-x_3 t_1} \\ \vdots & \vdots & \vdots \\ 1 & e^{-x_3 t_m} & -x_2 t_m e^{-x_3 t_m} \end{pmatrix}$$

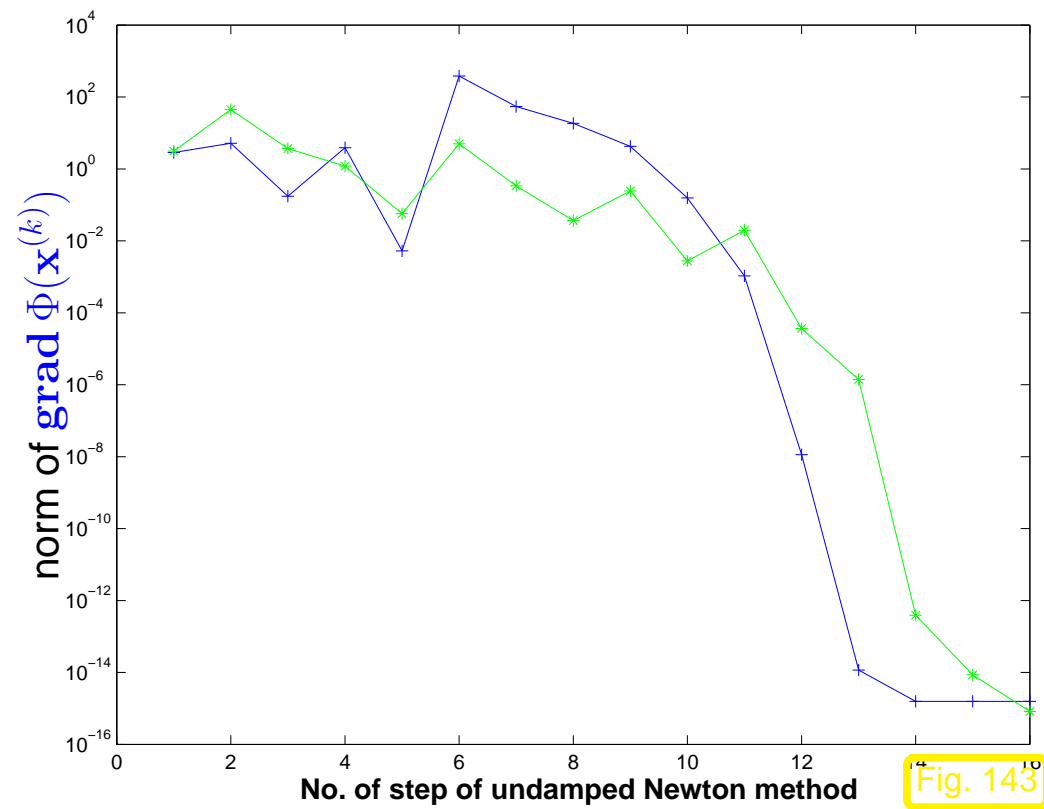
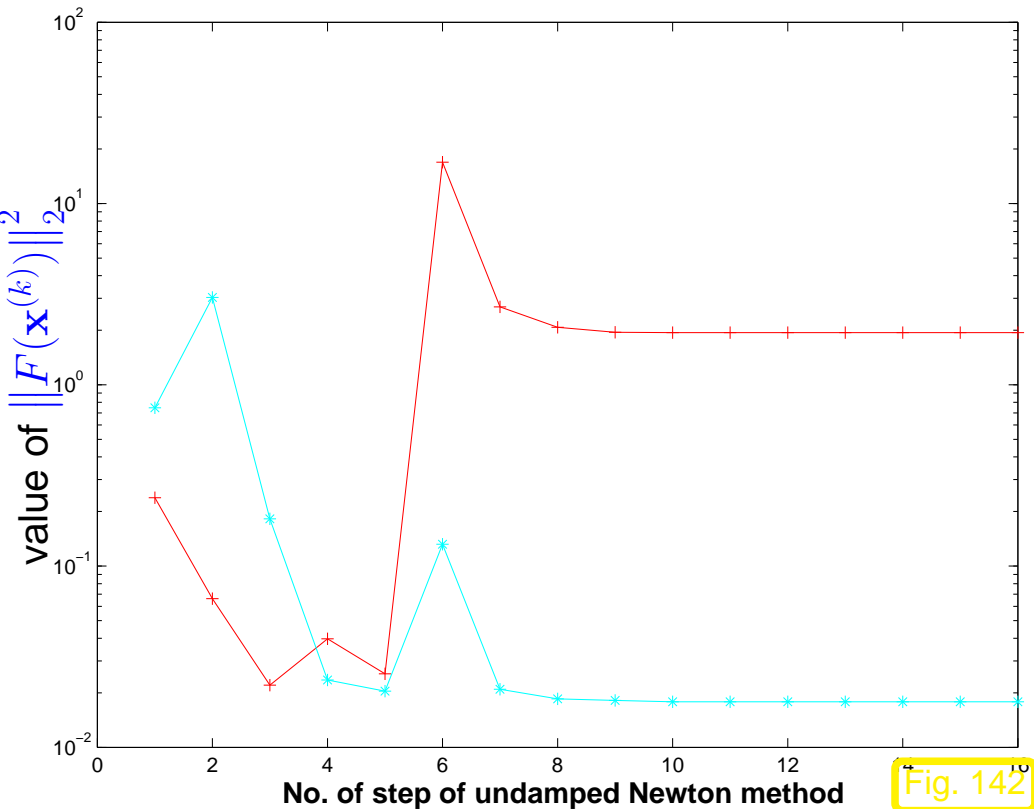
Numerical experiment:

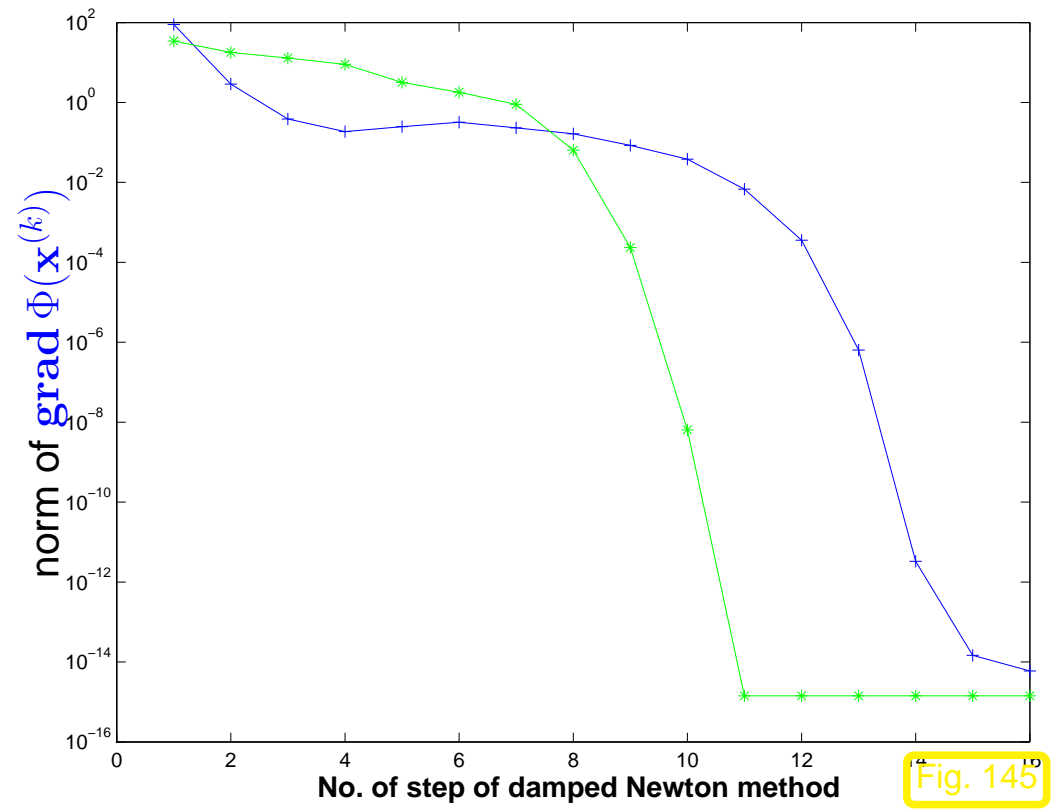
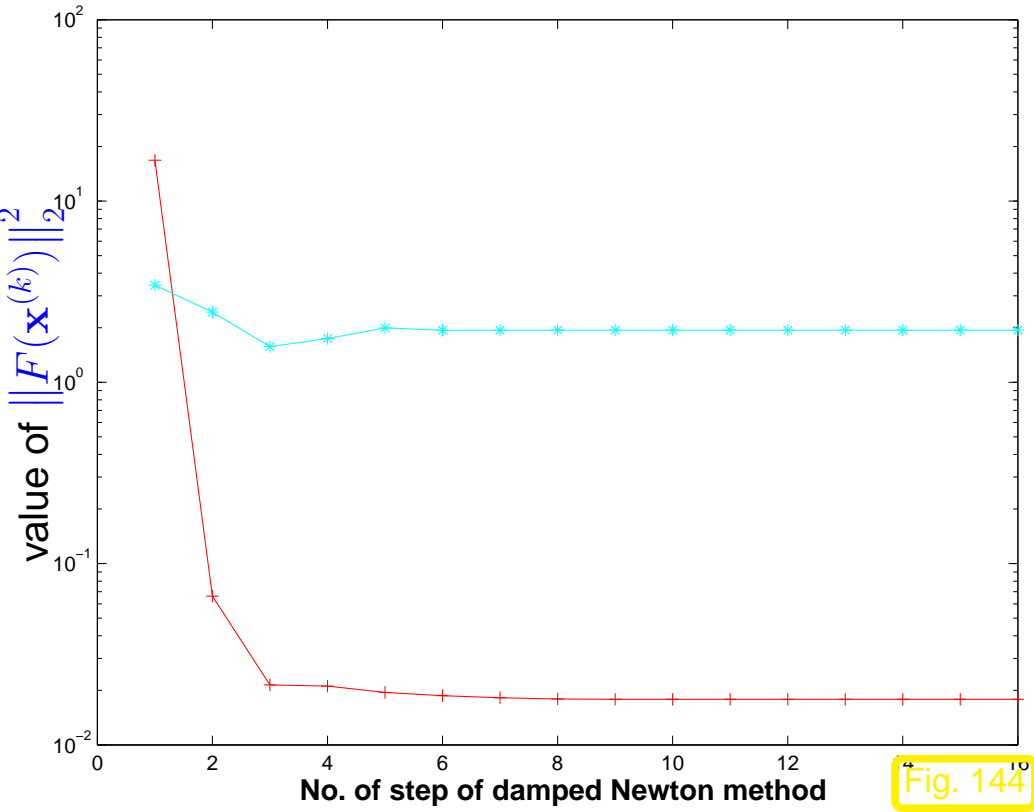
convergence of the Newton method,
damped Newton method (\rightarrow Section
4.4.4) and Gauss-Newton method for
different initial values

```

rand('seed',0);
t = (1:0.3:7)';
y = x(1) + x(2)*exp(-x(3)*t);
y = y+0.1*(rand(length(y),1)-0.5);
    
```

- 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)





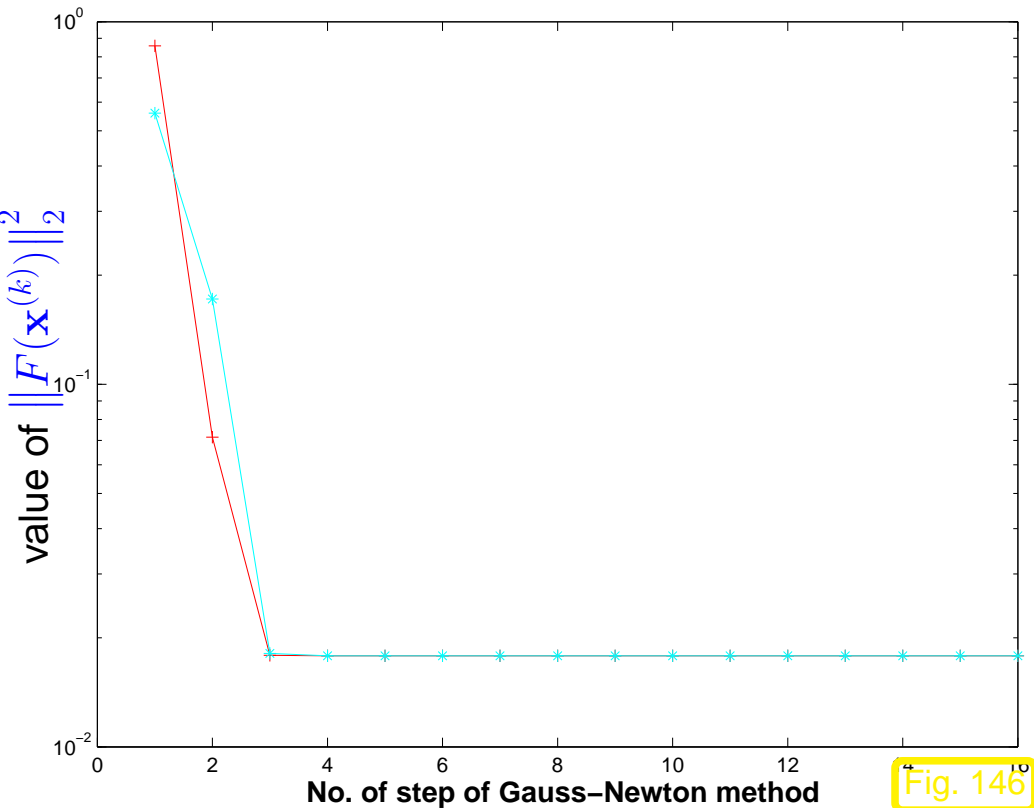


Fig. 146

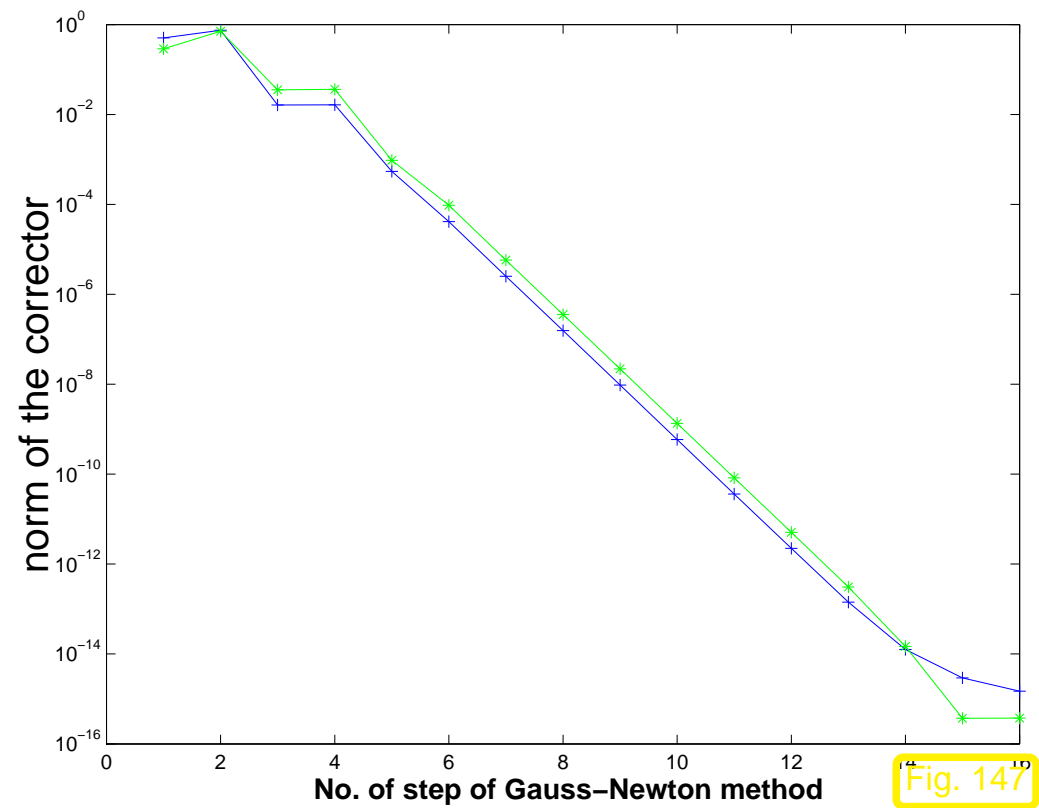


Fig. 147



7.5.3 Trust region method (Levenberg-Marquardt method)

Idea: damping of the Gauss-Newton correction in (7.5.9) using a **penalty term**

instead of $\|F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})\mathbf{s}\|_2^2$ minimize $\|F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})\mathbf{s}\|_2^2 + \lambda \|\mathbf{s}\|_2^2$.

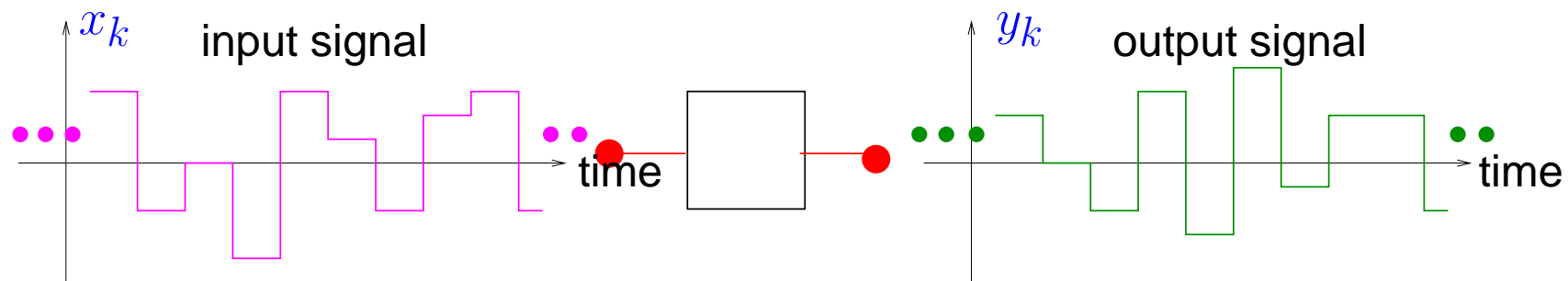
$\lambda > 0 \hat{=}$ penalty parameter (how to choose it ? \rightarrow heuristic)

8

Filtering Algorithms

8.1 Discrete convolutions

Example 8.1.1 (Discrete finite linear time-invariant causal channel (filter)).



Impulse response of channel (filter):

$$\mathbf{h} = (h_0, \dots, h_{n-1})^T$$

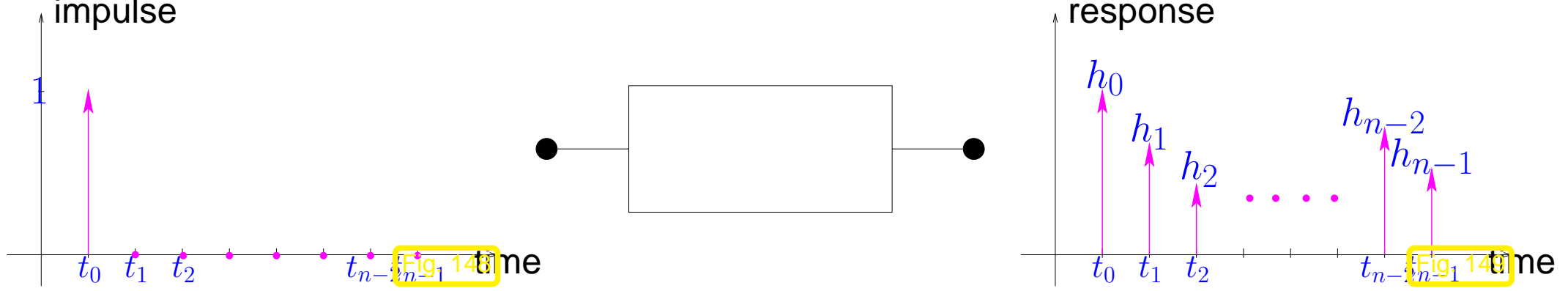


Fig. 149

Fig. 149

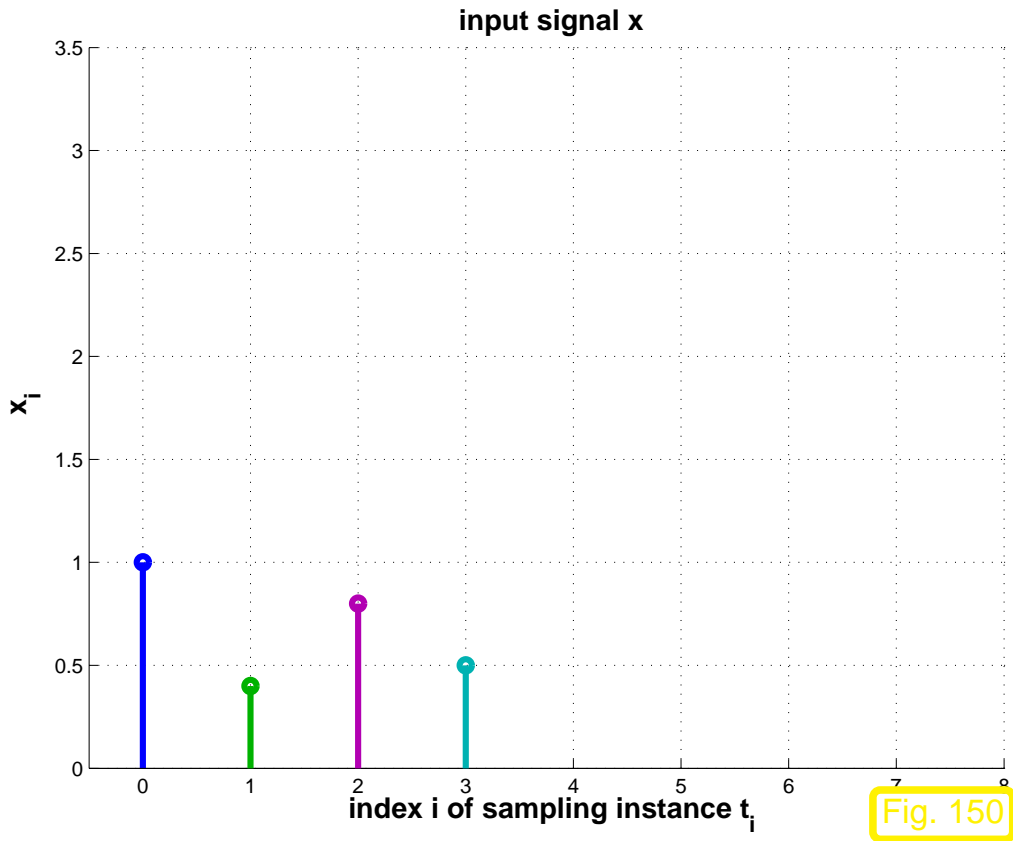


Fig. 150

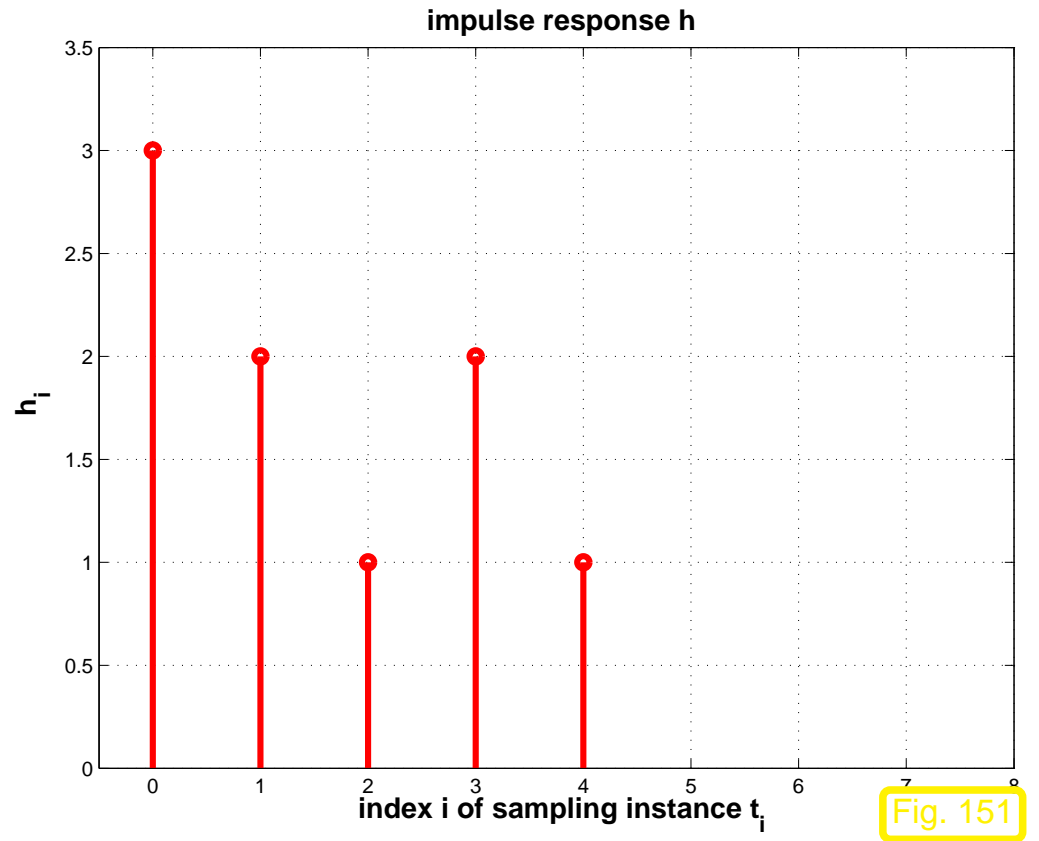
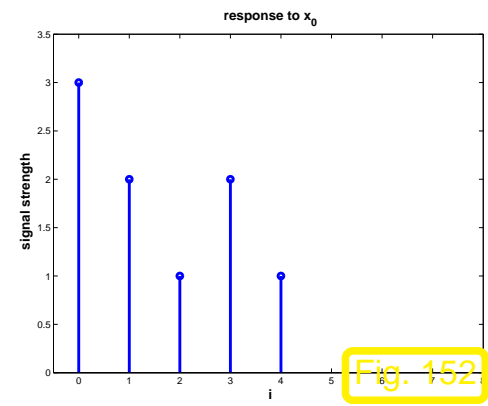
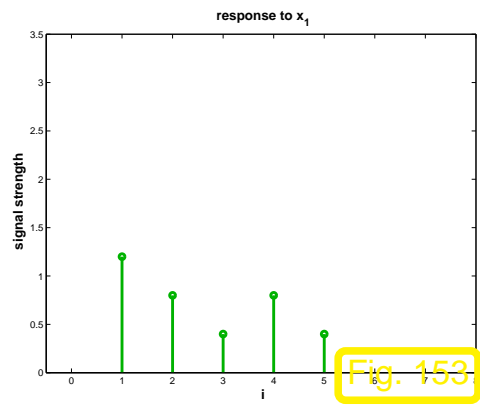


Fig. 151

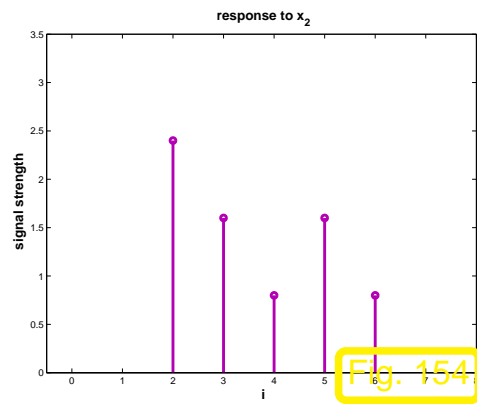
Output = linear superposition of impulse responses:



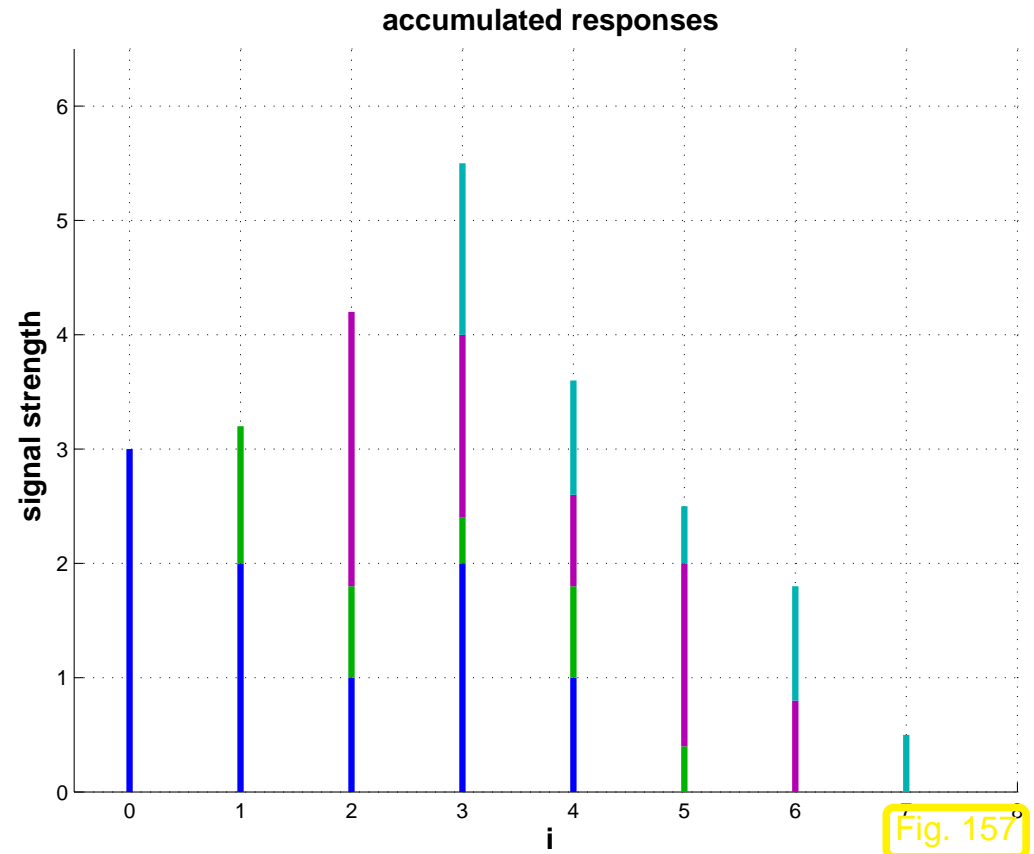
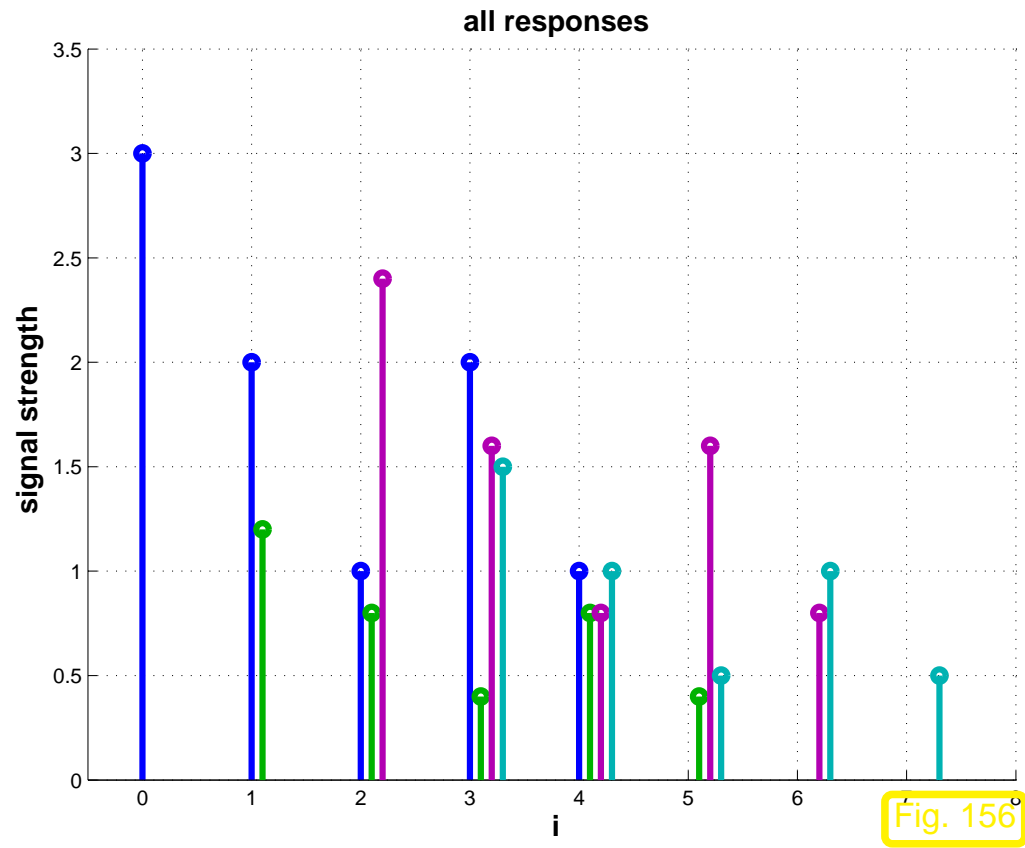
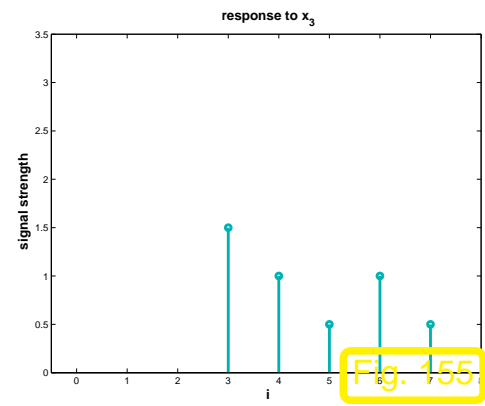
+



+



+



$$\begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \\ \vdots \\ y_{2n-3} \\ y_{2n-2} \end{pmatrix} = x_0 \begin{pmatrix} h_0 \\ \vdots \\ h_{n-1} \\ 0 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{pmatrix} + x_1 \begin{pmatrix} 0 \\ h_0 \\ \vdots \\ h_{n-1} \\ 0 \\ \vdots \\ \vdots \\ 0 \end{pmatrix} + x_2 \begin{pmatrix} 0 \\ 0 \\ h_0 \\ \vdots \\ h_{n-1} \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \cdots + x_{n-1} \begin{pmatrix} 0 \\ \vdots \\ \vdots \\ \vdots \\ 0 \\ h_0 \\ \vdots \\ h_{n-1} \end{pmatrix} .$$

channel is causal!

$$y_k = \sum_{j=0}^{n-1} h_{k-j} x_j, \quad k = 0, \dots, 2n-2 \quad (h_j := 0 \text{ for } j < 0 \text{ and } j \geq n) . \quad (8.1.3)$$

$\mathbf{x} = (x_0, \dots, x_{n-1})^\top \in \mathbb{R}^n \hat{=}$ input signal $\mapsto \mathbf{y} = (y_0, \dots, y_{2n-2})^\top \in \mathbb{R}^{2n-1} \hat{=}$ output signal.

Matrix notation of (8.1.3):

$$\begin{pmatrix} y_0 \\ \vdots \\ y_{2n-2} \end{pmatrix} = \begin{pmatrix} h_0 & 0 & \cdots & 0 \\ h_1 & & \cdots & 0 \\ & \ddots & \ddots & \\ h_{n-1} & \cdots & h_1 & h_0 \\ 0 & & & \\ & \ddots & \ddots & \\ 0 & \cdots & 0 & h_{n-1} \end{pmatrix} \begin{pmatrix} x_0 \\ \vdots \\ x_{n-1} \end{pmatrix} . \tag{8.1.4}$$



Example 8.1.5 (Multiplication of polynomials).



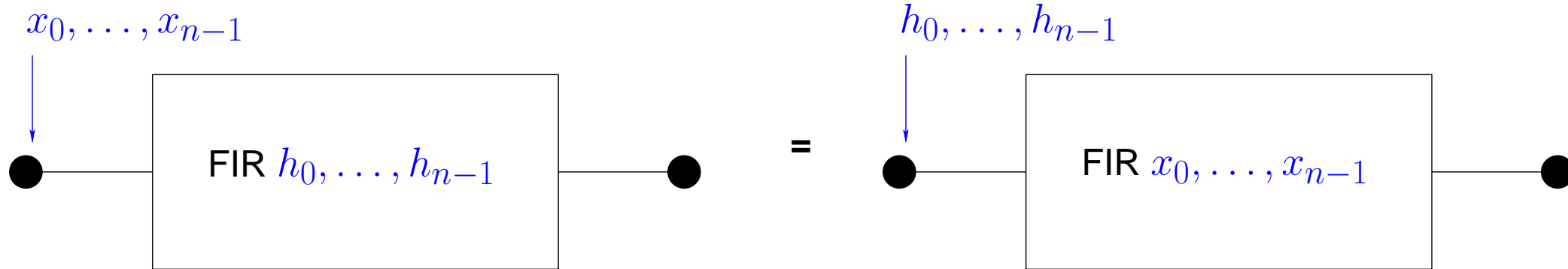
Definition 8.1.7 (Discrete convolution).

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** (ger.: *diskrete Faltung*) is the vector $\mathbf{y} \in \mathbb{K}^{2n-1}$ with components

$$y_k = \sum_{j=0}^{n-1} h_{k-j} x_j, \quad k = 0, \dots, 2n-2 \quad (h_j := 0 \text{ for } j < 0). \quad (8.1.8)$$

✎ Notation for discrete convolution (8.1.8):

$$\mathbf{y} = \mathbf{h} * \mathbf{x}.$$



Example 8.1.10 (Linear filtering of periodic signals).

n -periodic signal ($n \in \mathbb{N}$) = sequence $(x_j)_{j \in \mathbb{Z}}$ with $x_{j+n} = x_j \quad \forall j \in \mathbb{Z}$

➤ n -periodic signal $(x_j)_{j \in \mathbb{Z}}$ fixed by $x_0, \dots, x_{n-1} \leftrightarrow$ vector $\mathbf{x} = (x_0, \dots, x_{n-1})^\top \in \mathbb{R}^n$.

$$y_k = \sum_{j=0}^{n-1} p_{k-j} x_j \quad \text{for some } p_0, \dots, p_{n-1} \in \mathbb{R}, \quad p_k := p_{k-n} \text{ for all } k \in \mathbb{Z}. \quad (8.1.11)$$

Matrix notation:

$$\begin{pmatrix} y_0 \\ \vdots \\ y_{n-1} \end{pmatrix} = \underbrace{\begin{pmatrix} p_0 & p_{n-1} & p_{n-2} & \cdots & \cdots & p_1 \\ p_1 & p_0 & p_{n-1} & & & \vdots \\ p_2 & p_1 & p_0 & \ddots & & \\ \vdots & & \ddots & \ddots & \ddots & \\ \vdots & & & \ddots & \ddots & \\ p_{n-1} & \cdots & & & \cdots & \ddots & p_{n-1} \\ & & & & p_1 & p_0 \end{pmatrix}}{=:P} \begin{pmatrix} x_0 \\ \vdots \\ x_{n-1} \end{pmatrix}. \quad (8.1.12)$$



Definition 8.1.13 (Discrete periodic convolution).

The *discrete periodic convolution* of two n -periodic sequences $(x_k)_{k \in \mathbb{Z}}$, $(y_k)_{k \in \mathbb{Z}}$ yields the n -periodic sequence

$$(z_k) := (x_k) *_{n} (y_k) \quad , \quad z_k := \sum_{j=0}^{n-1} x_{k-j} y_j = \sum_{j=0}^{n-1} y_{k-j} x_j \quad , \quad k \in \mathbb{Z} .$$

Definition 8.1.17 (Circulant matrix). \rightarrow [35, Sect. 54]

A matrix $\mathbf{C} = (c_{ij})_{i,j=1}^n \in \mathbb{K}^{n,n}$ is *circulant* (ger.: zirkulant)

$$:\Leftrightarrow \quad \exists (u_k)_{k \in \mathbb{Z}} \text{ } n\text{-periodic sequence: } c_{ij} = u_{j-i}, 1 \leq i, j \leq n.$$

Structure of circulant matrix \triangleright

$$\begin{pmatrix} u_0 & u_1 & u_2 & \cdots & & \cdots & u_{n-1} \\ u_{n-1} & u_0 & & & & & u_{n-2} \\ u_{n-2} & & & & & & \vdots \\ \vdots & & & & & & \vdots \\ \vdots & & & & & & \vdots \\ u_2 & & & & & & u_1 \\ u_1 & u_2 & \cdots & & \cdots & u_{n-1} & u_0 \end{pmatrix}$$

Remark 8.1.19 (Reduction to periodic convolution).

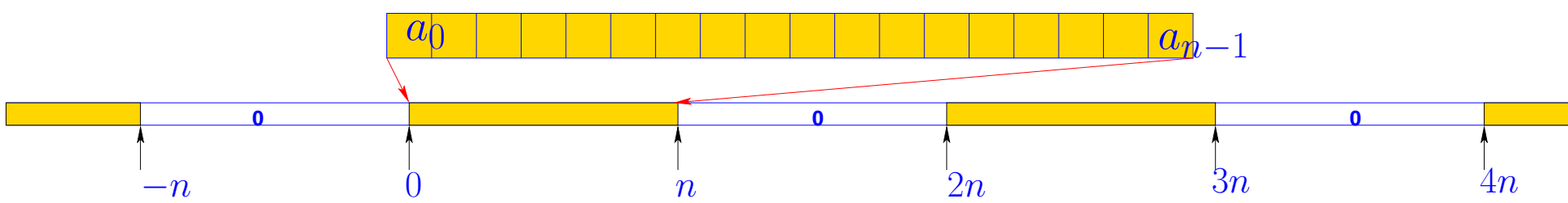
Recall discrete convolution (\rightarrow Def. 8.1.7) of $\mathbf{a} = (a_0, \dots, a_{n-1})^\top \in \mathbb{K}^n$, $\mathbf{b} = (b_0, \dots, b_{n-1})^\top \in \mathbb{K}^n$:

$$(\mathbf{a} * \mathbf{b})_k = \sum_{j=0}^{n-1} a_j b_{k-j}, \quad k = 0, \dots, 2n - 2.$$

Expand a_0, \dots, a_{n-1} and b_0, \dots, b_{n-1} to $2n - 1$ -periodic sequences by **zero padding**:

$$x_k := \begin{cases} a_k & , \text{ if } 0 \leq k < n, \\ 0 & , \text{ if } n \leq k < 2n - 1 \end{cases}, \quad y_k := \begin{cases} b_k & , \text{ if } 0 \leq k < n, \\ 0 & , \text{ if } n \leq k < 2n - 1, \end{cases} \quad (8.1.20)$$

and periodic extension: $x_k = x_{2n-1+k}$, $y_k = y_{2n-1+k}$ for all $k \in \mathbb{Z}$.



► $(\mathbf{a} * \mathbf{b})_k = (\mathbf{x} *_{2n-1} \mathbf{y})_k, \quad k = 0, \dots, 2n - 2. \quad (8.1.21)$

$$\begin{pmatrix} y_0 \\ \vdots \\ \vdots \\ y_{2n-2} \end{pmatrix} = \underbrace{\begin{pmatrix} h_0 & 0 & \cdots & 0 & h_{n-1} & \cdots & h_1 & h_0 \\ h_1 & & & 0 & & & & \\ \vdots & & & \vdots & & & & \\ h_{n-1} & \cdots & h_1 & h_0 & 0 & \cdots & 0 & h_{n-1} \\ 0 & & & h_1 & h_0 & & & \\ \vdots & & & \vdots & & & & \\ 0 & \cdots & 0 & h_{n-1} & \cdots & h_1 & h_0 & \end{pmatrix}}_{\text{a } (2n-1) \times (2n-1) \text{ circulant matrix!}} \begin{pmatrix} x_0 \\ \vdots \\ \vdots \\ x_{n-1} \\ 0 \\ \vdots \\ 0 \end{pmatrix}. \tag{8.1.22}$$



8.2 Discrete Fourier Transform (DFT)

Example 8.2.1 (Eigenvectors of circulant matrices).

Random 8×8 circulant matrices C_1, C_2 (\rightarrow Def. 8.1.17)

eigenvalues (real part) \triangleright

Generated by MATLAB-command:

```
C = gallery('circul',rand(n,1));
```

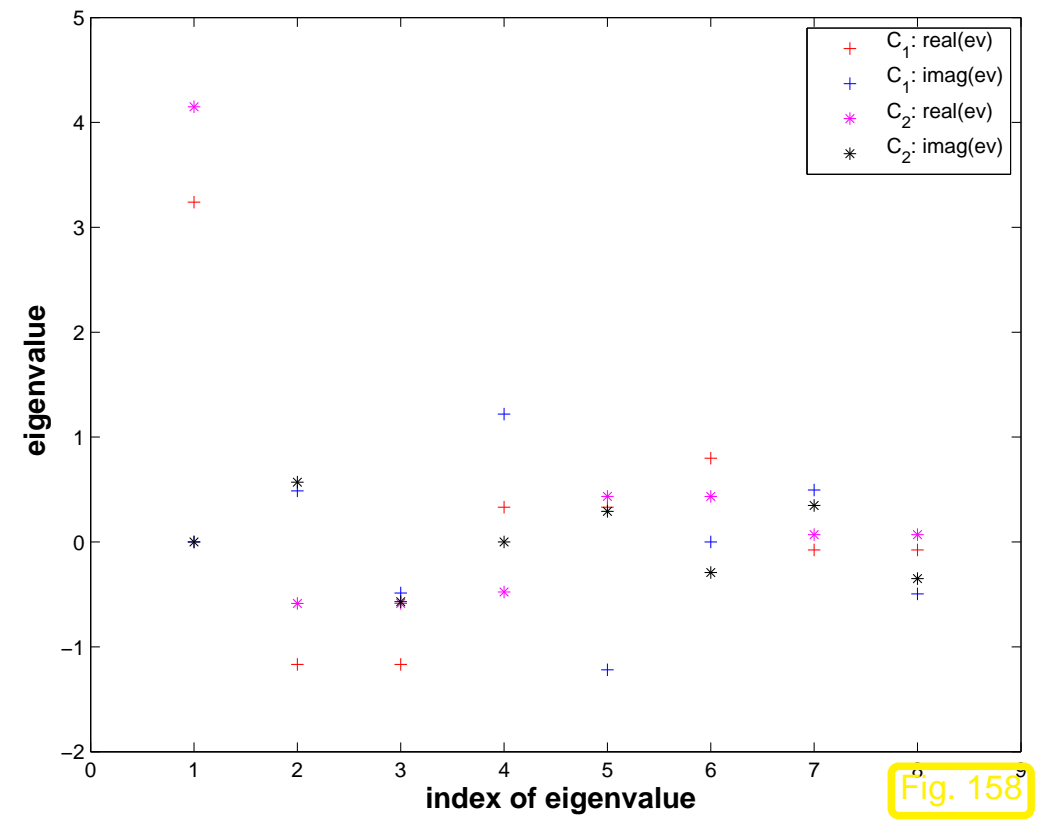
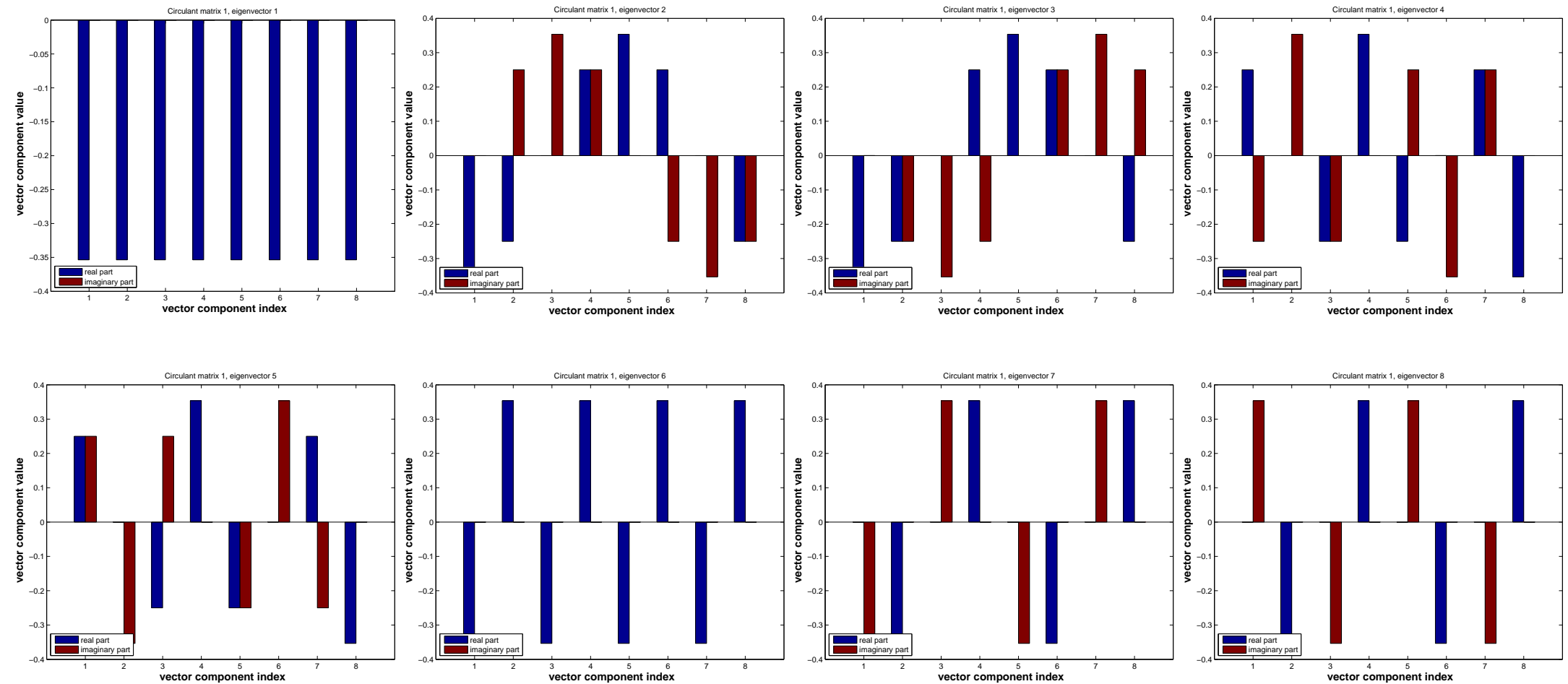
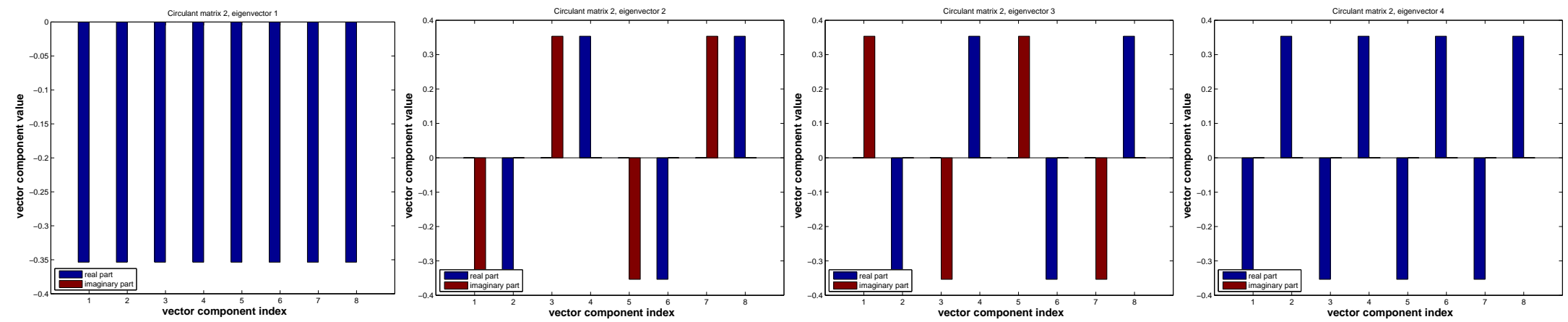


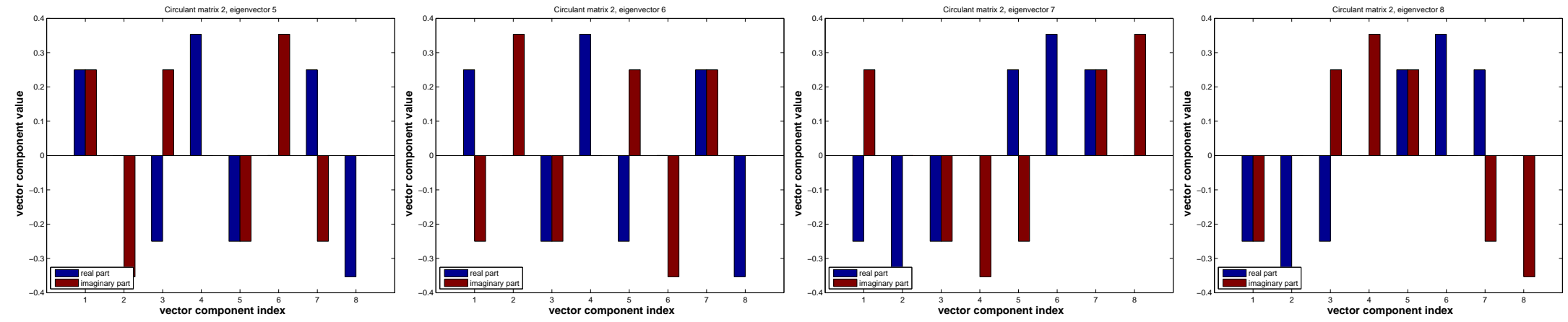
Fig. 158

Eigenvectors of matrix C_1 :

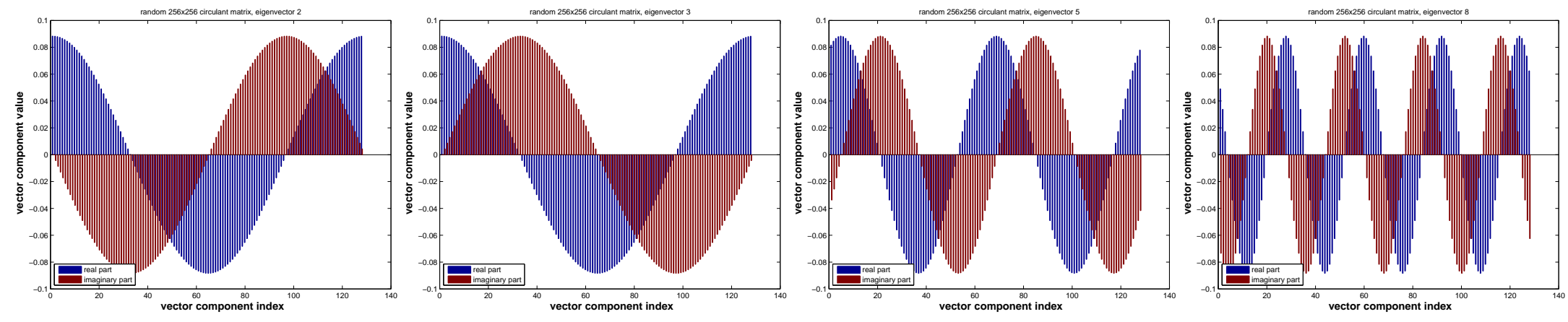


Eigenvectors of matrix C_2





Eigenvectors of $C = \text{gallery}(\text{'circul'}, (1:128)); :$



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 $\bar{\omega}_n = \omega_n^{-1}$, $\omega_n^n = 1$, $\omega_n^{n/2} = -1$, $\omega_n^k = \omega_n^{k+n} \quad \forall k \in \mathbb{Z}$, (8.2.3)

$$\sum_{k=0}^{n-1} \omega_n^{kj} = \begin{cases} n & , \text{ if } j = 0 \pmod n , \\ 0 & , \text{ if } j \neq 0 \pmod n . \end{cases} \tag{8.2.4}$$

Orthogonal trigonometric basis of \mathbb{C}^n = eigenvector basis for circulant matrices

$$\left\{ \begin{pmatrix} \omega_n^0 \\ \vdots \\ \omega_n^0 \end{pmatrix}, \begin{pmatrix} \omega_n^0 \\ \omega_n^1 \\ \vdots \\ \omega_n^{n-1} \end{pmatrix}, \dots, \begin{pmatrix} \omega_n^0 \\ \omega_n^{n-2} \\ \omega_n^{2(n-2)} \\ \vdots \\ \omega_n^{(n-1)(n-2)} \end{pmatrix}, \begin{pmatrix} \omega_n^0 \\ \omega_n^{n-1} \\ \omega_n^{2(n-1)} \\ \vdots \\ \omega_n^{(n-1)^2} \end{pmatrix} \right\} .$$

Matrix of change of basis trigonometrical basis \rightarrow standard basis: **Fourier-matrix**

$$\mathbf{F}_n = \begin{pmatrix} \omega_n^0 & \omega_n^0 & \dots & \omega_n^0 \\ \omega_n^0 & \omega_n^1 & \dots & \omega_n^{n-1} \\ \omega_n^0 & \omega_n^2 & \dots & \omega_n^{2n-2} \\ \vdots & \vdots & & \vdots \\ \omega_n^0 & \omega_n^{n-1} & \dots & \omega_n^{(n-1)^2} \end{pmatrix} = \left(\omega_n^{lj} \right)_{l,j=0}^{n-1} \in \mathbb{C}^{n,n} . \tag{8.2.8}$$

Lemma 8.2.9 (Properties of Fourier matrix).

The scaled Fourier-matrix $\frac{1}{\sqrt{n}}\mathbf{F}_n$ is unitary (\rightarrow Def. 2.8.5): $\mathbf{F}_n^{-1} = \frac{1}{n}\mathbf{F}_n^H = \frac{1}{n}\overline{\mathbf{F}}_n$.

Lemma 8.2.10 (Diagonalization of circulant matrices (\rightarrow Def. 8.1.17)).

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

$$\mathbf{C}\overline{\mathbf{F}}_n = \overline{\mathbf{F}}_n \text{diag}(d_1, \dots, d_n) \quad , \quad \mathbf{d} = \mathbf{F}_n(u_0, \dots, u_{n-1})^\top .$$

Conclusion (from $\overline{\mathbf{F}}_n = n\mathbf{F}_n^{-1}$):

$$\mathbf{C} = \mathbf{F}_n^{-1} \text{diag}(d_1, \dots, d_n) \mathbf{F}_n . \quad (8.2.11)$$

Definition 8.2.12 (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 $\mathbf{c} := \mathcal{F}_n(\mathbf{y})$

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{kj} \quad , \quad k = 0, \dots, n-1 . \quad (8.2.13)$$

Terminology: $\mathbf{c} = \mathbf{F}_n \mathbf{y}$ is also called the (discrete) Fourier transform of \mathbf{y}

MATLAB-functions for discrete Fourier transform (and its inverse):

DFT: $\mathbf{c} = \text{fft}(\mathbf{y}) \leftrightarrow$ inverse DFT: $\mathbf{y} = \text{ifft}(\mathbf{c}) ;$

8.2.1 Discrete convolution via DFT

Code 8.2.16: discrete periodic convolution:
straightforward implementation

```
1 function z=pconv(u,x)
2 n = length(x); z = zeros(n,1);
3 for i=1:n, z(i)=dot(conj(u),
   x([i:-1:1,n:-1:i+1]));
4 end
```

Code 8.2.18: discrete periodic convolution:
DFT implementation

```
1 function z=pconvfft(u,x)
2 % Implementation of (8.2.13), cf.
   % Lemma 8.2.12
3 z = ifft(fft(u) .* fft(x));
```

Implementation of
discrete convolution
Def. 8.1.7) based
periodic discrete convolution

(\rightarrow)
on
 \triangleright

Built-in MATLAB-function:

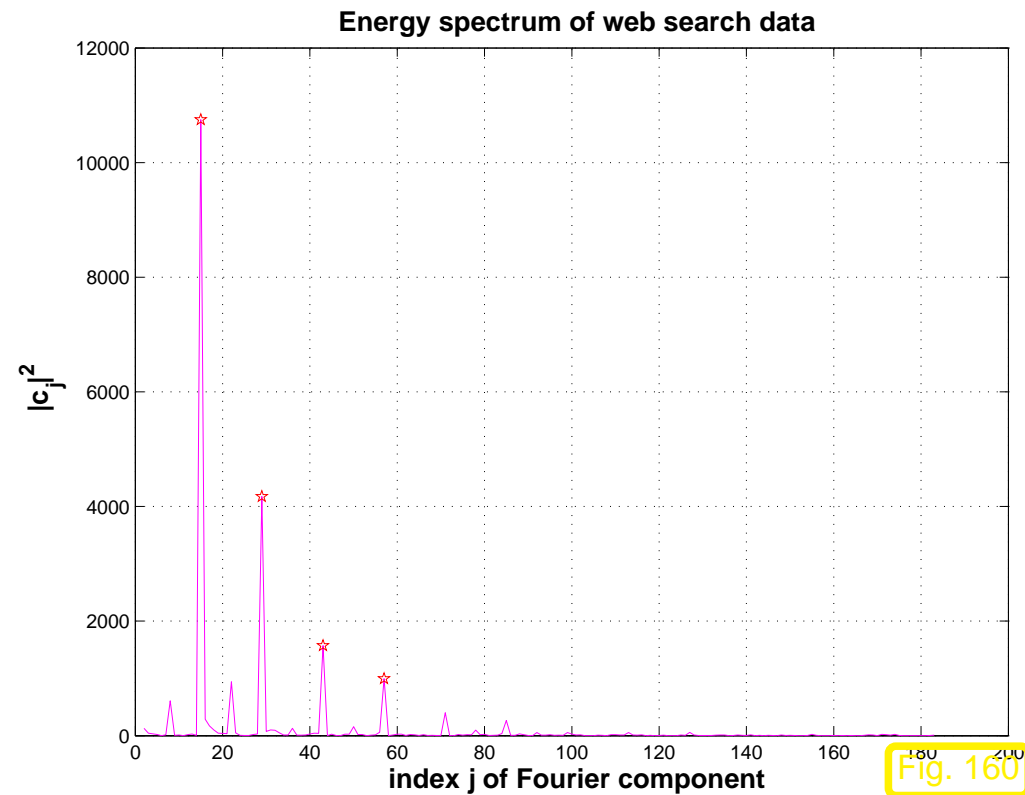
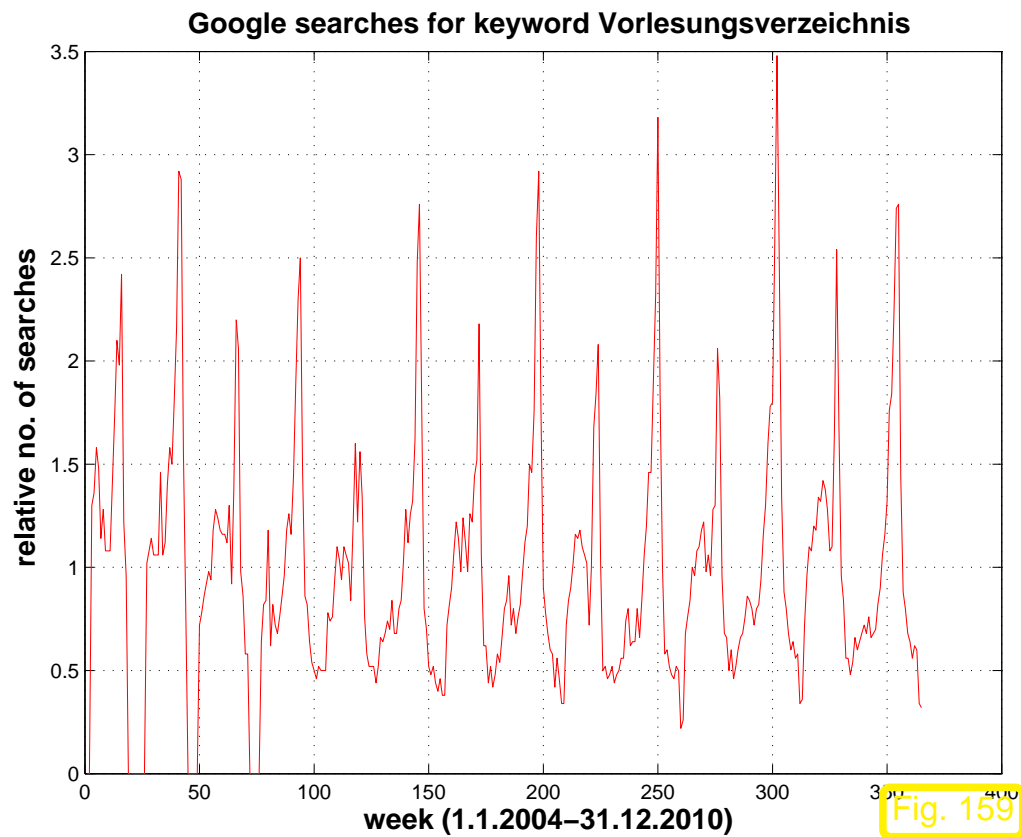
```
y = conv(h,x);
```

Code 8.2.19: discrete convolution: DFT implementation

```
1 function y = myconv(h,x)
2 n = length(h);
3 % Zero padding, cf. (8.1.20)
4 h = [h;zeros(n-1,1)];
5 x = [x;zeros(n-1,1)];
6 % Periodic discrete convolution of length 2n-1
7 y = pconvfft(h,x);
```

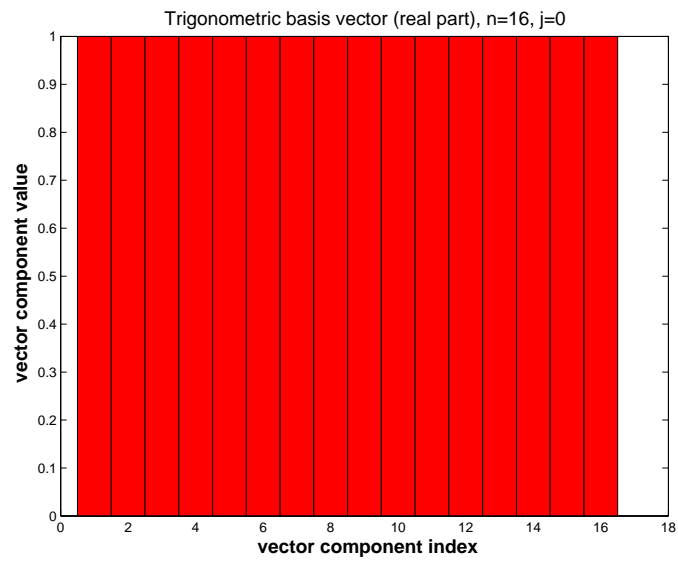
8.2.2 Frequency filtering via DFT

Example 8.2.21 (Detecting periodicity in data).

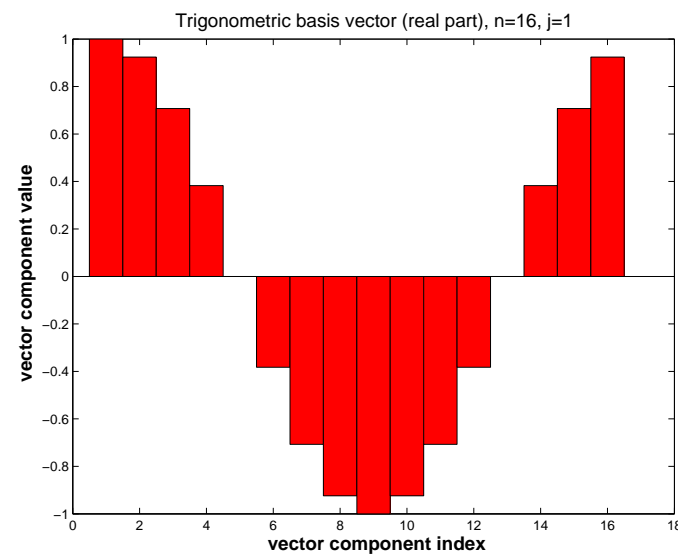


Remark 8.2.23 (“Low” and “high” frequencies).

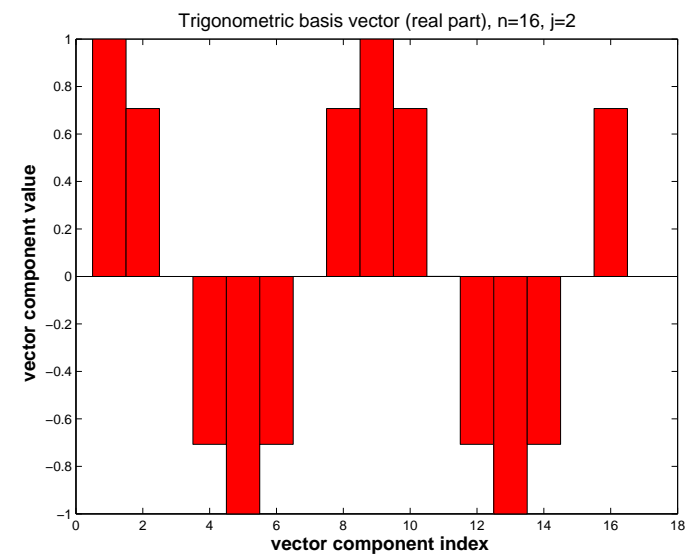
Plots of real parts of trigonometric basis vectors $(\mathbf{F}_n)_{:,j}$ (= columns of Fourier matrix), $n = 16$.



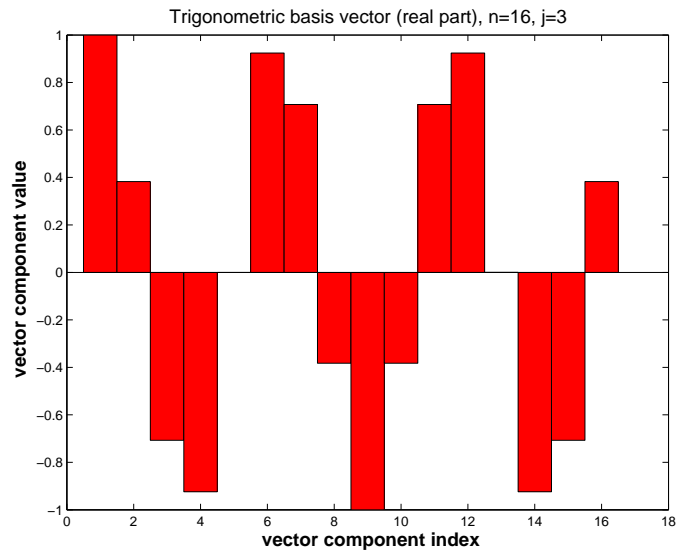
$\text{Re}(\mathbf{F}_{16})_{:,0}$



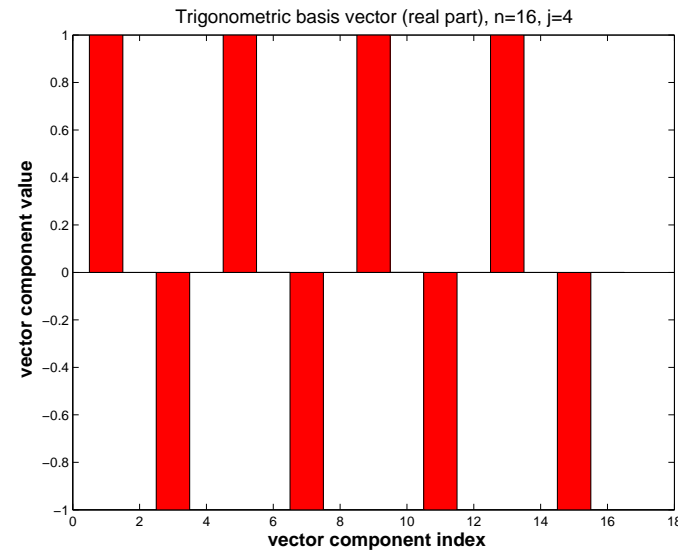
$\text{Re}(\mathbf{F}_{16})_{:,1}$



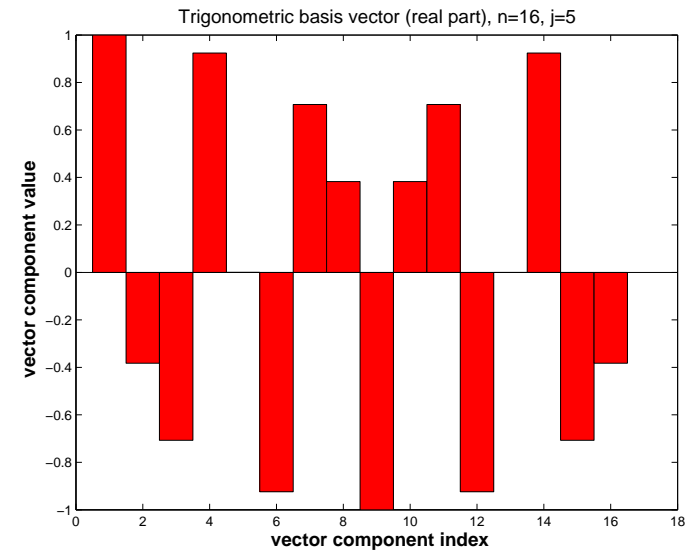
$\text{Re}(\mathbf{F}_{16})_{:,2}$



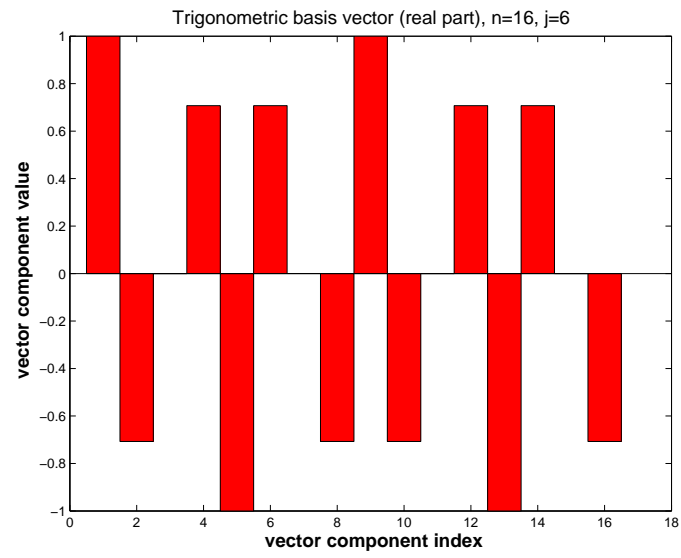
$\text{Re}(\mathbf{F}_{16})_{:,3}$



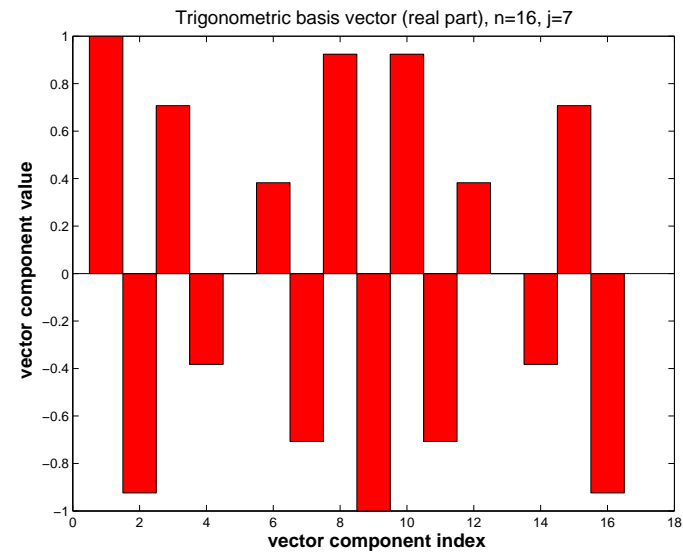
$\text{Re}(\mathbf{F}_{16})_{:,4}$



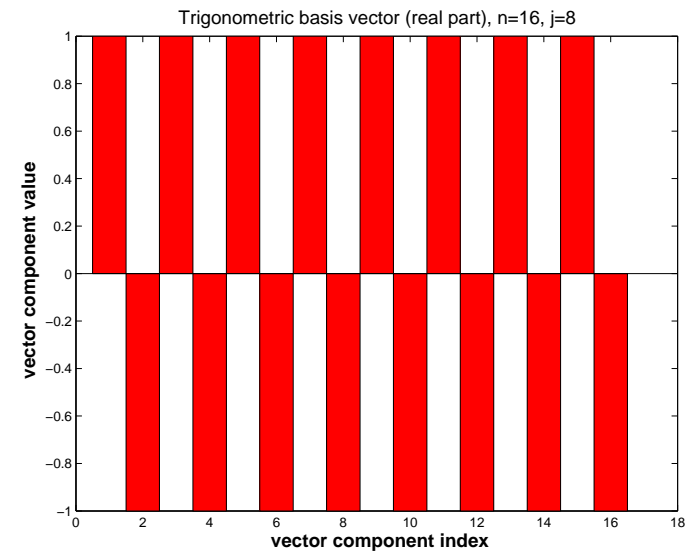
$\text{Re}(\mathbf{F}_{16})_{:,5}$



$\text{Re}(\mathbf{F}_{16})_{:,6}$



$\text{Re}(\mathbf{F}_{16})_{:,7}$



$\text{Re}(\mathbf{F}_{16})_{:,8}$



► Frequency filtering of real discrete periodic signals by suppressing certain “**Fourier coefficients**”.

Code 8.2.24: DFT-based frequency filtering

```

1 function [low,high] =
   freqfilter(y,k)
2 m = length(y)/2; c = fft(y);
3 clow = c; clow(m+1-k:m+1+k) = 0;
4 chigh = c-clow;
5 low = real(ifft(clow));
6 high = real(ifft(chigh));

```

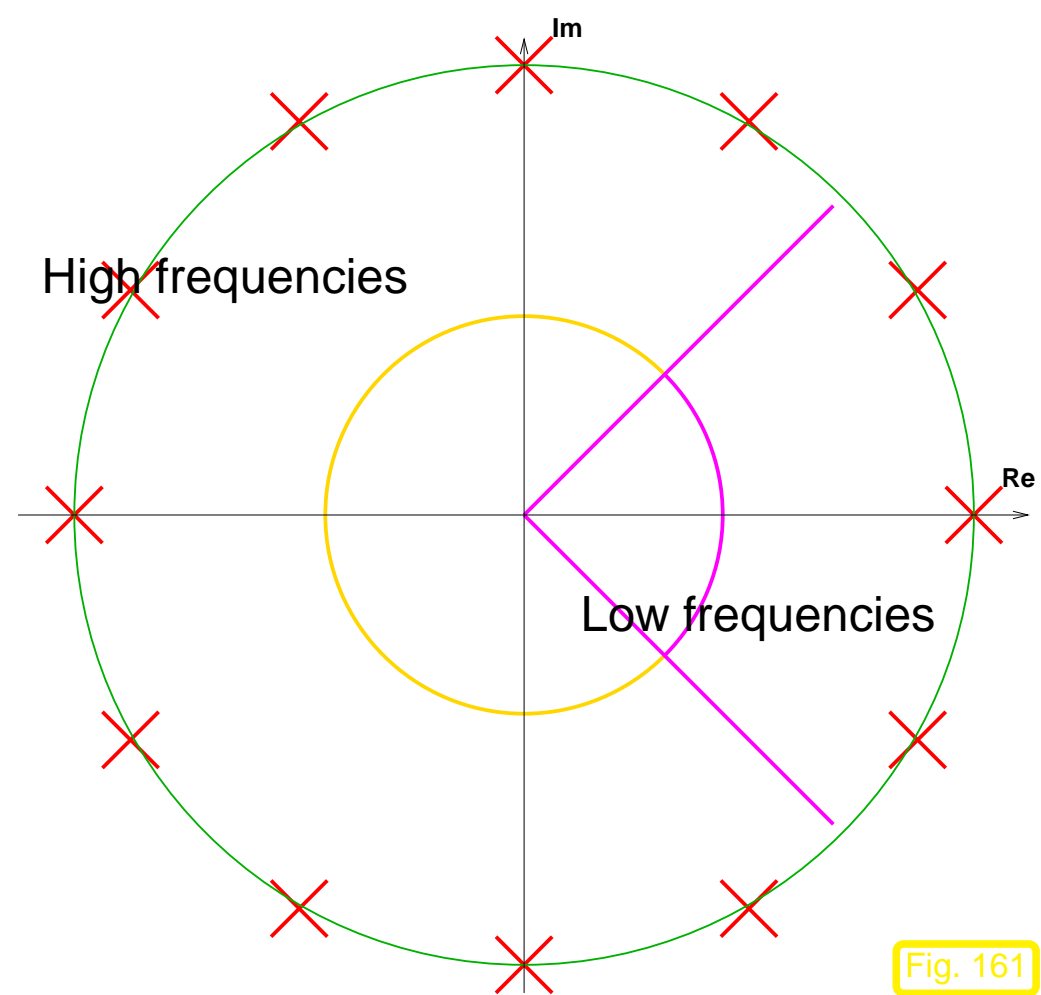


Fig. 161

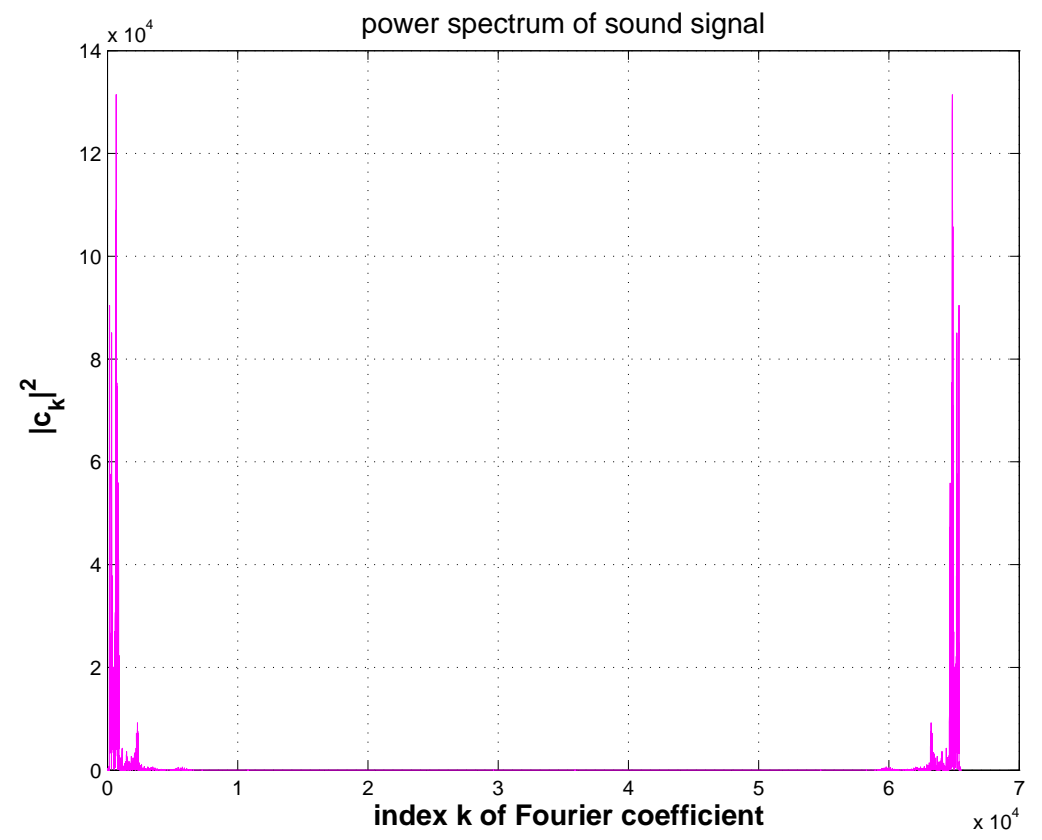
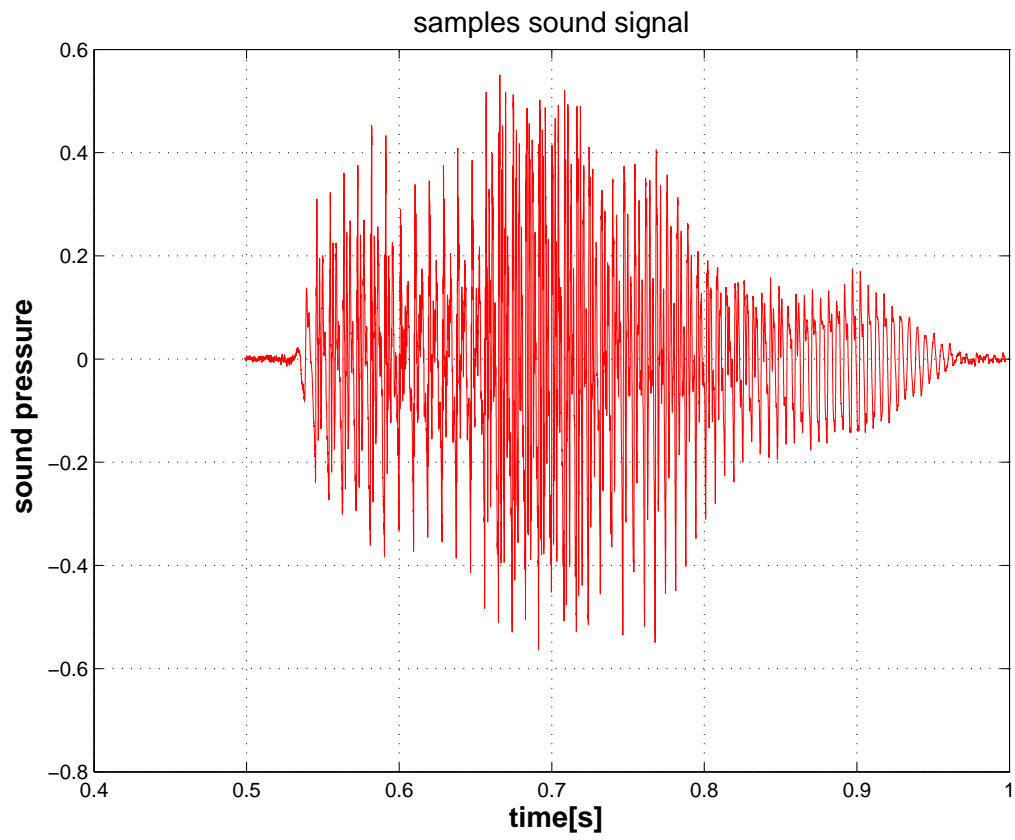
Example 8.2.26 (Sound filtering by DFT).

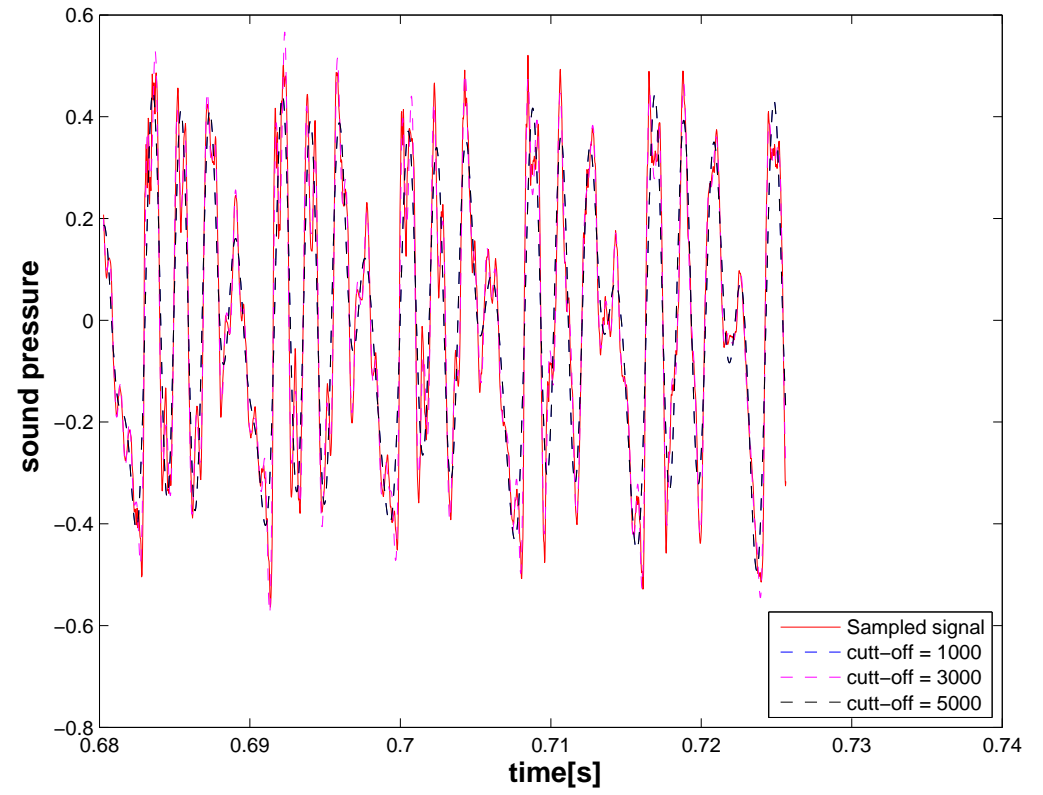
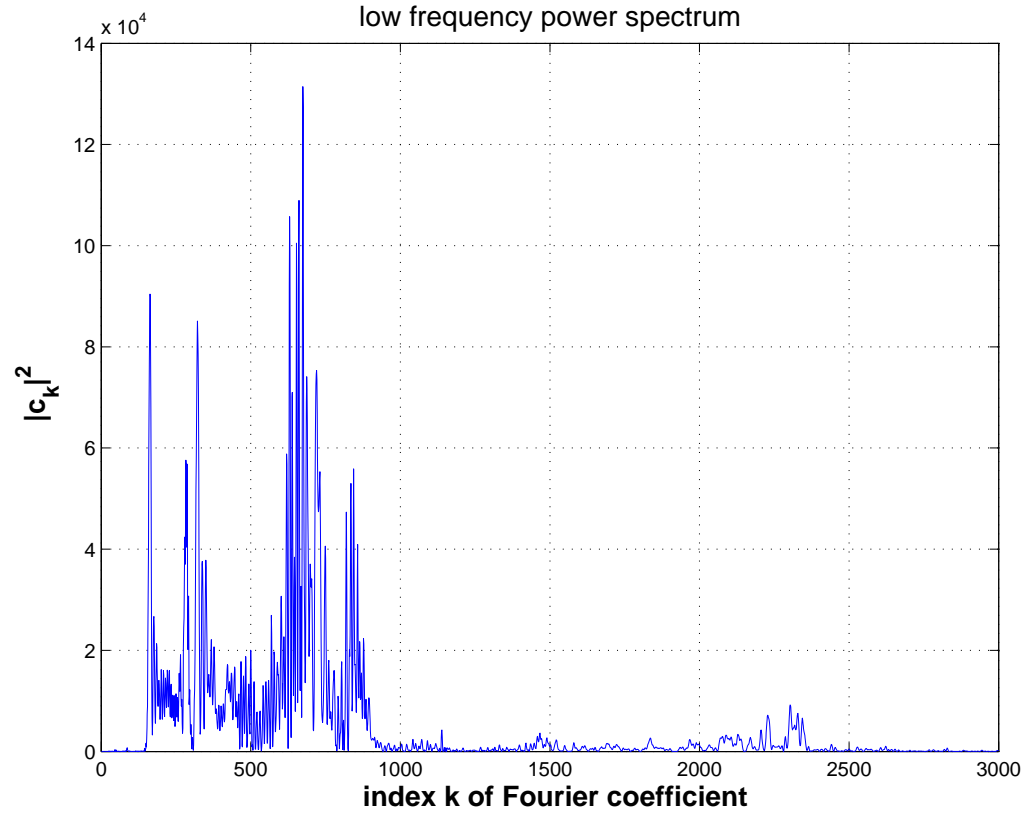
DFT based low pass frequency filtering of sound

```

[y,sf,nb] = wavread('hello.wav');
c = fft(y); c(m+1:end-m) = 0;
wavwrite(ifft(c),sf,nb,'filtered.wav');

```





$$h_k = \sum_{j=0}^{m-1} (y_{2j} + iy_{2j+1}) \omega_m^{jk} = \boxed{\sum_{j=0}^{m-1} y_{2j} \omega_m^{jk}} + i \cdot \boxed{\sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}}, \quad (8.2.28)$$

$$\bar{h}_{m-k} = \sum_{j=0}^{m-1} \overline{y_{2j} + iy_{2j+1}} \bar{\omega}_m^{j(m-k)} = \boxed{\sum_{j=0}^{m-1} y_{2j} \omega_m^{jk}} - i \cdot \boxed{\sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}}. \quad (8.2.29)$$

$$\Rightarrow \boxed{\sum_{j=0}^{m-1} y_{2j} \omega_m^{jk}} = \frac{1}{2}(h_k + \bar{h}_{m-k}), \quad \boxed{\sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}} = -\frac{1}{2}i(h_k - \bar{h}_{m-k}).$$

Use simple identities for roots of unity:

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{jk} = \boxed{\sum_{j=0}^{m-1} y_{2j} \omega_m^{jk}} + \omega_n^k \cdot \boxed{\sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}}. \quad (8.2.30)$$

$$\Rightarrow \begin{cases} c_k = \frac{1}{2}(h_k + \bar{h}_{m-k}) - \frac{1}{2}i\omega_n^k(h_k - \bar{h}_{m-k}), & k = 0, \dots, m-1, \\ c_m = \operatorname{Re}\{h_0\} - \operatorname{Im}\{h_0\}, \\ c_k = \bar{c}_{n-k}, & k = m+1, \dots, n-1. \end{cases} \quad (8.2.31)$$

```

1 function c = fftreal(y)
2 n = length(y); m = n/2;
3 if (mod(n,2) ~= 0), error('n must be even');
   end
4 y = y(1:2:n)+i*y(2:2:n); h = fft(y); h =
   [h;h(1)];
5 c = 0.5*(h+conj(h(m+1:-1:1))) - ...
6     (0.5*i*exp(-2*pi*i/n).^((0:m)')).*...
7     (h-conj(h(m+1:-1:1)));
8 c = [c;conj(c(m:-1:2))];

```

MATLAB-Implementation
(by a DFT of length $n/2$):

8.2.4 Two-dimensional DFT

Two-dimensional trigonometric basis of $\mathbb{C}^{m,n}$:

$$\text{tensor product matrices } \left\{ (\mathbf{F}_m)_{:,j} (\mathbf{F}_n)_{:,l}^\top, 1 \leq j \leq m, 1 \leq l \leq n \right\}. \quad (8.2.33)$$

Basis transform: for $y_{j_1,j_2} \in \mathbb{C}$, $0 \leq j_1 < m$, $0 \leq j_2 < n$ compute (nested DFTs !)

$$c_{k_1,k_2} = \sum_{j_1=0}^{m-1} \sum_{j_2=0}^{n-1} y_{j_1,j_2} \omega_m^{j_1 k_1} \omega_n^{j_2 k_2}, \quad 0 \leq k_1 < m, 0 \leq k_2 < n.$$

Two-dimensional DFT by *nested one-dimensional DFTs* (8.2.15):

$$\text{fft2}(Y) = \text{fft}(\text{fft}(Y) \cdot') \cdot'$$

Example 8.2.34 (Deblurring by DFT).

Gray-scale pixel image $\mathbf{P} \in \mathbb{R}^{m,n}$, actually $\mathbf{P} \in \{0, \dots, 255\}^{m,n}$, see Ex. 6.3.22.

$(p_{l,k})_{l,k \in \mathbb{Z}} \hat{=}$ 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)

$$c_{l,j} = \sum_{k=-L}^L \sum_{q=-L}^L s_{k,q} p_{l+k,j+q}, \quad \begin{array}{l} 0 \leq l < m, \\ 0 \leq j < n, \end{array} \quad L \in \{1, \dots, \min\{m, n\}\} . \quad (8.2.35)$$

blurred image

point spread function

Usually: L small, $s_{k,m} \geq 0$, $\sum_{k=-L}^L \sum_{q=-L}^L s_{k,q} = 1$ (an averaging)

Used in test calculations: $L = 5$

$$s_{k,q} = \frac{1}{1 + k^2 + q^2}.$$

Code 8.2.36: point spread function

```

1 function S = psf(L)
2 [X,Y] = meshgrid(-L:1:L,-L:1:L);
3 S = 1./(1+X.*X+Y.*Y);
4 S = S/sum(sum(S));

```

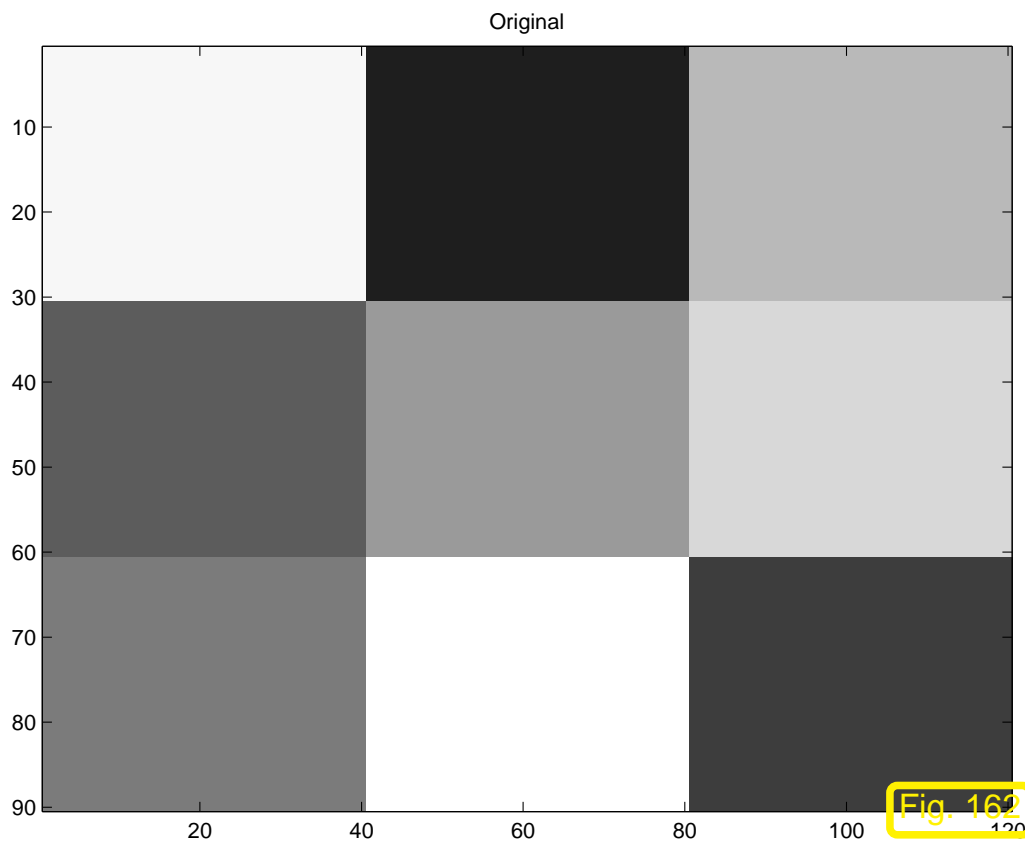


Fig. 162

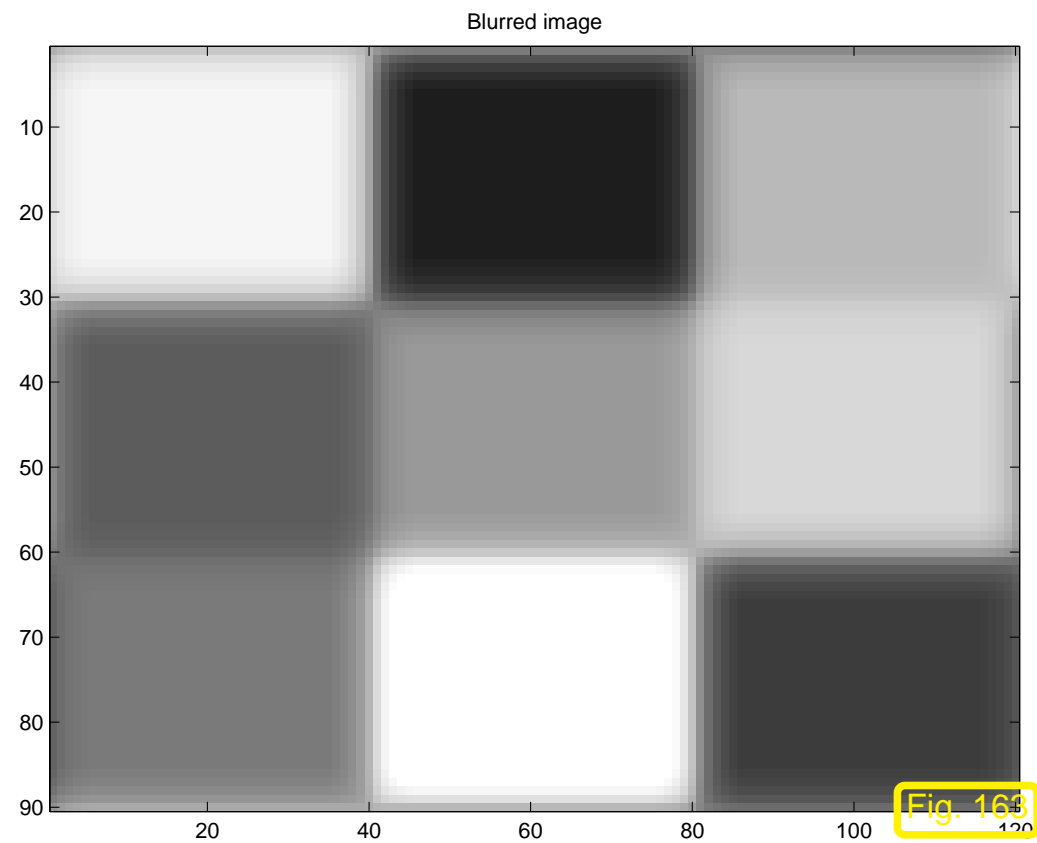


Fig. 163

Code 8.2.38: blurring operator

NumCSE,
autumn
2010

```

1 function C = blur(P,S)
2 [m,n] = size(P); [M,N] = size(S);
3 if (M ~= N), error('S not quadratic'); end
4 L = (M-1)/2; C = zeros(m,n);
5 for l=1:m, for j=1:n
6     s = 0;
7     for k=1:(2*L+1), for q=1:(2*L+1)
8         kl = l+k-L-1;
9         if (kl < 1), kl = kl + m; end
10        if (kl > m), kl = kl - m; end
11        jm = j+q-L-1;
12        if (jm < 1), jm = jm + n; end
13        if (jm > n), jm = jm - n; end
14        s = s+P(kl,jm)*S(k,q);
15        end, end
16        C(l,j) = s;
17 end, end

```

Note:

(8.2.37) defines a linear
operator

$$\mathcal{B} : \mathbb{R}^{m,n} \mapsto \mathbb{R}^{m,n}$$

(“blurring operator”)

Note: more efficient im-
plementation via MAT-
LAB function `conv2(P,S)`

R. Hiptmair
rev 38355,
November
17, 2011

Code 8.2.40: DFT based deblurring

```

1 function D = deblur(C,S,tol)
2 [m,n] = size(C); [M,N] = size(S);
3 if (M ~= N), error('S not quadratic'); end
4 L = (M-1)/2; Spad = zeros(m,n);
5 % Zero padding
6 Spad(1:L+1,1:L+1) = S(L+1:end,L+1:end);
7 Spad(m-L+1:m,n-L+1:n) = S(1:L,1:L);
8 Spad(1:L+1,n-L+1:n) = S(L+1:end,1:L);
9 Spad(m-L+1:m,1:L+1) = S(1:L,L+1:end);
10 % Inverse of blurring operator
11 SF = fft2(Spad);
12 % Test for invertibility
13 if (nargin < 3), tol = 1E-3; end
14 if (min(min(abs(SF))) < tol*max(max(abs(SF))))
15     error('Deblurring impossible');
16 end
17 % DFT based deblurring
18 D = fft2(ifft2(C)./SF);

```

Inversion of blurring
operator



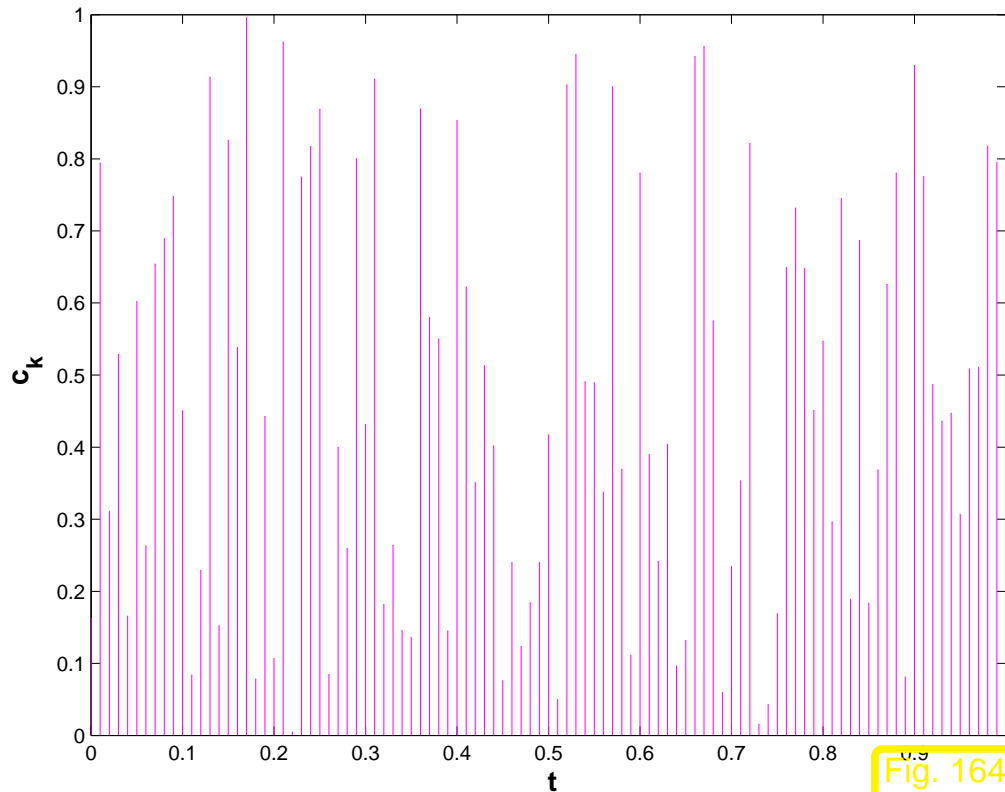
componentwise scaling
in “Fourier domain”



8.2.5 Semi-discrete Fourier Transform [51, Sect. 10.11]

$(y_k)_{k \in \mathbb{Z}}$: n -periodic sequence (signal), $n = 2m + 1$, $m \in \mathbb{N}$:

$$\text{DFT: } c_k = \sum_{j=-m}^m y_j \exp(-2\pi i \frac{kj}{n}), \quad k = 0, \dots, n-1. \quad (8.2.41)$$



◁ “Squeezing” a vector $\in \mathbb{R}^n$ into $[0, 1[$.

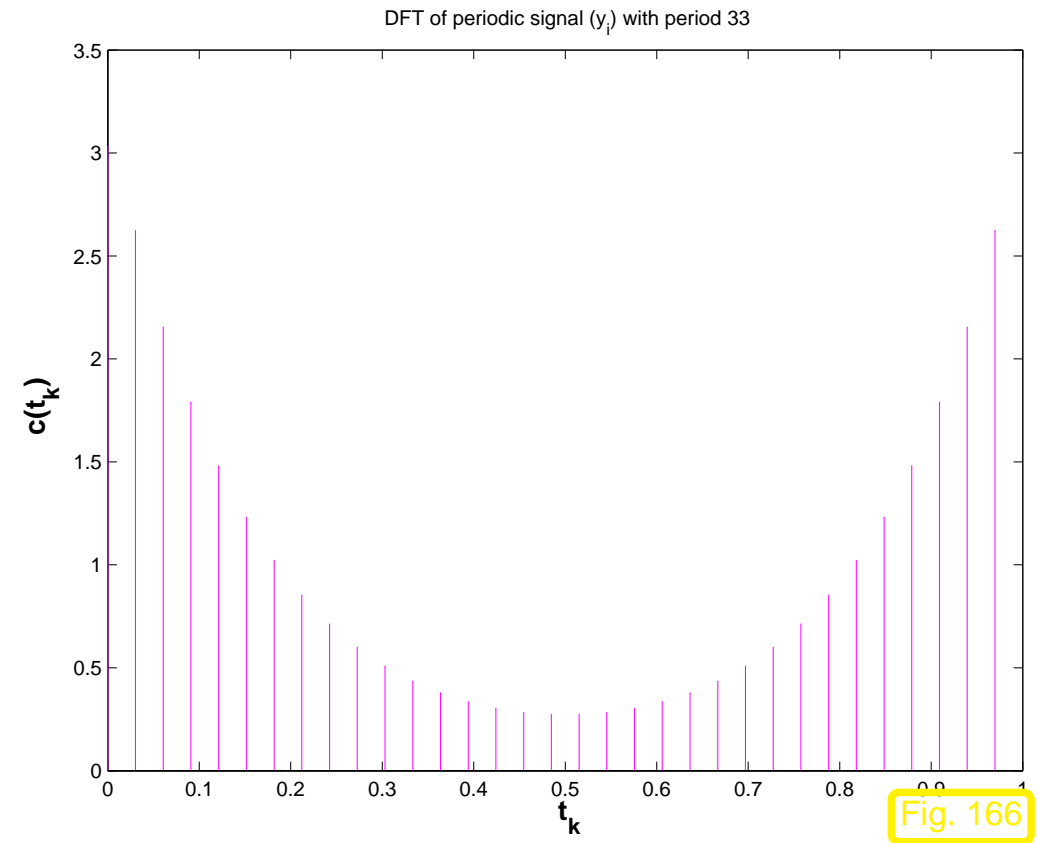
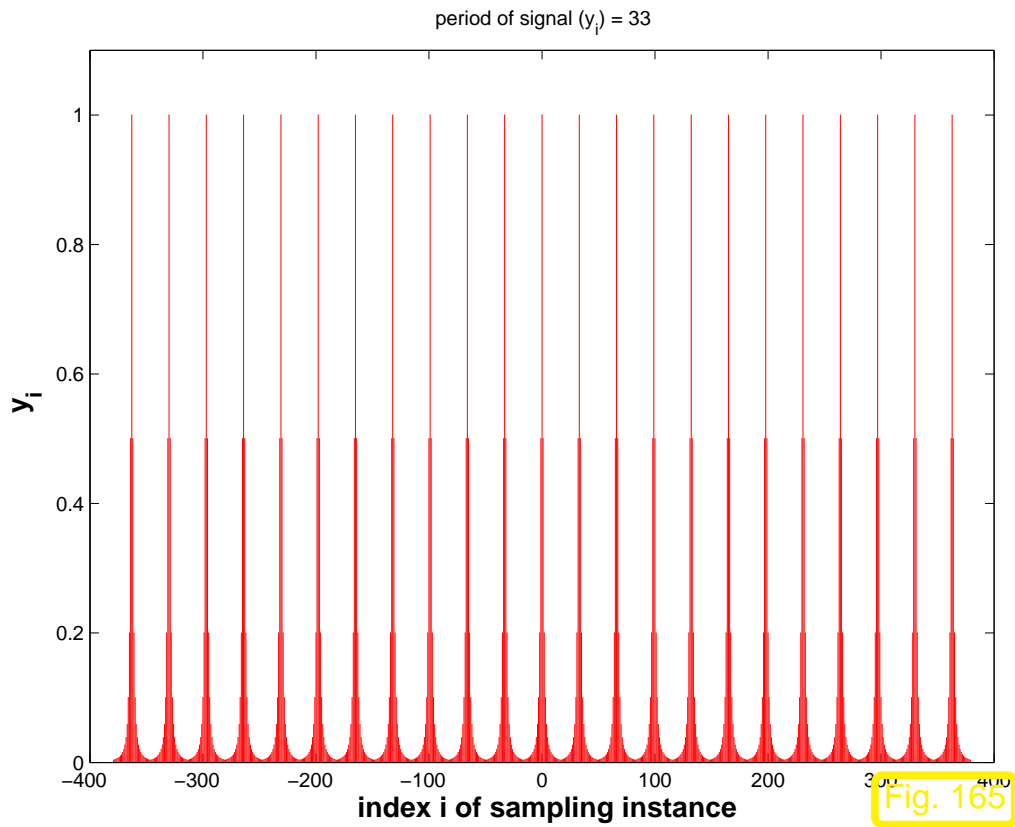
$$c_k \leftrightarrow c(t_k); ,$$

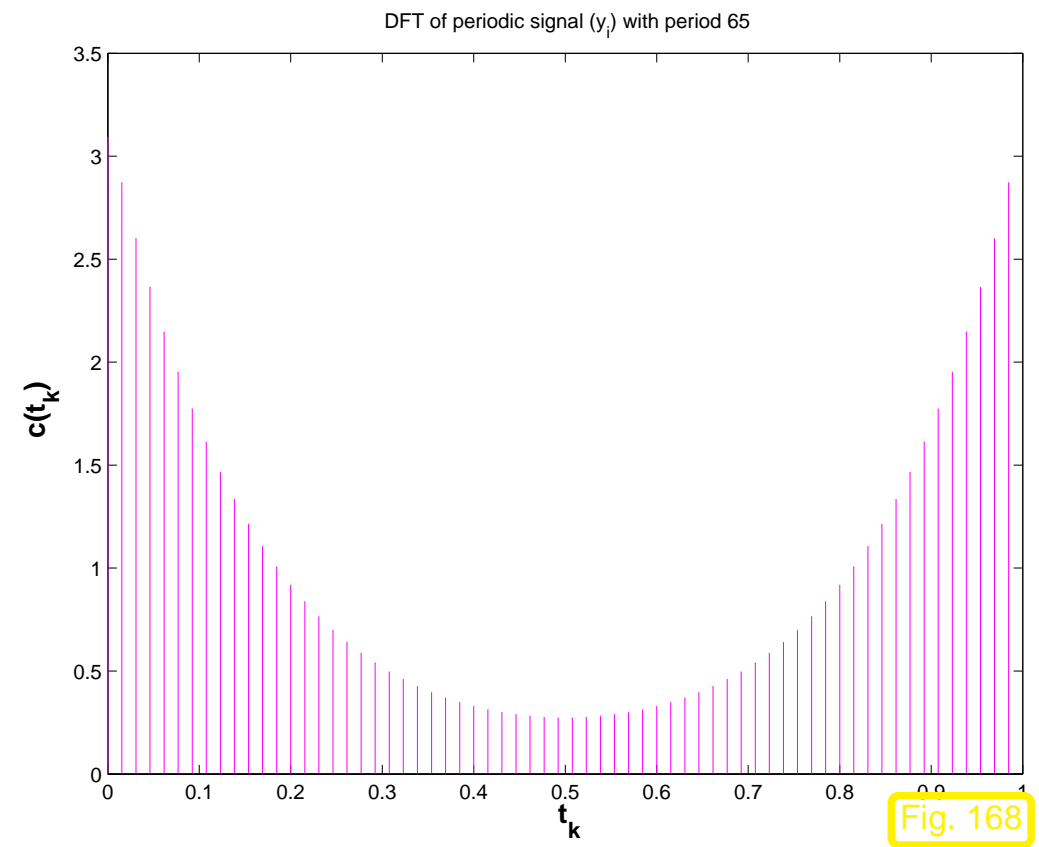
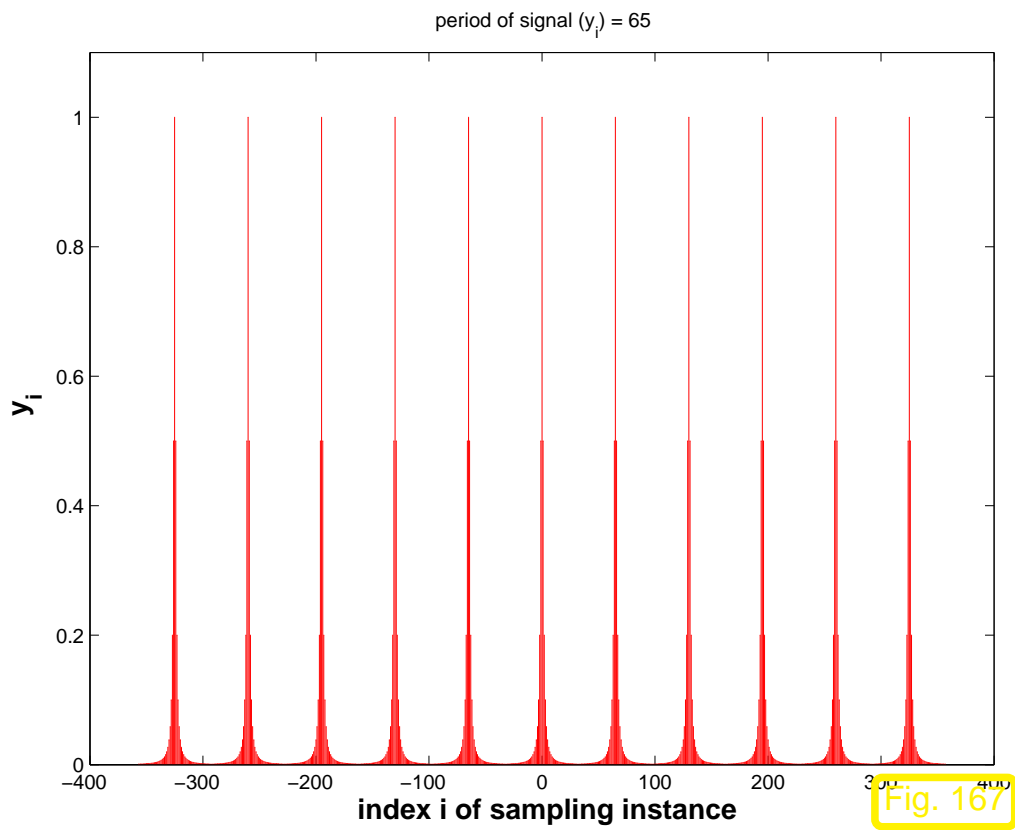
$$t_k = \frac{k}{n}, \quad k = 0, \dots, n-1.$$

Example 8.2.44 (“Squeezed” DFT of a periodically truncated signal).

Bi-infinite discrete signal, “concentrated around 0”:

$$y_j = \frac{1}{1 + j^2}, \quad j \in \mathbb{Z}.$$





period of signal (y_i) = 129

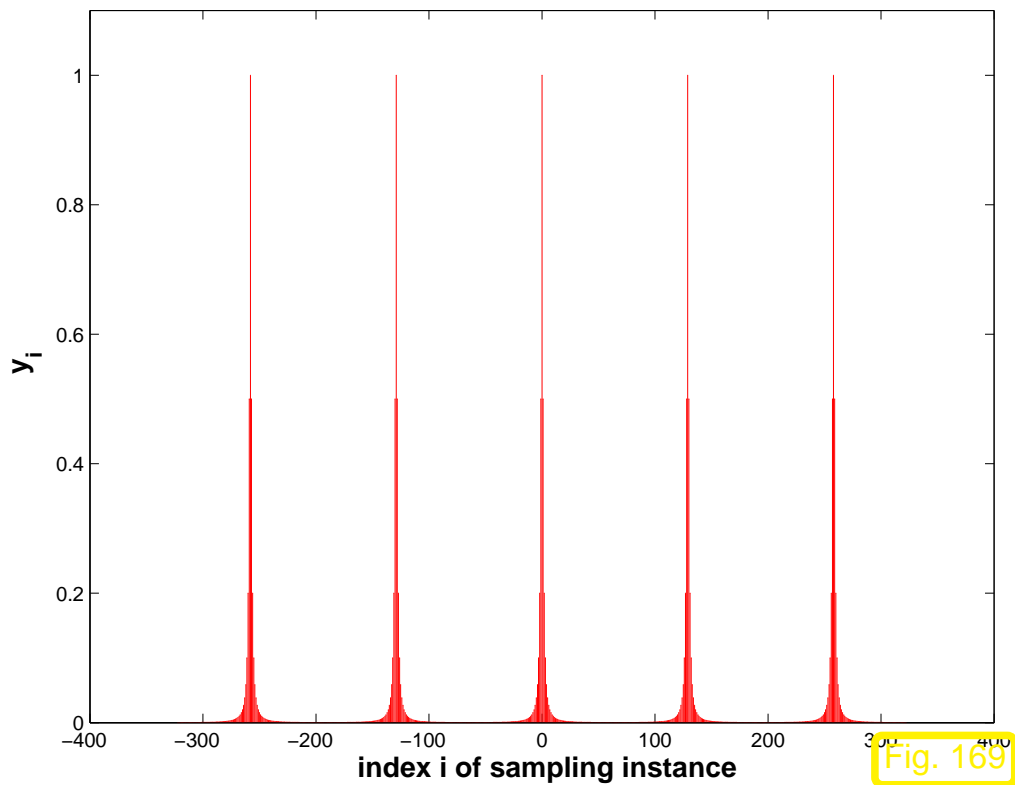


Fig. 169

DFT of periodic signal (y_i) with period 129

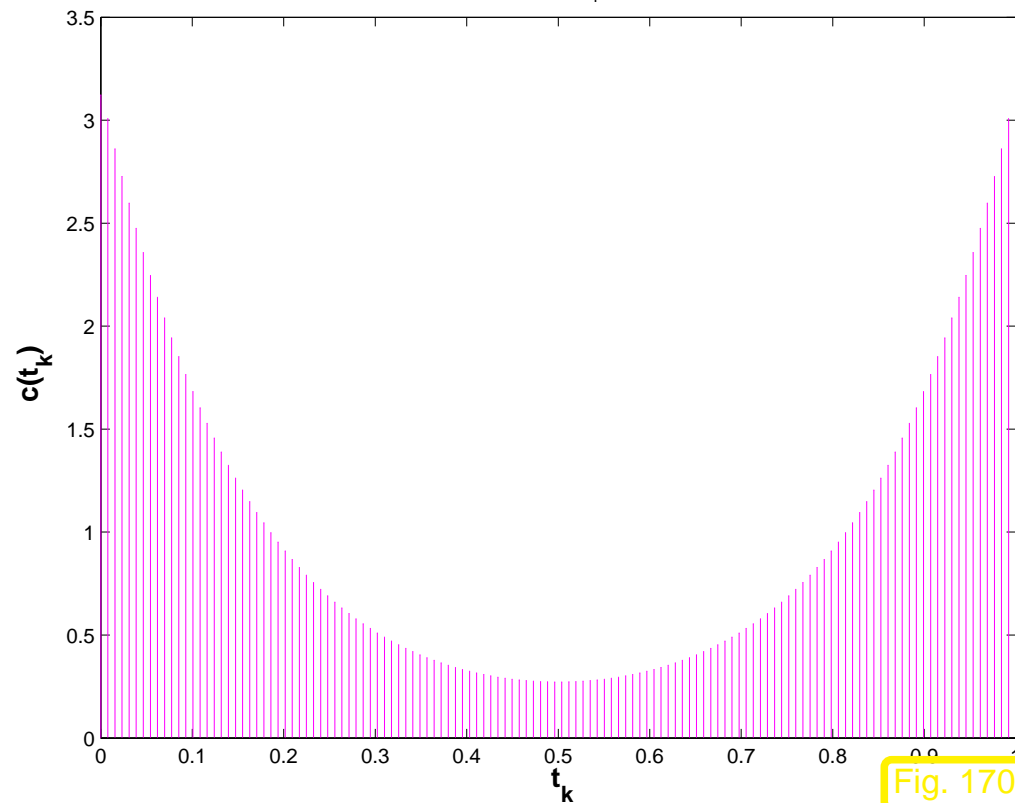


Fig. 170

period of signal (y_i) = 257

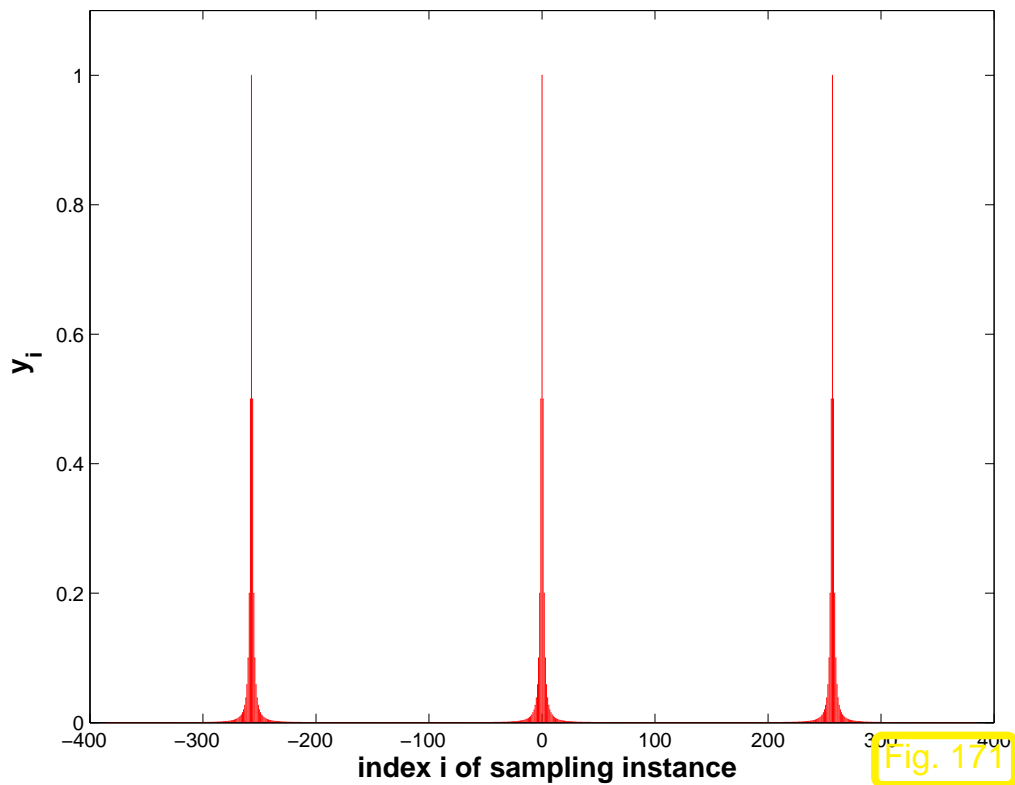


Fig. 171

DFT of periodic signal (y_i) with period 257

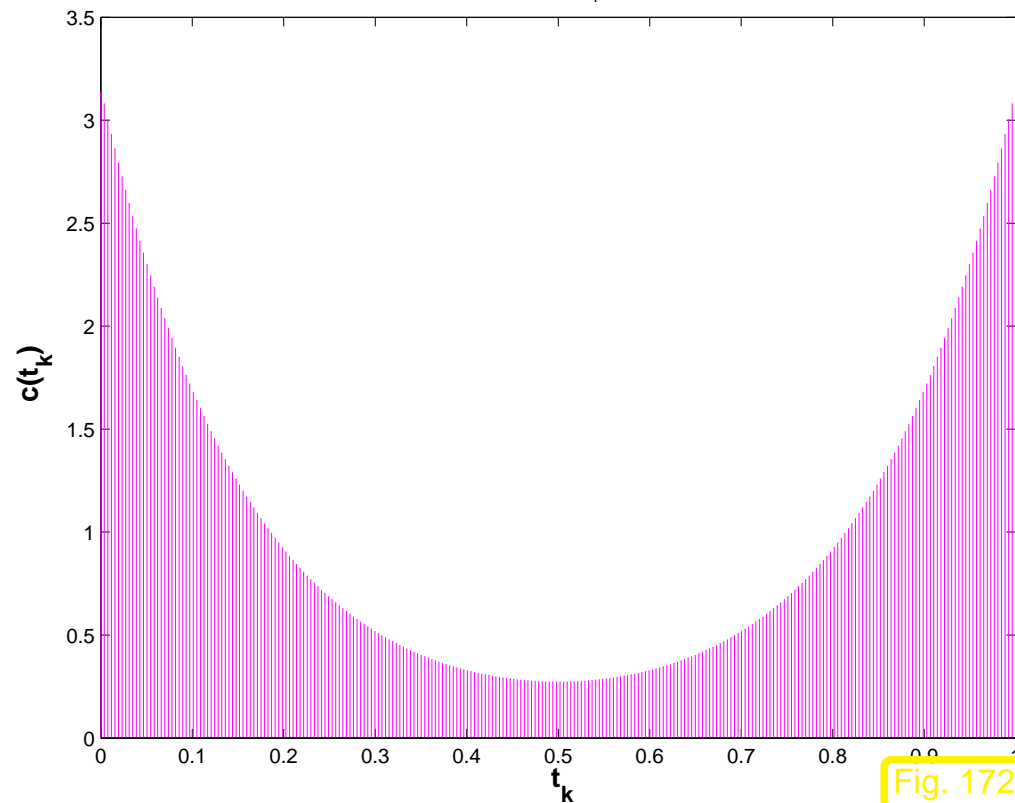
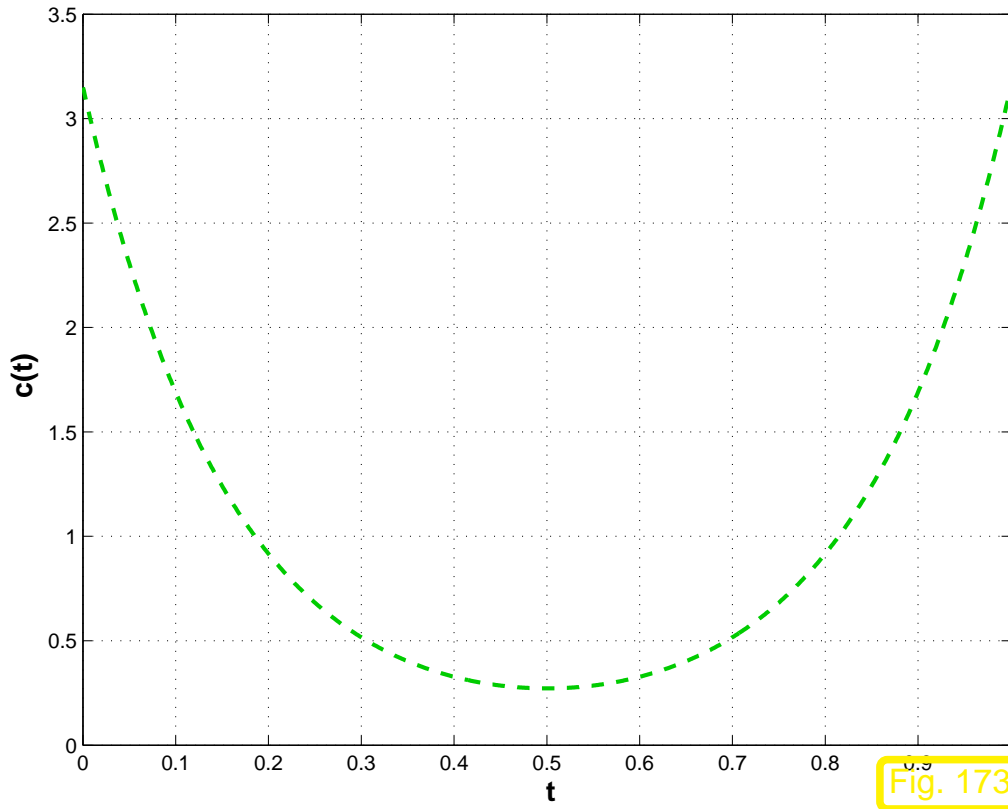


Fig. 172

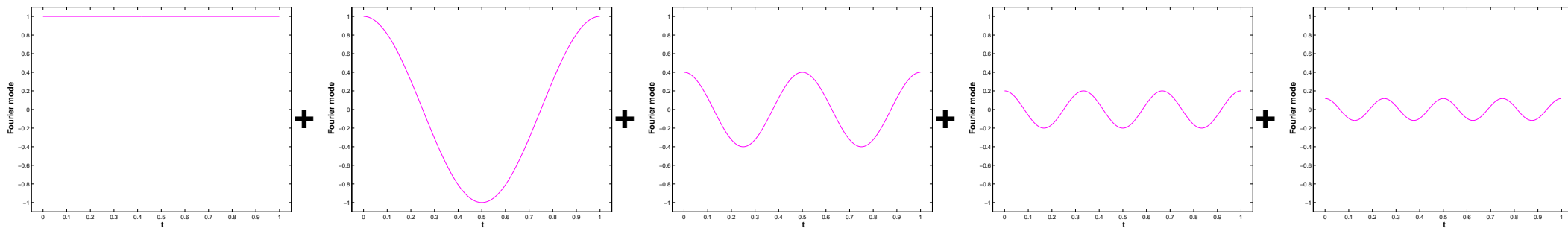


Fourier transform of $(1/1+k^2)_k$



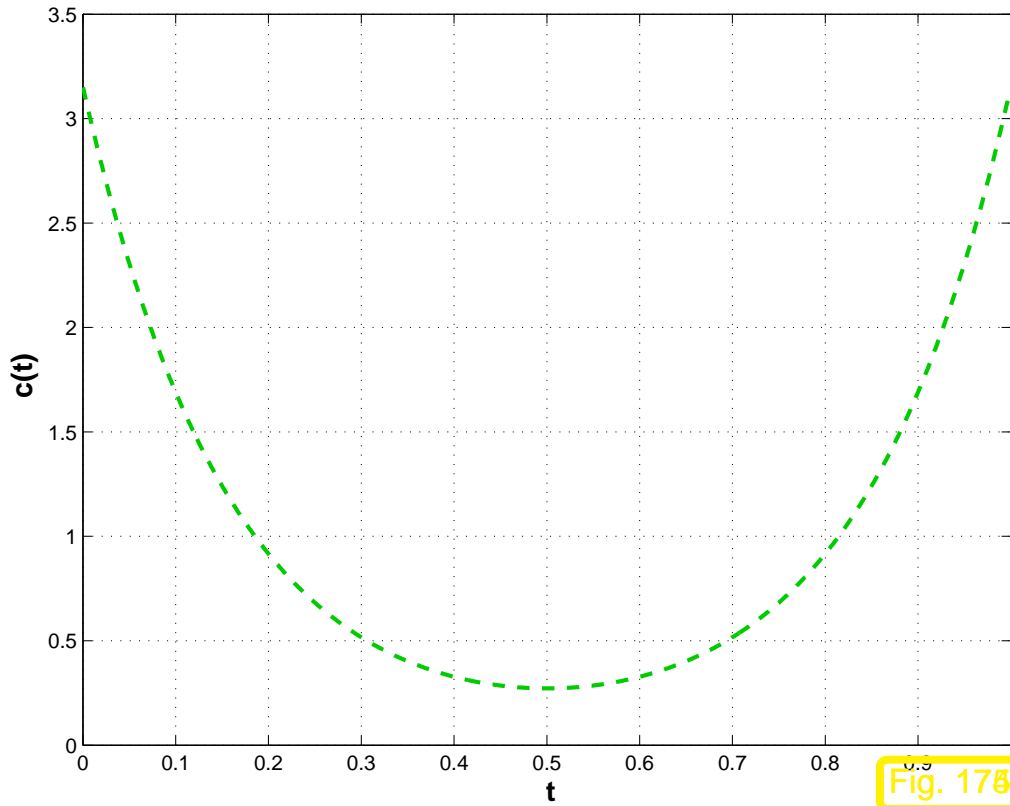
◁ Fourier transform of $y_k := \frac{1}{1+k^2}$

Fourier transform
=
sum of Fourier modes

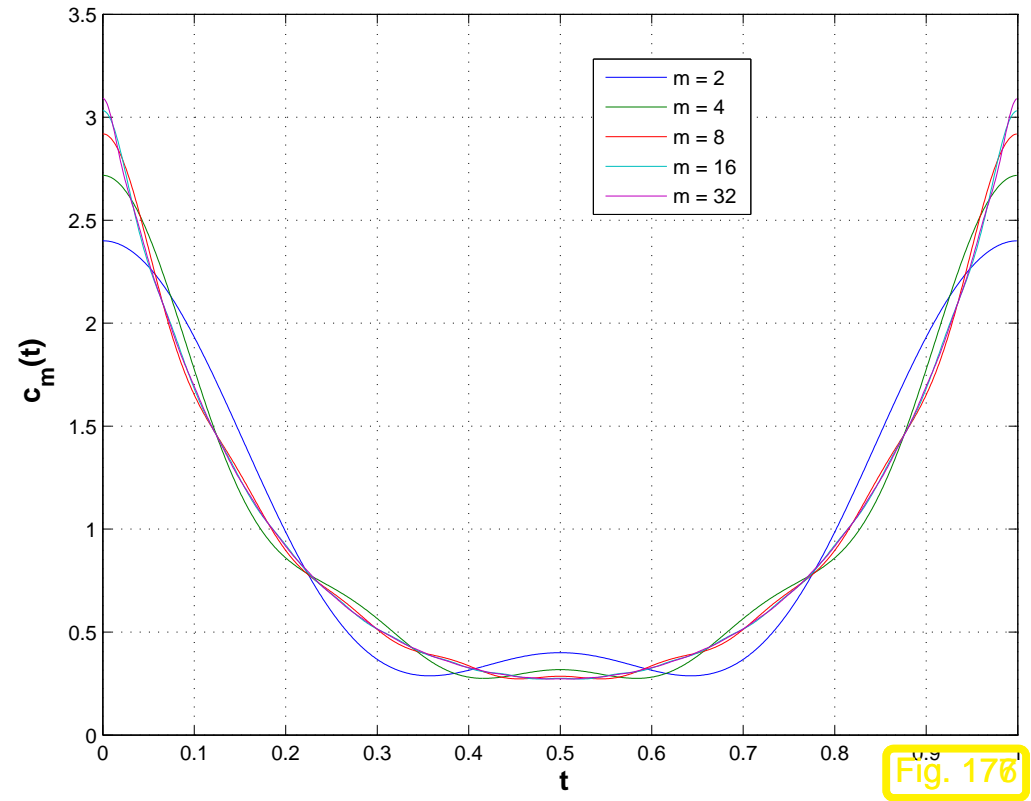


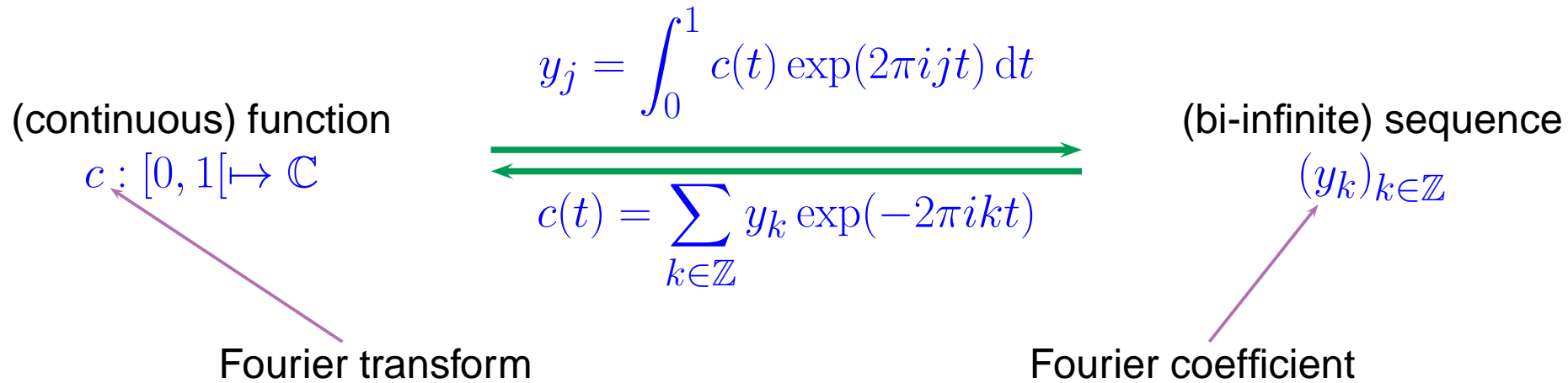
Example 8.2.55 (Convergence of Fourier sums).

Fourier transform of $(1/1+k^2)_k$



Fourier sum approximations with $2m+1$ terms, $y_k = 1/(1+k^2)$





Remark 8.2.59 (Filtering in Fourier domain).



Theorem 8.2.65 (Isometry property of Fourier transform).

If $\sum_{k \in \mathbb{Z}} |y_j|^2 < \infty$, then

$$c(t) = \sum_{k \in \mathbb{Z}} y_k \exp(-2\pi i k t) \Rightarrow \int_0^1 |c(t)|^2 dt = \sum_{k \in \mathbb{Z}} |y_j|^2 .$$

8.3 Fast Fourier Transform (FFT) [13, Sect. 8.7.3], [35, Sect. 53], [51, Sect. 10.9.2]

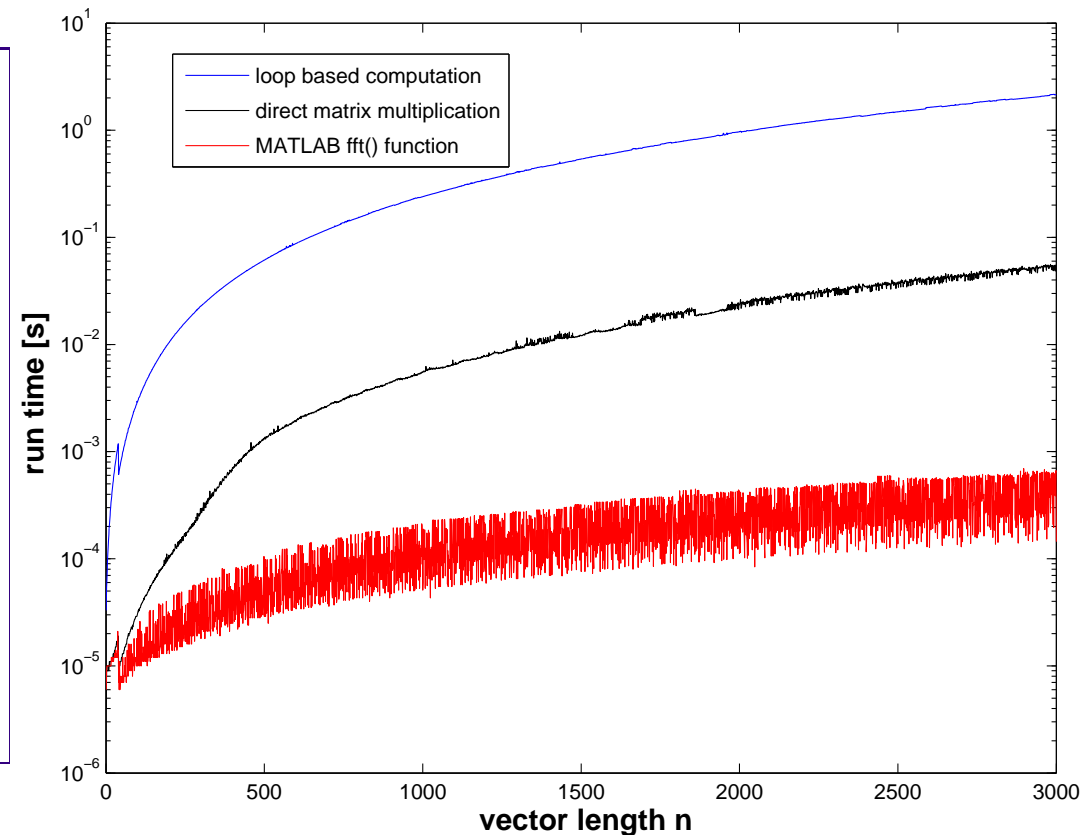
Example 8.3.1 (Efficiency of `fft`).

`tic-toc`-timing in MATLAB: compare `fft`, loop based implementation, and direct matrix multiplication

(MATLAB V6.5, Linux, Mobile Intel Pentium 4 - M CPU 2.40GHz, minimum over 5 runs)

MATLAB-CODE naive DFT-implementation

```
c = zeros(n,1);
omega = exp(-2*pi*i/n);
c(1) = sum(y); s = omega;
for j=2:n
    c(j) = y(n);
    for k=n-1:-1:1
        c(j) = c(j)*s+y(k);
    end
    s = s*omega;
end
```



The secret of MATLAB's `fft()`:

the **Fast Fourier Transform** algorithm [18]

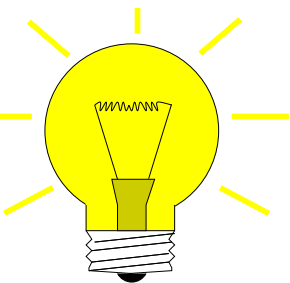
(discovered by C.F. Gauss in 1805, rediscovered by Cooley & Tuckey in 1965,
one of the “top ten algorithms of the century”).

Idea:

divide & conquer recursion

(for DFT of length $n = 2^L$)

FFT-algorithm



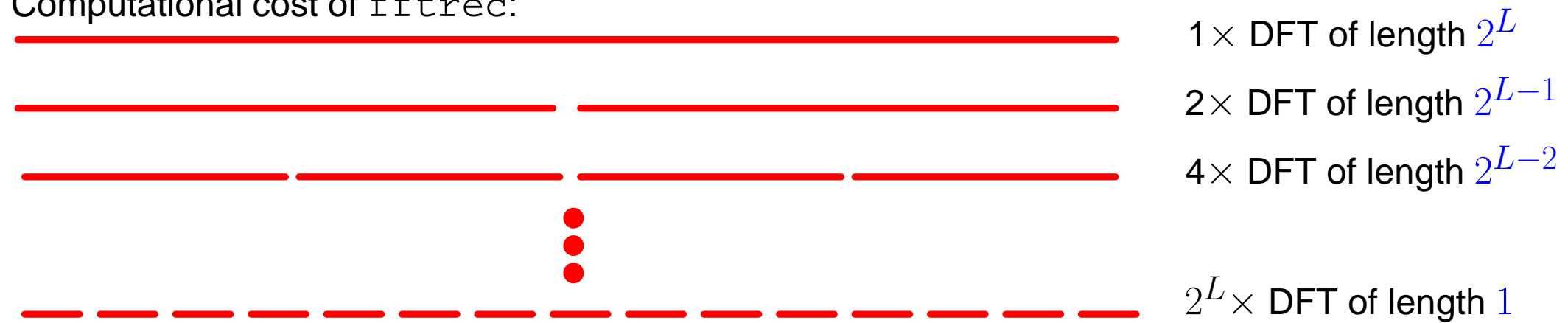
Code 8.3.3: Recursive FFT

```

1 function c = fftrec(y)
2 n = length(y);
3 if (n == 1), c = y; return;
4 else
5     c1 = fftrec(y(1:2:n));
6     c2 = fftrec(y(2:2:n));
7     c = [c1;c1] +
          (exp(-2*pi*i/n).^((0:n-1)'))
          .*[c2;c2];
8 end

```

Computational cost of `fftrec`:



Code 8.3.3: each level of the recursion requires $O(2^L)$ elementary operations.

Asymptotic complexity of FFT algorithm, $n = 2^L$: $O(L2^L) = O(n \log_2 n)$

(MATLAB `fft`-function: cost $\approx 5n \log_2 n$).


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 \rightarrow Ex. 8.3.1.

Asymptotic complexity of `c=fft(y)` for $\mathbf{y} \in \mathbb{C}^n = O(n \log n)$.

 \leftarrow Sect. 8.2.1

Asymptotic complexity of discrete periodic convolution/multiplication with circulant matrix, see Code 8.2.18:

$$\text{Cost}(\mathbf{z} = \text{pconvfft}(\mathbf{u}, \mathbf{x}), \mathbf{u}, \mathbf{x} \in \mathbb{C}^n) = O(n \log n).$$

Asymptotic complexity of discrete convolution, see Code 8.2.20:

$$\text{Cost}(\mathbf{z} = \text{myconv}(\mathbf{h}, \mathbf{x}), \mathbf{h}, \mathbf{x} \in \mathbb{C}^n) = O(n \log n).$$

8.4 Trigonometric transformations [35, Sect. 55]

8.4.1 Sine transform

Another trigonometric basis transform in \mathbb{R}^{n-1} , $n \in \mathbb{N}$:

$$\left\{ \begin{array}{c} \left(\begin{array}{c} 1 \\ 0 \\ \vdots \\ 0 \end{array} \right) \left(\begin{array}{c} 0 \\ 1 \\ \vdots \\ 0 \end{array} \right) \dots \left(\begin{array}{c} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \end{array} \right) \left(\begin{array}{c} 0 \\ \vdots \\ 0 \\ 0 \\ 1 \end{array} \right) \end{array} \right\} \leftarrow \left\{ \begin{array}{c} \left(\begin{array}{c} \sin(\frac{\pi}{n}) \\ \sin(\frac{2\pi}{n}) \\ \vdots \\ \sin(\frac{(n-1)\pi}{n}) \end{array} \right) \left(\begin{array}{c} \sin(\frac{2\pi}{n}) \\ \sin(\frac{4\pi}{n}) \\ \vdots \\ \sin(\frac{2(n-1)\pi}{n}) \end{array} \right) \dots \left(\begin{array}{c} \sin(\frac{(n-1)\pi}{n}) \\ \sin(\frac{2(n-1)\pi}{n}) \\ \vdots \\ \sin(\frac{(n-1)^2\pi}{n}) \end{array} \right) \end{array} \right\}$$

“Sine basis”

Basis transform matrix (sine basis \rightarrow standard basis): $\mathbf{S}_n := (\sin(jk\pi/n))_{j,k=1}^{n-1} \in \mathbb{R}^{n-1, n-1}$.

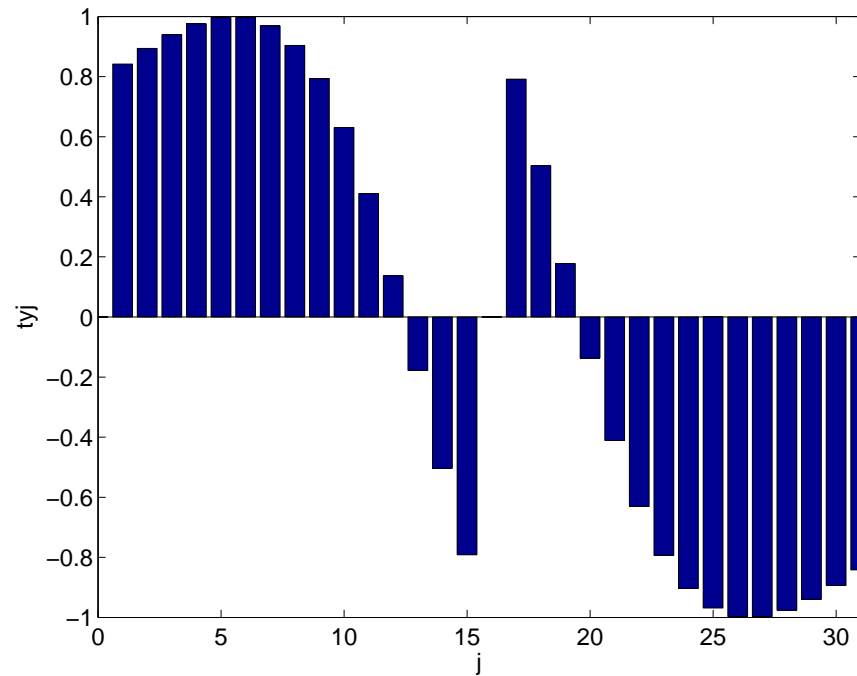
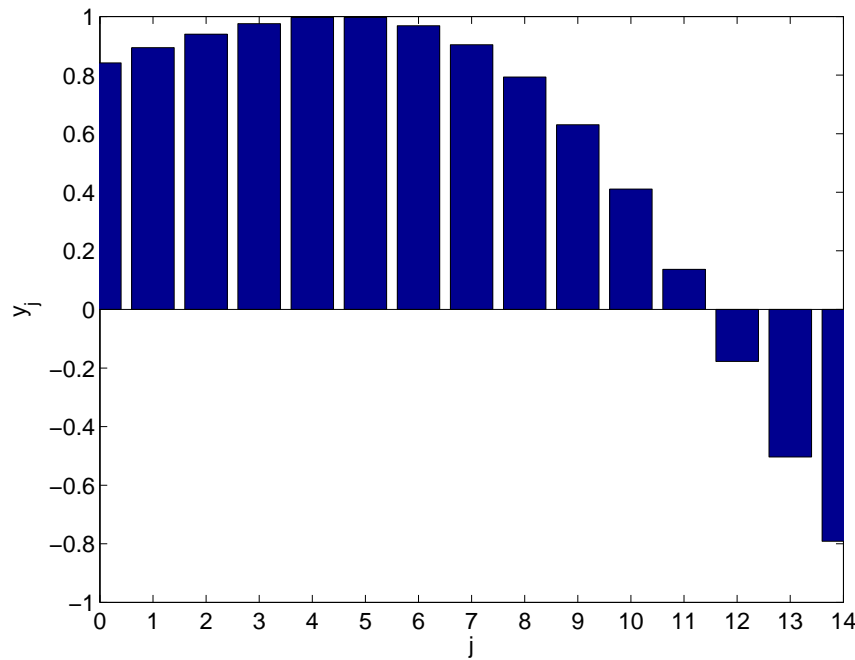
Lemma 8.4.1 (Properties of the sine matrix).

$\sqrt{2/n} \mathbf{S}_n$ is real, symmetric and orthogonal (\rightarrow Def. 2.8.5)

Sine transform: $s_k = \sum_{j=1}^{n-1} y_j \sin(\pi jk/n)$, $k = 1, \dots, n-1$. (8.4.2)

DFT-based algorithm for the sine transform ($\hat{=} \mathbf{S}_n \times \text{vector}$):

“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”)



$$\begin{aligned}
 (\mathbf{F}_{2n}\tilde{\mathbf{y}})_k & \stackrel{(8.2.15)}{=} \sum_{j=1}^{2n-1} \tilde{y}_j e^{-\frac{2\pi}{2n}kj} \\
 & = \sum_{j=1}^{n-1} y_j e^{-\frac{\pi}{n}kj} - \sum_{j=n+1}^{2n-1} y_{2n-j} e^{-\frac{\pi}{n}kj} \\
 & = \sum_{j=1}^{n-1} y_j (e^{-\frac{\pi}{n}kj} - e^{\frac{\pi}{n}kj}) \\
 & = -2i (\mathbf{S}_n \mathbf{y})_k, \quad k = 1, \dots, n-1.
 \end{aligned}$$

```

Wrap-around implementation
function c = sinetrans(y)
n = length(y)+1;
yt = [0,y,0,-y(end:-1:1)];
ct = fft(yt);
c = -ct(2:n)/(2*i);
MATLAB-CODE sine transform
    
```

Remark 8.4.3 (Sine transform via DFT of half length).

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 (\rightarrow Sect. 8.2.3) of $(\tilde{y}_0, \dots, \tilde{y}_{n-1}) \in \mathbb{R}^n$:
$$c_k := \sum_{j=0}^{n-1} \tilde{y}_j e^{-\frac{2\pi i}{n}jk}$$

Hence

$$\begin{aligned} \operatorname{Re}\{c_k\} &= \sum_{j=0}^{n-1} \tilde{y}_j \cos\left(-\frac{2\pi i}{n}jk\right) = \sum_{j=1}^{n-1} (y_j + y_{n-j}) \sin\left(\frac{\pi j}{n}\right) \cos\left(\frac{2\pi i}{n}jk\right) \\ &= \sum_{j=0}^{n-1} 2y_j \sin\left(\frac{\pi j}{n}\right) \cos\left(\frac{2\pi i}{n}jk\right) = \sum_{j=0}^{n-1} y_j \left(\sin\left(\frac{2k+1}{n}\pi j\right) - \sin\left(\frac{2k-1}{n}\pi j\right) \right) \\ &= s_{2k+1} - s_{2k-1}. \\ \operatorname{Im}\{c_k\} &= \sum_{j=0}^{n-1} \tilde{y}_j \sin\left(-\frac{2\pi i}{n}jk\right) = -\sum_{j=1}^{n-1} \frac{1}{2}(y_j - y_{n-j}) \sin\left(\frac{2\pi i}{n}jk\right) = -\sum_{j=1}^{n-1} y_j \sin\left(\frac{2\pi i}{n}jk\right) \\ &= -s_{2k}. \end{aligned}$$

Step ③: extraction of s_k

$$s_{2k+1}, \quad k = 0, \dots, \frac{n}{2} - 1 \quad \blacktriangleright \quad \text{from recursion } s_{2k+1} - s_{2k-1} = \operatorname{Re}\{c_k\}, \quad s_1 = \sum_{j=1}^{n-1} y_j \sin(\pi j/n),$$

$$s_{2k}, \quad k = 1, \dots, \frac{n}{2} - 2 \quad \blacktriangleright \quad s_{2k} = -\operatorname{Im}\{c_k\}.$$

MATLAB-Implementation (via a `fft` of length $n/2$):

MATLAB-CODE Sine transform

```
function s = sinetrans(y)
n = length(y)+1;
sinevals = imag(exp(i*pi/n).^(1:n-1));
yt = [0 (sinevals.*(y+y(end:-1:1)) + 0.5*(y-y(end:-1:1)))]';
c = fftreal(yt);
s(1) = dot(sinevals,y);
for k=2:N-1
if (mod(k,2) == 0), s(k) = -imag(c(k/2+1));
else, s(k) = s(k-2) + real(c((k-1)/2+1)); end
end
```

5-points-stencil-operator on $\mathbb{R}^{n,n}$, $n \in \mathbb{N}$, in grid representation:

$$T : \mathbb{R}^{n,n} \mapsto \mathbb{R}^{n,n}, \quad \mathbf{X} \mapsto T(\mathbf{X})$$

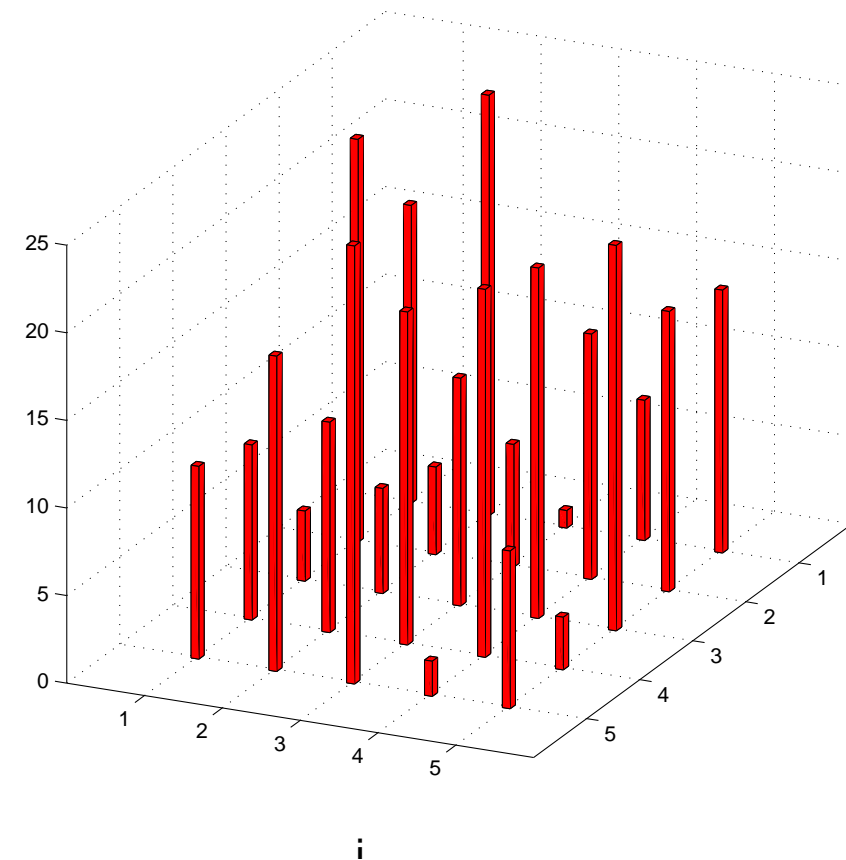
$$(T(\mathbf{X}))_{ij} := cx_{ij} + cyx_{i,j+1} + cyx_{i,j-1} + cx_{i+1,j} + cx_{i-1,j}$$

with $c, cy, cx \in \mathbb{R}$, convention: $x_{ij} := 0$ for $(i, j) \notin \{1, \dots, n\}^2$.

$$\mathbf{X} \in \mathbb{R}^{n,n}$$

$$\updownarrow$$

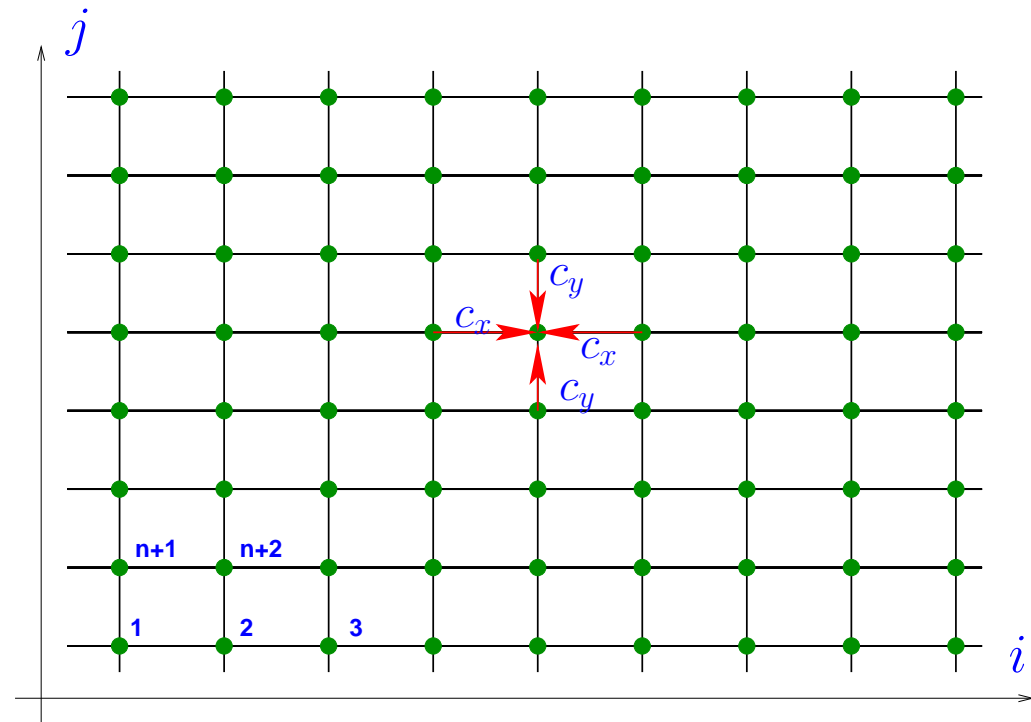
grid function $\in \{1, \dots, n\}^2 \mapsto \mathbb{R}$



Identification $\mathbb{R}^{n,n} \cong \mathbb{R}^{n^2}$, $x_{ij} \sim \tilde{x}_{(j-1)n+i}$ gives matrix representation $\mathbf{T} \in \mathbb{R}^{n^2,n^2}$ of T :

$$\mathbf{T} = \begin{pmatrix} \mathbf{C} & c_y \mathbf{I} & 0 & \dots & \dots & 0 \\ c_y \mathbf{I} & \mathbf{C} & c_y \mathbf{I} & & & \vdots \\ 0 & \ddots & \ddots & \ddots & & \\ \vdots & & & c_y \mathbf{I} & \mathbf{C} & c_y \mathbf{I} \\ 0 & \dots & \dots & 0 & c_y \mathbf{I} & \mathbf{C} \end{pmatrix} \in \mathbb{R}^{n^2,n^2},$$

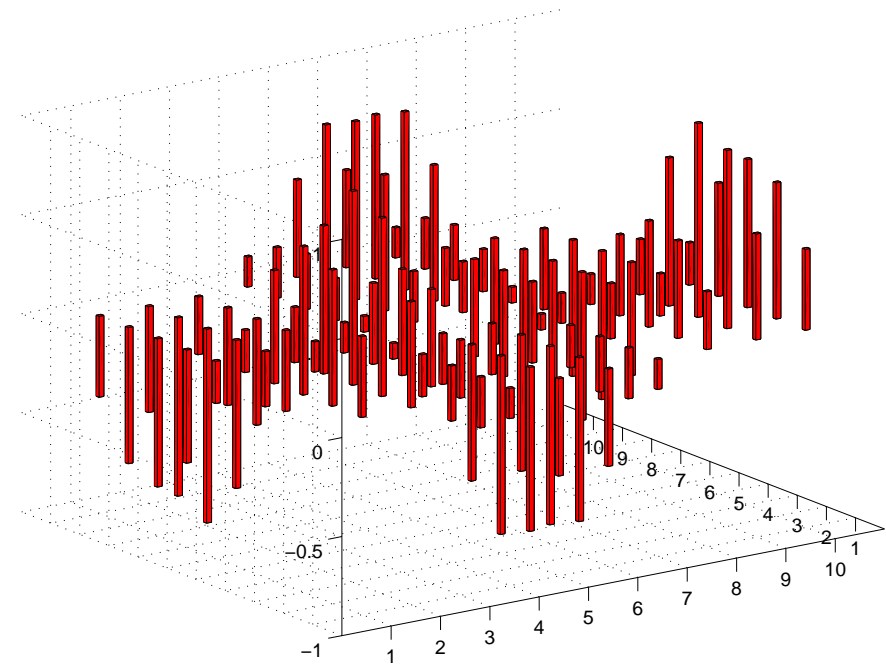
$$\mathbf{C} = \begin{pmatrix} c & c_x & 0 & \dots & \dots & 0 \\ c_x & c & c_x & & & \vdots \\ 0 & \ddots & \ddots & \ddots & & \\ \vdots & & & c_x & c & c_x \\ 0 & \dots & \dots & 0 & c_x & c \end{pmatrix} \in \mathbb{R}^{n,n}.$$



Sine basis of $\mathbb{R}^{n,n}$:

$$\mathbf{B}^{kl} = \left(\sin\left(\frac{\pi}{n+1}ki\right) \sin\left(\frac{\pi}{n+1}lj\right) \right)_{i,j=1}^n. \quad (8.4.4)$$

$n = 10$: grid function $\mathbf{B}^{2,3}$



$$\begin{aligned} (T(\mathbf{B}^{kl}))_{ij} &= c \sin\left(\frac{\pi}{n}ki\right) \sin\left(\frac{\pi}{n}lj\right) + c_y \sin\left(\frac{\pi}{n}ki\right) \left(\sin\left(\frac{\pi}{n+1}l(j-1)\right) + \sin\left(\frac{\pi}{n+1}l(j+1)\right) \right) + \\ &\quad c_x \sin\left(\frac{\pi}{n}lj\right) \left(\sin\left(\frac{\pi}{n+1}k(i-1)\right) + \sin\left(\frac{\pi}{n+1}k(i+1)\right) \right) \\ &= \sin\left(\frac{\pi}{n}ki\right) \sin\left(\frac{\pi}{n}lj\right) \left(c + 2c_y \cos\left(\frac{\pi}{n+1}l\right) + 2c_x \cos\left(\frac{\pi}{n+1}k\right) \right) \end{aligned}$$

Hence \mathbf{B}^{kl} is **eigenvector** of $T \leftrightarrow \mathbf{T}$ corresponding to eigenvalue $c + 2c_y \cos\left(\frac{\pi}{n+1}l\right) + 2c_x \cos\left(\frac{\pi}{n+1}k\right)$.

Algorithm for basis transform:

$$\mathbf{X} = \sum_{k=1}^n \sum_{l=1}^n y_{kl} \mathbf{B}^{kl} \Rightarrow x_{ij} = \sum_{k=1}^n \sin\left(\frac{\pi}{n+1}ki\right) \sum_{l=1}^n y_{kl} \sin\left(\frac{\pi}{n+1}lj\right).$$

MATLAB-CODE two dimensional sine tr.

```
function C = sinft2d(Y)
[m,n] = size(Y);
C = fft([zeros(1,n); Y;...
        zeros(1,n);...
        -Y(end:-1:1,:)]);
C = i*C(2:m+1,:)/2;
C = fft([zeros(1,m); C;...
        zeros(1,m);...
        -C(end:-1:1,:)]);
C = i*C(2:n+1,:)/2;
```

Hence nested sine transforms (\rightarrow Sect. 8.2.4)
for rows/columns of $\mathbf{Y} = (y_{kl})_{k,l=1}^n$.

Here: implementation of sine transform (8.4.2)
with “wrapping”-technique.

MATLAB-CODE FFT-based solution of local

```
function X = fftsolve(B,c,cx,cy)
[m,n] = size(B);
[I,J] = meshgrid(1:m,1:n);
X = 4*sinft2d(sinft2d(B)...
    ./((c+2*cx*cos(pi/(n+1)*I)+...
        2*cy*cos(pi/(m+1)*J)))...
    /((m+1)*(n+1));
```

translation invariant linear operators

Diagonalization of \mathbf{T}
via 2D sine transform



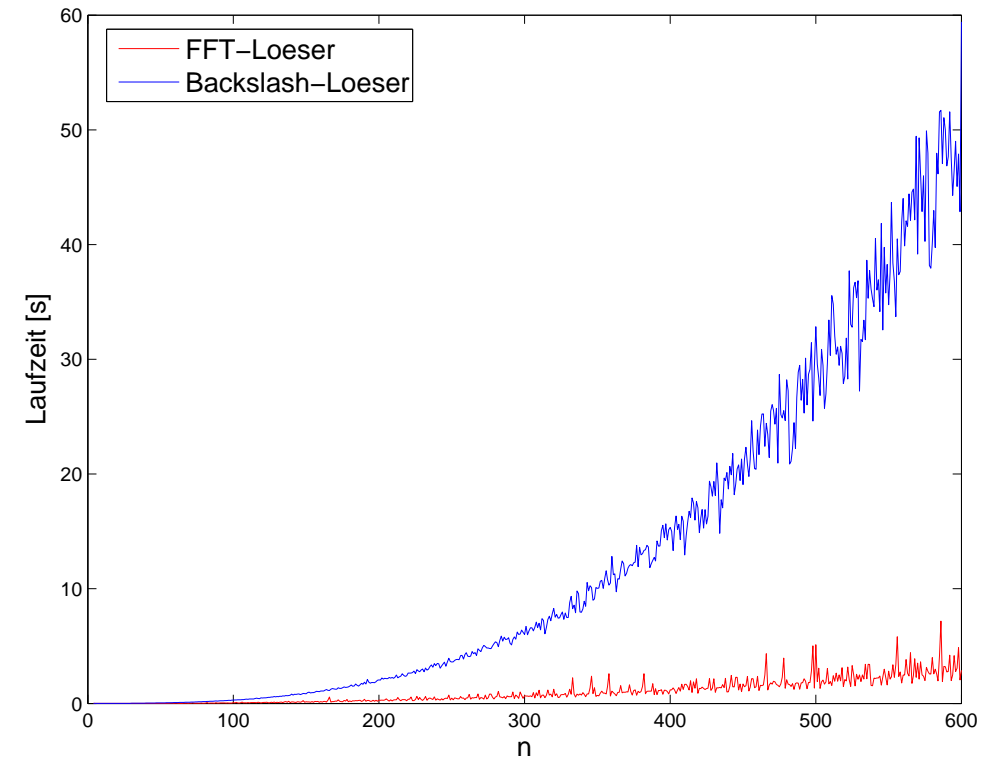
efficient algorithm

for solving linear system of equations $\mathbf{T}(\mathbf{X}) = \mathbf{B}$

computational cost $O(n^2 \log n)$.

Example 8.4.5 (Efficiency of FFT-based LSE-solver).


```
A = gallery('poisson',n);  
B = magic(n);  
b = reshape(B,n*n,1);  
tic;  
C = fftsolve(B,4,-1,-1);  
t1 = toc;  
tic; x = A\b; t2 = toc;
```



8.4.2 Cosine transform

Another trigonometric basis transform in \mathbb{R}^n , $n \in \mathbb{N}$:

standard basis of \mathbb{R}^n

“cosine basis”

$$\left\{ \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \\ 1 \end{pmatrix} \right\} \leftarrow \left\{ \begin{pmatrix} 2^{-1/2} \\ \cos(\frac{\pi}{2n}) \\ \cos(\frac{2\pi}{2n}) \\ \vdots \\ \cos(\frac{(n-1)\pi}{2n}) \end{pmatrix}, \begin{pmatrix} 2^{-1/2} \\ \cos(\frac{3\pi}{2n}) \\ \cos(\frac{6\pi}{2n}) \\ \vdots \\ \cos(\frac{3(n-1)\pi}{2n}) \end{pmatrix}, \dots, \begin{pmatrix} 2^{-1/2} \\ \cos(\frac{(2n-1)\pi}{2n}) \\ \cos(\frac{2(2n-1)\pi}{2n}) \\ \vdots \\ \cos(\frac{(n-1)(2n-1)\pi}{2n}) \end{pmatrix} \right\}.$$

R. Hiptmair
rev 38355,
November
17, 2011

Basis transform matrix (cosine basis \rightarrow standard basis):

$$\mathbf{C}_n = (c_{ij}) \in \mathbb{R}^{n,n} \quad \text{with} \quad c_{ij} = \begin{cases} 2^{-1/2} & , \text{ if } i = 1 , \\ \cos((i - 1)\frac{2j-1}{2n}\pi) & , \text{ if } i > 1 . \end{cases}$$

Lemma 8.4.6 (Properties of cosine matrix).

$\sqrt{2/n} \mathbf{C}_n$ is real and orthogonal (\rightarrow Def. 2.8.5).

cosine transform:

$$c_k = \sum_{j=0}^{n-1} y_j \cos\left(k \frac{2j+1}{2n} \pi\right), \quad k = 1, \dots, n-1, \quad (8.4.7)$$

$$c_0 = \frac{1}{\sqrt{2}} \sum_{j=0}^{n-1} y_j.$$

MATLAB-implementation of \mathbf{C}_y ("wrapping"-technique):

MATLAB-CODE cosine transform

```
function c = costrans(y)
n = length(y);
z = fft([y,y(end:-1:1)]);
c = real([z(1)/(2*sqrt(2)), ...
         0.5*(exp(-i*pi/(2*n)).^(1:n-1)).*z(2:n)]);
```

MATLAB-implementation of $\mathbf{C}_n^{-1}\mathbf{y}$ ("Wrapping"-technique):

```
function y=icostrans(c)
n = length(c);
y = [sqrt(2)*c(1), (exp(i*pi/(2*n)).^(1:n-1)).*c(2:end)];
y = ifft([y,0,conj(y(end:-1:2))]);
y = 2*y(1:n);
```

8.5 Toeplitz matrix techniques

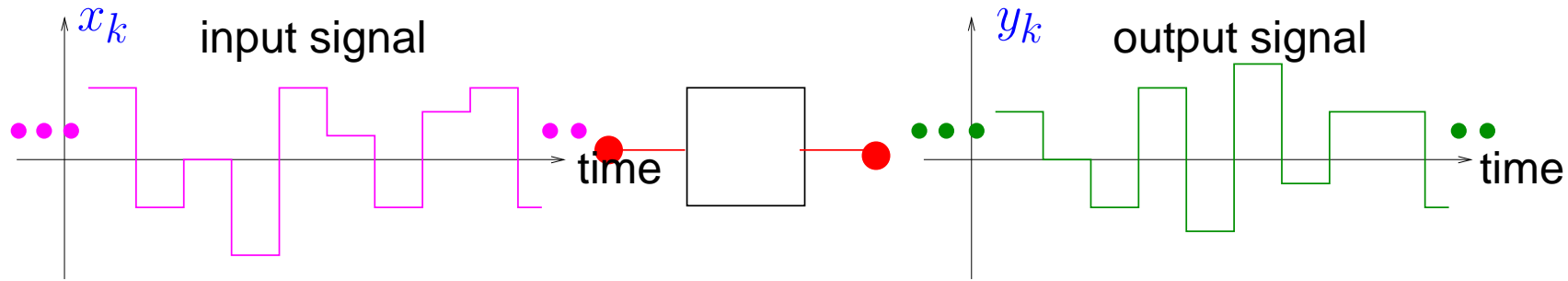
Example 8.5.1 (Parameter identification for linear time-invariant filters).

- $(x_k)_{k \in \mathbb{Z}}$ m -periodic discrete signal = *known* input
- $(y_k)_{k \in \mathbb{Z}}$ m -periodic *measured*^(*) output signal of a **linear time-invariant filter**, see Ex. 8.1.1.
- ^(*) \rightarrow measurement errors !
- Known: impulse response of filter has maximal duration $n\Delta t$, $n \in \mathbb{N}$, $n \leq m$

cf. (8.1.3)



$$\exists \mathbf{h} = (h_0, \dots, h_{n-1})^\top \in \mathbb{R}^n, \quad n \leq m : \quad y_k = \sum_{j=0}^{n-1} h_j x_{k-j} . \quad (8.5.2)$$



Parameter identification problem: seek $\mathbf{h} = (h_0, \dots, h_{n-1})^\top \in \mathbb{R}^n$ with

$$\|\mathbf{A}\mathbf{h} - \mathbf{y}\|_2 = \left\| \begin{pmatrix} x_0 & x_{-1} & \cdots & \cdots & x_{1-n} \\ x_1 & x_0 & x_{-1} & & \vdots \\ \vdots & x_1 & x_0 & \ddots & \\ & & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & x_{-1} \\ x_{n-1} & & & x_1 & x_0 \\ x_n & x_{n-1} & & & x_1 \\ \vdots & & & & \vdots \\ x_{m-1} & \cdots & & \cdots & x_{m-n} \end{pmatrix} \begin{pmatrix} h_0 \\ \vdots \\ h_{n-1} \end{pmatrix} - \begin{pmatrix} y_0 \\ \vdots \\ y_{m-1} \end{pmatrix} \right\|_2 \rightarrow \min .$$

➤ **Linear least squares problem**, → Ch. 7 with Toeplitz matrix \mathbf{A} : $(\mathbf{A})_{ij} = x_{i-j}$.



Definition 8.5.5 (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 $t_{ij} = u_{j-i}$, $1 \leq i \leq m$, $1 \leq j \leq n$.

$$\mathbf{T} = \begin{pmatrix} u_0 & u_1 & \cdots & \cdots & u_{n-1} \\ u_{-1} & u_0 & u_1 & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & u_1 \\ u_{1-m} & \cdots & & \cdots & u_{-1} & u_0 \end{pmatrix}$$


8.5.1 Toeplitz matrix arithmetic

$$\mathbf{C} = \begin{pmatrix} \mathbf{T} & \mathbf{S} \\ \mathbf{S} & \mathbf{T} \end{pmatrix} = \left(\begin{array}{cccc|cccc}
 u_0 & u_1 & \cdots & \cdots & u_{n-1} & 0 & u_{1-n} & \cdots & \cdots & u_{-1} \\
 u_{-1} & u_0 & u_1 & & \vdots & u_{n-1} & 0 & \ddots & & \vdots \\
 \vdots & \ddots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots & & \vdots \\
 \vdots & & \ddots & \ddots & \vdots & \vdots & & & \ddots & \vdots \\
 \vdots & & & \ddots & u_1 & \vdots & & \ddots & \ddots & u_{1-n} \\
 \hline
 u_{1-n} & \cdots & & \cdots & u_{-1} & u_1 & & & u_{n-1} & 0 \\
 0 & u_{1-n} & \cdots & \cdots & u_{-1} & u_0 & u_1 & \cdots & \cdots & u_{n-1} \\
 u_{n-1} & 0 & \ddots & & \vdots & u_{-1} & u_0 & u_1 & & \vdots \\
 \vdots & \ddots & \ddots & & \vdots & \vdots & \ddots & \ddots & & \vdots \\
 & & & \ddots & \vdots & \vdots & & \ddots & \ddots & \vdots \\
 \vdots & & & \cdots & \ddots & u_{1-n} & & \ddots & \ddots & u_1 \\
 u_1 & & & & u_{n-1} & 0 & u_{1-n} & \cdots & \cdots & u_{-1} & u_0
 \end{array} \right)$$

In general:

$$\mathbf{T} = \text{toeplitz}(u(0:-1:1-m), u(0:n-1));$$

$$\mathbf{S} = \text{toeplitz}([0, u(n-1:-1:n-m+1)], [0, u(1-m:1:-1)]);$$

zero padding  $\mathbf{C} \begin{pmatrix} \mathbf{x} \\ 0 \end{pmatrix} = \begin{pmatrix} \mathbf{T}\mathbf{x} \\ \mathbf{S}\mathbf{x} \end{pmatrix}$

▶ Computational effort $O(n \log n)$ for computing $\mathbf{T}\mathbf{x}$ (FFT based, Sect. 8.3)

8.5.2 The Levinson algorithm

Levinson algorithm

(recursive, u_{n+1} not used!)

Linear recursion:

Computational cost $\sim (n-k)$ on level k , $k = 0, \dots, n-1$

➤ Asymptotic complexity $O(n^2)$



Code 8.5.9: Levinson algorithm

```
1 function [x,y] = levinson(u,b)
2 k = length(u)-1;
3 if (k == 0), x=b(1); y = u(1);
   return; end
4 [xk,yk] =
   levinson(u(1:k),b(1:k));
5 sigma = 1-dot(u(1:k),yk);
6 t =
   (b(k+1)-dot(u(k:-1:1),xk))/sigma;
7 x = [ xk-t*yk(k:-1:1);t];
8 s =
   (u(k+1)-dot(u(k:-1:1),yk))/sigma;
9 y = [yk-s*yk(k:-1:1); s];
```


9

Approximation of Functions in 1D

9.1 Error estimates for polynomial interpolation

We consider Lagrangian polynomial interpolation on node set

$$\mathcal{T} := \{t_0, \dots, t_n\} \subset I, \quad I \subset \mathbb{R}, \quad \text{interval of length } |I|.$$

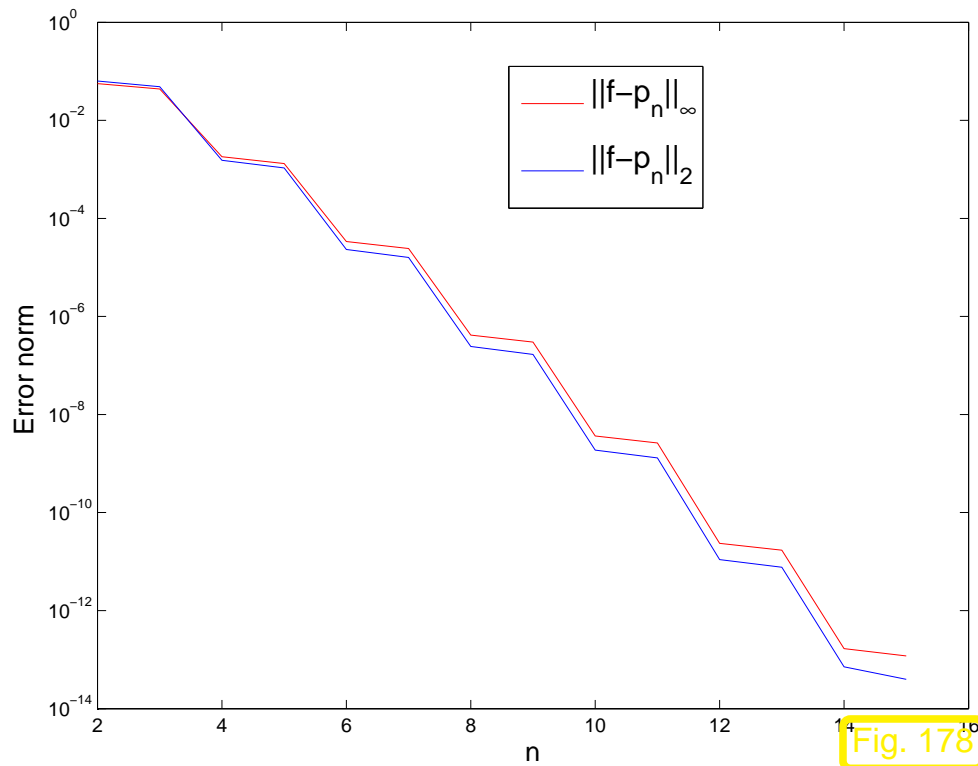
Goal: estimate of the **interpolation error** norm $\|f - I_{\mathcal{T}}f\|$ (for some norm on $C(I)$).

Focus: **asymptotic** behavior of interpolation error for $n \rightarrow \infty$

Example 9.1.2 (Asymptotic behavior of polynomial interpolation error).

Interpolation of $f(t) = \sin t$ on equispaced nodes in $I = [0, \pi]$: $\mathcal{T} = \{j\pi/n\}_{j=0}^n$.

Interpolating polynomial $p := I_{\mathcal{T}}f \in \mathcal{P}_n$.



$$\exists C \neq C(n) > 0: \quad \|f - I_{\mathcal{T}}f\| \leq C T(n) \quad \text{for } n \rightarrow \infty. \quad (9.1.3)$$

Classification (best bound for $T(n)$):

$$\begin{aligned} \exists p > 0: \quad T(n) \leq n^{-p} & : \text{ algebraic convergence, with rate } p > 0, \quad \forall n \in \mathbb{N}. \\ \exists 0 < q < 1: \quad T(n) \leq q^n & : \text{ exponential convergence,} \end{aligned}$$

Example 9.1.4 (Runge's example). \rightarrow Ex. 3.5.1

Polynomial interpolation of $f(t) = \frac{1}{1+t^2}$ with equispaced nodes:

$$\mathcal{T} := \left\{ t_j := -5 + \frac{10}{n} j \right\}_{j=0}^n, \quad y_j = \frac{1}{1+t_j^2}, j = 0, \dots, n.$$

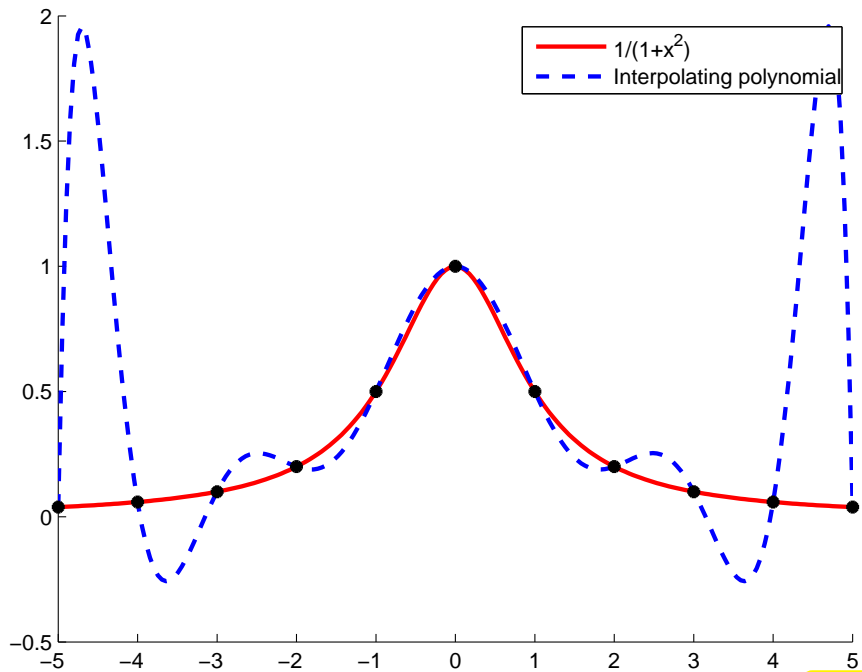


Fig. 179

Interpolating polynomial, $n = 10$

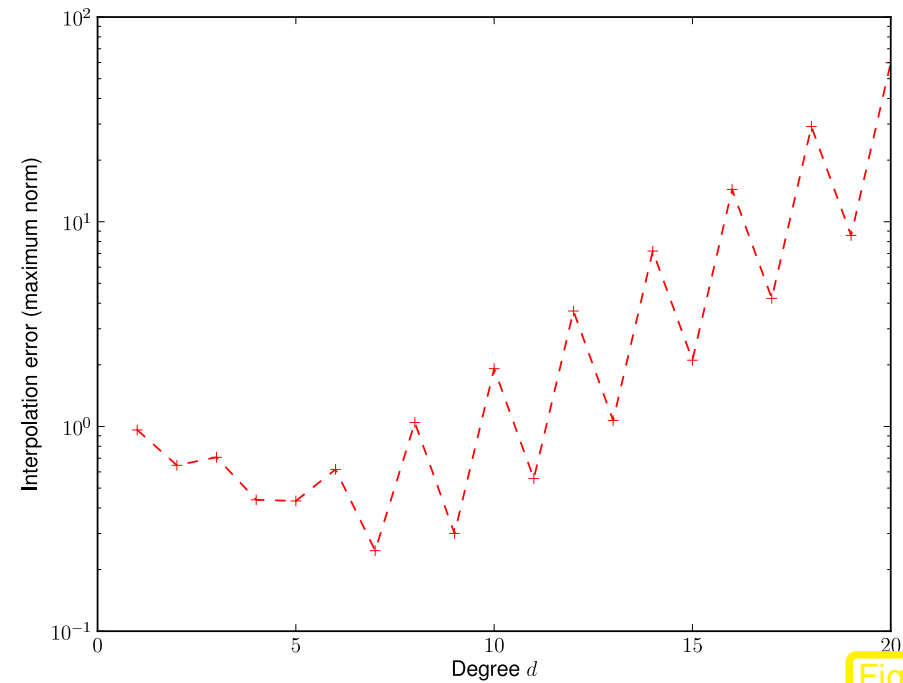


Fig. 180

Approximate $\|f - \mathcal{I}_{\mathcal{T}}f\|_{\infty}$ on $[-5, 5]$



Theorem 9.1.6 (Representation of interpolation error). [13, Thm. 8.22], [35, Thm. 37.4]

$f \in C^{n+1}(I)$: $\forall t \in I$: $\exists \tau_t \in]\min\{t, t_0, \dots, t_n\}, \max\{t, t_0, \dots, t_n\}[$:

$$f(t) - \mathcal{I}_{\mathcal{T}}(f)(t) = \frac{f^{(n+1)}(\tau_t)}{(n+1)!} \cdot \prod_{j=0}^n (t - t_j). \quad (9.1.7)$$

$$\text{Thm. 9.1.7} \quad \Rightarrow \quad \|f - I_{\mathcal{T}}f\|_{L^\infty(I)} \leq \frac{\|f^{(n+1)}\|_{L^\infty(I)}}{(n+1)!} \max_{t \in I} |(t - t_0) \cdots (t - t_n)|. \quad (9.1.11)$$

Interpolation error estimate requires smoothness!

9.2 Chebychev Interpolation

9.2.1 Motivation and definition

Setting: Mesh of nodes: $\mathcal{T} := \{t_0 < t_1 < \cdots < t_{n-1} < t_n\}$, $n \in \mathbb{N}$,
function $f : I \rightarrow \mathbb{R}$ continuous; without loss of generality $I = [-1, 1]$.

Thm. 9.1.7:

$$\|f - p\|_{L^\infty(I)} \leq \frac{1}{(n+1)!} \|f^{(n+1)}\|_{L^\infty(I)} \|w\|_{L^\infty(I)},$$

$$w(t) := (t - t_0) \cdots (t - t_n).$$

Idea: choose nodes t_0, \dots, t_n such that $\|w\|_{L^\infty(I)}$ is minimal!

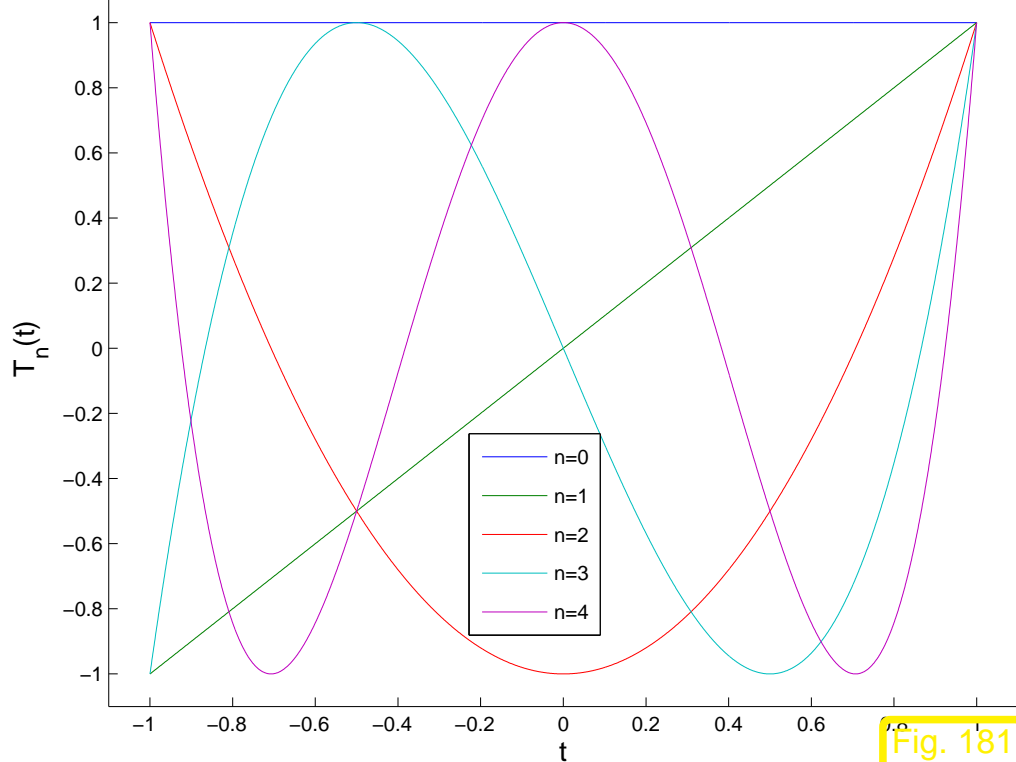
Equivalent to finding $q \in \mathcal{P}_{n+1}$, with

- leading coefficient = 1,
- such that $\|q\|_{L^\infty(I)}$ is minimal.

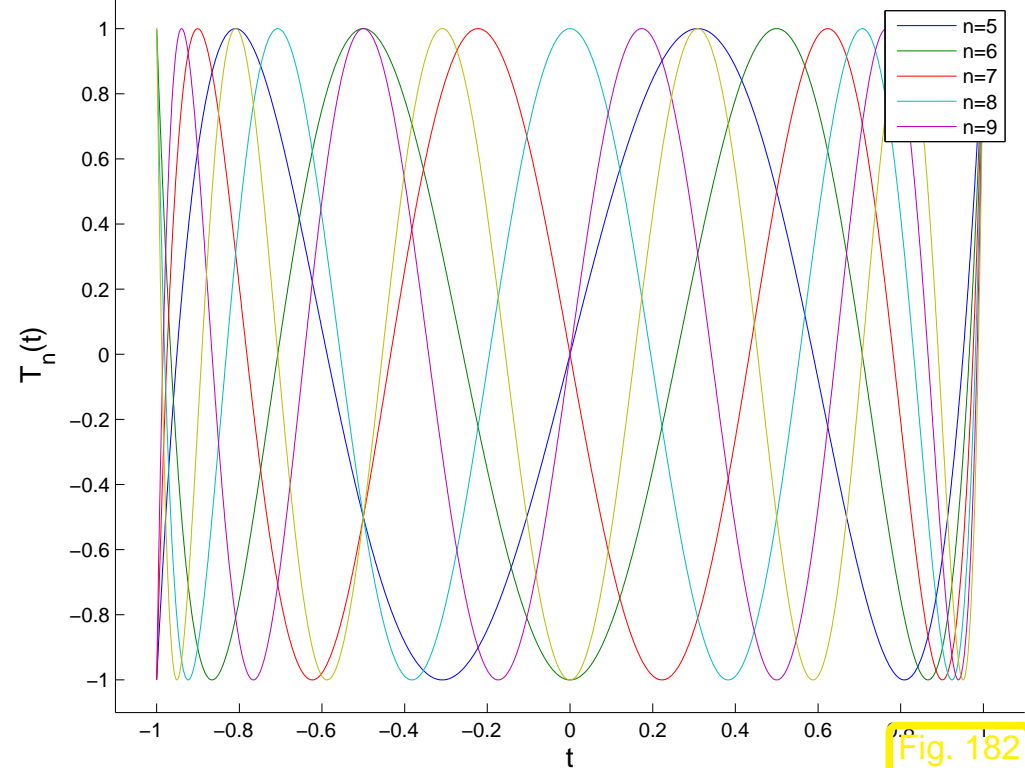
► Choice of $t_0, \dots, t_n =$ zeros of q (caution: t_j must belong to I).

Definition 9.2.1 (Chebychev polynomial). $\rightarrow [35, \text{Ch. } 32]$

The n^{th} *Chebychev polynomial* is $T_n(t) := \cos(n \arccos t)$, $-1 \leq t \leq 1$.



Chebyshev polynomials T_0, \dots, T_4



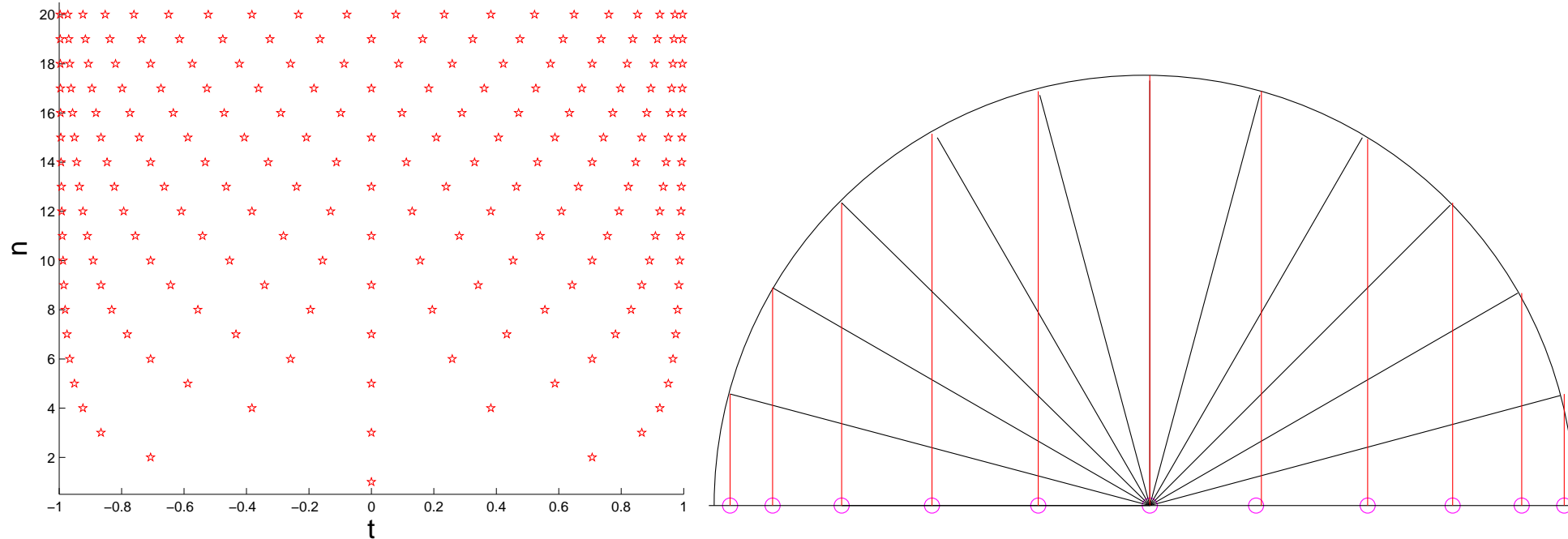
Chebyshev polynomials T_5, \dots, T_9

Theorem 9.2.7 (Minimax property of the Chebyshev polynomials). [17, Section 7.1.4.], [35, Thm. 32.2]

$$\|T_n\|_{L^\infty([-1,1])} = \inf\{\|p\|_{L^\infty([-1,1])} : p \in \mathcal{P}_n, p(t) = 2^{n-1}t^n + \dots\}, \quad \forall n \in \mathbb{N}.$$

Zeros of T_n : $t_k = \cos\left(\frac{2k-1}{2n}\pi\right), \quad k = 1, \dots, n.$ (9.2.8)

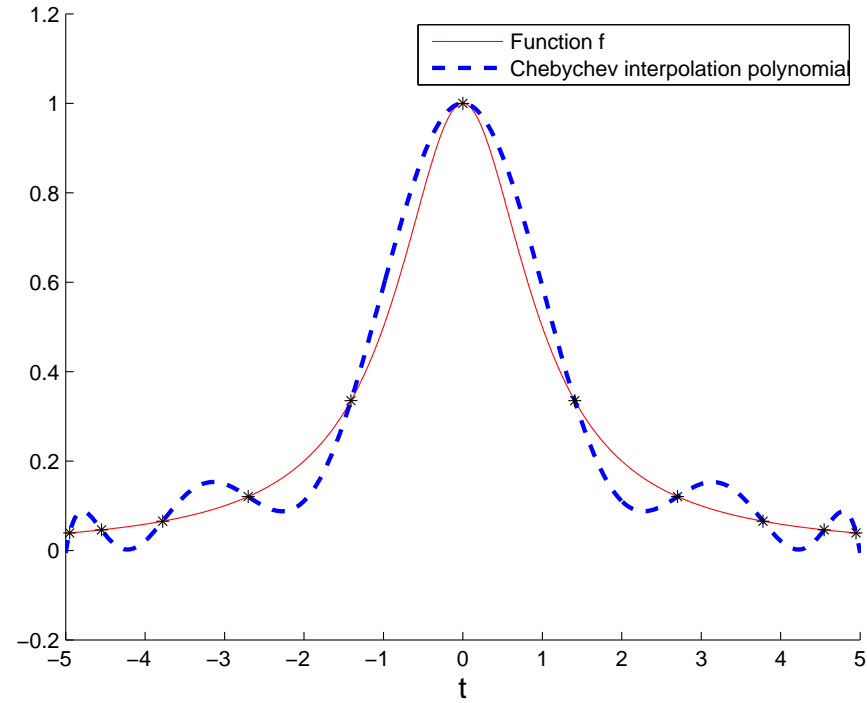
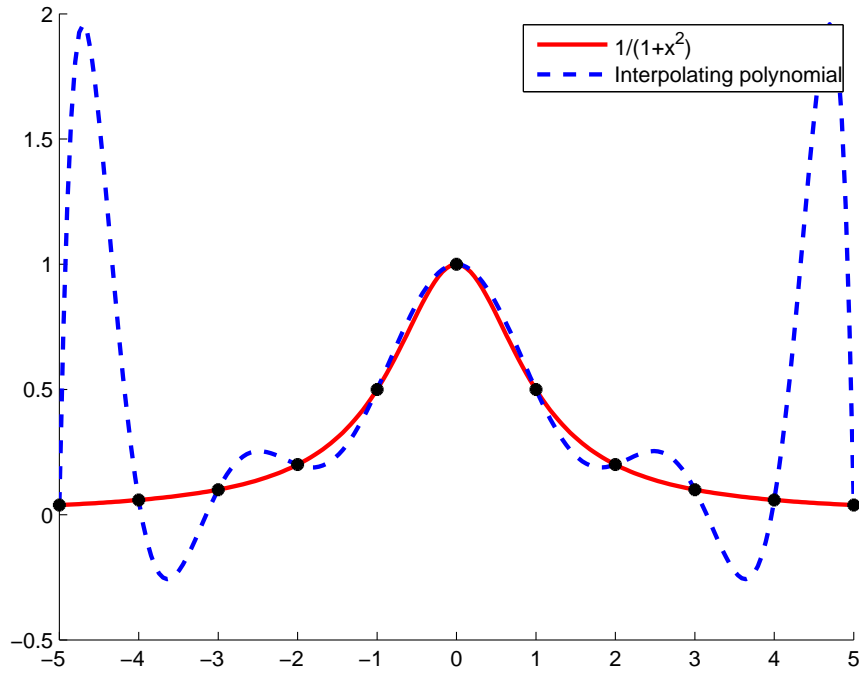
Location of **Chebyshev nodes** t_k from (9.2.8):



9.2.2 Chebyshev interpolation error estimates

Example 9.2.12 (Polynomial interpolation: Chebyshev nodes versus equidistant nodes).

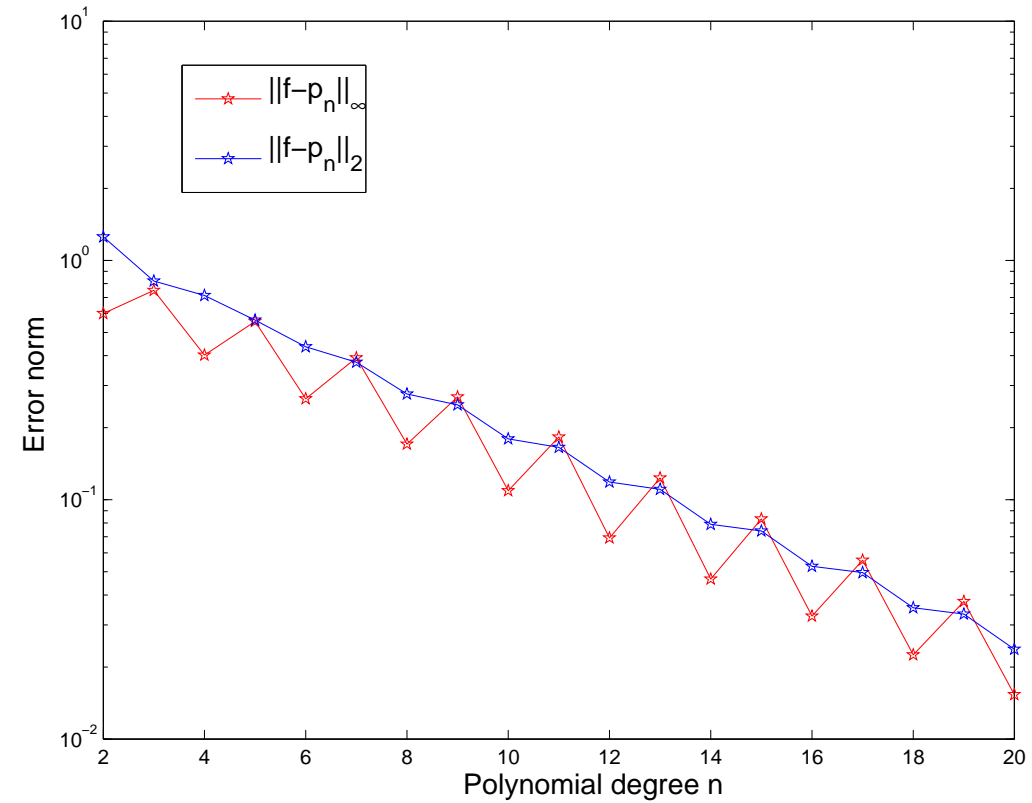
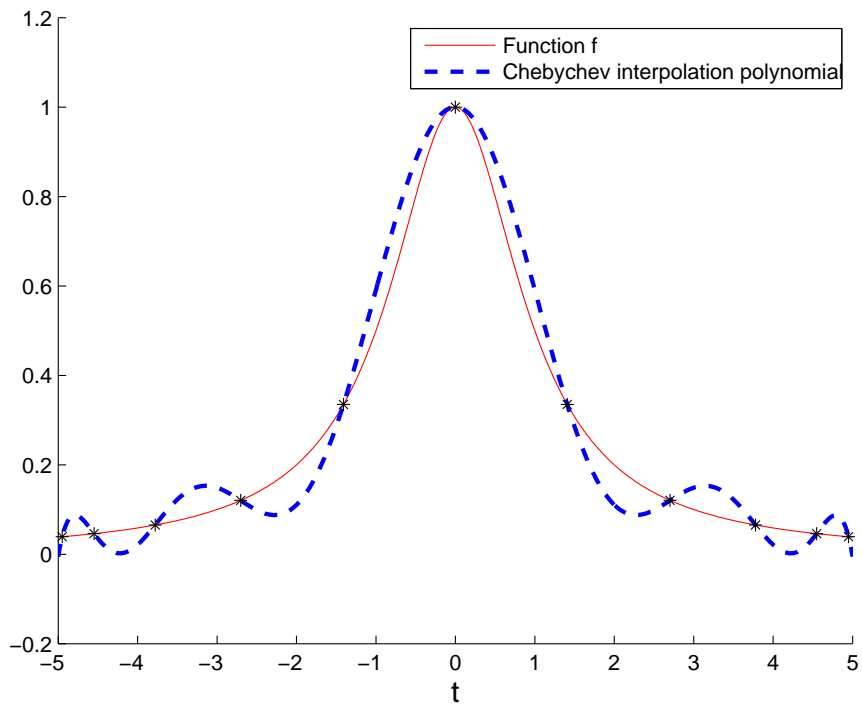
Runge's function $f(t) = \frac{1}{1+t^2}$, see Ex. 9.1.5, polynomial interpolation based on uniformly spaced nodes and Chebychev nodes:



Example 9.2.14 (Chebyshev interpolation error).

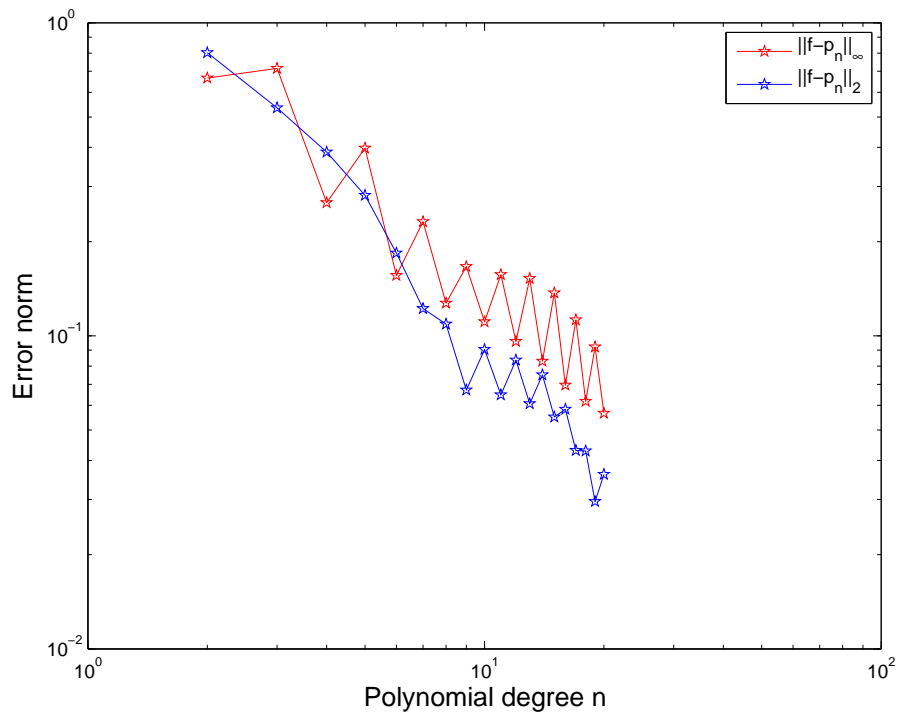
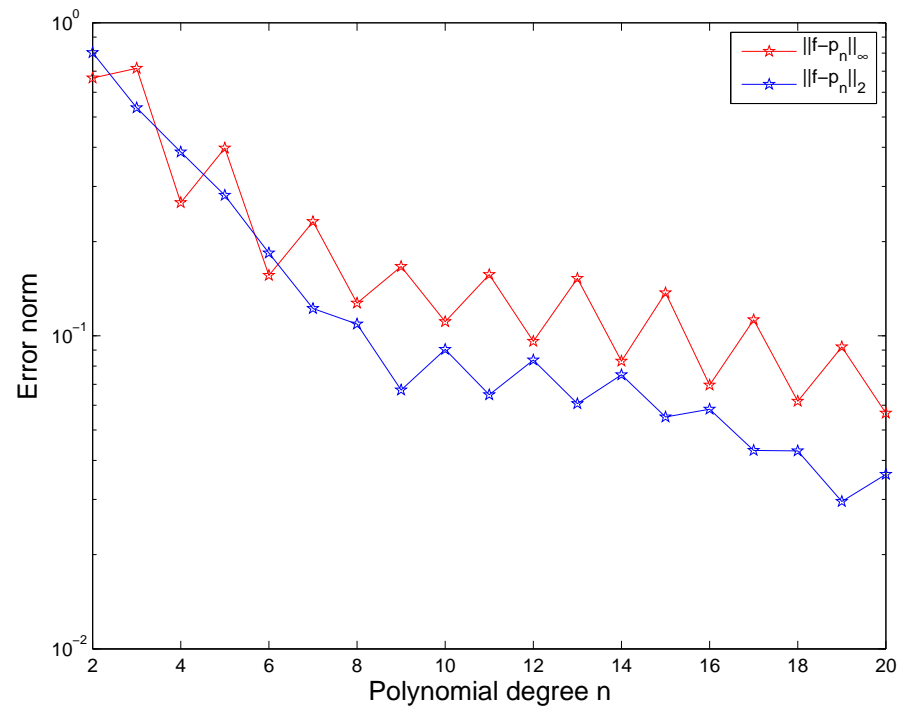
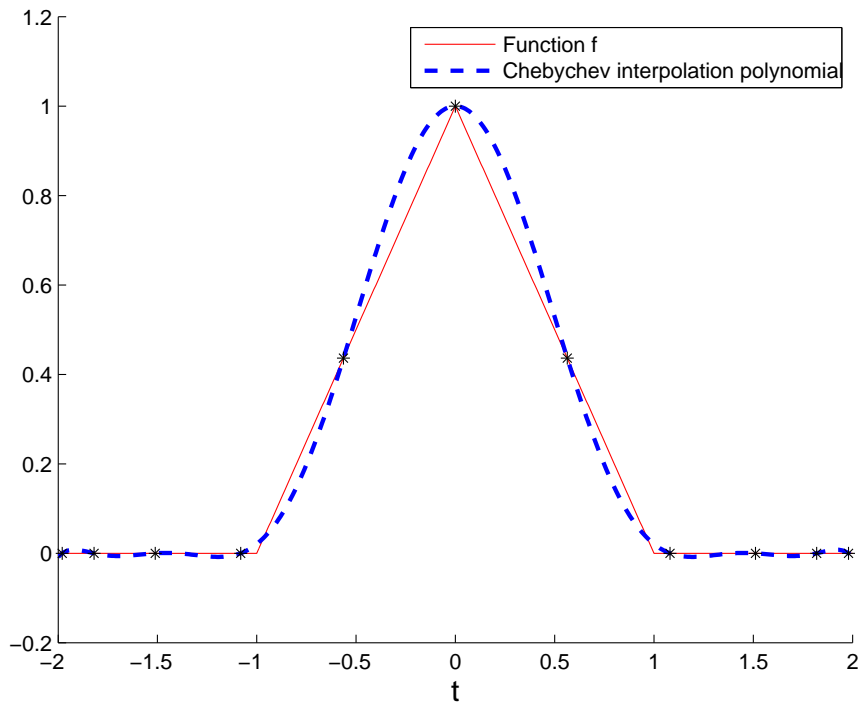
① $f(t) = (1 + t^2)^{-1}$, $I = [-5, 5]$ (see Ex. 9.1.5)

Interpolation with $n = 10$ Chebyshev nodes (plot on the left).

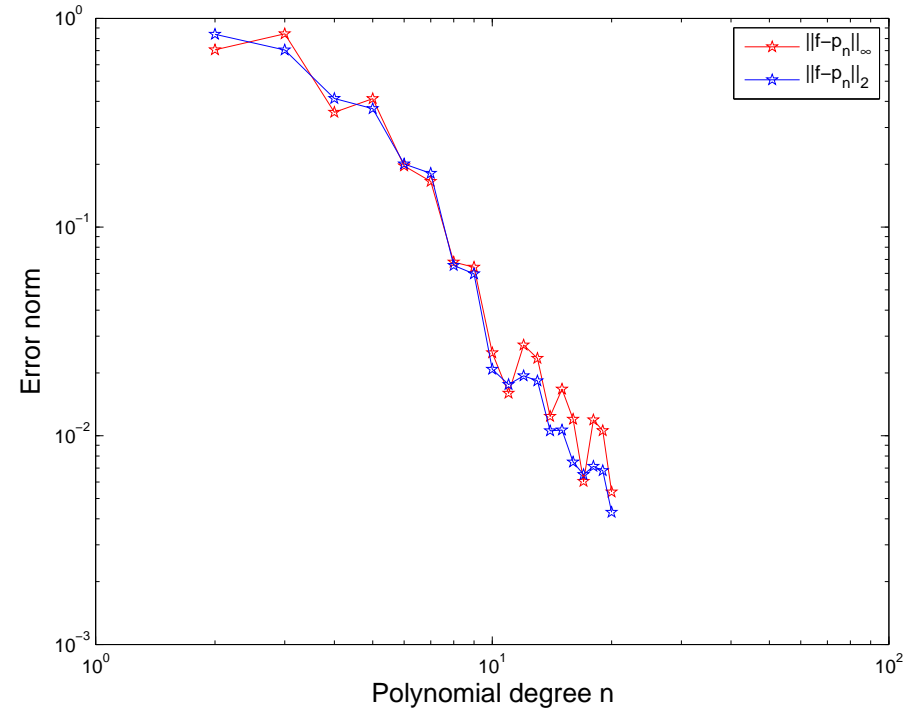
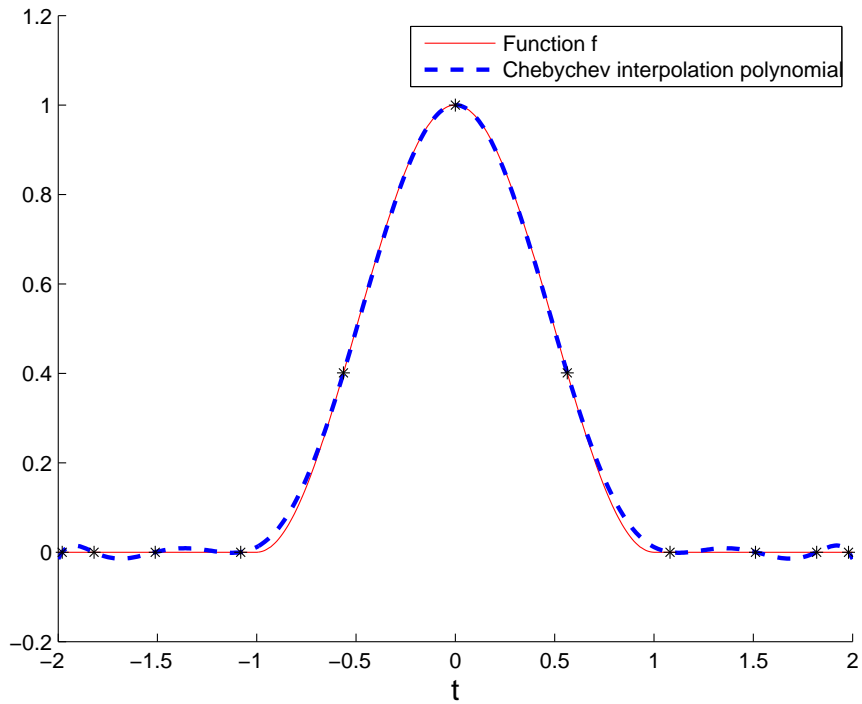


② $f(t) = \max\{1 - |t|, 0\}$, $I = [-2, 2]$, $n = 10$ nodes (plot on the left).

Now $f \in C^0(I)$ but $f \notin C^1(I)$.



③ $f(t) = \begin{cases} \frac{1}{2}(1 + \cos \pi t) & |t| < 1 \\ 0 & 1 \leq |t| \leq 2 \end{cases}$ $I = [-2, 2]$, $n = 10$ (plot on the left).



9.2.3 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 interpolating polynomial $p \in \mathcal{P}_n$ in Chebychev nodes (9.2.12) on $[-1, 1]$

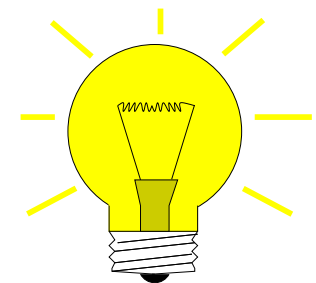
```

1 class ChebInterp {
2   private :
3   // various internal data describing Chebychev interpolating polynomial p
4   public :
5   // Constructor taking function f and degree n as arguments
6   PolyInterp(const Function &f, unsigned int n);
7   // Evaluation operator:  $y_j = p(x_j)$ ,  $j = 1, \dots, m$  ( $m$  "large")
8   double eval(const vector<double> &x, vector <double> &y) const;
9 };
```

R. Hiptmair
rev 38355,
October 27,
2011

Trick: expand p into Chebychev polynomials

$$p = \sum_{j=0}^n \alpha_j T_j, \quad \alpha_j \in \mathbb{R}. \quad (9.2.20)$$



Remark 9.2.21 (Fast evaluation of Chebychev expansion). → [35, Alg. 32.1]


Code 9.2.23: Clenshaw algorithm for evaluation of Chebychev expansion (9.2.22)

```

1 function y = clenshaw(a,x)
2 % Clenshaw algorithm for evaluating  $p = \sum_{j=1}^{n+1} a_j T_{j-1}$  at points passed in vector x
3 n = length(a)-1; % degree of polynomial
4 d = repmat(reshape(a,n+1,1),1,length(x));
5 for j=n:-1:2
6     d(j,:) = d(j,:) + 2*x.*d(j+1,:); % see (9.2.24)
7     d(j-1,:) = d(j-1,:) - d(j+1,:);
8 end
9 y = d(1,:) + x.*d(2,:);

```

R. Hiptmair


 rev 38355,
October 27,
2011

Code 9.2.29: Efficient computation of Chebychev expansion coefficient of Chebychev interpolant

```

1 function a = chebexp(y)
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  $t_k$ ,  $k = 0, \dots, n$ . These values are

```

```

5 % passed in the row vector y.
6 n = length(y)-1; % degree of polynomial
7 % create vector z by wrapping and componentwise scaling
8 z = exp(-pi*i*n/(n+1)*(0:2*n+1)).*[y,y(end:-1:1)]; % r.h.s. vector
9 c = ifft(z); % Solve linear system (9.2.31) with effort O(n log n)
0 b = real(exp(0.5*pi*i/(n+1)*(-n:n+1)).*c); % recover beta_j, see (9.2.31)
1 a = [b(n+1),2*b(n+2:2*n+1)]; % recover alpha_j, see (9.2.28)

```

9.3 Trigonometric interpolation [13, Sect. 8.5]

Trigonometric interpolation [13, pp. 304]

Given $t_0 < t_1 < \dots < t_{2n}$, $t_k \in [0, 1]$, and $y_k \in \mathbb{R}$, $k = 0, \dots, 2n$ find

$$q \in \mathcal{P}_{2n}^T := \text{Span} \{t \mapsto \cos(2\pi jt), t \mapsto \sin(2\pi jt)\}_{j=0}^n, \quad (9.3.1)$$

$$\text{with } q(t_k) = y_k \text{ for all } k = 0, \dots, 2n. \quad (9.3.2)$$

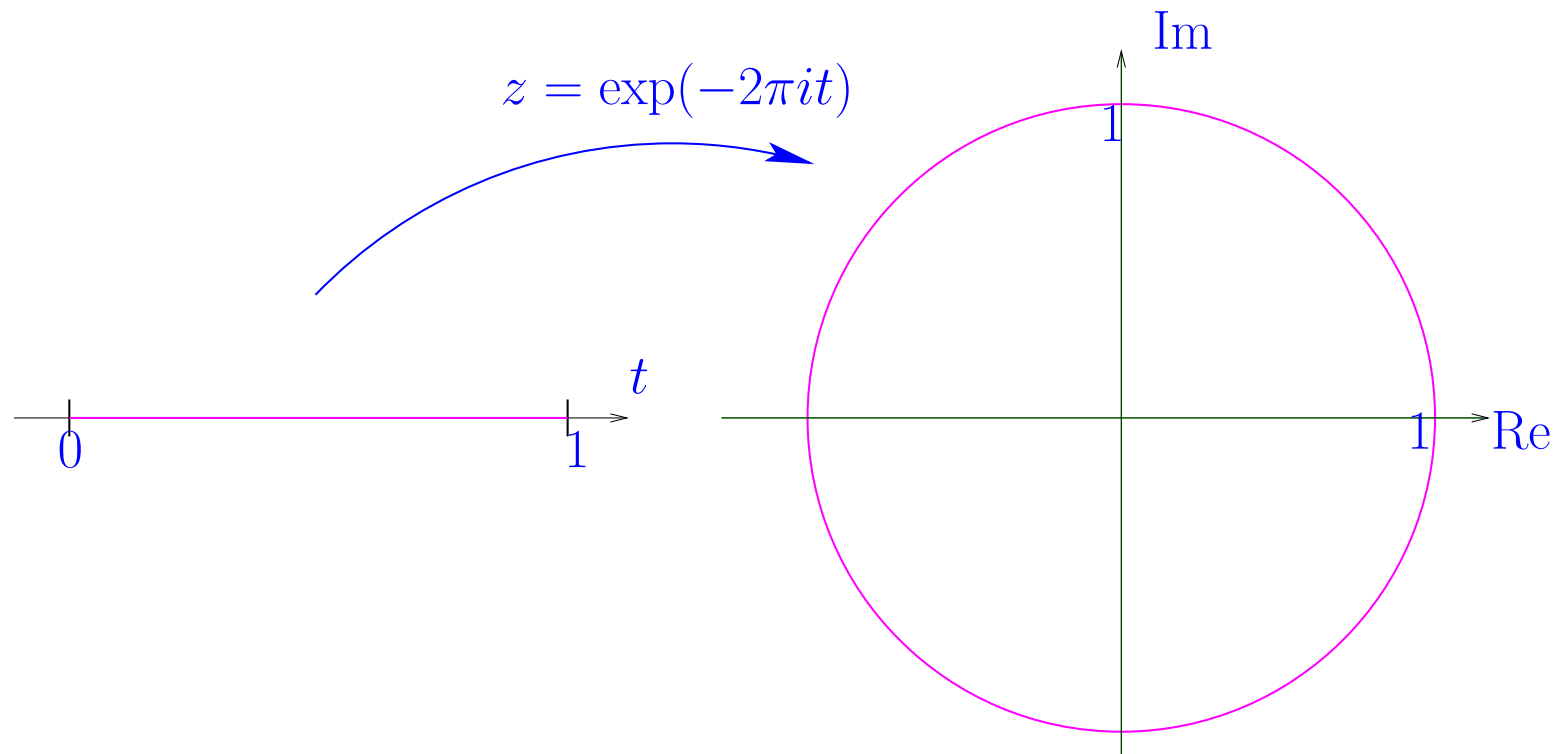
Terminology: $\mathcal{P}_{2n}^T \hat{=}$ space of **trigonometric polynomials** of degree $2n$.

$$q \in \mathcal{P}_{2n+1}^T \Rightarrow q(t) = e^{-2\pi i n t} \cdot p(e^{2\pi i t}) \quad \text{with} \quad p(z) = \sum_{j=0}^{2n} \gamma_j z^j \in \mathcal{P}_{2n}, \quad (9.3.6)$$

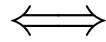
a polynomial !

and γ_j from (9.3.5).

$t \rightarrow \exp(2\pi i n t)q(t)$ is a polynomial $p \in \mathcal{P}_{2n}$ restricted to the unit circle \mathbb{S}^1 in \mathbb{C} .



Trigonometric interpolation
through data points (t_k, y_k)



Polynomial interpolation
through data points $(e^{2\pi i t_k}, y_k)$

Code 9.3.7: Evaluation of trigonometric interpolation polynomial in many points

```

1 function q = trigpolyval(t,y,x)
2 % Evaluation of trigonometric interpolant at numerous points
3 % t: row vector of nodes  $t_0, \dots, t_n \in [0, 1[$ 
4 % y: row vector of data  $y_0, \dots, y_n$ 
5 % x: row vector of evaluation points  $x_1, \dots, x_N$ 
6 N = length(y); if (mod(N,2)~=1), error('#pts odd required'); end
7 n = (N-1)/2;
8 tc = exp(2*pi*i*t); % Interpolation nodes on unit circle
9 z = exp(2*pi*i*n*t).*y; % Rescaled values, according to  $q(t) = e^{-2\pi i n t} \cdot p(e^{2\pi i t})$ 
10 % Evaluation of interpolating polynomial on unit circle, see Code 3.4.1
11 p = intpolyval(tc,z,exp(2*pi*i*x));
12 q = exp(-2*pi*i*n*x).*p; % Undo the scaling, see (9.3.6)

```

► $(2n + 1) \times (2n + 1)$ linear system of equations:

$$\sum_{j=0}^{2n} \gamma_j \exp\left(2\pi i \frac{jk}{2n+1}\right) = z_k := \exp\left(2\pi i \frac{nk}{2n+1}\right) y_k, \quad k = 0, \dots, 2n.$$

$$\Updownarrow$$

$$\bar{\mathbf{F}}_{2n+1} \mathbf{c} = \mathbf{z}, \quad \mathbf{c} = (\gamma_0, \dots, \gamma_{2n})^T \quad \xRightarrow{\text{Lemma 8.2.10}} \quad \mathbf{c} = \frac{1}{2n+1} \mathbf{F}_{2n} \mathbf{z}. \quad (9.3.8)$$

$(2n+1) \times (2n+1)$ (conjugate) **Fourier matrix**, see (8.2.9)

Code 9.3.10: Efficient computation of coefficient of trigonometric interpolation polynomial (*equidistant nodes*)

```

1 function [a,b] = trigipequid(y)
2 % Efficient computation of coefficients in expansion (9.3.3) for a trigonometric
3 % interpolation polynomial in equidistant points  $(\frac{j}{2n+1}, y_j)$ ,  $j = 0, \dots, 2n$ 
4 % y has to be a row vector of odd length, return values are column vectors
5 N = length(y); if (mod(N,2)~=1), error('#pts odd!'); end;
6 n = (N-1)/2;
7 c = fft(exp(2*pi*i*(n/N)*(0:2*n)).*y)/N; % see (9.3.8)
8 % From (9.3.5):  $\alpha_j = \frac{1}{2}(\gamma_{n-j} + \gamma_{n+j})$  and  $\beta_j = \frac{1}{2i}(\gamma_{n-j} - \gamma_{n+j})$ ,  $j = 1, \dots, n$ ,  $\alpha_0 = \gamma_n$ 
9 a = transpose([c(n+1), c(n:-1:1)+c(n+2:N)]);
10 b = transpose(-i*[c(n:-1:1)-c(n+2:N)]);

```

Code 9.3.11: Computation of coefficients of trigonometric interpolation polynomial, general nodes

```

1 function [a,b] = trigpolycoeff(t,y)

```

```
2 % Computes expansion coefficients of trigonometric polynomials (9.3.3)
3 % t: row vector of nodes  $t_0, \dots, t_n \in [0, 1[$ 
4 % y: row vector of data  $y_0, \dots, y_n$ 
5 % return values are column vectors of expansion coefficients  $\alpha_j, \beta_j$ 
6 N = length(y); if (mod(N,2)~=1), error('#pts odd!'); end
7 n = (N-1)/2;
8 M = [cos(2*pi*t'*(0:n)), sin(2*pi*t'*(1:n))];
9 c = M\y';
10 a = c(1:n+1); b = c(n+2:end);
```

Example 9.3.12 (Runtime comparison for computation of coefficient of trigonometric interpolation polynomials).

tic-toc-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

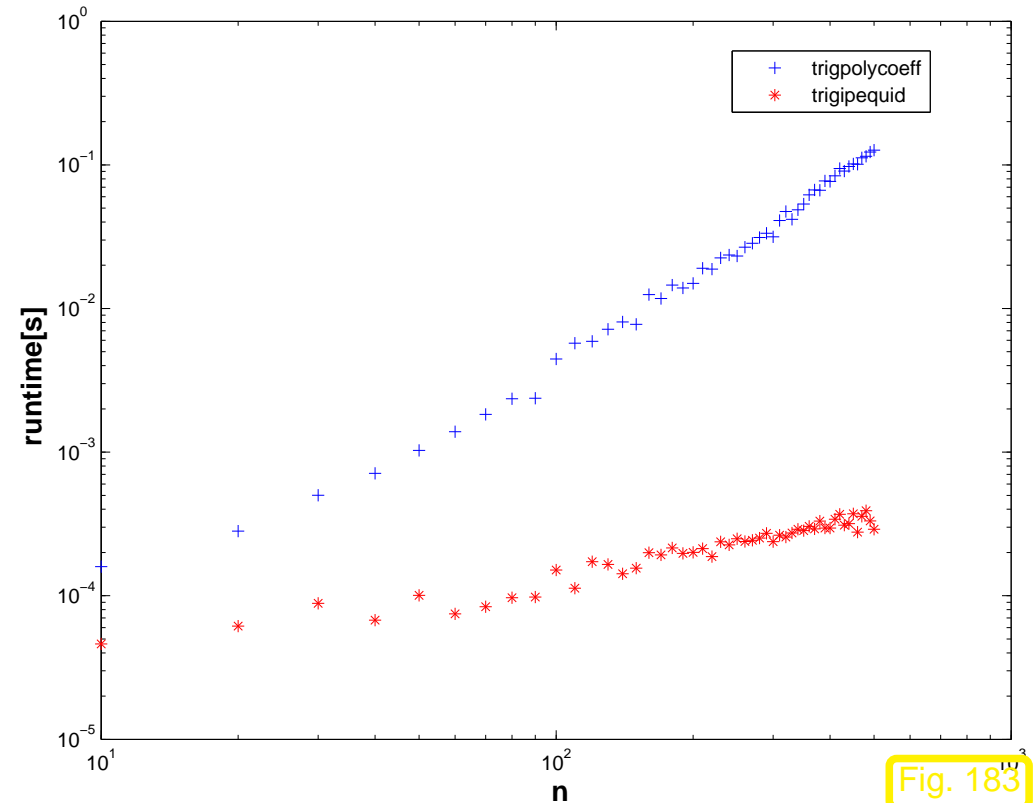


Fig. 183

*Remark 9.3.14* (Efficient evaluation of trigonometric interpolation polynomials).Code 9.3.16: Fast evaluation of trigonometric polynomial at *equidistant* points

```

1 function q = trigipequidcomp(a,b,N)
2 % Efficient evaluation of trigonometric polynomial at equidistant points
3 % column vectors a and b pass coefficients  $\alpha_j, \beta_j$  in
4 % representation (9.3.3)
5 n = length(a)-1; if (N < (2*n-1)), error('N too small'); end;
6 gamma = transpose(0.5*[a(end:-1:2)+i*b(end:-1:1)];...
```

```

7         2*a(1);a(2:end)-i*b(1:end)]);
8 ch = [gamma, zeros(1,N-(2*n+1))]; % zero padding
9 v = conj(fft(conj(ch))); % Multiplication with conjugate Fourier matrix
10 q = exp(-2*pi*i*n*(0:N-1)/N).*v; % undo rescaling

```

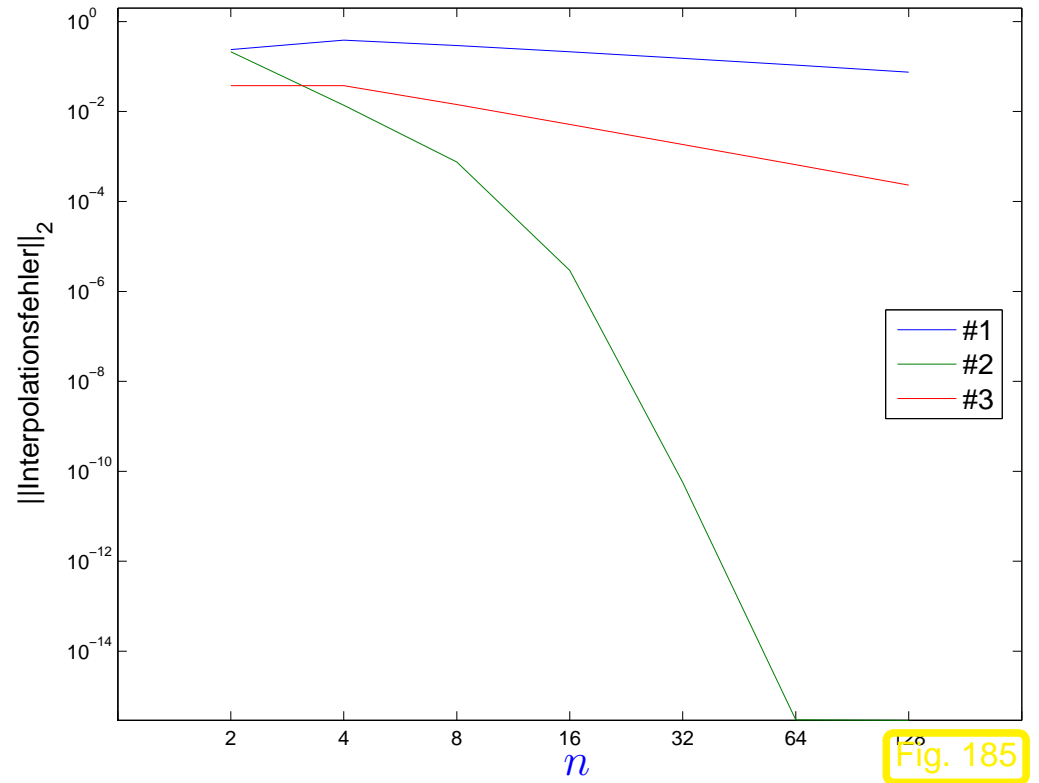
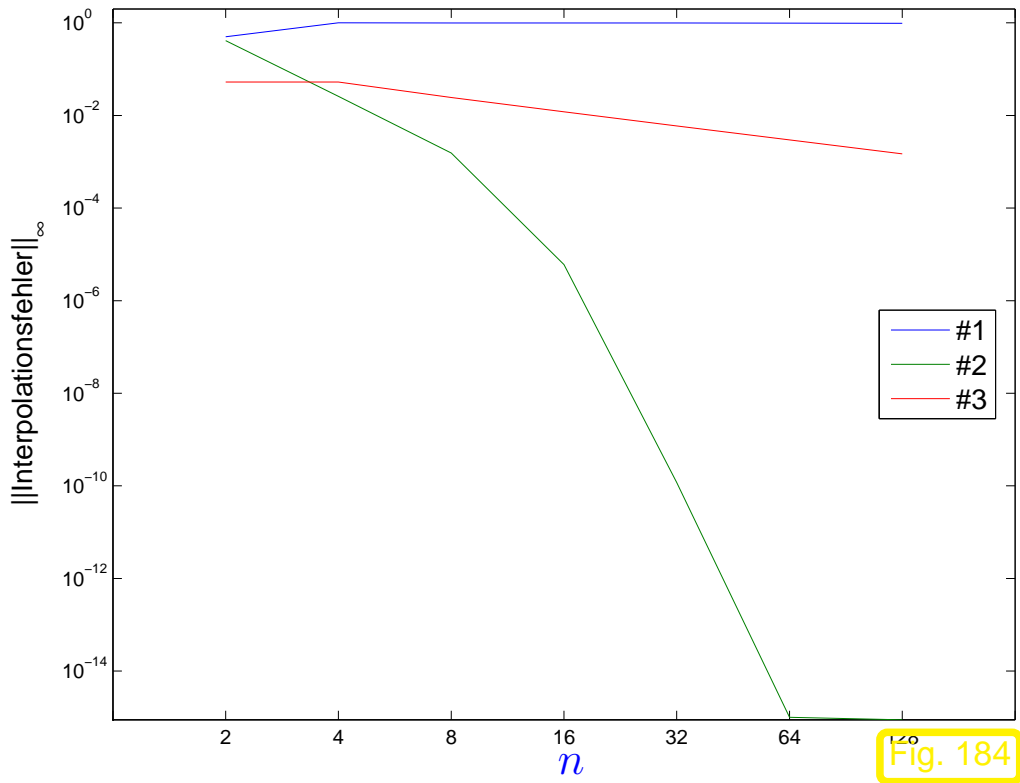


Example 9.3.18 (Interpolation error: trigonometric interpolation).

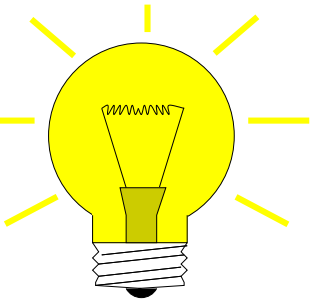
#1 Step function: $f(t) = 0$ for $|t - \frac{1}{2}| > \frac{1}{4}$, $f(t) = 1$ for $|t - \frac{1}{2}| \leq \frac{1}{4}$

#2 C^∞ periodic function: $f(t) = \frac{1}{\sqrt{1 + \frac{1}{2} \sin(2\pi t)}}$.

#3 “wedge function”: $f(t) = |t - \frac{1}{2}|$



9.4 Approximation by piecewise polynomials



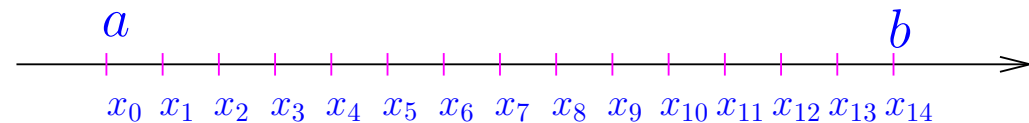
Idea: use **piecewise polynomials** with respect to a **grid/mesh**

$$\mathcal{M} := \{a = x_0 < x_1 < \dots < x_{m-1} < x_m = b\} \quad (9.4.1)$$

to approximate function $f : [a, b] \mapsto \mathbb{R}$, $a < b$.

Terminology:

- $x_j \hat{=}$ **nodes** of the mesh \mathcal{M} ,
- $[x_{j-1}, x_j[\hat{=}$ **intervals/cells** of the mesh,
- $h_{\mathcal{M}} := \max_j |x_j - x_{j-1}| \hat{=}$ **mesh width**,
- If $x_j = a + jh \hat{=}$ **equidistant** (uniform) mesh with meshwidth $h > 0$



9.4.1 Piecewise Lagrange interpolation

Focus: **asymptotic** behavior of (some norm of) interpolation error

$$\|f - I_{\mathcal{M}}f\| \leq CT(N) \quad \text{for } N \rightarrow \infty, \quad (9.4.4)$$

where
$$N := \sum_{j=1}^m (n_j + 1).$$

Example 9.4.5 (*h*-convergence of piecewise polynomial interpolation).

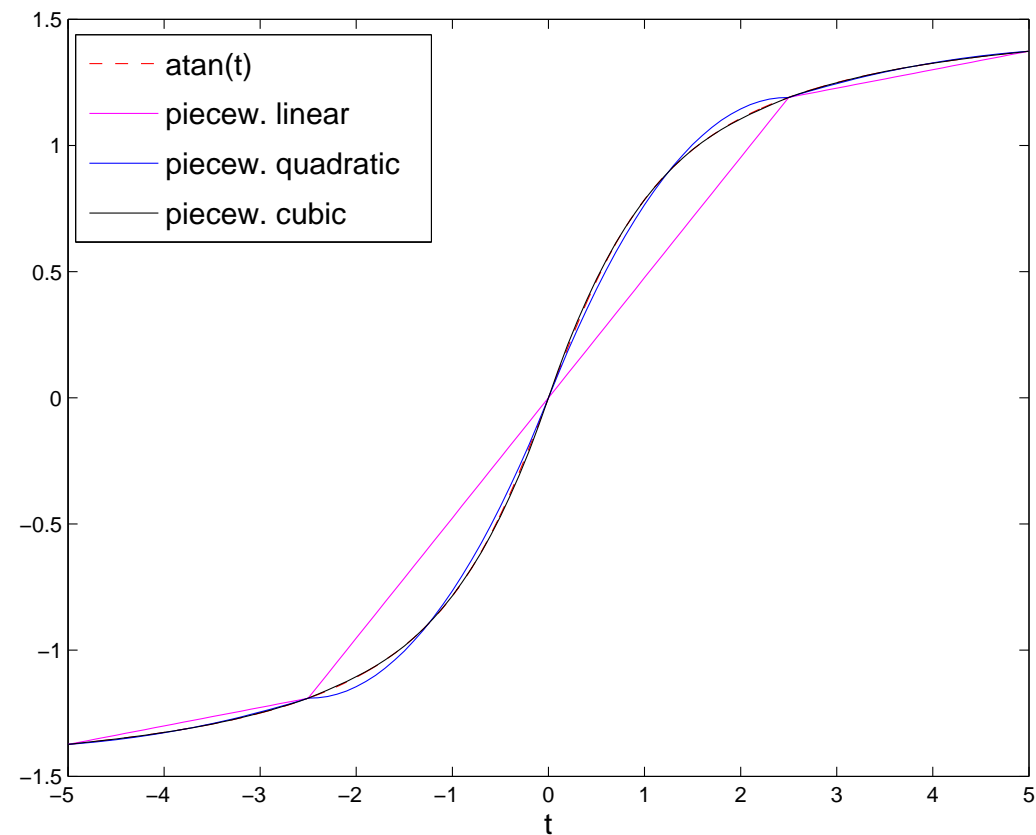
Compare Ex. 3.6.5:

$$f(t) = \arctan t, I = [-5, 5]$$

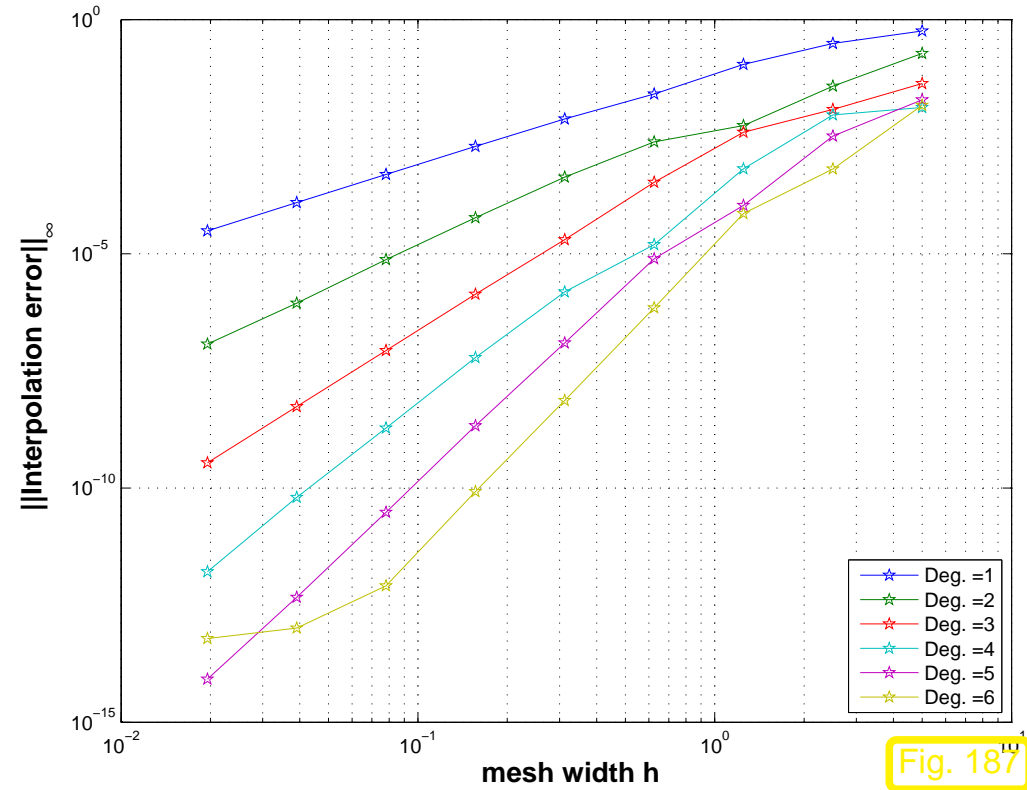
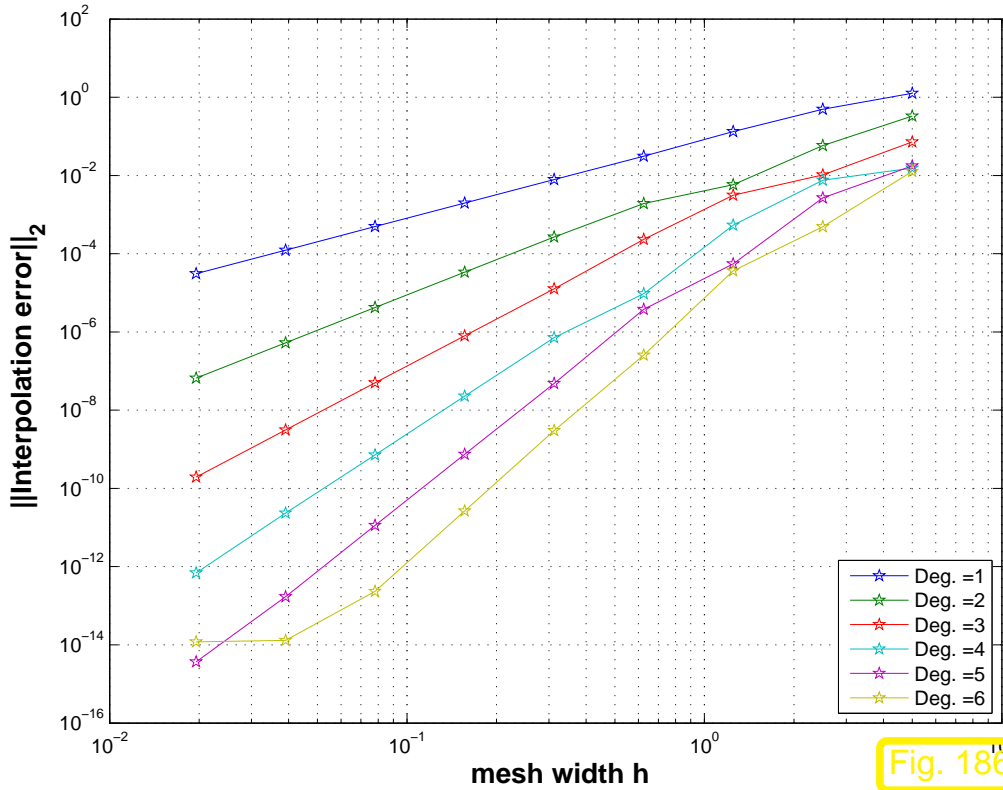
$$\text{Grid } \mathcal{M} := \left\{-5, -\frac{5}{2}, 0, \frac{5}{2}, 5\right\}$$

\mathcal{T}^j equidistant in I_j .

Plots of the piecewise linear, quadratic and cubic
polynomial interpolants \rightarrow



- Sequence of (equidistant) meshes: $\mathcal{M}_i := \left\{-5 + j 2^{-i} 10\right\}_{j=0}^{2^i}, \quad i = 1, \dots, 6.$
- Equidistant local interpolation nodes (endpoints of grid intervals included).



Rates α of algebraic convergence $O(h_M^\alpha)$ of norms of interpolation error:

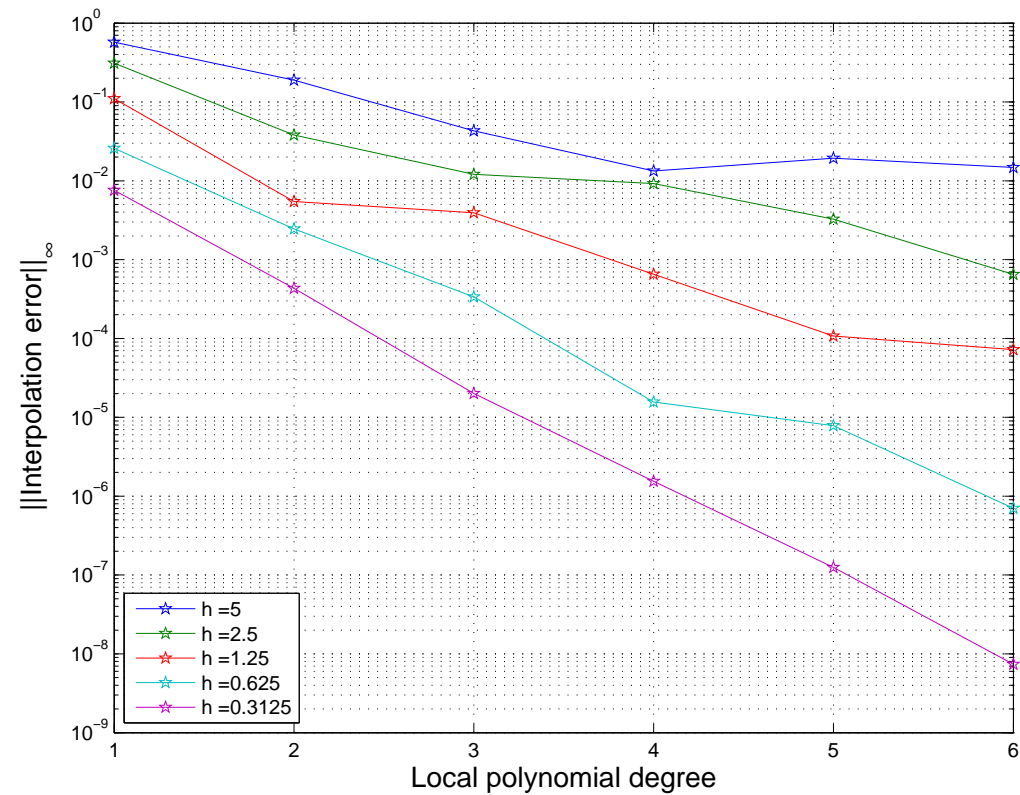
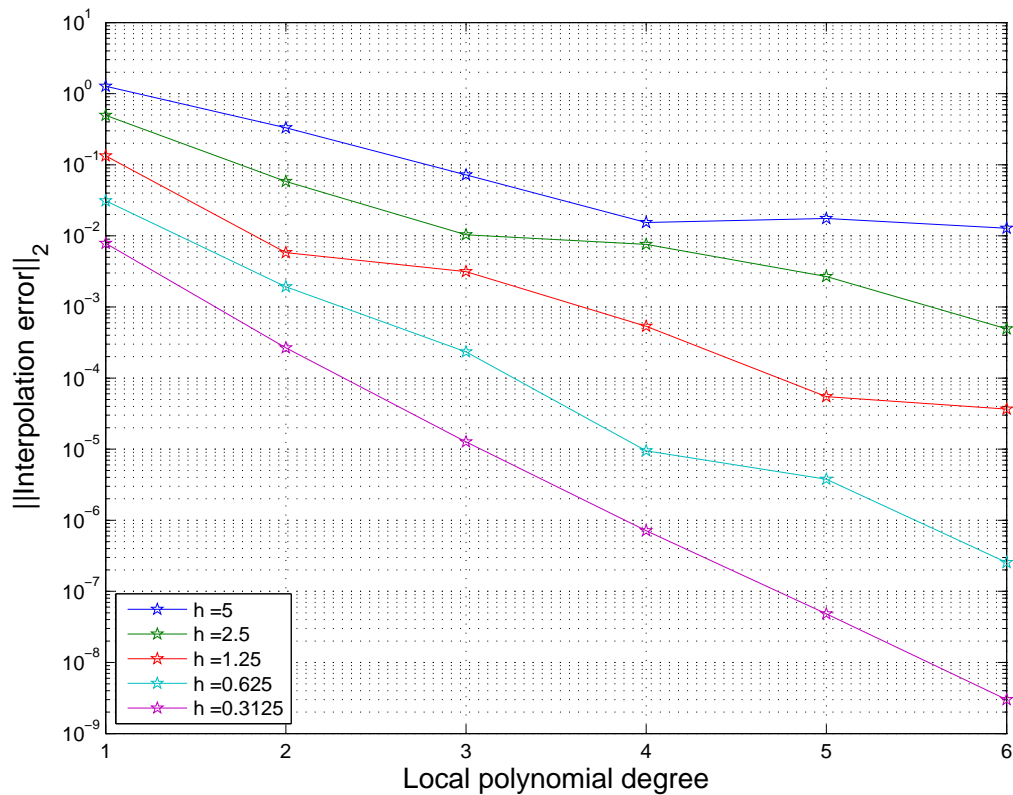
| n | 1 | 2 | 3 | 4 | 5 | 6 |
|-------------------------|--------|--------|--------|--------|--------|--------|
| w.r.t. L^2 -norm | 1.9957 | 2.9747 | 4.0256 | 4.8070 | 6.0013 | 5.2012 |
| w.r.t. L^∞ -norm | 1.9529 | 2.8989 | 3.9712 | 4.7057 | 5.9801 | 4.9228 |



$$(9.1.12) \Rightarrow \|f - s\|_{L^\infty([x_0, x_m])} \leq \frac{h^{n+1}}{(n+1)!} \|f^{(n+1)}\|_{L^\infty([x_0, x_m])}, \quad (9.4.6)$$

with **mesh width** $h := \max\{|x_j - x_{j-1}| : j = 1, \dots, m\}$.

Example 9.4.7 (p -convergence of piecewise polynomial interpolation).



9.4.2 Cubic Hermite interpolation: error estimates

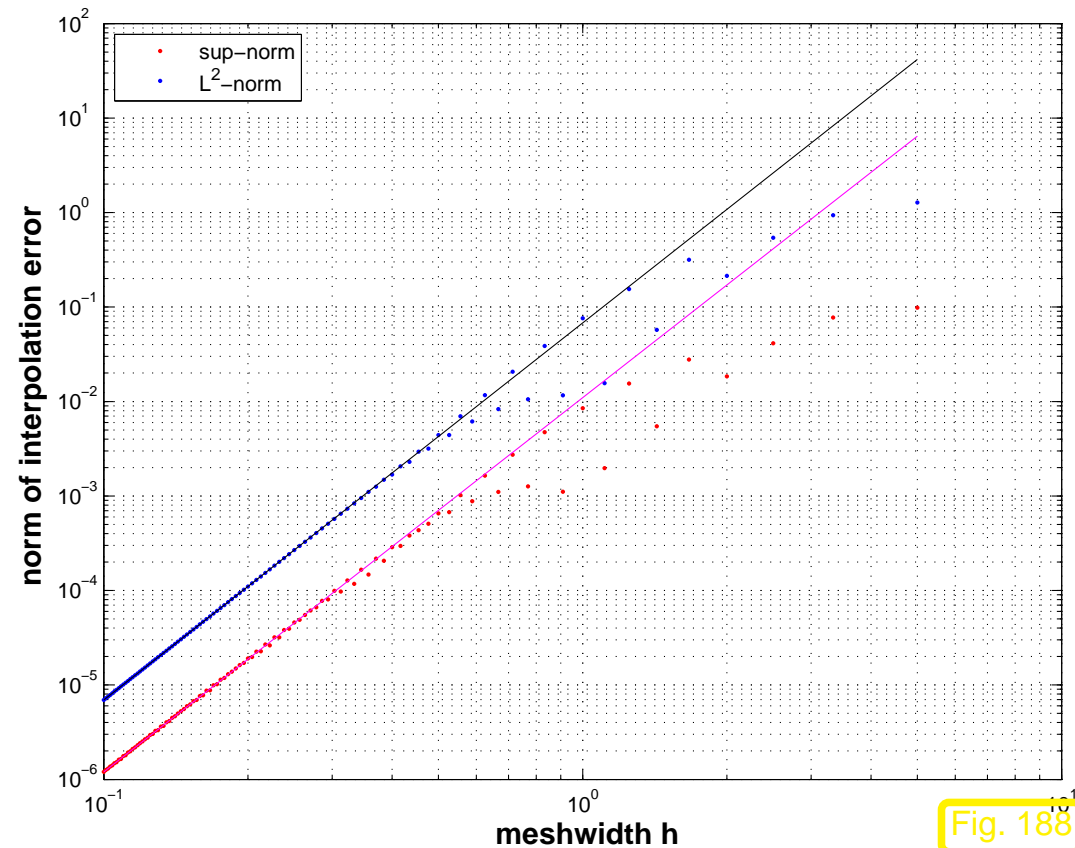
Example 9.4.8 (Convergence of Hermite interpolation with exact slopes).

Piecewise cubic Hermite interpolation of

$$f(x) = \arctan(x) .$$

- domain: $I = (-5, 5)$
- mesh $\mathcal{T} = \{-5 + hj\}_{j=0}^n \subset I$, $h = \frac{10}{n}$,
- exact slopes $c_i = f'(t_i)$, $i = 0, \dots, n$

► algebraic convergence $O(h^4)$



Example 9.4.10 (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 Ex. 9.4.8,
- averaged local slopes, see (3.7.4)

▶ algebraic convergence $O(h^3)$ in meshwidth
(see Code 9.4.10)

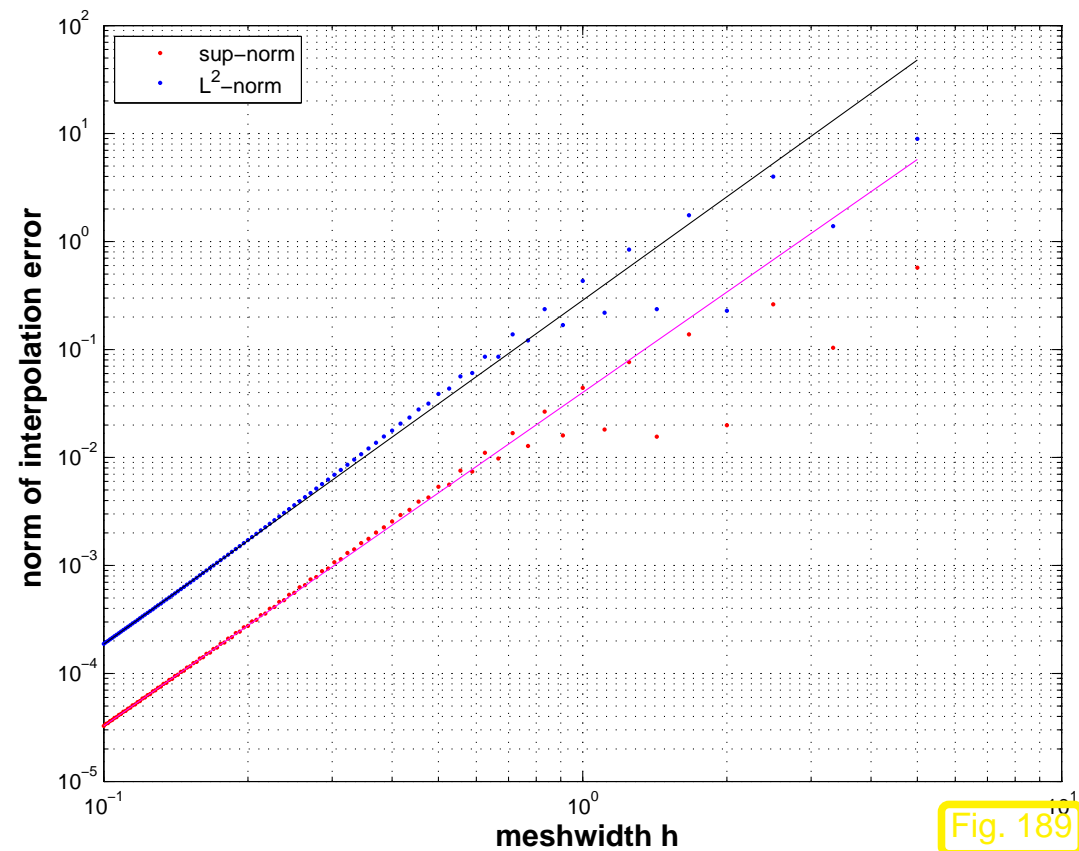


Fig. 189

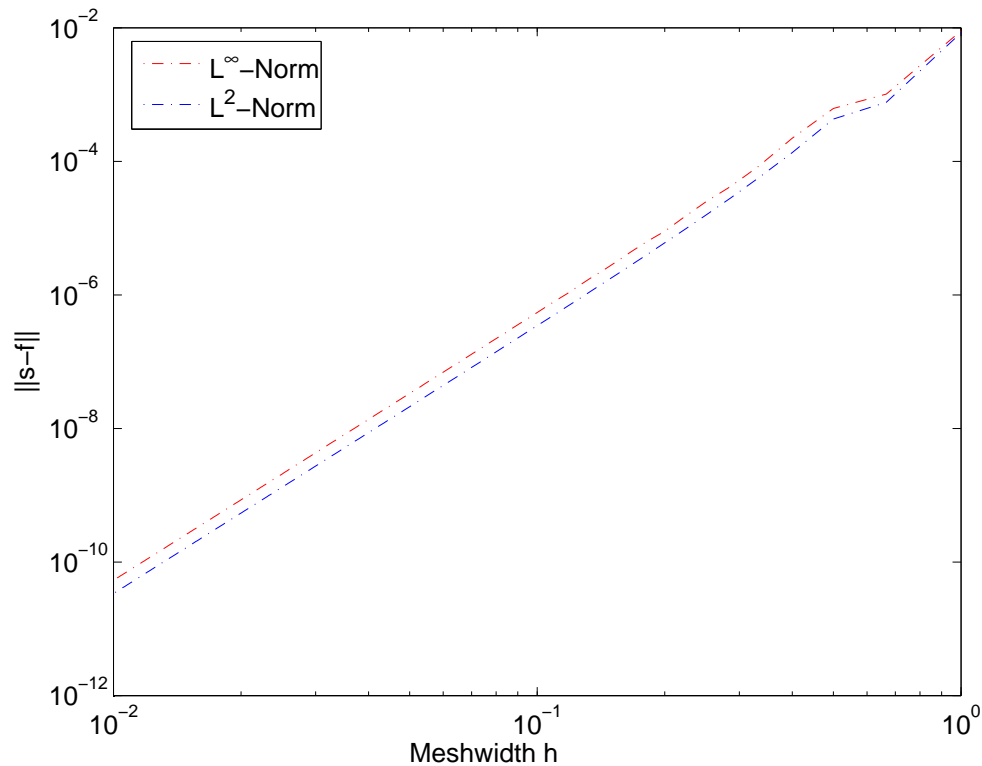


9.4.3 Cubic spline interpolation: error estimates [35, Ch. 47]

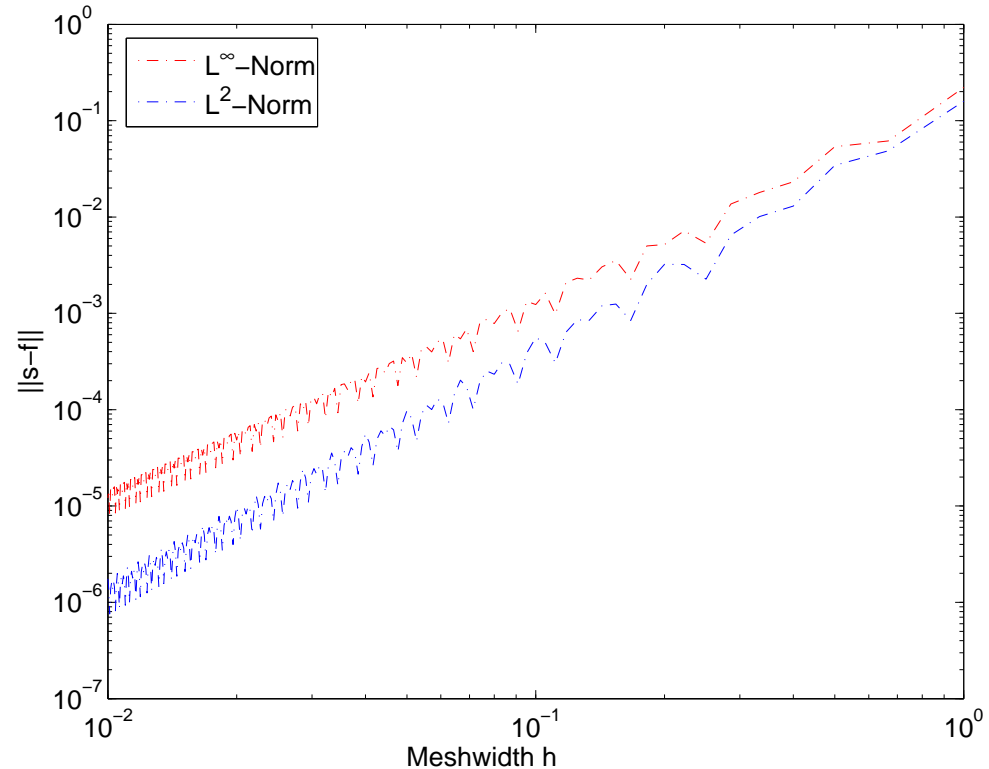
Example 9.4.12 (Approximation by complete cubic spline interpolants).

Grid (knot set) $\mathcal{M} := \{-1 + \frac{2}{n}j\}_{j=0}^n, n \in \mathbb{N} \rightarrow$ meshwidth $h = 2/n, I = [-1, 1]$

$$f_1(t) = \frac{1}{1 + e^{-2t}} \in C^\infty(I) \quad , \quad f_2(t) = \begin{cases} 0 & , \text{if } t < -\frac{2}{5} , \\ \frac{1}{2}(1 + \cos(\pi(t - \frac{3}{5}))) & , \text{if } -\frac{2}{5} < t < \frac{3}{5} , \\ 1 & \text{otherwise.} \end{cases} \in C^1(I) .$$



$$\|f_1 - s\|_{L^\infty([-1,1])} = O(h^4)$$



$$\|f_2 - s\|_{L^\infty([-1,1])} = O(h^2)$$

Algebraic order of convergence in $h = \min \{ 1 + \text{regularity of } f, 4 \}$.

Theory [34]: $f \in C^4([t_0, t_n]) \rightarrow \|f - s\|_{L^\infty([t_0, t_n])} \leq \frac{5}{384} h^4 \left\| f^{(4)} \right\|_{L^\infty([t_0, t_n])}$

Numerical Quadrature [35, VII], [13, Ch. 10]

10

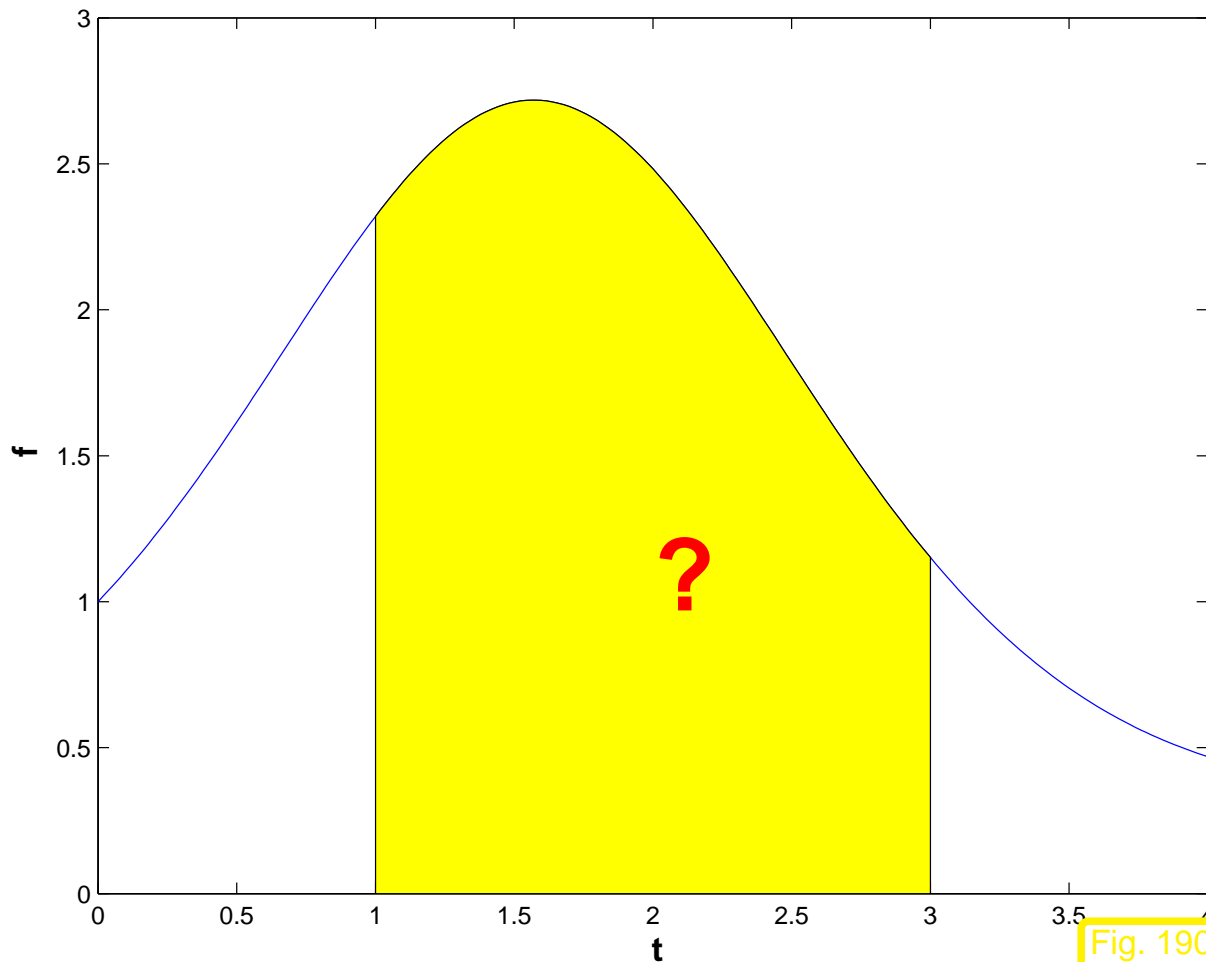


Fig. 190

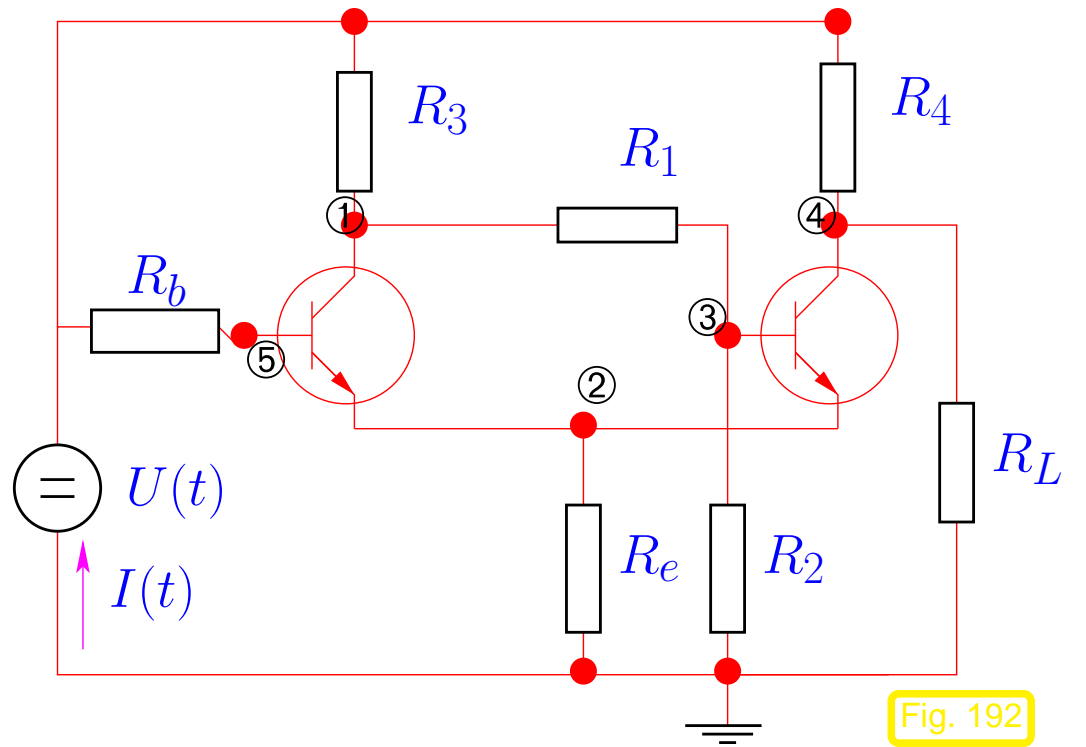
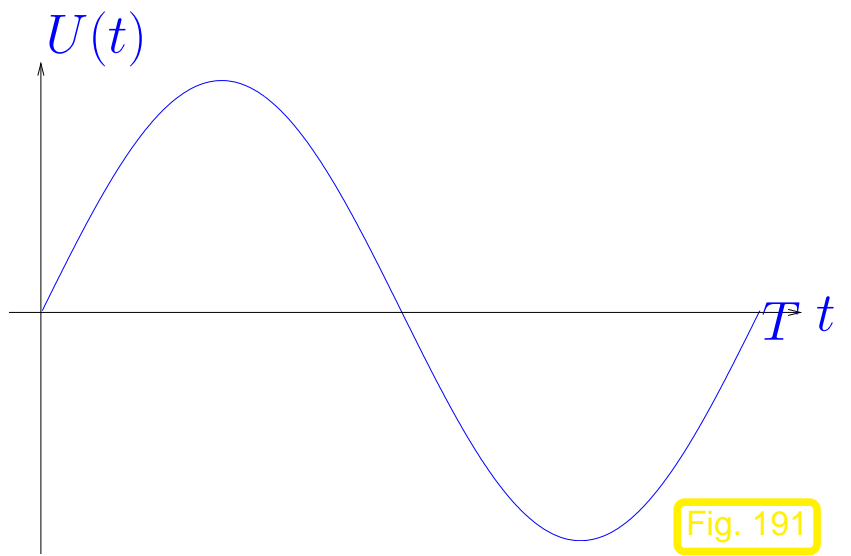
Numerical quadrature methods

approximate

$$\int_a^b f(t) dt$$

Example 10.0.2 (Heating production in electrical circuits).

Time-harmonic excitation:



Integrating power $P = UI$ over period $[0, T]$ yields heat production per period:

$$W_{\text{therm}} = \int_0^T U(t)I(t) dt, \text{ where } I = I(U).$$

function $I = \text{current}(U)$ involves solving non-linear system of equations, see Ex. 4.0.1!



10.1 Quadrature Formulas

n -point **quadrature formula** on $[a, b]$:
(n -point quadrature rule)

$$\int_a^b f(t) dt \approx Q_n(f) := \sum_{j=1}^n w_j^n f(c_j^n). \quad (10.1.1)$$

w_j^n : **quadrature weights** $\in \mathbb{R}$ (ger.: Quadraturgewichte)
 c_j^n : **quadrature nodes** $\in [a, b]$ (ger.: Quadraturknoten)

Remark 10.1.3 (Transformation of quadrature rules).



Our focus (cf. interpolation error estimates, Sect. 9.1):

given **families** of quadrature rules $\{Q_n\}_n$ 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

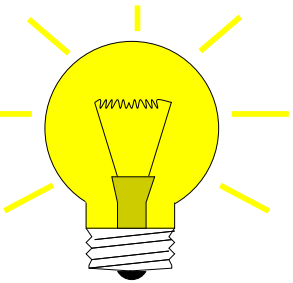
study the *asymptotic* behavior of the **quadrature error** $E(n)$ for $n \rightarrow \infty$

Qualitative distinction, see (9.1.3):

- ▷ algebraic convergence $E(n) = O(n^{-p}), p > 0$
- ▷ exponential convergence $E(n) = O(q^n), 0 \leq q < 1$

10.2 Polynomial Quadrature Formulas [13, Sect. 10.2]

Idea: replace integrand f with $p_{n-1} \in \mathcal{P}_{n-1}$ = polynomial interpolant of f for given interpolation nodes $\{t_0, \dots, t_{n-1}\} \subset [a, b]$

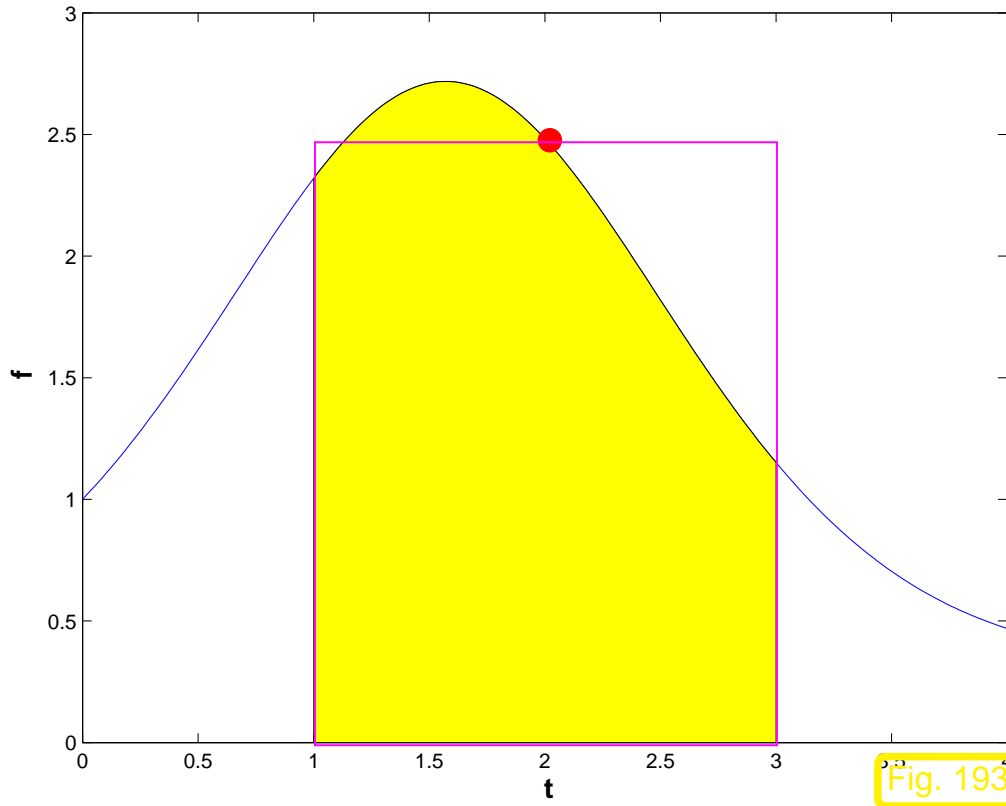


$$\blacktriangleright \int_a^b f(t) dt \approx Q_n(f) := \int_a^b p_{n-1}(t) dt . \quad (10.2.3)$$

Lagrange polynomials: $L_i(t) := \prod_{\substack{j=0 \\ j \neq i}}^{n-1} \frac{t - t_j}{t_i - t_j}, \quad i = 0, \dots, n-1 \quad (3.3.4) \quad \blacktriangleright \quad p_{n-1}(t) = \sum_{i=0}^{n-1} f(t_i) L_i(t) .$

$$\int_a^b p_{n-1}(t) dt = \sum_{i=0}^{n-1} f(t_i) \int_a^b L_i(t) dt \quad \blacktriangleright \quad \begin{array}{l} \text{nodes } c_i = t_{i-1} , \\ \text{weights } w_i := \int_a^b L_{i-1}(t) dt . \end{array} \quad (10.2.4)$$

Example 10.2.5 (Midpoint rule).



1-point quadrature formula:

$$\int_a^b f(t) dt \approx Q_{\text{mp}}(f) = (b - a) f\left(\frac{1}{2}(a + b)\right).$$

“midpoint”



Example 10.2.6 (Newton-Cotes formulas). [35, Ch. 38]

Equidistant quadrature nodes $t_j := a + hj, \quad h := \frac{b - a}{n}, \quad j = 0, \dots, n:$

Symbolic computation of quadrature formulas on $[0, 1]$ using MAPLE:

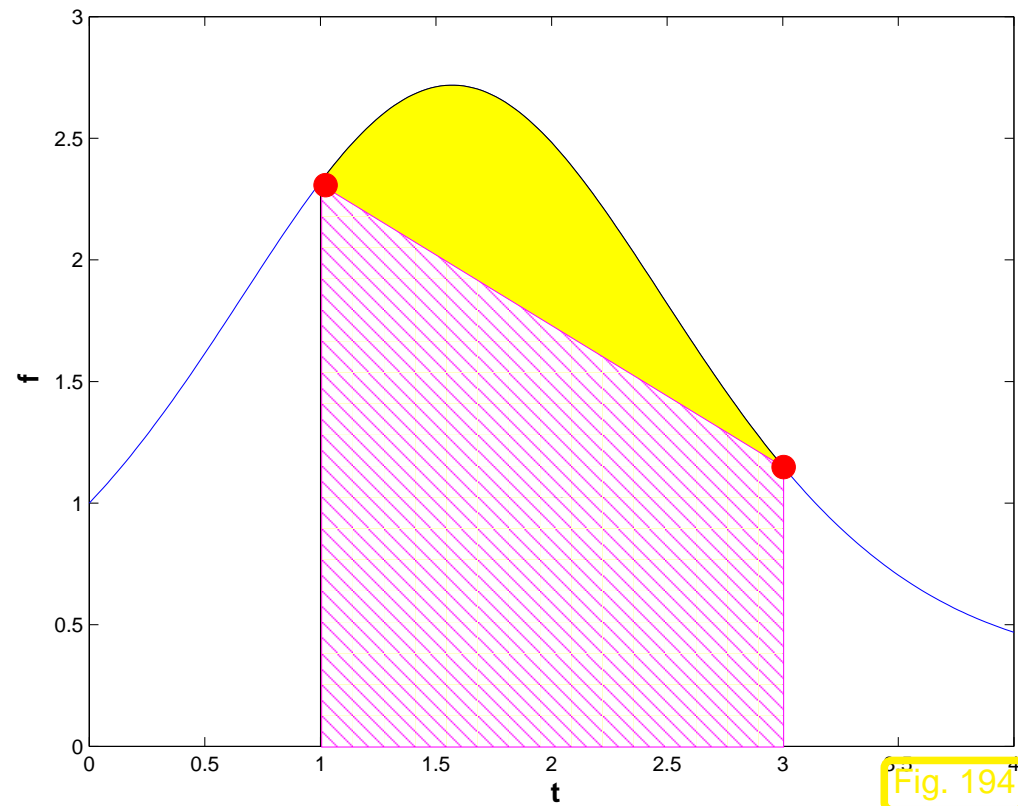
```
> newtoncotes := n -> factor(int(interp([seq(i/n, i=0..n)],
    [seq(f(i/n), i=0..n)], z), z=0..1)):
```

- $n = 1$: Trapezoidal rule

```
> trapez := newtoncotes(1);
```

$$\hat{Q}_{\text{trp}}(f) := \frac{1}{2} (f(0) + f(1)) \quad (10.2.7)$$

$$\left(\int_a^b f(t) dt \approx \frac{b-a}{2} (f(a) + f(b)) \right)$$



- $n = 2$: Simpson rule

```
> simpson := newtoncotes(2);
```

$$\frac{h}{6} \left(f(0) + 4 f\left(\frac{1}{2}\right) + f(1) \right) \left(\int_a^b f(t) dt \approx \frac{b-a}{6} \left(f(a) + 4 f\left(\frac{a+b}{2}\right) + f(b) \right) \right) \quad (10.2.8)$$

- $n \geq 8$: quadrature formulas with *negative* weights

> newtoncotes(8);

$$\frac{1}{28350} h \left(989 f(0) + 5888 f\left(\frac{1}{8}\right) - 928 f\left(\frac{1}{4}\right) + 10496 f\left(\frac{3}{8}\right) \right. \\ \left. - 4540 f\left(\frac{1}{2}\right) + 10496 f\left(\frac{5}{8}\right) - 928 f\left(\frac{3}{4}\right) + 5888 f\left(\frac{7}{8}\right) + 989 f(1) \right)$$



Negative weights compromise numerical stability (\rightarrow Def. 2.5.11) !

10.3 Composite Quadrature

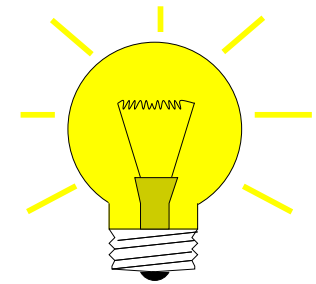
Idea: • Partition integration domain $[a, b]$ by **mesh** (grid, \rightarrow Sect.9.4)

$$\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$$

• Apply quadrature formulas from Sects. 10.2, 10.4 **locally** on mesh intervals

$I_j := [x_{j-1}, x_j], j = 1, \dots, m$, and sum up.

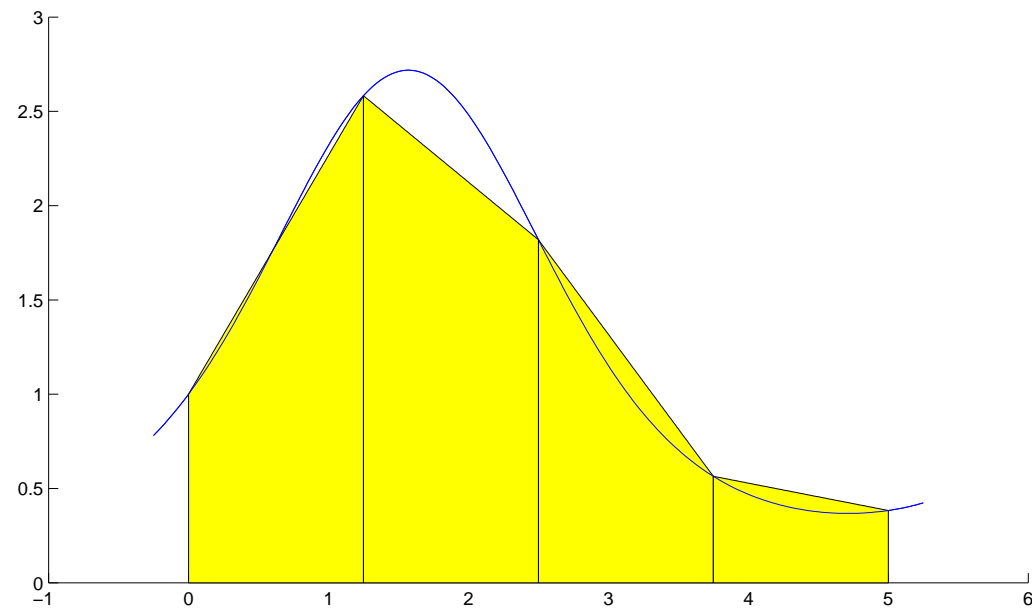
composite quadrature rule



Example 10.3.2 (Simple composite polynomial quadrature rules).

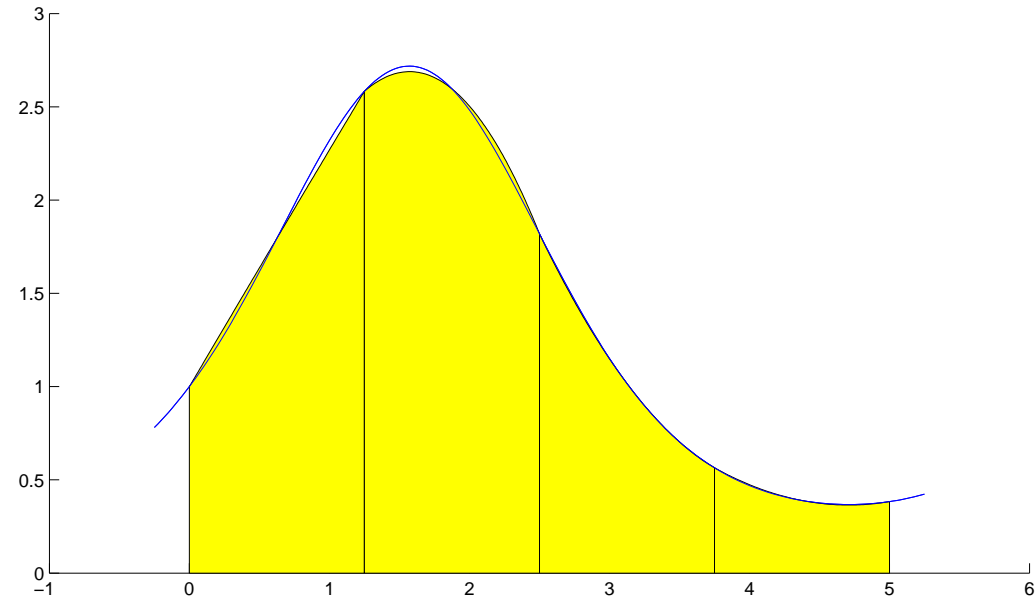
Composite trapezoidal rule, *cf.* (10.2.7)

$$\int_a^b f(t) dt = \frac{1}{2}(x_1 - x_0)f(a) + \sum_{j=1}^{m-1} \frac{1}{2}(x_{j+1} - x_{j-1})f(x_j) + \frac{1}{2}(x_m - x_{m-1})f(b) . \quad (10.3.3)$$



Composite Simpson rule, cf. (10.2.8)

$$\begin{aligned}
 \int_a^b f(t) dt = & \\
 & \frac{1}{6}(x_1 - x_0)f(a) + \\
 & \sum_{j=1}^{m-1} \frac{1}{6}(x_{j+1} - x_{j-1})f(x_j) + \\
 & \sum_{j=1}^m \frac{2}{3}(x_j - x_{j-1})f\left(\frac{1}{2}(x_j + x_{j-1})\right) + \\
 & \frac{1}{6}(x_m - x_{m-1})f(b) .
 \end{aligned}
 \tag{10.3.4}$$

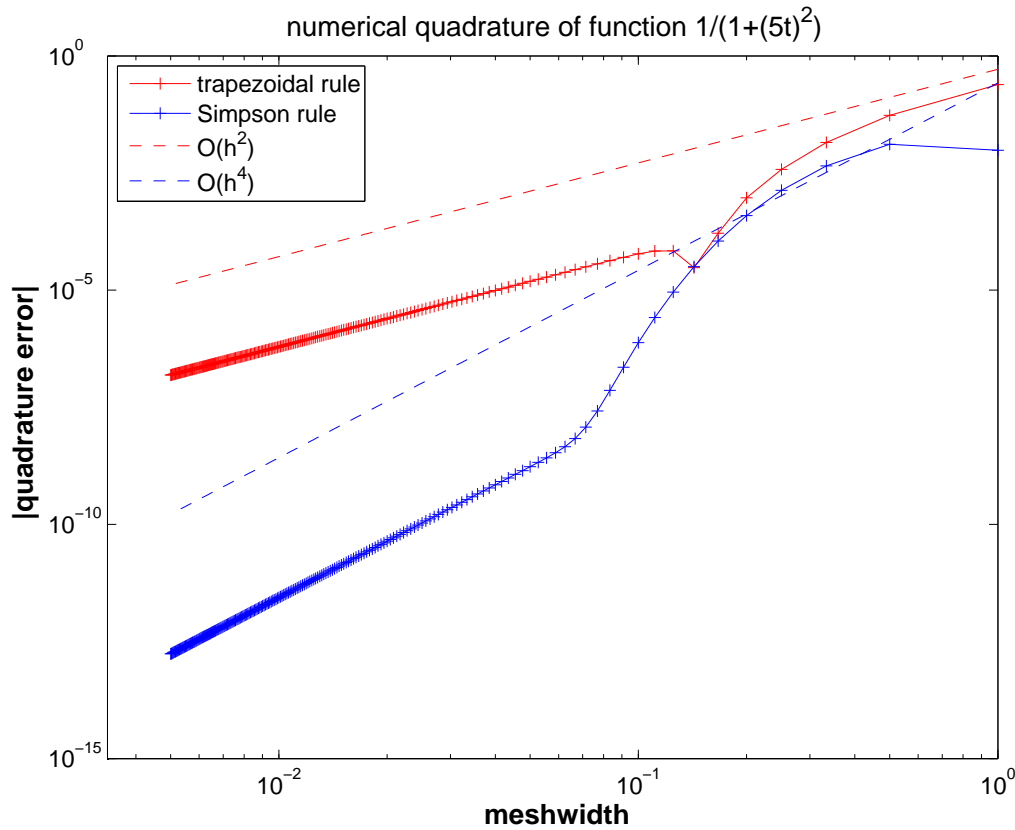


Example 10.3.5 (Quadrature errors for composite quadrature rules).

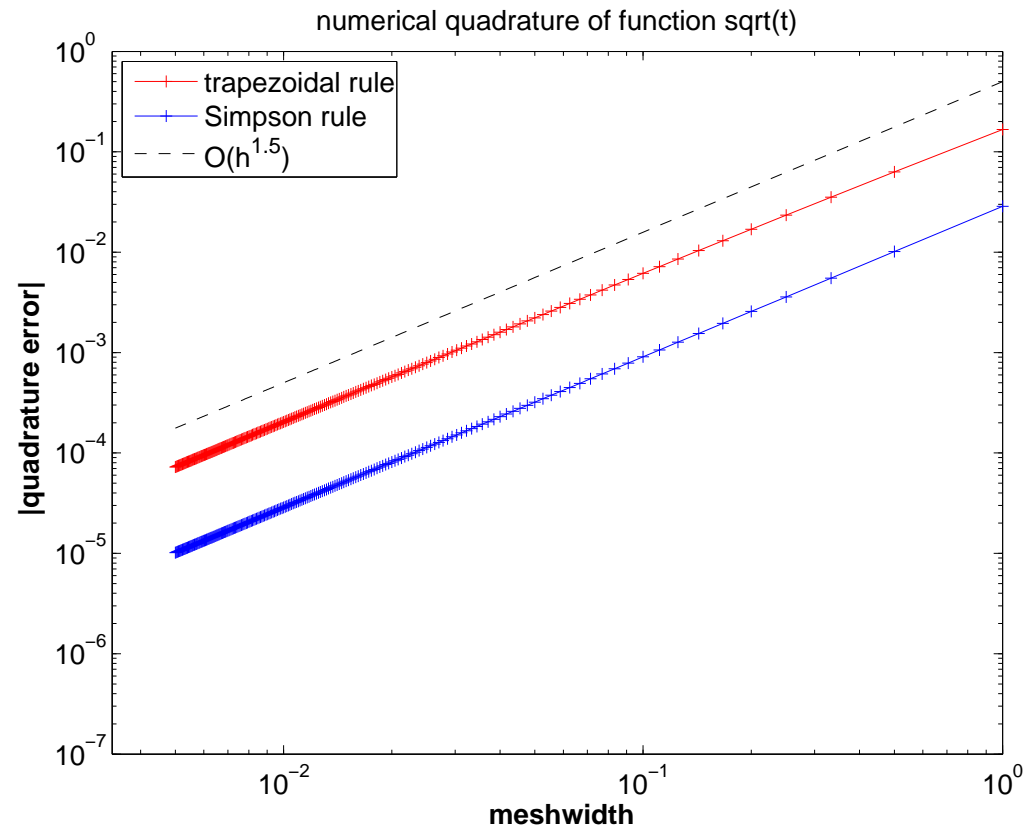
Composite quadrature rules based on

- trapezoidal rule (10.2.7) \triangleright local order 2 (exact for linear functions),
- Simpson rule (10.2.8) \triangleright local order 3 (exact for quadratic polynomials)

on equidistant mesh $\mathcal{M} := \{jh\}_{j=0}^n$, $h = 1/n$, $n \in \mathbb{N}$.



quadrature error, $f_1(t) := \frac{1}{1+(5t)^2}$ on $[0, 1]$



quadrature error, $f_2(t) := \sqrt{t}$ on $[0, 1]$

Composite Simpson rule: rate = 4 ? investigate *local* quadrature error on $[0, h]$ with 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);
```

$$err := \left(\frac{1}{90} \left(D^{(4)} \right) (f) (0) h^5 + O \left(h^6 \right), h, 6 \right)$$



➤ Composite Simpson rule converges with rate 4, indeed !

Gauge for “quality” of a quadrature formula Q_n :

$$\text{Order}(Q_n) := \max\{n \in \mathbb{N}_0: Q_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_n\} + 1$$

Focus: *asymptotic* behavior of quadrature error for

$$\text{mesh width} \quad h := \max_{j=1, \dots, m} |x_j - x_{j-1}| \rightarrow 0$$

Theorem 10.3.10 (Convergence of composite quadrature formulas).

For a composite quadrature formula Q based on a local quadrature formula of order $p \in \mathbb{N}$ holds

$$\exists C > 0: \left| \int_I f(t) dt - Q(f) \right| \leq Ch^p \left\| f^{(p)} \right\|_{L^\infty(I)} \quad \forall f \in C^p(I), \forall \mathcal{M}.$$



Remark 10.3.11 (Removing a singularity by transformation).

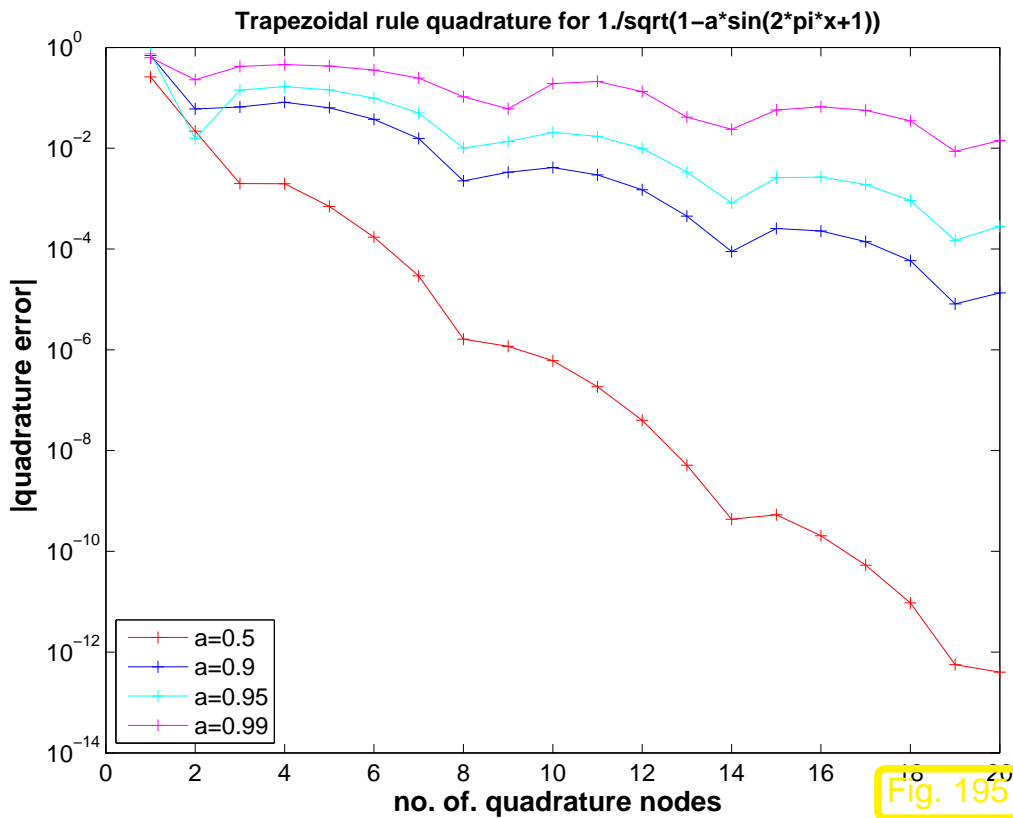
Example 10.3.13 (Convergence of equidistant trapezoidal rule).

Equidistant trapezoidal rule (order 2), see (10.3.3)

$$\int_a^b f(t) dt \approx T_m(f) := h \left(\frac{1}{2}f(a) + \sum_{k=1}^{m-1} f(kh) + \frac{1}{2}f(b) \right), \quad h := \frac{b-a}{m}. \quad (10.3.14)$$

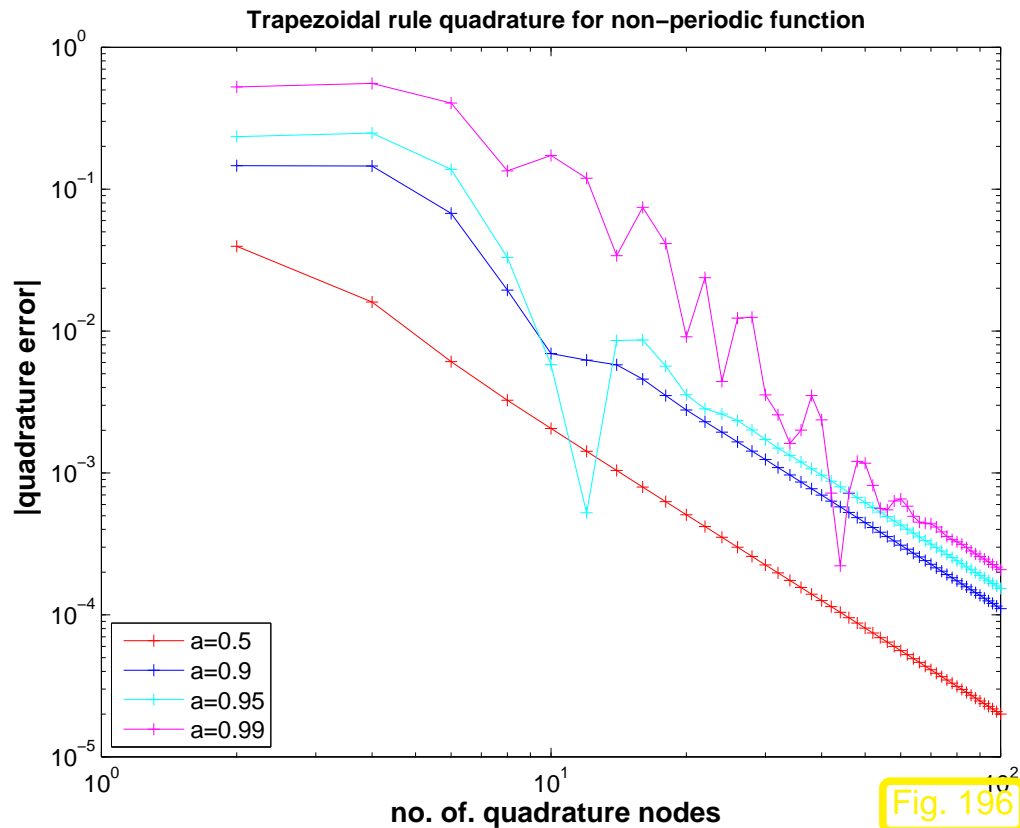
1-periodic smooth (**analytic**) integrand

$$f(t) = \frac{1}{\sqrt{1 - a \sin(2\pi t - 1)}}, \quad 0 < a < 1.$$



quadrature error for $T_n(f)$ on $[0, 1]$

exponential convergence !!



quadrature error for $T_n(f)$ on $[0, \frac{1}{2}]$

merely algebraic convergence



Remark 10.3.15 (Approximate computation of Fourier coefficients).

Code 10.3.16: DFT-based approximate computation of Fourier coefficients

```

1 function y = fourcoeffcomp(c,m,ovsmp1)
2 % Compute the Fourier coefficients  $y_{-m}, \dots, y_m$  of the function
    
```

```
3 % c: [0,1[ → ℂ using an oversampling factor ovsmpl.
4 % c must be a handle to a function @(t), which accepts row
5 % vector arguments
6 if ( nargin < 3 ), ovsmpl = 2; else ovsmpl = ceil(ovsmpl); end
7 N = (2*m+1)*ovsmpl; h = 1/N; % Number of quadrature points
8 % (Inverse) discrete Fourier transform
9 y = ifft(c(0:h:1-h));
0 % Undo oversampling and wrapping of Fourier coefficient array
1 y = [y(N-m+1:N), y(1:m+1)];
```



Remark 10.3.17 (Choice of (local) quadrature weights).




Lemma 10.3.20 (Bound for order of quadrature formula).

There is no n -point quadrature formula of order $2n + 1$

10.4 Gauss Quadrature [35, Ch. 40-41], [13, Sect.10.3]

Example 10.4.1 (2-point quadrature rule of order 4).

Necessary & sufficient conditions for order 4 , *cf.* (10.3.21):

$$Q_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_3 \quad \Leftrightarrow \quad Q_n(t^q) = \frac{1}{q+1}(b^{q+1} - a^{q+1}), \quad q = 0, 1, 2, 3 .$$


4 equations for weights ω_j and nodes $\xi_j, j = 1, 2$ ($a = -1, b = 1$), *cf.* Rem. 10.3.19

$$\begin{aligned} \int_{-1}^1 1 dt = 2 = 1\omega_1 + 1\omega_2 & , & \int_{-1}^1 t dt = 0 = \xi_1\omega_1 + \xi_2\omega_2 \\ \int_{-1}^1 t^2 dt = \frac{2}{3} = \xi_1^2\omega_1 + \xi_2^2\omega_2 & , & \int_{-1}^1 t^3 dt = 0 = \xi_1^3\omega_1 + \xi_2^3\omega_2 . \end{aligned} \tag{10.4.2}$$

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: $\left\{ \omega_2 = 1, \omega_1 = 1, \xi_1 = 1/3 \sqrt{3}, \xi_2 = -1/3 \sqrt{3} \right\}$

▶ quadrature formula:
$$\int_{-1}^1 f(x) dx \approx f\left(\frac{1}{\sqrt{3}}\right) + f\left(-\frac{1}{\sqrt{3}}\right) \quad (10.4.3)$$



Theorem 10.4.8 (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 $\{\xi_j^n\}_{j=1}^m$, $m \leq n$, of real zeros of \bar{P}_n is contained in $[-1, 1]$.

Then

$$Q_n(f) := \sum_{j=1}^m \omega_j^n f(\xi_j^n)$$

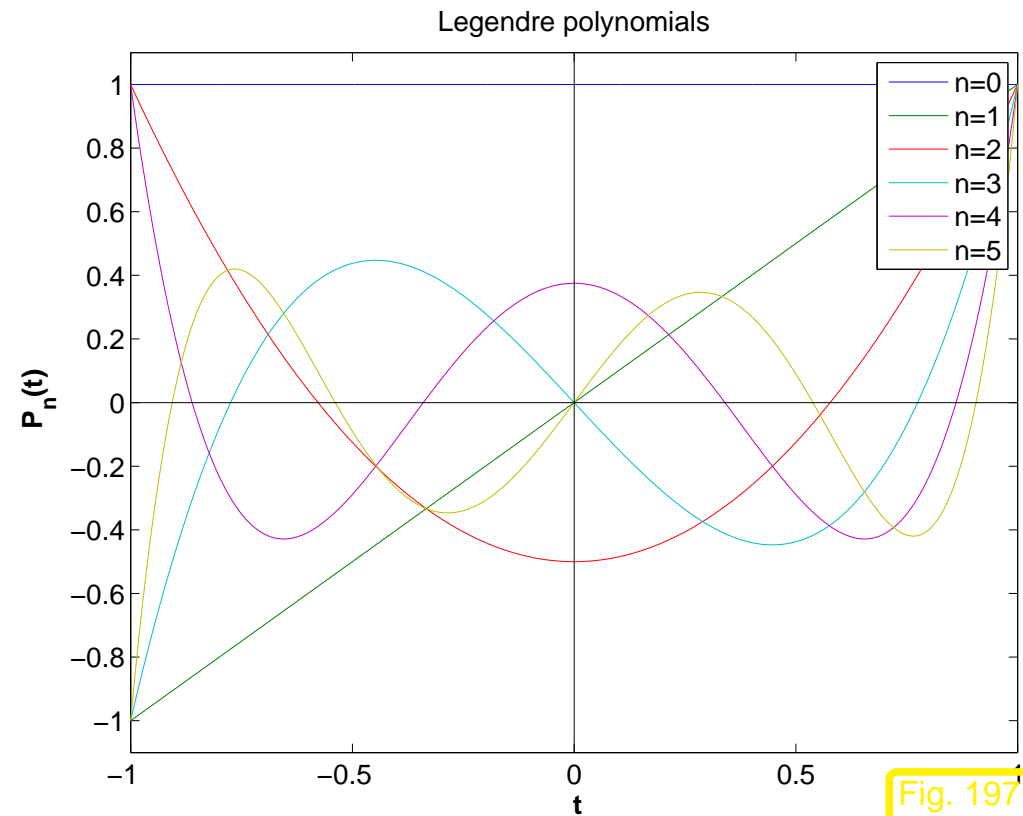
with weights chosen according to Rem. 10.3.19 provides a quadrature formula of order $2n$ on $[-1, 1]$.

Definition 10.4.12 (Legendre polynomials).

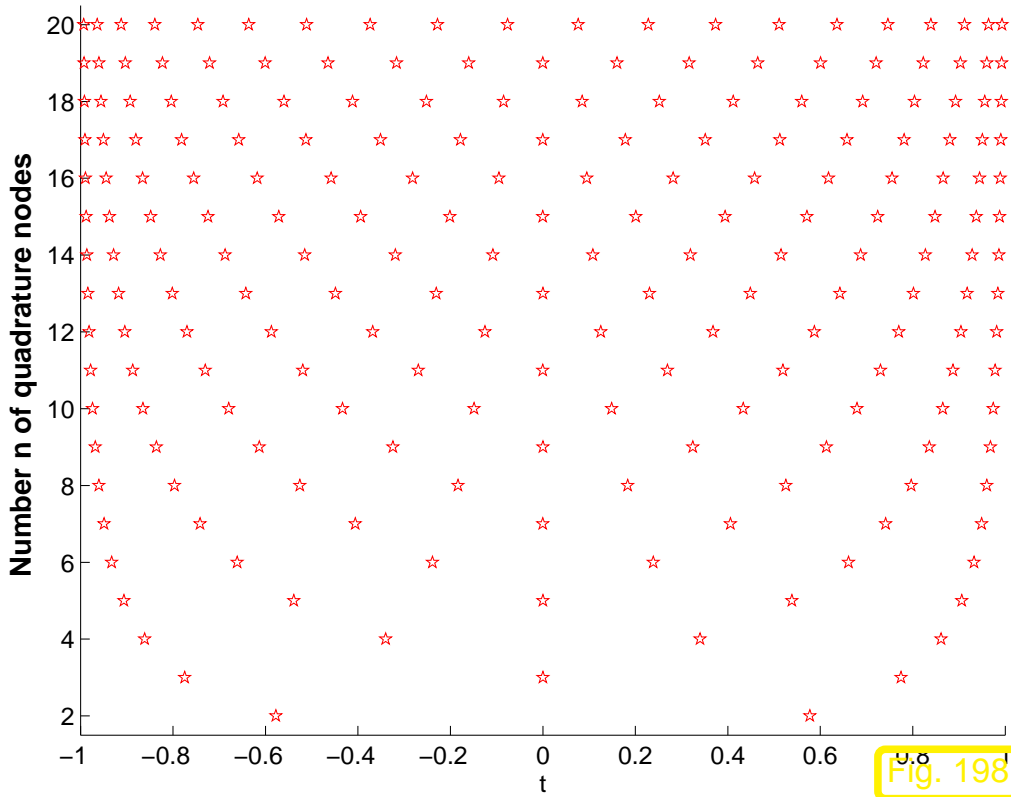
The n -th **Legendre polynomial** P_n is defined by

- $P_n \in \mathcal{P}_n$,
- $\int_{-1}^1 P_n(t)q(t) dt = 0 \quad \forall q \in \mathcal{P}_{n-1}$,
- $P_n(1) = 1$.

Legendre polynomials P_0, \dots, P_5



Zeros of Legendre polynomials in $[-1,1]$



◁ Obviously:

Lemma 10.4.13 (Zeros of Legendre polynomials).

P_n has n distinct zeros in $] - 1, 1[$.

Zeros of Legendre polynomials = **Gauss points**

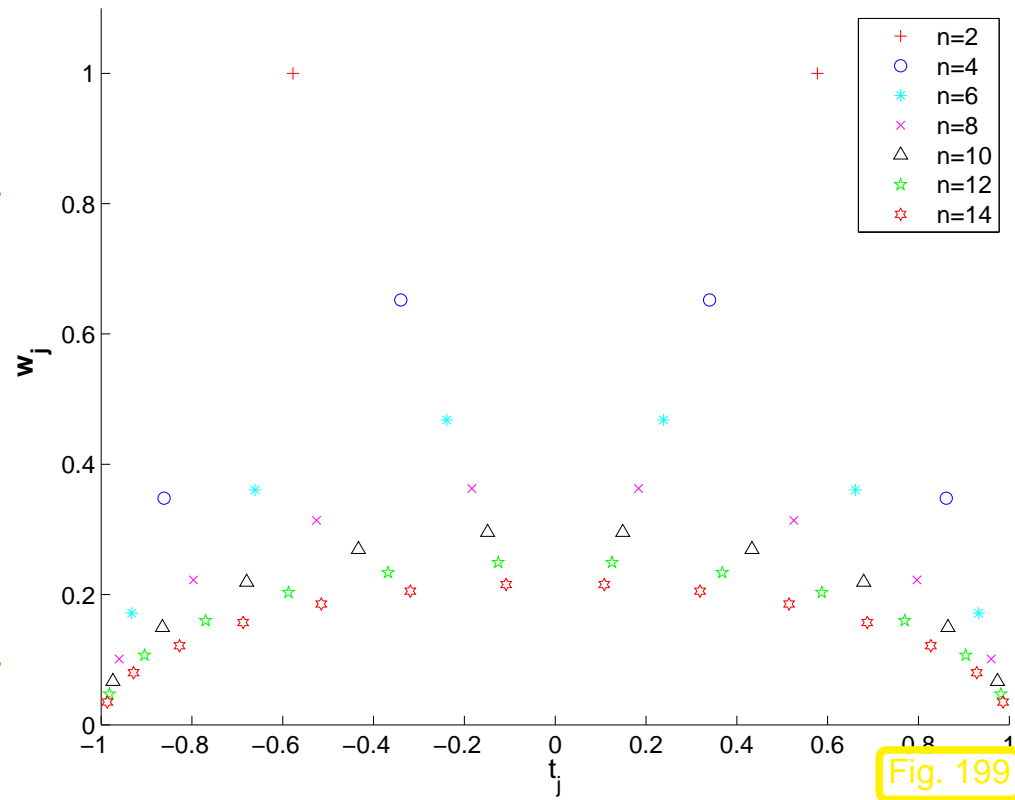
Quadrature formula from Thm. 10.4.8: **Gauss-Legendre quadrature**
(nodes ξ_j^n = Gauss points)

Obviously ▷

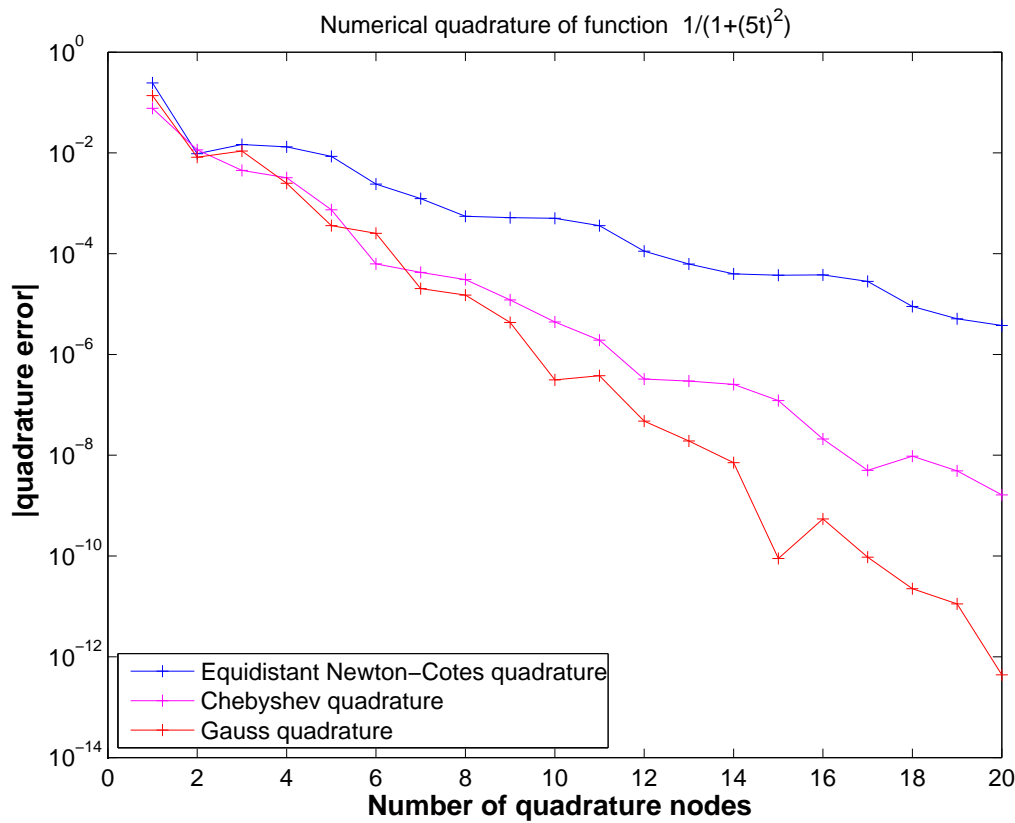
Lemma 10.4.14. *(Positivity of Gauss-Legendre quadrature weights)*

*The weights of
Gauss-Legendre quadrature formulas
are positive.*

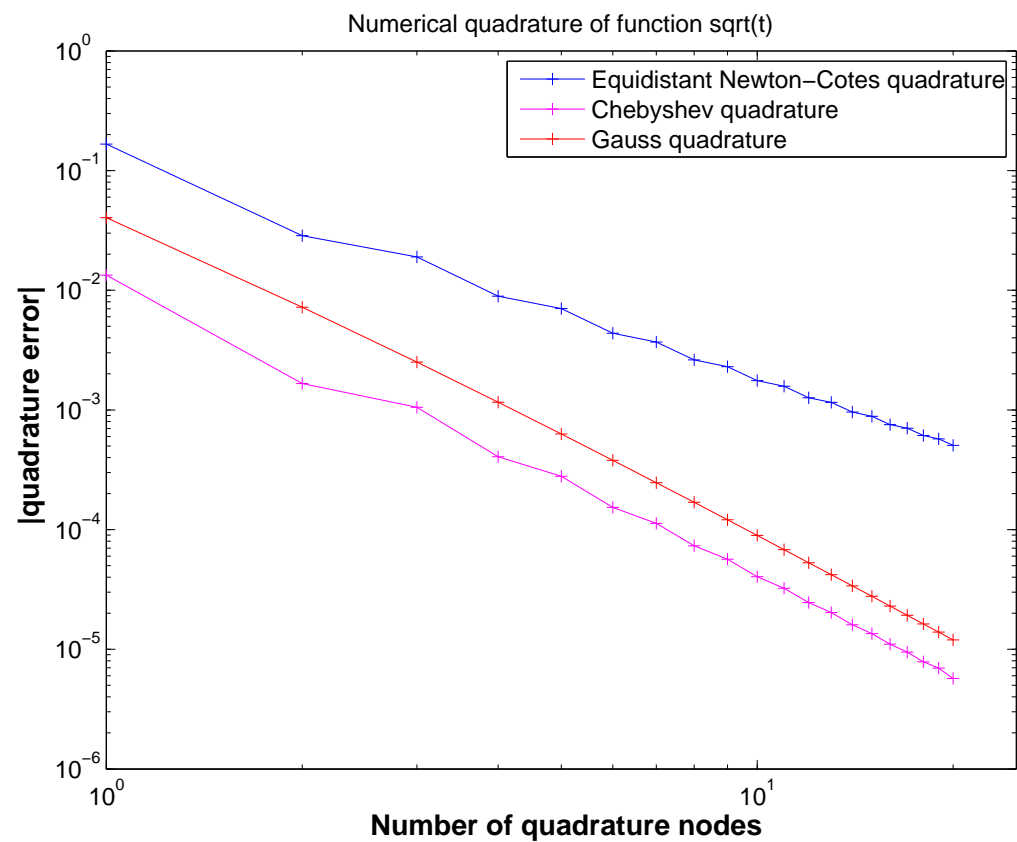
Gauss-Legendre weights for $[-1,1]$



Example 10.4.22 (Error of (non-composite) quadratures).



quadrature error, $f_1(t) := \frac{1}{1+(5t)^2}$ on $[0, 1]$



quadrature error, $f_2(t) := \sqrt{t}$ on $[0, 1]$



10.5 Adaptive Quadrature

Example 10.5.1 (Rationale for adaptive quadrature).

Consider composite trapezoidal rule (10.3.3) on mesh $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$:

Local quadrature error (for $f \in C^2([a, b])$):

$$\int_{x_{k-1}}^{x_k} f(t) dt - \frac{1}{2}(f(x_{k-1}) + f(x_k))$$

$$\leq (x_k - x_{k-1})^3 \|f''\|_{L^\infty([x_{k-1}, x_k])} \cdot$$

➤ Do not use equidistant mesh !

Refine \mathcal{M} , where $|f''|$ large !

Makes sense, e.g., for “spike function”

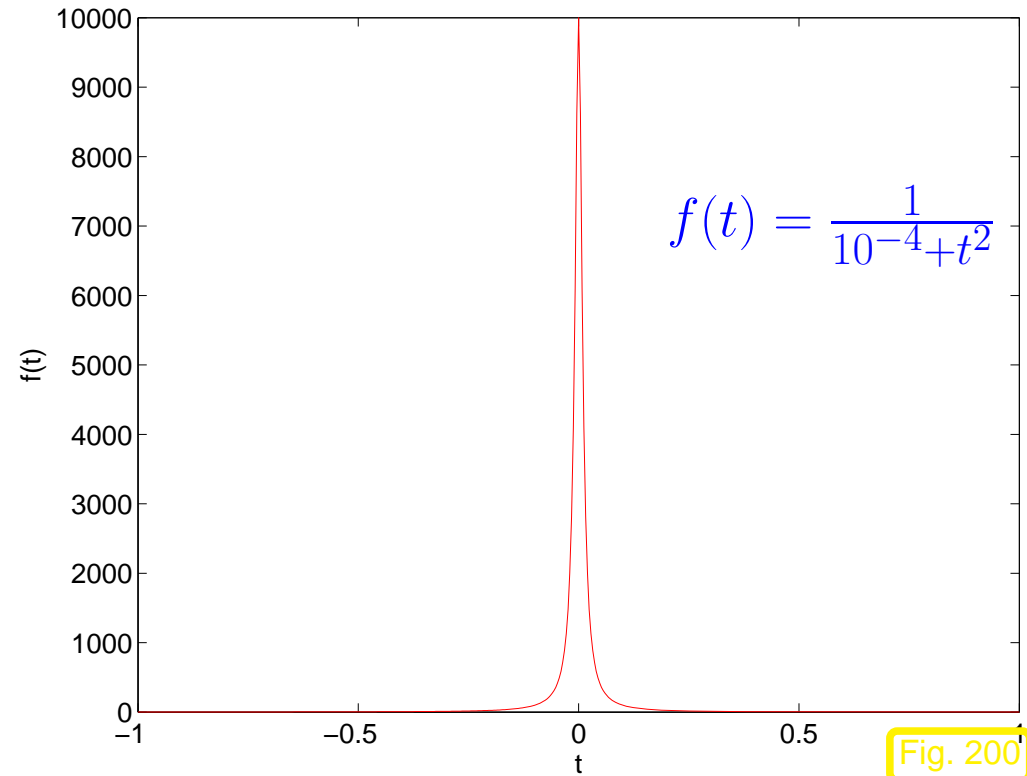


Fig. 200



Goal: *Equilibrate error contributions of all mesh intervals*

Tool: Local **a posteriori error estimation**
(Estimate contributions of mesh intervals from intermediate results)

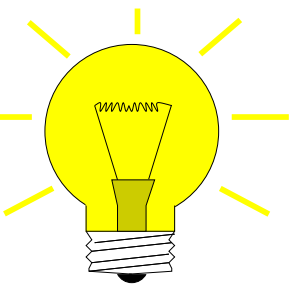
Policy: **Local mesh refinement**

Idea: local error estimation by comparing local results of two quadrature formulas

Q_1, Q_2 of *different* order \rightarrow local error estimates

heuristics: $\text{error}(Q_2) \ll \text{error}(Q_1) \Rightarrow \text{error}(Q_1) \approx Q_2(f) - Q_1(f)$.

Now: $Q_1 =$ trapezoidal rule (order 2) \leftrightarrow $Q_2 =$ Simpson rule (order 4)



Given: 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 $p_k := \frac{1}{2}(x_{k-1} + x_k)$)

$$\text{EST}_k := \left| \underbrace{\frac{h_k}{6}(f(x_{k-1}) + 4f(p_k) + f(x_k))}_{\text{Simpson rule}} - \underbrace{\frac{h_k}{4}(f(x_{k-1}) + 2f(p_k) + f(x_k))}_{\text{trapezoidal rule on split mesh interval}} \right|. \quad (10.5.2)$$

② (Termination)

Simpson rule on $\mathcal{M} \Rightarrow$ preliminary result I

$$\text{If } \sum_{k=1}^m \text{EST}_k \leq \text{RTOL} \cdot I \quad (\text{RTOL} := \text{prescribed tolerance}) \Rightarrow \text{STOP} \quad (10.5.3)$$

③ (local mesh refinement)

$$\mathcal{S} := \{k \in \{1, \dots, m\} : \text{EST}_k \geq \eta \cdot \frac{1}{m} \sum_{j=1}^m \text{EST}_j\}, \quad \eta \approx 0.9. \quad (10.5.4)$$

▶ new mesh: $\mathcal{M}^* := \mathcal{M} \cup \{p_k : k \in \mathcal{S}\}.$

Then continue with step ① and mesh $\mathcal{M} \leftarrow \mathcal{M}^*.$

Non-optimal recursive MATLAB implementation:

Code 10.5.5: h -adaptive numerical quadrature

```

1 function I = adaptquad(f,M,rtol,abstol)
2 h = diff(M); %
3 mp = 0.5*(M(1:end-1)+M(2:end)); %
4 fx = f(M); fm = f(mp); %
5 trp_loc = h.*(fx(1:end-1)+2*fm+fx(2:end))/4; %
6 simp_loc = h.*(fx(1:end-1)+4*fm+fx(2:end))/6; %
7 I = sum(simp_loc); %
8 est_loc = abs(simp_loc -trp_loc); %
9 err_tot = sum(est_loc); %
0 %
1 if ((err_tot > rtol*abs(I)) and (err_tot > abstol))

```

```

2  refcells = find(est_loc > 0.9*sum(est_loc)/length(h));
3  I = adaptquad(f, sort([M,mp(refcells)]),rtol,abstol); %
4  end

```

Example 10.5.6 (*h*-adaptive numerical quadrature).

• approximate $\int_0^1 \exp(6 \sin(2\pi t)) dt$, initial mesh $\mathcal{M}_0 = \{j/10\}_{j=0}^{10}$

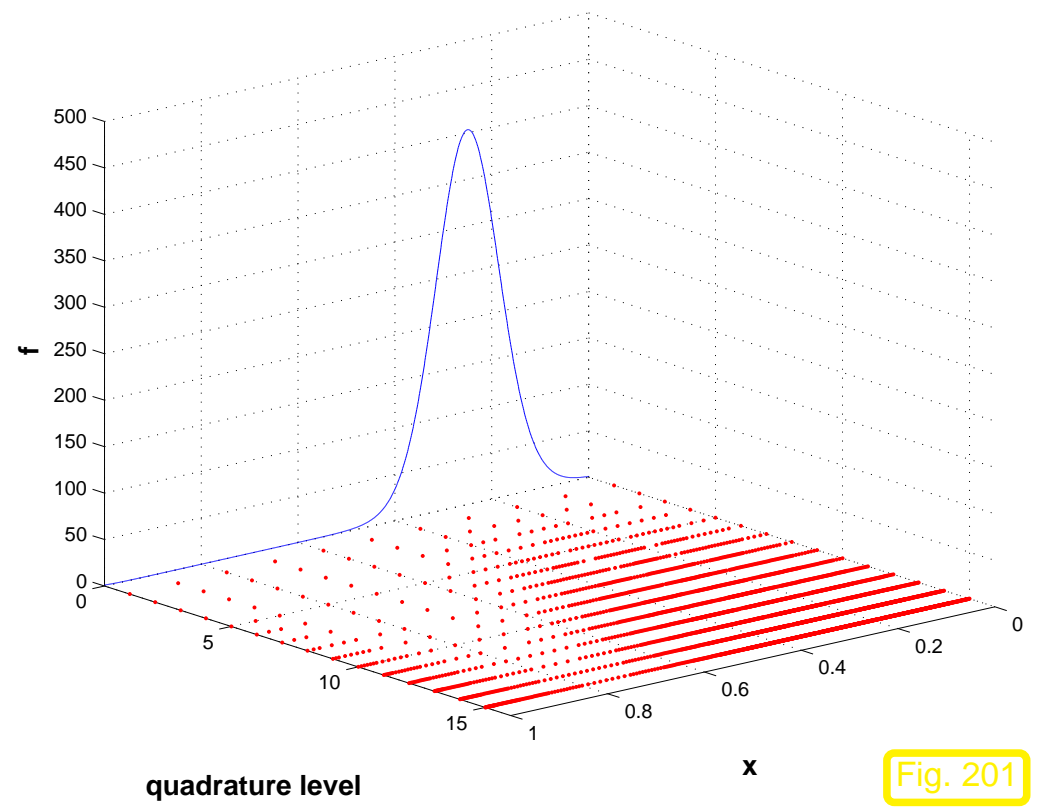


Fig. 201

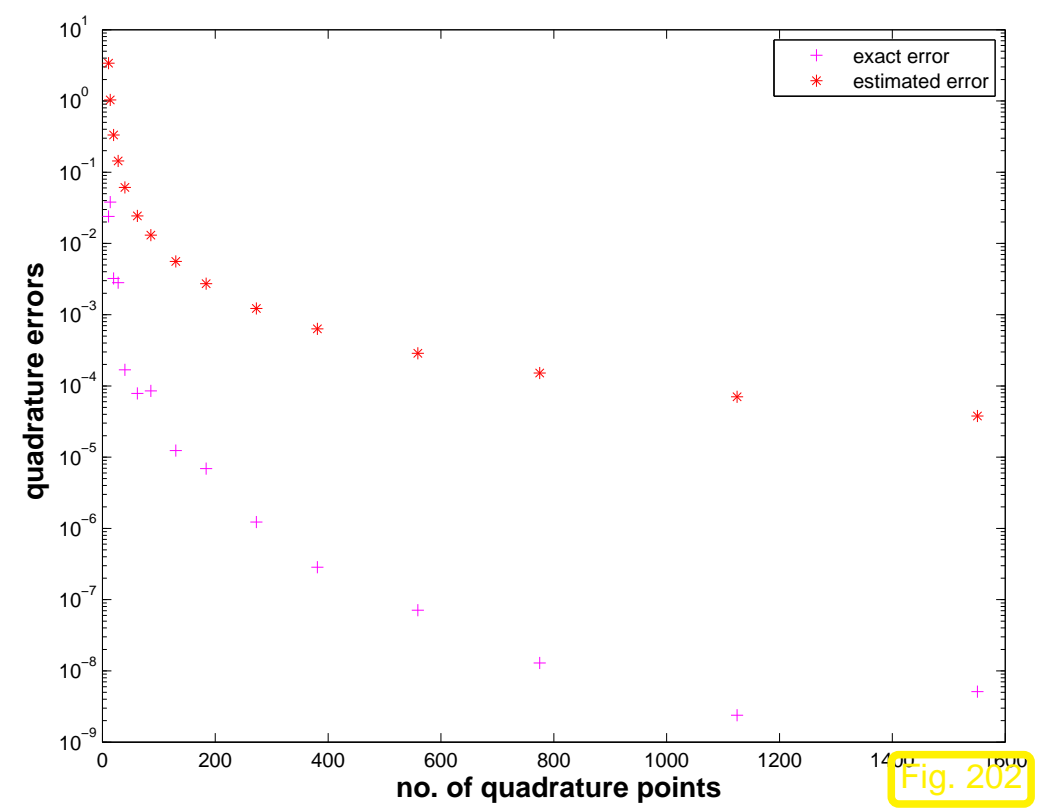


Fig. 202

approximate $\int_0^1 \min\{\exp(6 \sin(2\pi t)), 100\} dt$, initial mesh as above

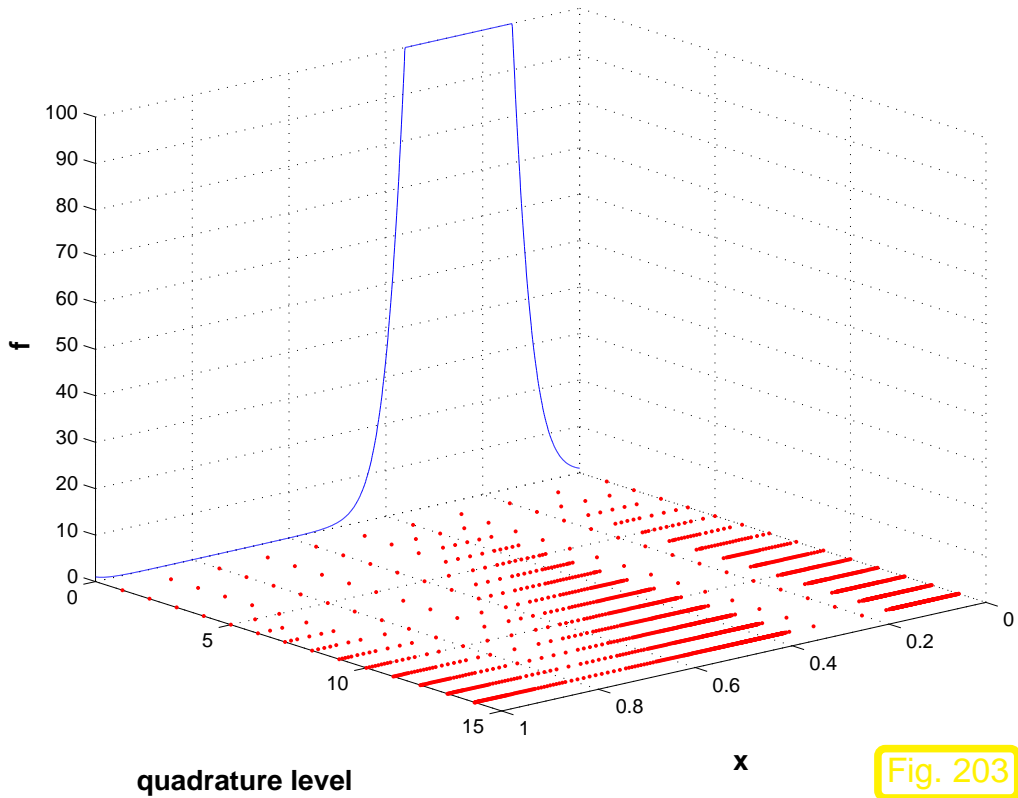


Fig. 203

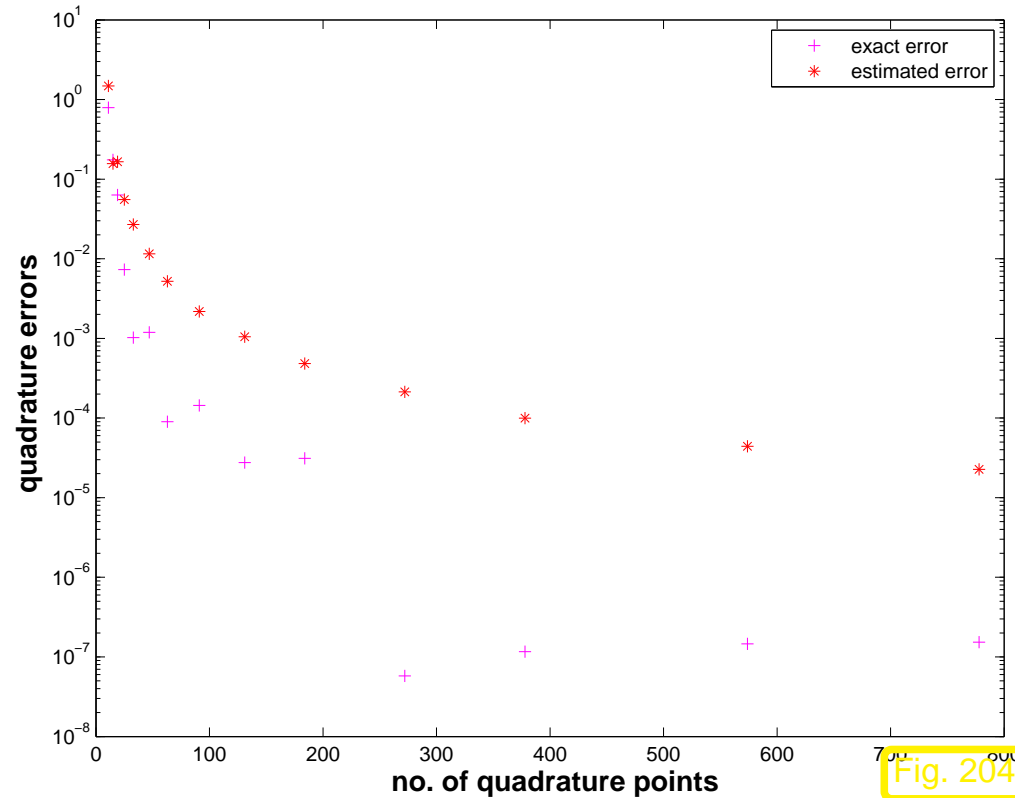


Fig. 204



Remark 10.5.7 (Adaptive quadrature in MATLAB).

`q = quad(fun, a, b, tol)`: adaptive multigrid quadrature
(local low order quadrature formulas)

`q = quadl(fun, a, b, tol)`: adaptive Gauss-Lobatto quadrature



11

Clustering Techniques

11.1 Kernel Matrices

TASK:

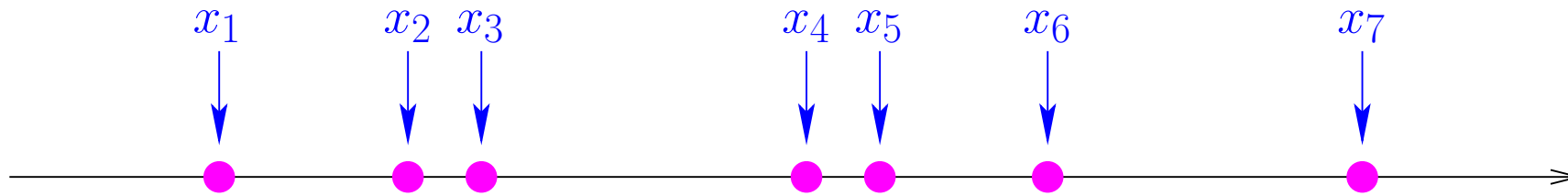
Given: **Kernel Function** $G : I \times J \mapsto \mathbb{C}$, $I, J \subset \mathbb{R}$ Interval: $G(x, y)$ smooth for $x \neq y$ Collocation Points $x_1 < x_2 < \dots < x_n$, $x_j \in I$, $y_1 < y_2 < \dots < y_m$, $y_j \in J$

Collocation Matrix: $\mathbf{M} \in \mathbb{C}^{n,m} \quad :\Leftrightarrow \quad (\mathbf{M})_{ij} := G(x_i, y_j), \quad 1 \leq i \leq n, 1 \leq j \leq m.$

(11.1.1)

We have to find: *Efficient* algorithms for **approximate** evaluation of $\mathbf{M} \times \text{Vector}$ (Note: Computational Effort $O(mn)$!)

Example 11.1.2 (Interaction calculations for many body systems).



n parallel wires with current flowing through them.

Wire j has current $c_j \in \mathbb{R}$, and is at position $x_j \in \mathbb{R}$

Our Aim : To compute magnetic force on *each* wire

- Force on wire j due to all wires: $f_j = \sum_{\substack{k=1 \\ k \neq j}}^n \frac{1}{|x_j - x_k|} c_k c_j$, $j = 1, \dots, n$.

- Force on every wire is given by vector $\mathbf{f} = \text{diag}(c_1, \dots, c_n) \mathbf{M} \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix}$,

$$\text{where } \mathbf{M} = (m_{ij})_{i,j=1}^n, \quad m_{ij} = \begin{cases} \frac{1}{|x_j - x_i|} & \text{for } i \neq j, \\ 0 & \text{for } i = j. \end{cases}$$

Collocation matrix \mathbf{M} will be formed using **kernel function** $G(x, y) = \frac{1}{|x - y|}$

Example 11.1.3 (Gravitational forces in galaxy).

Number of stars in galaxy n ($\approx 10^9$) with position $\mathbf{x}_i \in \mathbb{R}^3$ and mass m_i , $i = 1, \dots, n$.

Gravitational force on each star is required for simulation of dynamics of galaxy

Our Aim : To compute gravitational force on each star

- Gravitational force on star j :

$$f_j = \frac{G}{4\pi} \sum_{\substack{i \neq j \\ j \in \{1, \dots, n\}}} \frac{1}{\|\mathbf{x}_i - \mathbf{y}_j\|} m_i m_j ,$$



R. Hiptmair
rev 38355,
August 18,
2011

- Gravitational force on every star is given by

$$\mathbf{f} = \text{diag}(m_1, \dots, m_n) \mathbf{M} \begin{pmatrix} m_1 \\ \vdots \\ m_n \end{pmatrix} , \quad m_{ij} := \begin{cases} \frac{G}{4\pi} \frac{1}{\|\mathbf{x}_i - \mathbf{y}_j\|} & \text{for } i \neq j , \\ 0 & \text{for } i = j , \end{cases} \quad 1 \leq i, j \leq n .$$

The above example is a 3D generalization of our original task.

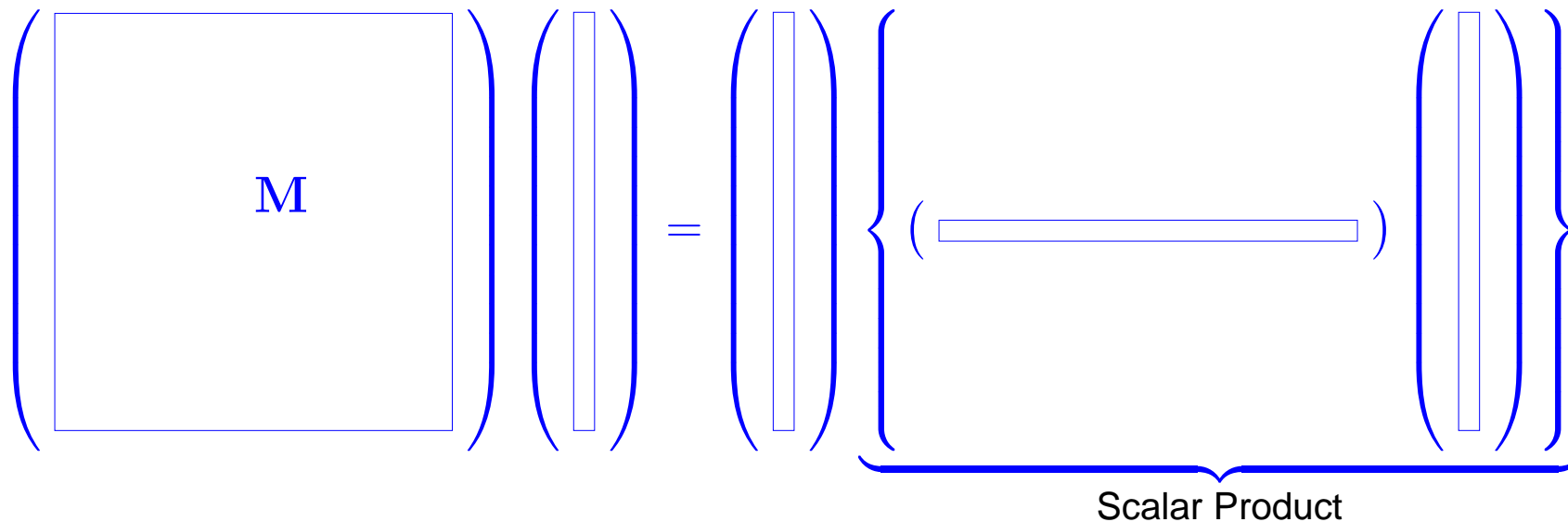


11.2 Local Separable Approximation

If kernel function is **separable** i.e. : $G(x, y) = g(x)h(y)$, $g : I \mapsto \mathbb{C}$, $h : J \mapsto \mathbb{C}$

$$\mathbf{M} = \left(g(x_j) \right)_{j=1}^n \cdot \left(h(y_j) \right)_{j=1}^m {}^T \quad \text{rank}(\mathbf{M}) = 1 .$$

Computational Effort ($\mathbf{M} \times \text{Vector}$) = $m + n$



Generalization: $G(x, y) = \sum_{j=1}^q g_j(x)h_j(y)$, $g_j : I \mapsto \mathbb{C}$, $h_j : J \mapsto \mathbb{C}$, $q \in \mathbb{N}$

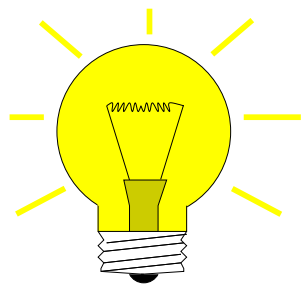
$$\mathbf{M} = \mathbf{U}\mathbf{V}^T, \quad \begin{array}{l} \mathbf{U} \in \mathbb{R}^{n,q}, \quad u_{ij} = g_j(x_i), \quad j \in \{1, \dots, q\}, i \in \{1, \dots, n\}, \\ \mathbf{V} \in \mathbb{R}^{q,m}, \quad v_{ij} = h_i(y_j), \quad i \in \{1, \dots, q\}, j \in \{1, \dots, m\}. \end{array}$$

If $\text{rank}(\mathbf{M}) = q$ then Computational Effort ($\mathbf{M} \times \text{Vector}$) = $q(m + n)$

$$\left(\begin{array}{|c|} \hline \mathbf{M} \\ \hline \end{array} \right) \left(\begin{array}{|c|} \hline \\ \hline \end{array} \right) = \left(\begin{array}{|c|} \hline \mathbf{U} \\ \hline \end{array} \right) \underbrace{\left\{ \left(\begin{array}{|c|} \hline \mathbf{V}^T \\ \hline \end{array} \right) \left(\begin{array}{|c|} \hline \\ \hline \end{array} \right) \right\}}_{q \text{ Scalar Product}}$$

Idea: Global approximation of G through sum of separated kernel function:

$$G(x, y) \approx \tilde{G}(x, y) := \sum_{l=0}^d \sum_{k=0}^d \kappa_{l,k} g_l(x) h_k(y), \quad \kappa_{l,j} \in \mathbb{C}, \quad (x, y) \in I \times J. \quad (11.2.1)$$



As G is approximated by \tilde{G} therefore \mathbf{M} is approximated by $\tilde{\mathbf{M}}$, where $\tilde{m}_{ij} = \tilde{G}(x_i, y_j)$

Remark 11.2.2 (Quality measure for kernel approximation).

$$\|\mathbf{M} - \widetilde{\mathbf{M}}\|_2^2 \leq \|\mathbf{M} - \widetilde{\mathbf{M}}\|_F^2 = \sum_{i,j} (m_{ij} - \tilde{m}_{ij})^2 = \sum_{i,j} (G(x_i, y_j) - \tilde{G}(x_i, y_j))^2 .$$

Normalizing quality measure:
$$\frac{1}{mn} \sum_{i,j} (G(x_i, y_j) - \tilde{G}(x_i, y_j))^2 . \quad (11.2.3)$$

△

Now, our aim is to find separable kernel approximation

Definition 11.2.4 (Tensor product interpolation polynomial). $L_j^x \in \mathcal{P}_n$ ($L_k^y \in \mathcal{P}_m$), $j = 0, \dots, n$ ($k = 0, \dots, m$), Lagrange polynomial on nodes of mesh $\mathcal{X} := \{x_j\}_{j=0}^n \subset I$ ($\mathcal{Y} := \{y_k\}_{k=0}^m \subset J$), $I, J \subset \mathbb{R}$ interval. Continuous $f : I \times J \mapsto \mathbb{C}$ defined by

$$(l_{\mathcal{X} \times \mathcal{Y}} f)(x, y) := \sum_{j=0}^n \sum_{k=0}^m f(x_j, y_k) L_j^x(x) L_k^y(y), \quad x, y \in \mathbb{R},$$

This is *tensor product interpolation polynomial*.

$$G(x, y) \approx \tilde{G}(x, y) := \sum_{j=0}^d \sum_{k=0}^d G(t_j^x, t_k^y) L_j^x(x) L_k^y(y), \quad (11.2.5)$$

$t_0^x, \dots, t_d^x / t_0^y, \dots, t_d^y$ Chebyshev-nodes in I/J , and L_j^x, L_k^y subordinate Lagrange-polynomial.

Example 11.2.6 (Global separable approximation by smooth kernel function).

- Smooth (even analytic) kernel function $G(x, y) = \frac{1}{1 + |x - y|^2}$, collocation points $x_i = y_i = \frac{i-1}{n}, i = 1, \dots, n$.
- \tilde{G} approximated by tensor product Chebyshev interpolation polynomial (11.2.10) of degree $d \in \mathbb{N}$.

$$\frac{1}{n^2} \sum_{i,j} (G(x_i, y_j) - \tilde{G}(x_i, y_j))^2$$

as function of d for $n \in \{100, 200, 400\}$

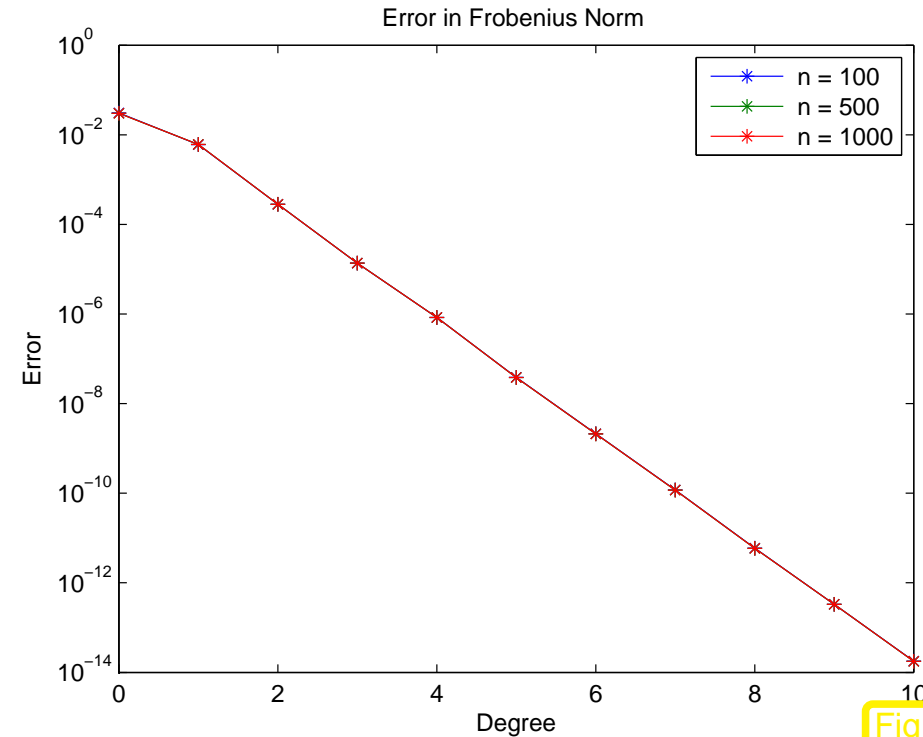


Fig. 205

In file `changeable.c`

- put the admissibility condition as return `true` in `admissible_ClusterPair`
- put kernel function as `1.0/(1.0 + pow(fabs(dX - dY),2))` in `kernel_function`
- put `NONE` as return values in `get_iOperation`
 - Run `main.cpp`, input as follows
Choice - 2, NumberofPts - 100 1000 0,
Number of Increments - 2, Increment 1 - 400, Increment 2 - 500
Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1
Degree - 0 10 1, Admissibility Coefficient - 0 0 1
Output File from C++: `Error.txt`
 - OutputFile from MATLAB : `Error_in_Frobenius_Norm_Pts.eps`
 - Run `error_plot_NumberofPts.m` from Matlab

Note: Exponential Convergence $\|M - \tilde{M}\|_F \rightarrow 0$ in dependence of d



Example 11.2.7 (Global separable approximation by non-smooth kernel function).

Analysis as in ex. 11.2.6, now for the kernel function

$$G(x, y) = \begin{cases} \frac{1}{|x-y|} & \text{if } x \neq y, \\ 0 & \text{otherwise.} \end{cases} \quad (11.2.8)$$

from ex. 11.1.2.

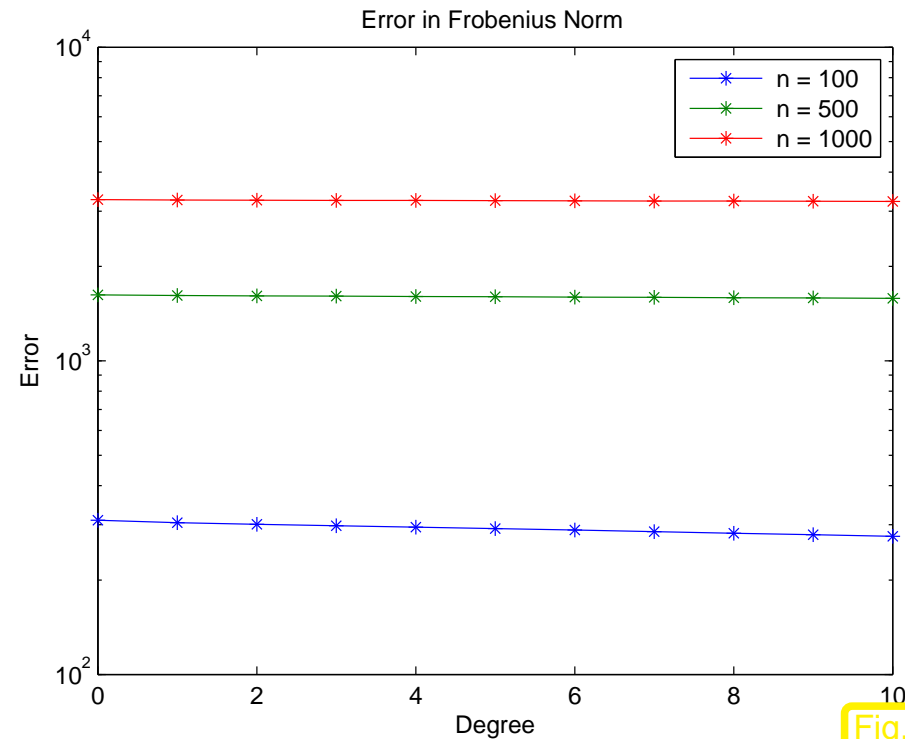


Fig. 206

In file changeable.c

- a) put the admissibility condition as return **true** in `admissible_ClusterPair`
- b) put kernel function as **0.0** for `fabs(dX-dY) < EPS` as **1.0/fabs(dX - dY)** otherwise in `kernel_function`
- c) put **NONE** as return values in `get_iOperation`
 - 1) Run `main.cpp`, input as follows
 - Choice - 2, NumberofPts - 100 1000 0,
 - Number of Increments - 2, Increment 1 - 400, Increment 2 - 500
 - Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1
 - Degree - 0 10 1, Admissibility Coefficient - 0 0 1
 - Output File from C++: `Error.txt`
 - OutputFile from MATLAB : `Error_in_Frobenius_Norm_Pts.eps`
 - 2) Run `error_plot_NumberofPts.m` from Matlab



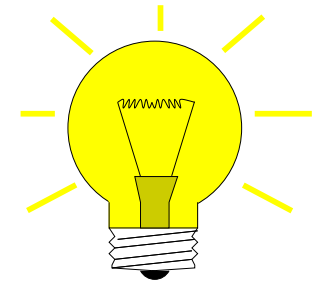
Note: (Virtually) No Convergence $\left\| \mathbf{M} - \widetilde{\mathbf{M}} \right\|_F \rightarrow 0$ for $d \rightarrow \infty$

Reason: Missing *global* smoothness

Poor approximation of $G(x, y) := |x - y|^{-1}$ in region of $\{(x, y) \in I \times J : x = y\}$.

However $G(x, y)$ from (11.2.8) is smooth (even analytic, Def. 9.2.20) at “large distances” from $\{(x, y) \in I \times J : x = y\}$.

Idea: **Local** approximation of G through sum of separated kernel function




$$G(x, y) \approx \sum_{l=0}^d \sum_{k=0}^d \kappa_{l,k} g_l(x) h_k(y), \quad \kappa_{l,j} \in \mathbb{C}, \tag{11.2.9}$$

$$(x, y) \in \widetilde{I} \times \widetilde{Y}, \quad \widetilde{I} \subset I, \quad \widetilde{J} \subset J,$$

with $\widetilde{I} \times \widetilde{J} \cap \{(x, y) : x = y\} = \emptyset$.

Actually:

Local approximation of G on rectangle $\tilde{I} \times \tilde{J}$ which is a *partition* of $I \times J$

 Possible partition for separable kernel approximation through local tensor product Chebyshev polynomial interpolation (11.2.9)
(Admissible rectangles)

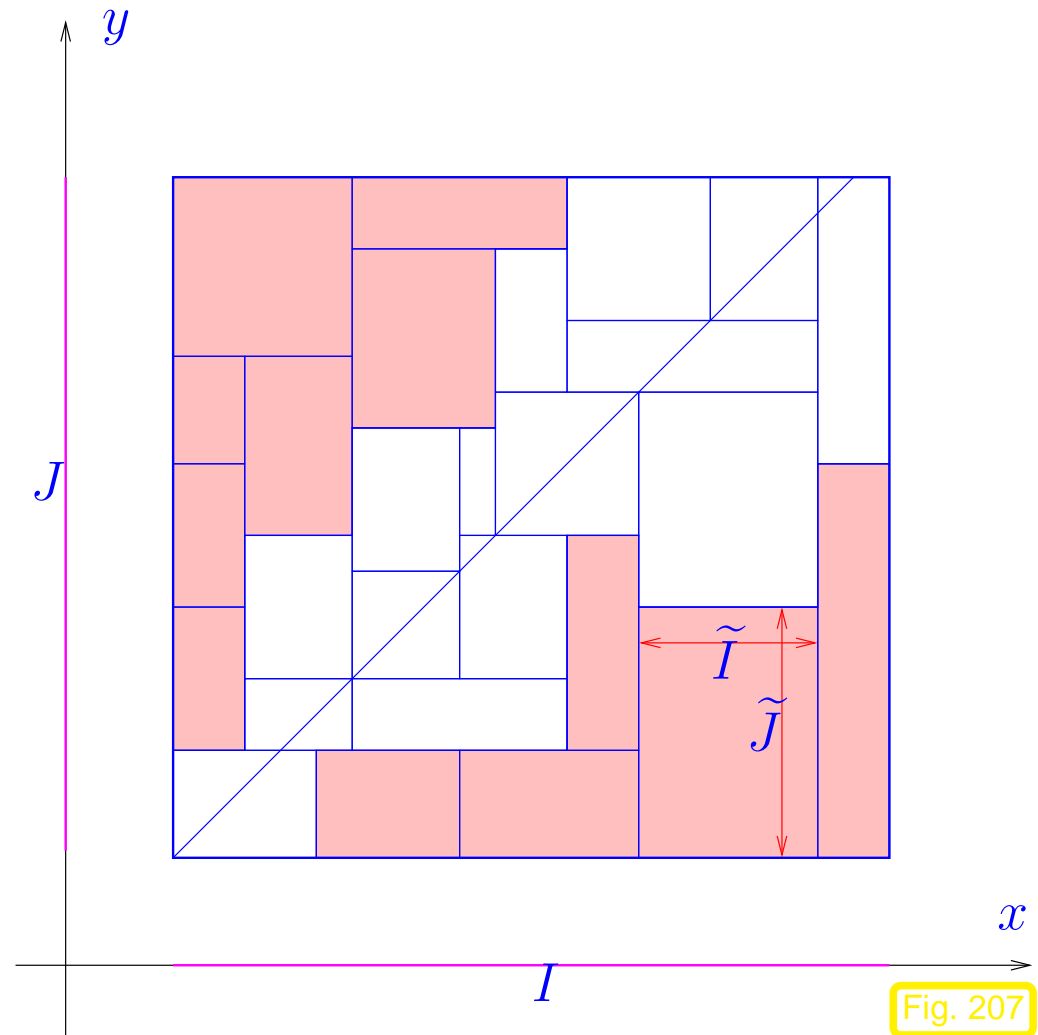
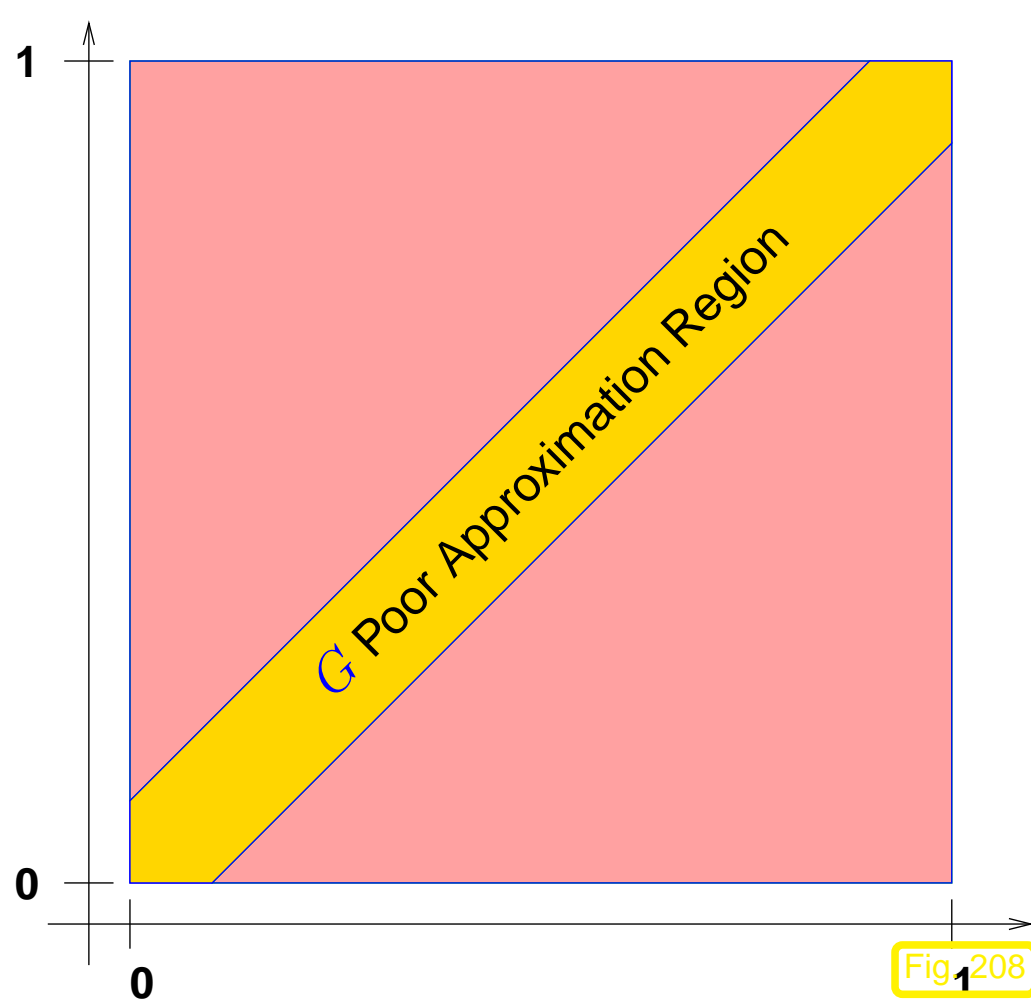


Fig. 207



Terminology:

 **Near Field:** G Poor approximation

 **Far Field:** G Good approximation

Fig 1 208

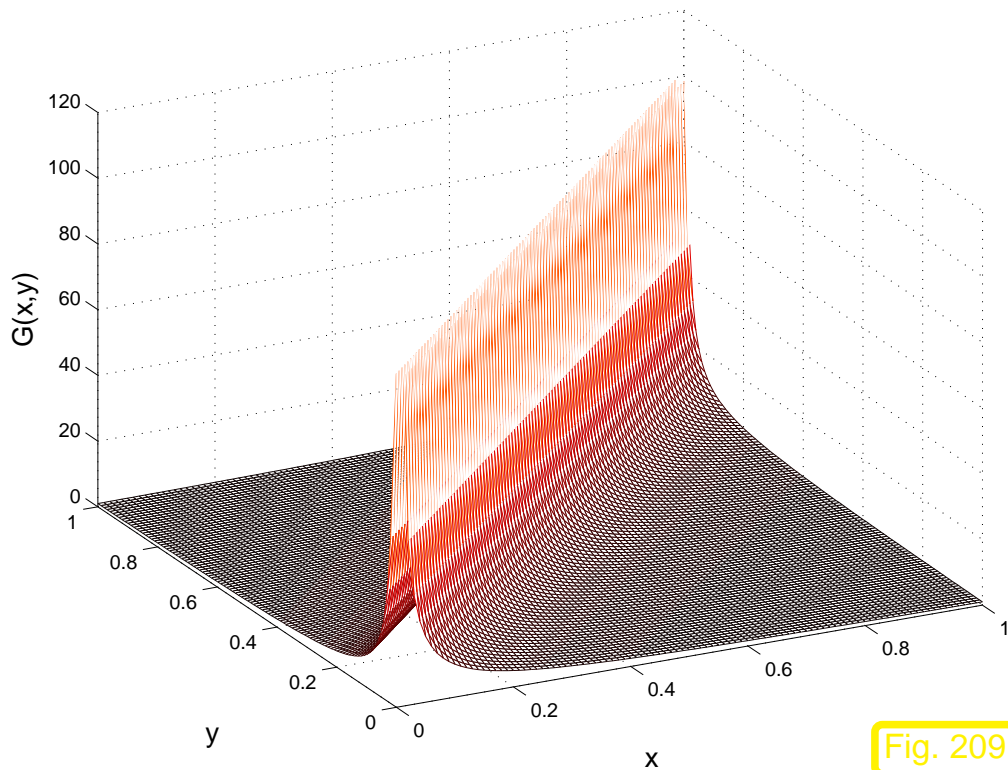


Fig. 209

Kernel function (11.2.8) for $[0, 1]^2$

$$G(x, y) = \frac{1}{|x - y|}$$

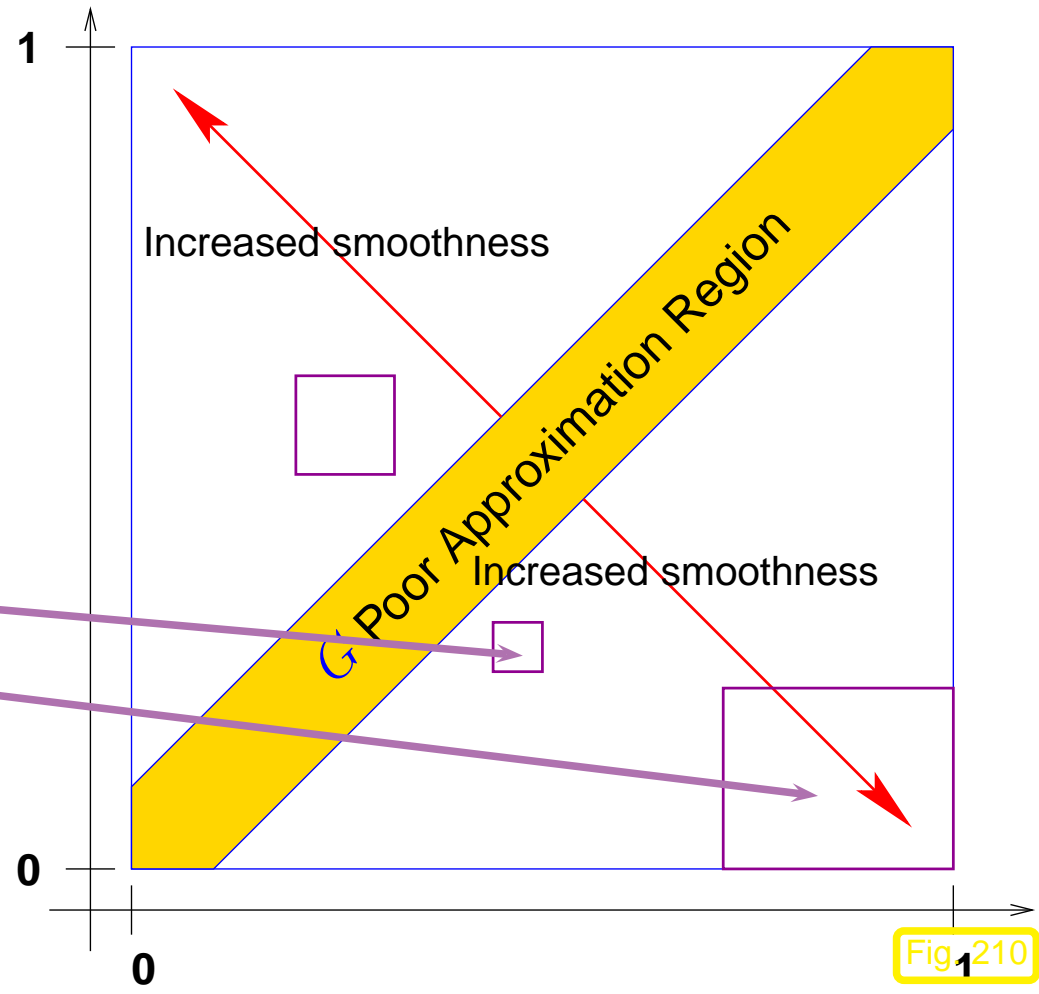
- Singular for $x = y$
- Analytic far from $\{(x, y): x = y\}$, analyticity (\rightarrow Chapter ??) increases with $|x - y|$

Similar kernel function: $G(x, y) = \log |x - y|$, $G(x, y) = \frac{\partial}{\partial x} \frac{1}{|x - y|}$, etc.

“Smoothness of $G(x, y) = \frac{1}{|x-y|}$ increases with growing distance from diagonal $\{x = y\}$ ”:

- No approximation, when $|x - y|$ is “small”
- Tensor product interpolation polynomial for
 - Small rectangles near diagonal
 - Large rectangles far from diagonal

Now, our aim is to find appropriate sizes of above mentioned rectangles



Example 11.2.10 (Tensor product Chebyshev interpolation on rectangles).

$I = J = [0, 1]^2$, $G(x, y) = |x - y|^{-1}$ from (11.2.8), (Approximate) Max norm of the error of local tensor product Chebyshev interpolation polynomial of G on rectangles of constant size, but with growing distance from $\{(x, y): x = y\}$

$$\tilde{I}_k = [0.55 + k \cdot 0.05, 0.75 + k \cdot 0.05] \quad , \quad \tilde{J}_k = [0.25 - k \cdot 0.05, 0.45 - k \cdot 0.05] \quad , \quad k \in \{0, \dots, 5\} .$$

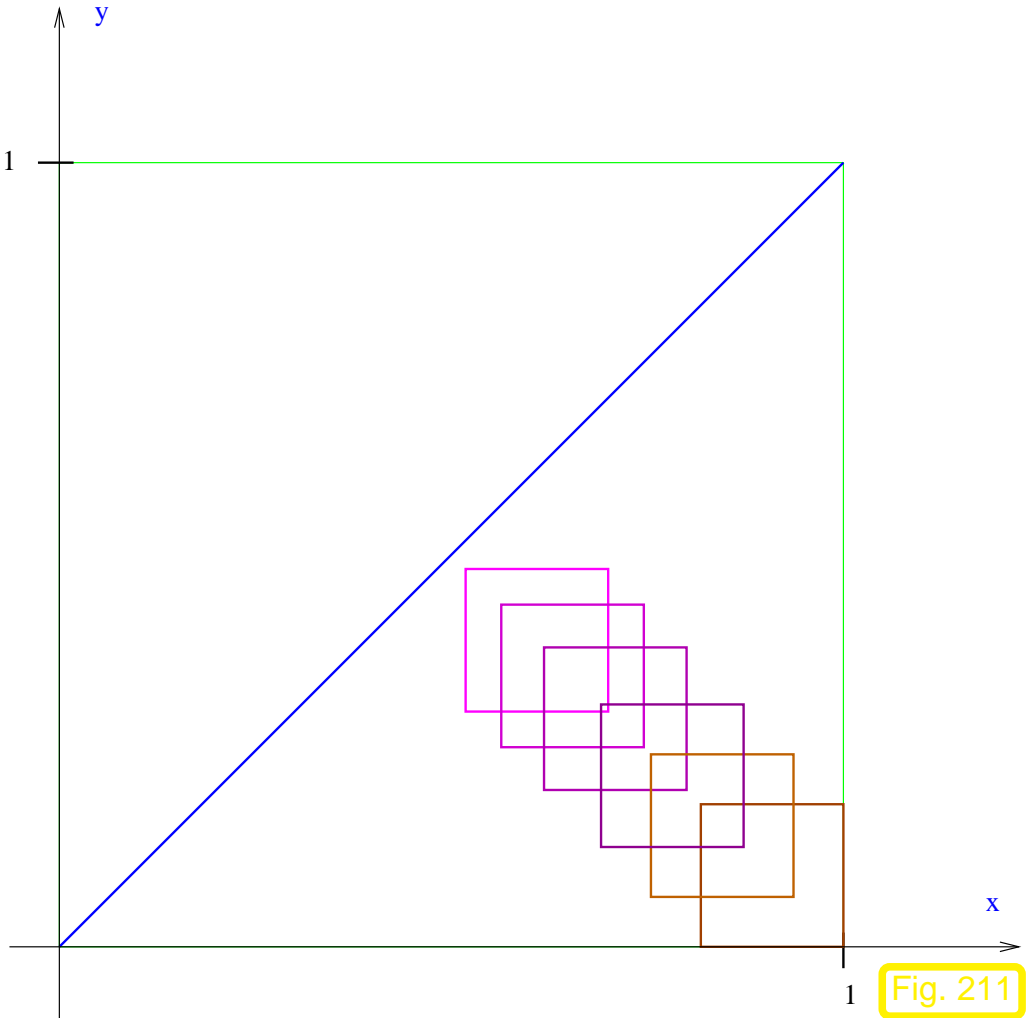


Fig. 211

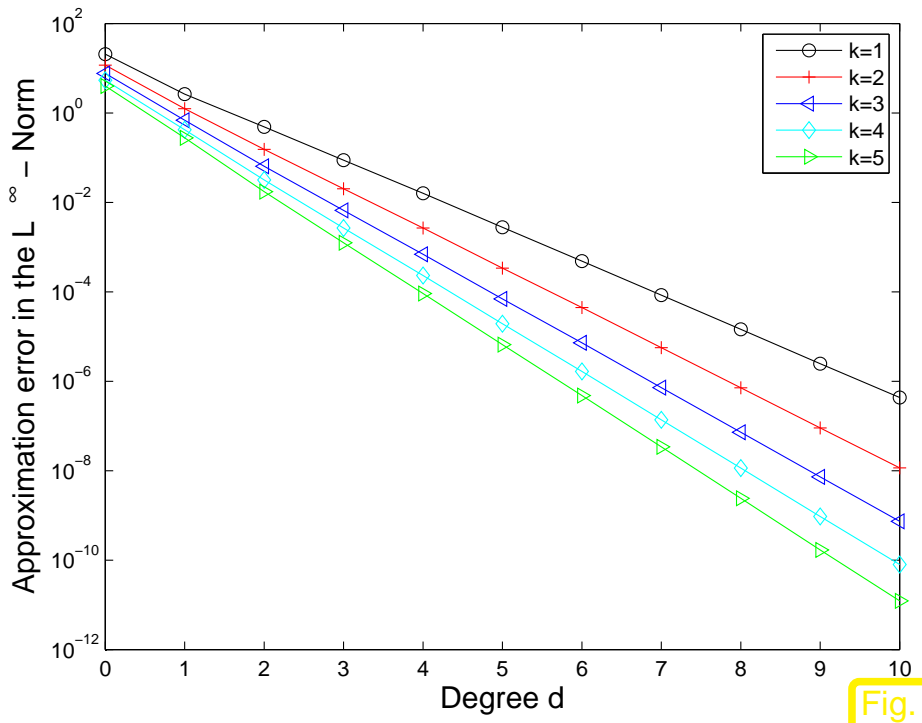


Fig. 212

NOTE: Decreasing interpolation errors with increasing distance from $\{(x, y): x = y\}$

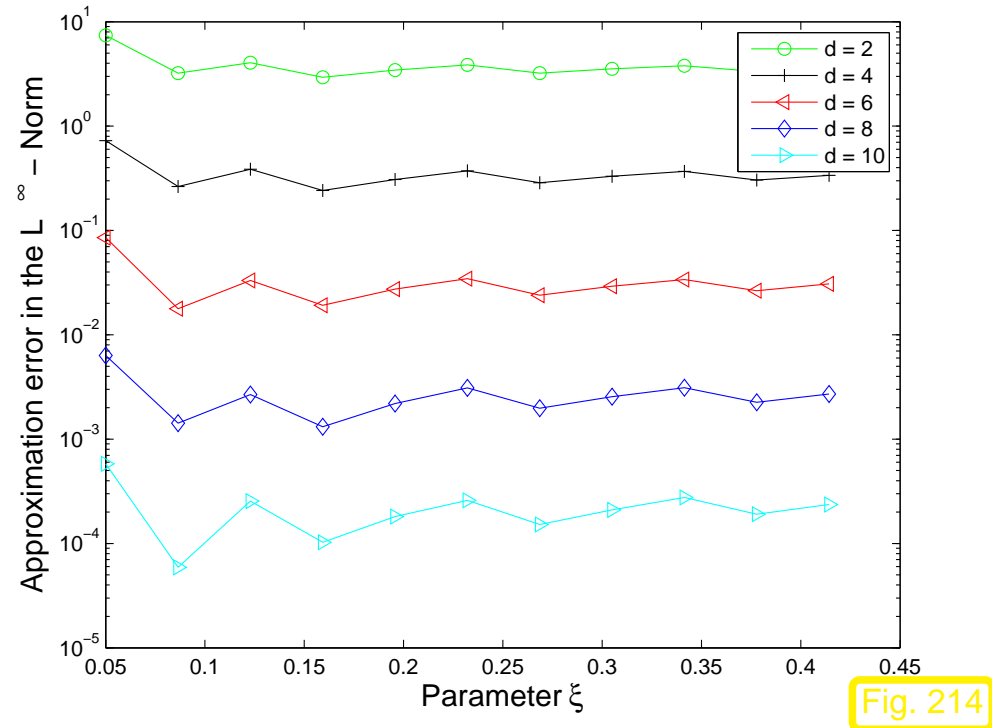
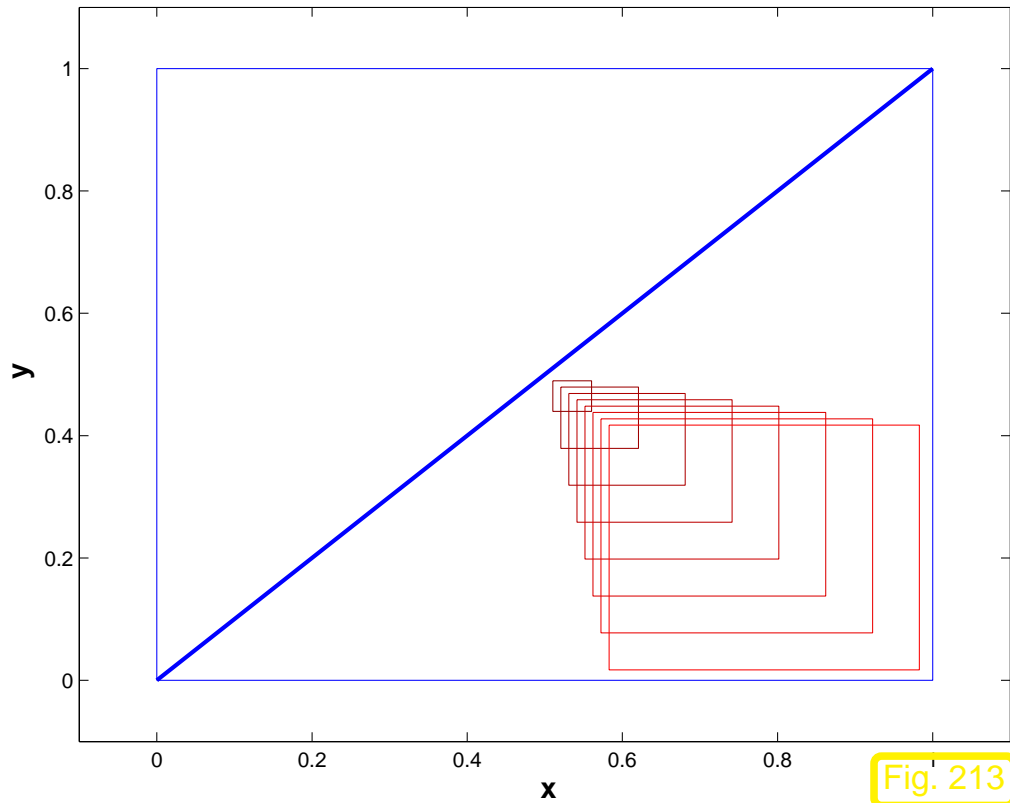


Example 11.2.11 (Tensor product Chebyshev interpolation for variable rectangle sizes).

Taking the variable rectangle sizes as

$$\left[\frac{1}{2}(\sqrt{2} - 1)\xi + \frac{1}{2}, \frac{1}{2}(\sqrt{2} + 1)\xi + \frac{1}{2}\right] \times \left[-\frac{1}{2}(\sqrt{2} - 1)\xi + \frac{1}{2}, -\frac{1}{2}(\sqrt{2} + 1)\xi + \frac{1}{2}\right], \quad 0.05 \leq \xi \leq \frac{1}{1 + \sqrt{2}}$$

that is, Size of rectangles \sim Distance from diagonal



NOTE: Controlled decrease of interpolation errors with increase in size of rectangles



Acceptability Criteria

$[a, b] \times [c, d]$ is called η -admissible, $\eta > 0$, if

$$\eta \text{dist}([a, b], [c, d]) \geq \max\{b - a, d - c\} . \quad (11.2.12)$$

($\text{dist}([a, b], [c, d]) := \min\{|x - y| : x \in [a, b], y \in [c, d]\}$)

11.3 Cluster Trees

TASK: Given $\eta > 0$, find (efficient algorithms) partitioning of $I \times J$ "far from diagonal" which are η -admissible rectangles.

Remark 11.3.1. Perspective:

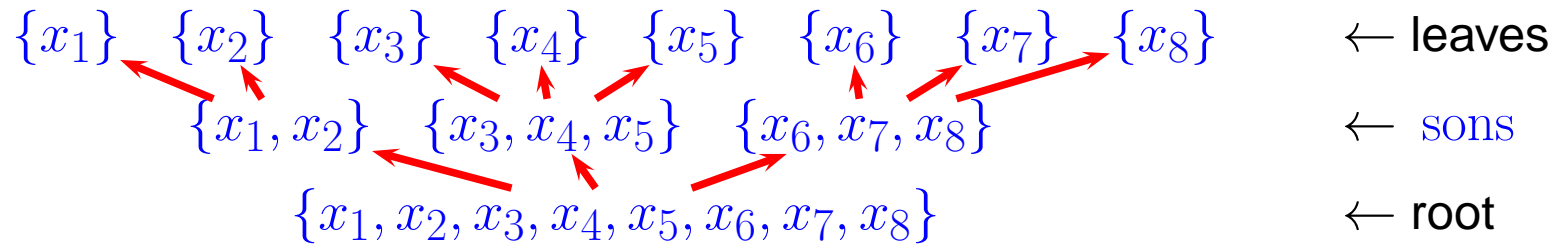
Partitioning of $I \times J \implies$ Partitioning of $\{x_i\}_{i=1}^n \times \{y_j\}_{j=1}^m \iff$ Partitioning of index mesh $\{1, \dots, n\} \times \{1, \dots, m\}$



Definition 11.3.2 (Cluster Tree). A Tree (\rightarrow Computer Science, Graph theory) T is called **Cluster Tree** on $\mathbb{P} := \{x_1, \dots, x_n\} \subset \mathbb{R} \iff$

- The nodes of the Tree T (= **Cluster**) are subset of \mathbb{P} .
- $\text{root}(T) = \mathbb{P}$.
- For every node σ of T : $\{\sigma' : \sigma' \in \text{sons}(\sigma)\}$ is Partitioning of σ .

Bounding Box of Cluster $\sigma \in T$: $\Gamma(\sigma) := [\min \{x\}_{x \in \sigma}, \max \{x\}_{x \in \sigma}]$



Terminology (\rightarrow Graph theory):

Cluster Tree T : $\sigma \in T: \text{sons}(\sigma) = \emptyset \iff \sigma$ leaf

Cluster Tree T :

Recursive Construction

MATLAB-Data Structure:

$N \times 6$ -Matrix, $N = \#T$,

Lines $\hat{=}$ Cluster

$T = [T; i, j, xl, xr, s1, s2]$

i, j : $\sigma = \{x_i, \dots, x_j\}$

xl : $xl = \min_{i \leq k \leq j} x_k$

xr : $xr = \max_{i \leq k \leq j} x_k$

$s1, s2$: The line indices of son

NOTE : $\#T \leq 2n$

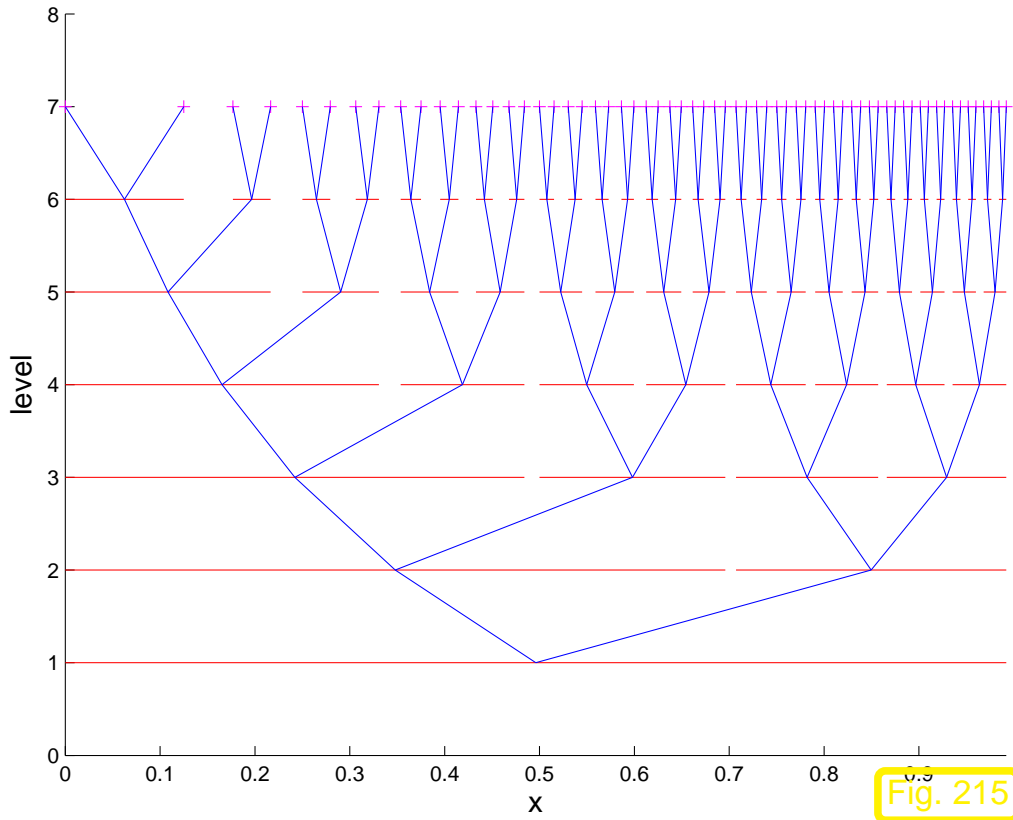
MATLAB-CODE : Construction of Cluster Trees

```
function [T,idx] = ct_rec(x,ofs,m,T)
N = length(x);
if (N <= m)
    T = [T;ofs,ofs+N-1,x(1),x(end),0,0];
else
    n = round(N/2);
    [T,s1] = ct_rec(x(1:n),ofs,m,T);
    [T,s2] = ct_rec(x(n+1:N),ofs+n,m,T);
    T = [T;ofs,ofs+N-1,x(1),x(end),s1,s2];
end
idx = size(T,1);
```

C++ code: Construction of Cluster Trees

```
1 void clsClusteringApproximation::build_ClusterTree_Linear(int iOffset, int
  iNumberOfPtsInClusterLinear, int* piDimension){
2   int iTempLeftChild, iTempRightChild, iStart, iEnd;
3   int iStartLeft, iPtsInLeft, iStartRight, iPtsInRight;
4
5   iStart = iOffset;
6   iEnd = iOffset + iNumberOfPtsInClusterLinear - 1;
7   if (iNumberOfPtsInClusterLinear == 1){
8     add_ClusterLinear(*piDimension, iStart, iEnd, -1, -1);
9   }
10  else{
11    iStartLeft = iOffset;
12    iPtsInLeft = iNumberOfPtsInClusterLinear - (iNumberOfPtsInClusterLinear/2);
13    build_ClusterTree_Linear(iStartLeft, iPtsInLeft, piDimension);
14    iTempLeftChild = vec2pclClusterLinear[*piDimension].size() - 1;
15
16    iStartRight = iOffset + iPtsInLeft;
17    iPtsInRight = iNumberOfPtsInClusterLinear/2;
18    build_ClusterTree_Linear(iStartRight, iPtsInRight, piDimension);
19    iTempRightChild = vec2pclClusterLinear[*piDimension].size() - 1;
20
21    add_ClusterLinear(*piDimension, iStart, iEnd, iTempLeftChild, iTempRightChild);
22  }
23 }
```

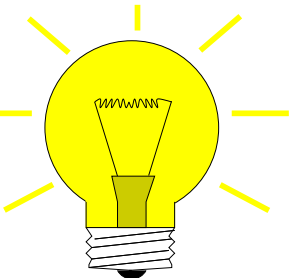
Example 11.3.3 (Cluster Tree).



← Binary cluster tree:

```
x = sqrt(0:1/64:(1-1/64));
```

- 2^{l-1} Cluster on Level l with 2^{7-l} points.
(balanced Cluster tree)



Idea: Choose approximation rectangles = Tensor product of Bounding Box of Cluster $\in T$

Cluster σ, μ admissible $:\Leftrightarrow$

$\Gamma(\sigma), \Gamma(\mu)$ admissible

Recursive construction for
admissible Partitioning:

P_{near} = **Near field**:

The area near the diagonal
(\rightarrow No approximation)

P_{far} = **Far field**:

Partitioning with η -admissible
rectangles

MATLAB-CODE : Recursive construction for admissible partitioning

```
function [Pnear,Pfar] = ...
    divide(Tx,Ty,sigma,mu,Pnear,Pfar,eta)
cls = Tx(sigma,:); clm = Ty(mu,:);
if (sigma == leaf | mu == leaf)
    Pnear = [Pnear; cls(2),cls(3),clm(2),clm(3)];
elseif ((sigma,mu) admissible)
    Pfar = [Pfar; cls(2),cls(3), clm(2),clm(3)];
else
    for s1 = cls(6:7), for s2 = clm(6:7)
        [Pnear,Pfar] = divide(Tx,Ty,s1,s2,Pnear,Pfar,eta);
    end; end
end
```

MATLAB-CODE : Parent Code

```
function [Pnear,Pfar] = partition(Tx,Ty,eta)
Pnear = []; Pfar = [];
sigma = find(Tx(:,1) == min(Tx(:,1)));
mu = find(Ty(:,1) == min(Ty(:,1)));
[Pnear,Pfar] = divide(Tx,Ty,sigma,mu,Pnear,Pfar,eta);
```

Note: clusters of a far field cluster pair have a common level!


```
1 void clsClusteringApproximation::build_ClusterPair (int iIndex1,  
2     int iIndex2){  
3     clsClusterLinear *pclsClusterLinear1, *pclsClusterLinear2;  
4  
5     pclsClusterLinear1 = vec2pclsClusterLinear[0][iIndex1];  
6     pclsClusterLinear2 = vec2pclsClusterLinear[1][iIndex2];  
7     if (leaf(0,iIndex1) || leaf(1,iIndex2)){  
8         add_ClusterPairNear(pclsClusterLinear1,  
9             pclsClusterLinear2);  
10        (pclsClusterLinear1->get_ptr_AppearsIn())->push_back(-vecp  
11        (pclsClusterLinear2->get_ptr_AppearsIn())->push_back(-vecp  
12    }  
13    else if (admissible_ClusterPair(iIndex1, iIndex2)){  
14        add_ClusterPairFar(pclsClusterLinear1,  
15            pclsClusterLinear2);  
16        (pclsClusterLinear1->get_ptr_AppearsIn())->push_back(vecpc  
17        (pclsClusterLinear2->get_ptr_AppearsIn())->push_back(vecpc  
18    }  
19    else {  
20        tree_traverse(pclsClusterLinear1, pclsClusterLinear2);  
    }
```

```
1 void clsClusteringApproximation::preprocess(){
2     int iLoopVari, iLoopVarj;
3     for (iLoopVari = 0; iLoopVari < iDimension; iLoopVari++){
4         build_ClusterTree_Linear(0, veciNumberofPts[iLoopVari],
5             &iLoopVari);
6     }
7     build_ClusterPair(vec2pclsClusterLinear[0].size() - 1,
8         vec2pclsClusterLinear[1].size() - 1);
9 }
```

Cluster Pairs Formed for Number of Pts (64,64) and Admissibility Coefficient : 0.5

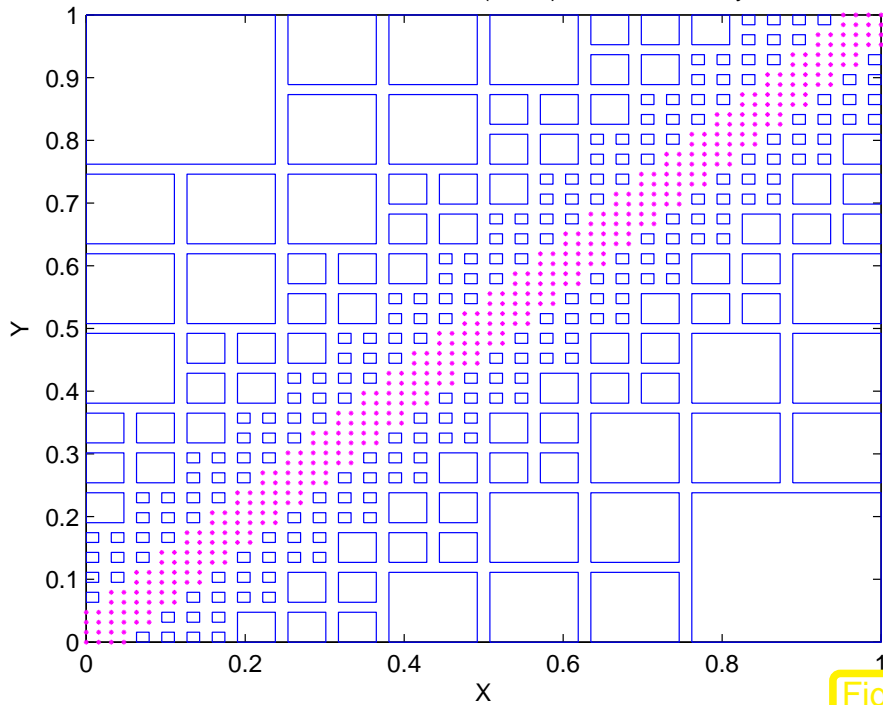


Fig. 216

$$x = (0 : 1/64 : (1 - 1/64)) ;$$

In file changeable.c

- a) put the admissibility condition in `admissible_ClusterPair`
- b) put kernel function in `kernel_function`
- c) put **NONE** as return values in `get_iOperation`

1) Run `main.cpp`, input as follows

Choice - 2, NumberofPts - 64 64 1,

Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1

Degree - 0 0 1, Admissibility Coefficient - 0.5 0.5 1

Output File from C++: `ClusterLinear1_AA.txt`, `ClusterLinear2_AA.txt`,

`ClusterFar_AA.txt`, `ClusterLinearNear_AA.txt`,

OutputFile from MATLAB : `ClusterFigure_64_64_0.5.eps`

2) Run `cluster_pair_indiv.m` from Matlab, arguments:

0 0 (argument here is for the Point Number

and Admissibility Coefficient. As 64 Number of Pts is the

first in the iteration over Number of Pts,

therefore argument is given as 0

same holds for Admissibility Coefficient)

Cluster Pairs Formed for Number of Pts (64,64) and Admissibility Coefficient : 0.5

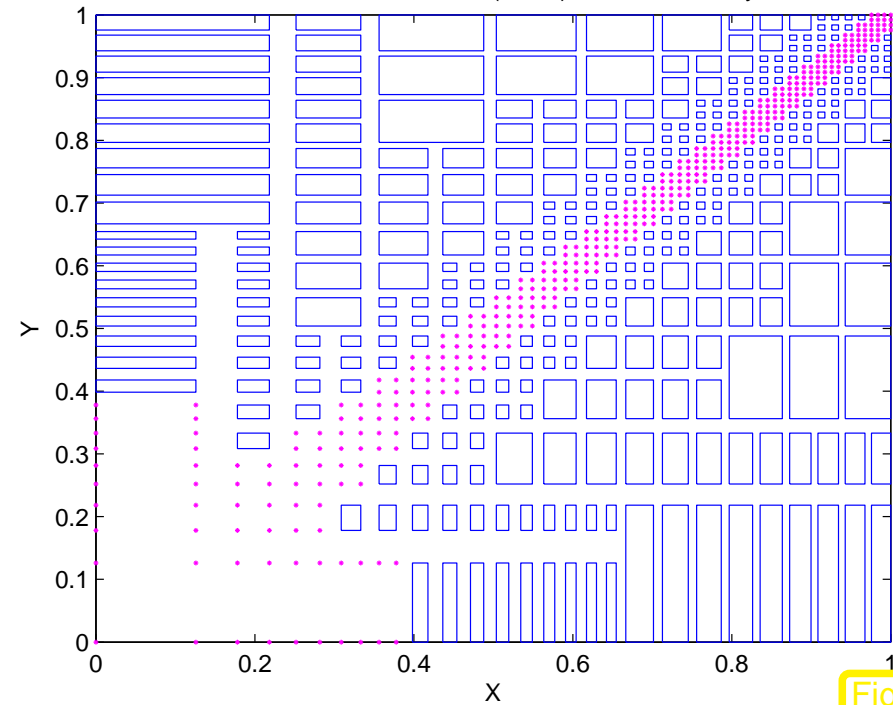


Fig. 217

$$x = \text{sqrt}(0 : 1/64 : (1 - 1/64)) ;$$

In file changeable.c

- a) put the admissibility condition in `admissible_ClusterPair`
- b) put kernel function in `kernel_function`
- c) put **SQRT** as return values in `get_iOperation`

1) Run `main.cpp`, input as follows

Choice - 2, NumberofPts - 64 64 1,

Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1

Degree - 0 0 1, Admissibility Coefficient - 0.5 0.5 1

Output File from C++: `ClusterLinear1_AA.txt`, `ClusterLinear2_AA.txt`,

`ClusterFar_AA.txt`, `ClusterLinearNear_AA.txt`,

OutputFile from MATLAB : `ClusterFigure_64_64_0.5.eps`

2) Run `cluster_pair_indiv.m` from Matlab, arguments:

0 0 (argument here is for the Point Number

and Admissibility Coefficient. As 64 Number of Pts is the

first in the iteration over Number of Pts,

therefore argument is given as 0

same holds for Admissibility Coefficient)

Points on mesh(n): * = Near field points , □ = Partition rectangles (Far field)

NOTE : Axis sections of the partition rectangles $\in \{\Gamma(\sigma) : \sigma \in T\}$, $T \hat{=} \text{Cluster Tree}$ (\rightarrow Def. 11.3.2)

11.4 Algorithm

Local kernel approximation of
Far field-Rectangles through Tensor
product Chebyshev interpolation
polynomial



Low rank approximation of \mathbf{M} on
blocks as (11.1.1)

see Figs. 245,246

R. Hiptmair
rev 38355,
October 28,
2009

Notation:

$$\begin{aligned} \sigma &:= \{i_1, \dots, i_2\} \subset \{1, \dots, n\} \\ \mu &:= \{j_1, \dots, j_2\} \subset \{1, \dots, m\} \end{aligned} \Rightarrow \mathbf{M}_{|\sigma \times \mu} := (m_{ij})_{\substack{i=i_1, \dots, i_2 \\ j=j_1, \dots, j_2}}.$$

Consider kernel approximation of $[x_{i_1}, x_{i_2}] \times [y_{j_1}, y_{j_2}]$ through Tensor product Chebyshev interpolation polynomial (see (11.2.10)):

$$G(x, y) \approx \tilde{G}(x, y) := \sum_{j=0}^d \sum_{k=0}^d G(t_j^x, t_k^y) L_j^x(x) L_k^y(y) ,$$

- $t_0^x, \dots, t_d^x / t_0^y, \dots, t_d^y \hat{=}$ Chebyshev nodes for $[x_{i_1}, x_{i_2}] / [y_{j_1}, y_{j_2}]$,
- $L_j^x, L_k^y \hat{=}$ associated Lagrange polynomial.

$$\tilde{G}(x, y) = \sum_{j=0}^d \sum_{k=0}^d G(t_j^x, t_k^y) L_j^x(x) L_k^y(y) , \quad x \in [x_{i_1}, x_{i_2}], y \in [y_{j_1}, y_{j_2}] .$$

$$\tilde{\mathbf{M}}_{|\sigma \times \mu} = \underbrace{\left(L_j^x(x_i) \right)_{\substack{i=i_1, \dots, i_2 \\ j=0, \dots, d}}}_{\in \mathbb{R}^{\#\sigma, d+1}} \underbrace{\left(G(t_j^x, t_k^y) \right)_{j, k=0, \dots, d}}_{\in \mathbb{R}^{d+1, d+1}} \underbrace{\left(L_k^y(y_j) \right)_{\substack{k=0, \dots, d \\ j=j_1, \dots, j_2}}}_{\in \mathbb{R}^{d+1, \#\mu}} .$$

rank $\leq d+1$

$$\tilde{\mathbf{M}}_{|\sigma \times \mu} = \left(\begin{array}{|c|} \hline \phantom{\text{matrix}} \\ \hline \end{array} \right) \left(\begin{array}{|c|} \hline \phantom{\text{matrix}} \\ \hline \end{array} \right) \left(\begin{array}{|c|} \hline \phantom{\text{matrix}} \\ \hline \end{array} \right) .$$

Application of partitioning

$$\{x_1, \dots, x_n\} \times \{y_1, \dots, y_m\} = \left(\bigcup_{(\sigma, \mu) \in P^{\text{near}}} \sigma \times \mu \right) \cup \left(\bigcup_{(\sigma, \mu) \in P^{\text{far}}} \sigma \times \mu \right) .$$

Algorithm: Analysis of $\mathbf{f} = \mathbf{M}\mathbf{c}$, \mathbf{M} as (11.1.1)

$$f_i = \sum_{j=1}^m G(x_i, y_j) c_j = \underbrace{\sum_{j \in P^{\text{near}}(i)} G(x_i, y_j) c_j}_{\text{Near field contribution}} + \underbrace{\sum_{\substack{\sigma \in T_x \\ x_i \in \sigma}} \sum_{\substack{\mu \in T_y \\ (\sigma, \mu) \in P^{\text{far}}}} \sum_{\substack{1 \leq j \leq m \\ y_j \in \mu}} G(x_i, y_j) c_j}_{\text{Far field contribution}} .$$

Near field coupling indices $P^{\text{near}}(i) := \{j \in \{1, \dots, m\} : (\{x_i\} \times \{y_j\}) \in P^{\text{near}}\}$

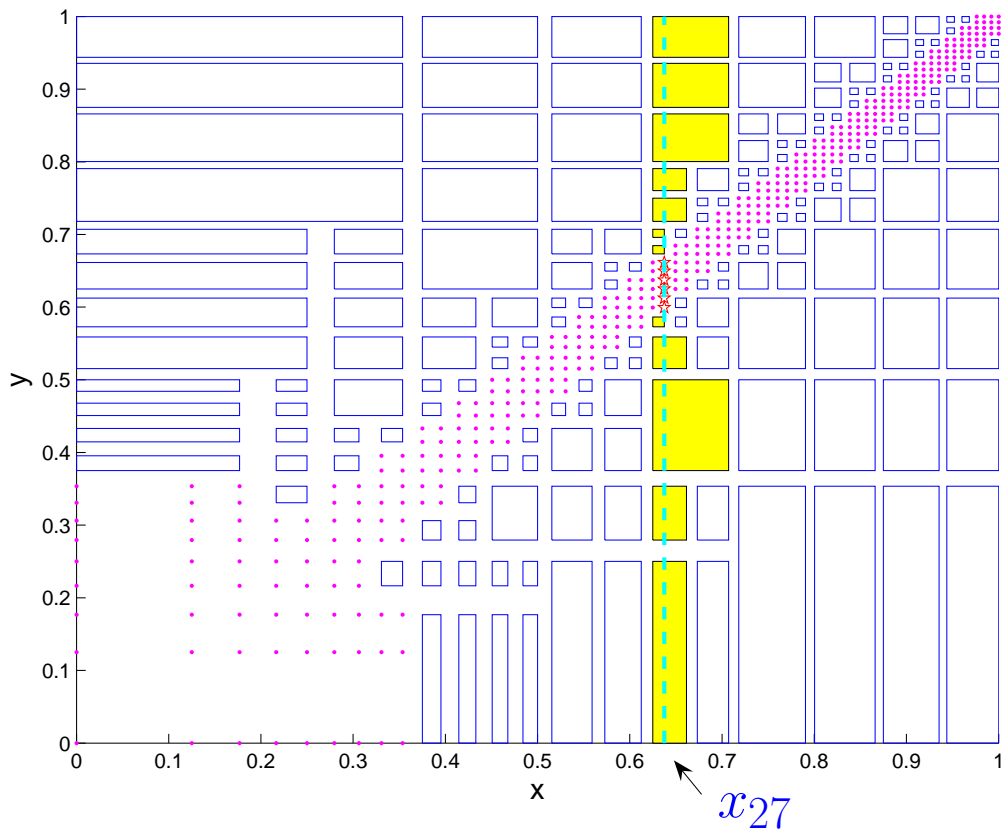


Illustration:

$\mathbf{x} = \text{sqrt}(0:1/64:1) ; i = 27, \eta = 1:$

■ → Far field clusters with contribution to f_i

★ → Near field clusters with contribution to f_i

$$\approx \sum_{j \in P^{\text{near}}(i)} G(x_i, y_j) c_j + \sum_{\substack{\sigma \in T_x \\ x_i \in \sigma}} \sum_{\substack{\mu \in T_y \\ (\sigma, \mu) \in P^{\text{far}}}} \sum_{\substack{1 \leq j \leq m \\ y_j \in \mu}} \sum_{l=0}^d \sum_{k=0}^d G(t_l^\sigma, t_k^\mu) L_l^\sigma(x_i) L_k^\mu(y_j) c_j$$

$$\approx \sum_{j \in P^{\text{near}}(i)} G(x_i, y_j) c_j + \sum_{\substack{\sigma \in T_x \\ x_i \in \sigma}} \sum_{\substack{\mu \in T_y \\ (\sigma, \mu) \in P^{\text{far}}}} (\mathbf{V}_\sigma \mathbf{X}_{\sigma, \mu} \mathbf{V}_\mu^T \mathbf{c}_{|\mu})_i,$$

with

$t_l^\sigma, t_k^\mu :=$ Chebyshev nodes (9.2.12) in $\Gamma(\sigma), \Gamma(\mu)$, $l = 0, \dots, d$,

$$\mathbf{X}_{\sigma, \mu} := (G(t_l^\sigma, t_k^\mu))_{l, k=0}^d \in \mathbb{R}^{d+1, d+1}, \quad (11.4.1)$$

$$\mathbf{V}_\sigma := (L_l^\sigma(x_i))_{\substack{i: x_i \in \sigma \\ l=0, \dots, d}} \in \mathbb{R}^{\#\sigma, d+1}, \quad (11.4.2)$$

$$\mathbf{V}_\mu := (L_k^\mu(y_j))_{\substack{j: y_j \in \mu \\ k=0, \dots, d}} \in \mathbb{R}^{\#\mu, d+1}, \quad (11.4.3)$$

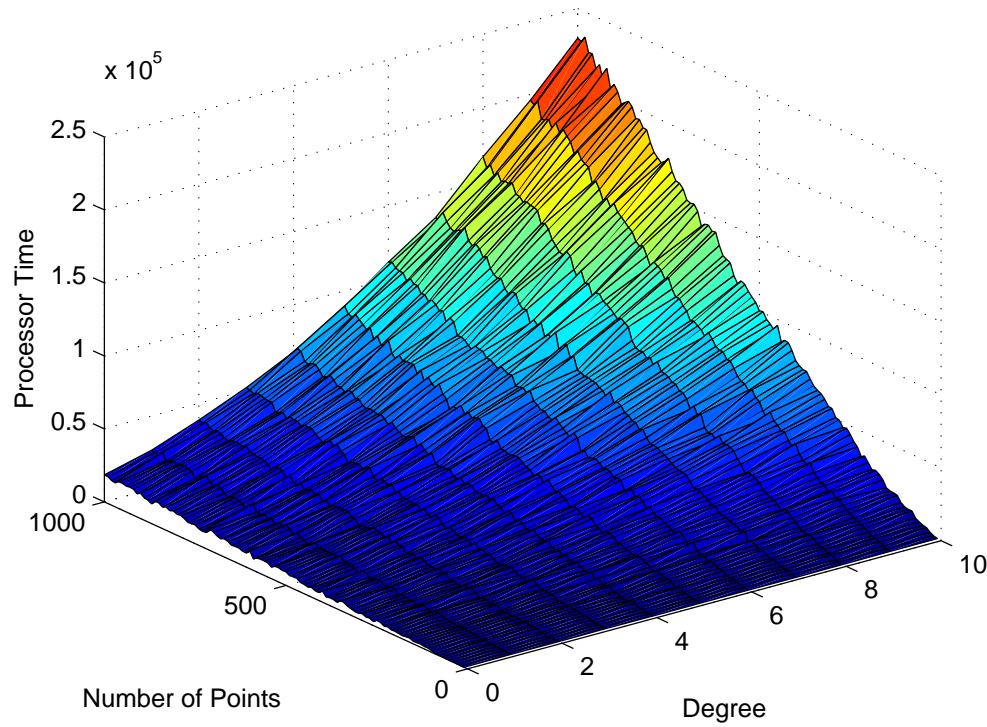
Analysis of Complexity (under the assumption $n \approx m$):

①: Trick: Calculate on above assumption

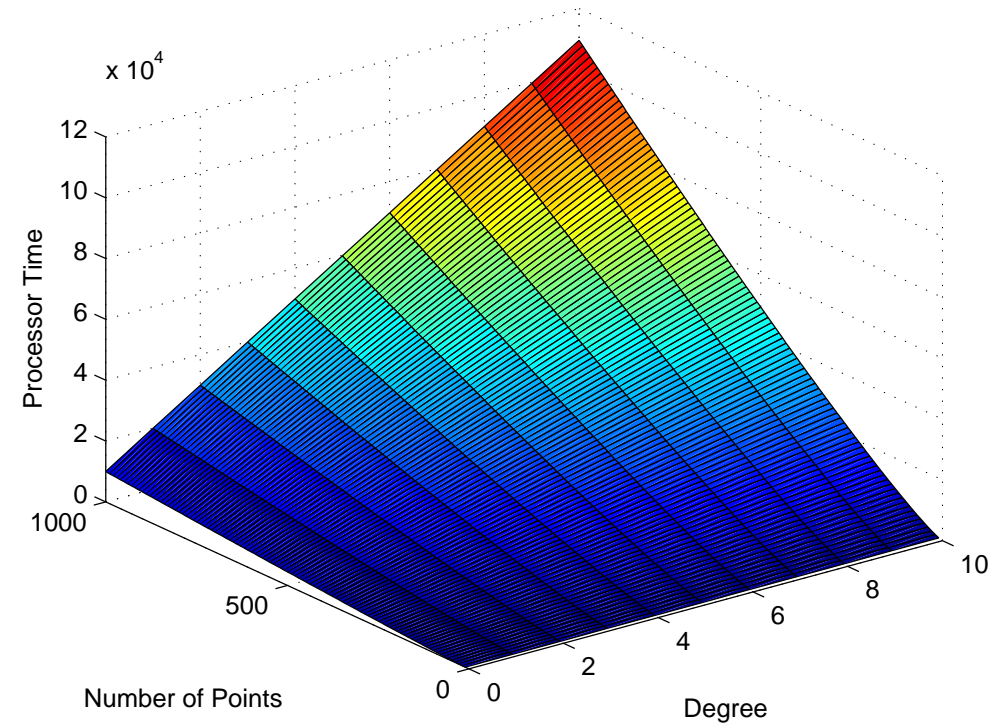
Compute $\mathbf{w}_\mu := \mathbf{V}_\mu^T \mathbf{c}|_\mu \in \mathbb{R}^{d+1}$, $\mu \in T_y \rightarrow O(n)$ Operations on every level of T_y

Computational Effort $O((d+1)n \log_2 n)$, Memory Required $O((d+1)n \log_2 n)$

Computational Time – Processor Time Vector W



Computational Time – Processor Time – Theoretical Vector W



In file changeable.c -

- a) put the admissibility condition in `admissible_ClusterPair`
- b) put kernel function in `kernel_function`
- c) put **NONE** as return values in `get_iOperation`
- d) put **100** as return values in `get_iNumber of Times`

Output File from C++:

OutputFile from MATLAB :

Time.txt
 Time_Vector_W.eps, Time_Theoretical_Vector_W.eps, Time_Chebyshev.eps,
 Time_Theoretical_Chebyshev.eps, Time_Far.eps, Time_Theoretical_Far.eps, Time_Near.eps,
 Time_Theoretical_Near.eps, Time_All.eps, Time_Theoretical_All.eps
 Choice - 2, NumberofPts - 10 1000 10, Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1
 Degree - 0 10 1, Admissibility Coefficient - **0.5 0.5 1.0**

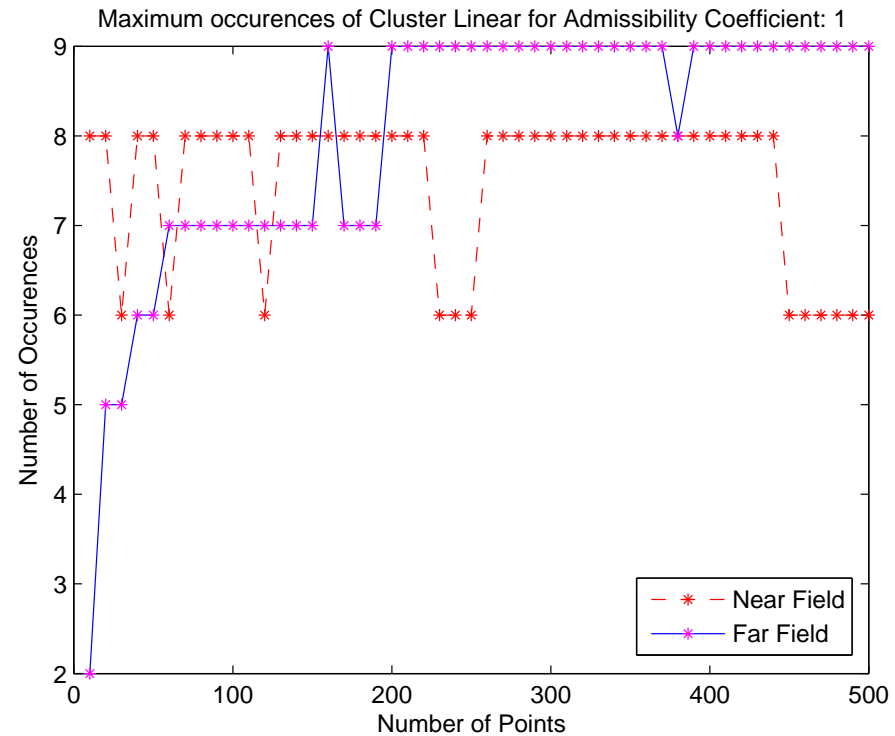
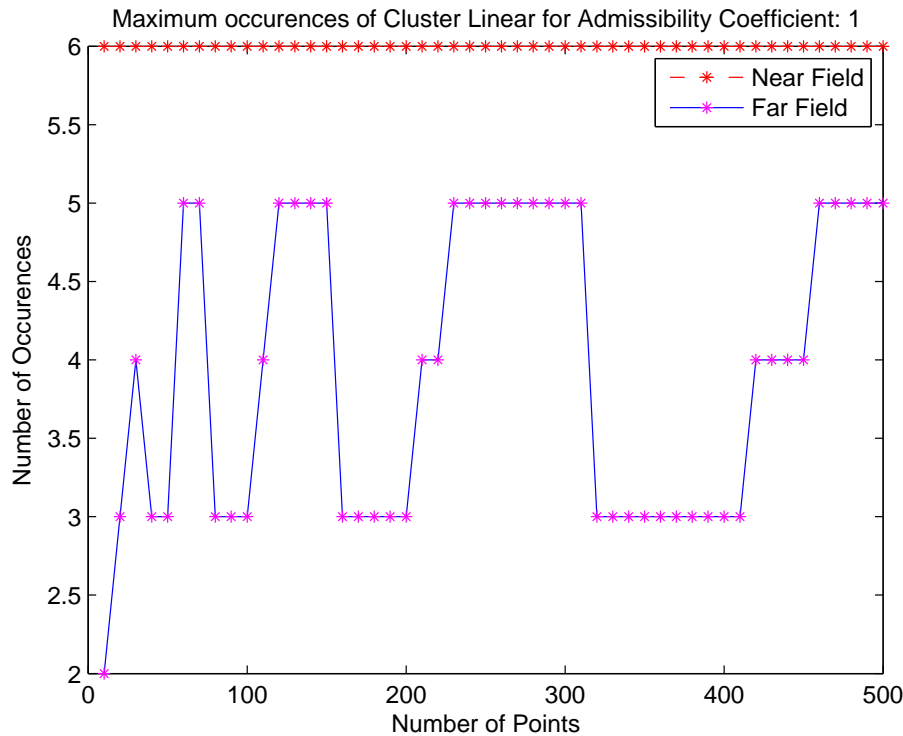
1) Run main.cpp, input as follows :

2) Run `time_mesh.m` and `time_mesh_theoretical` from Matlab

②: Compute $X_{\sigma,\mu}$ as in (11.4.1), $\sigma \in T_x$, $\mu \in T_y$, $(\sigma, \mu) \in P^{\text{far}}$.

Uniform distribution of points $\max\{\#\{\mu \in T_y: (\sigma, \mu) \in P^{\text{far}}\}, \sigma \in T_x\} = O(1)$

Example 11.4.4 (Occurrence of clusters in partition rectangles).



$x = (0:1/n:(1-1/n))$, $\eta = 1$
In file changeable.c

- put the admissibility condition in `admissible_ClusterPair`
- put kernel function in `kernel_function`
- put **NONE** as return values in `get_iOperation`
 - Run `main.cpp`, input as follows
Choice - 2, NumberofPts - 10 500 10,
Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1
Degree - 0 0 1, Admissibility Coefficient - 1 1 1
Output File from C++: `Count.txt`
OutputFile from MATLAB : `MaxCluster_1.eps`
 - Run `count_indiv.m` from Matlab, arguments:
0 (argument here is for specifying the index of Admissibility Coefficient.)

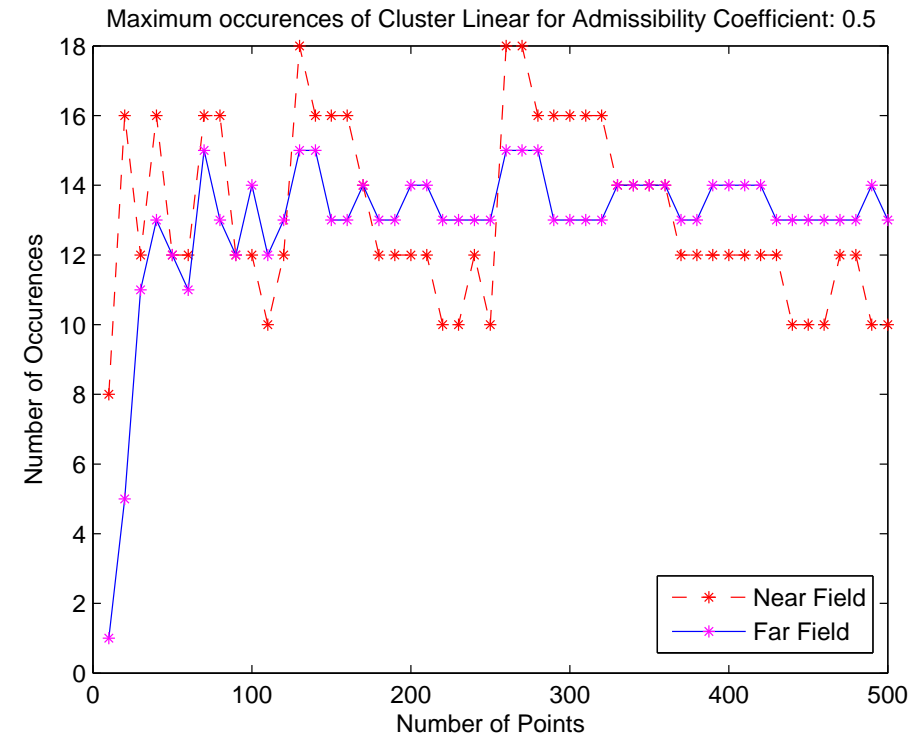
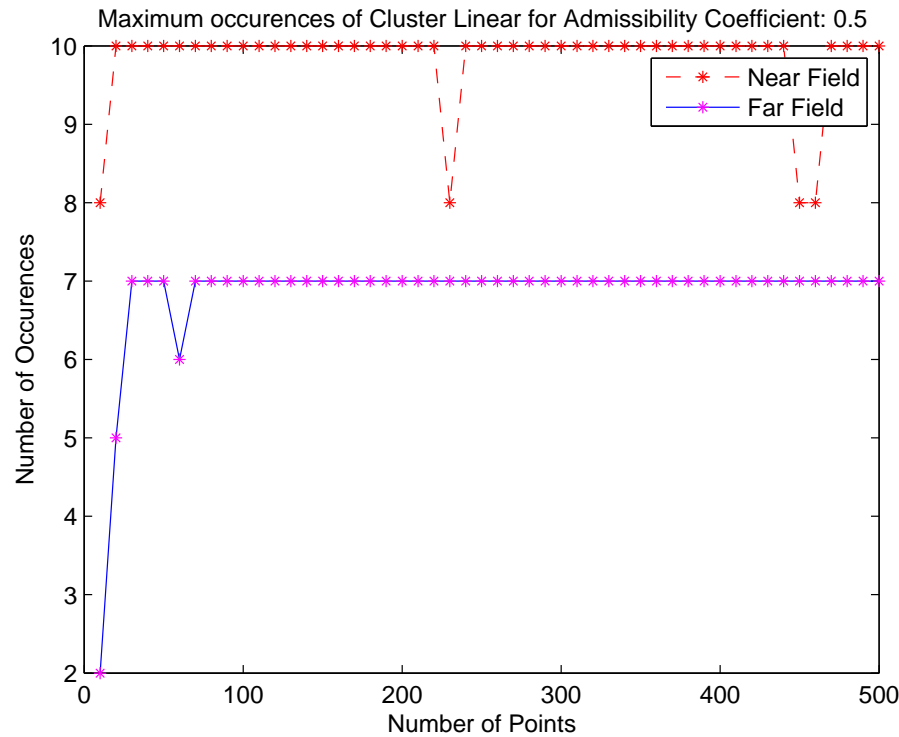
$x = \text{sqrt}(0:1/n:(1-1/n))$, $\eta = 1$
In file changeable.c

- put the admissibility condition in `admissible_ClusterPair`
- put kernel function in `kernel_function`
- put **SQRT** as return values in `get_iOperation`
 - Run `main.cpp`, input as follows
Choice - 2, NumberofPts - 10 500 10,
Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1
Degree - 0 0 1, Admissibility Coefficient - 1 1 1
Output File from C++: `Count.txt`
OutputFile from MATLAB : `MaxCluster_1.eps`
 - Run `count_indiv.m` from Matlab, arguments:
0 (argument here is for specifying the index of Admissibility Coefficient.)

Uniform distribution of x_i, y_j : for every $\sigma \in T_x$

$$\#\{\mu \in T_y: (\sigma, \mu) \in \text{Pfar}\} = O(1),$$

each cluster contributes only a small number of partitions of rectangle



$$x = (0:1/n:(1-1/n))i, \eta = 0.5$$

In file changeable.c

a) put the admissibility condition in `admissible_ClusterPair`

b) put kernel function in `kernel_function`

c) put **NONE** as return values in `get_iOperation`

1) Run `main.cpp`, input as follows

Choice - 2, NumberofPts - 10 500 10,

Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1

Degree - 0 0 1, Admissibility Coefficient - 0.5 0.5 1

Output File from C++: `Count.txt`

OutputFile from MATLAB : `MaxCluster_0.5.eps`

2) Run `count_indiv.m` from Matlab, arguments:

0 (argument here is for specifying the index of Admissibility Coefficient.)

$$x = \text{sqrt}(0:1/n:(1-1/n))i, \eta = 0.5$$

In file changeable.c

a) put the admissibility condition in `admissible_ClusterPair`

b) put kernel function in `kernel_function`

c) put **SQRT** as return values in `get_iOperation`

1) Run `main.cpp`, input as follows

Choice - 2, NumberofPts - 10 500 10,

Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1

Degree - 0 0 1, Admissibility Coefficient - 0.5 0.5 1

Output File from C++: `Count.txt`

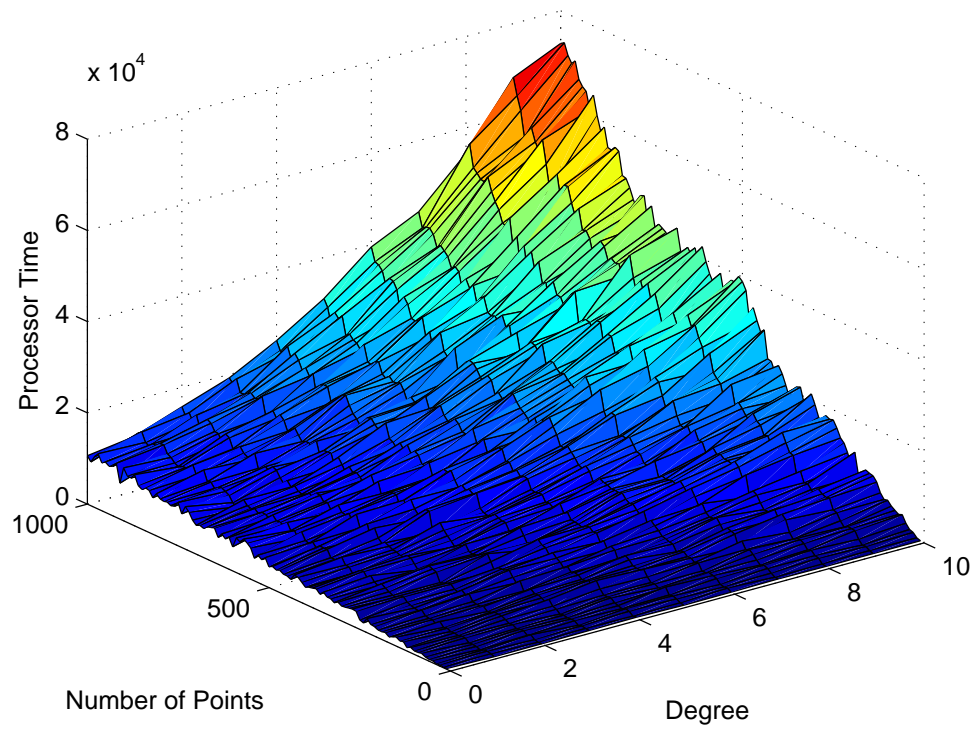
OutputFile from MATLAB : `MaxCluster_0.5.eps`

2) Run `count_indiv.m` from Matlab, arguments:

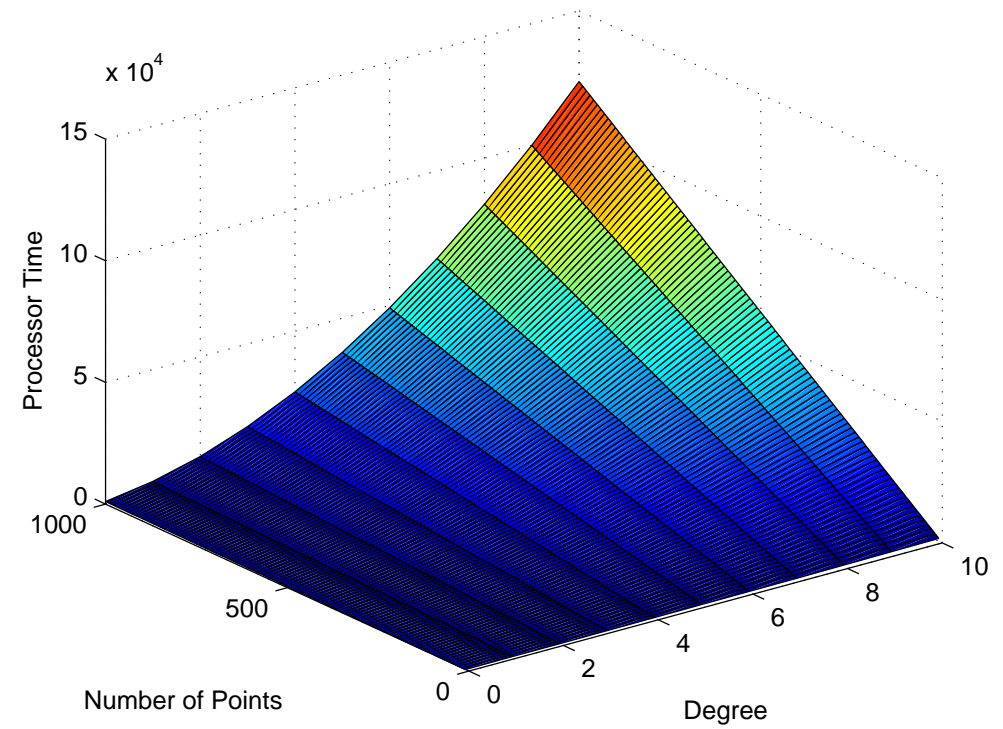
0 (argument here is for specifying the index of Admissibility Coefficient.)

Computational Effort for step ②: $O(n \log_2 n (d + 1)^2)$, Memory Required $O(n \log_2 n (d + 1)^2)$

Computational Time – Processor Time Chebyshev Nodes



Computational Time – Processor Time – Theoretical Chebyshev Nodes



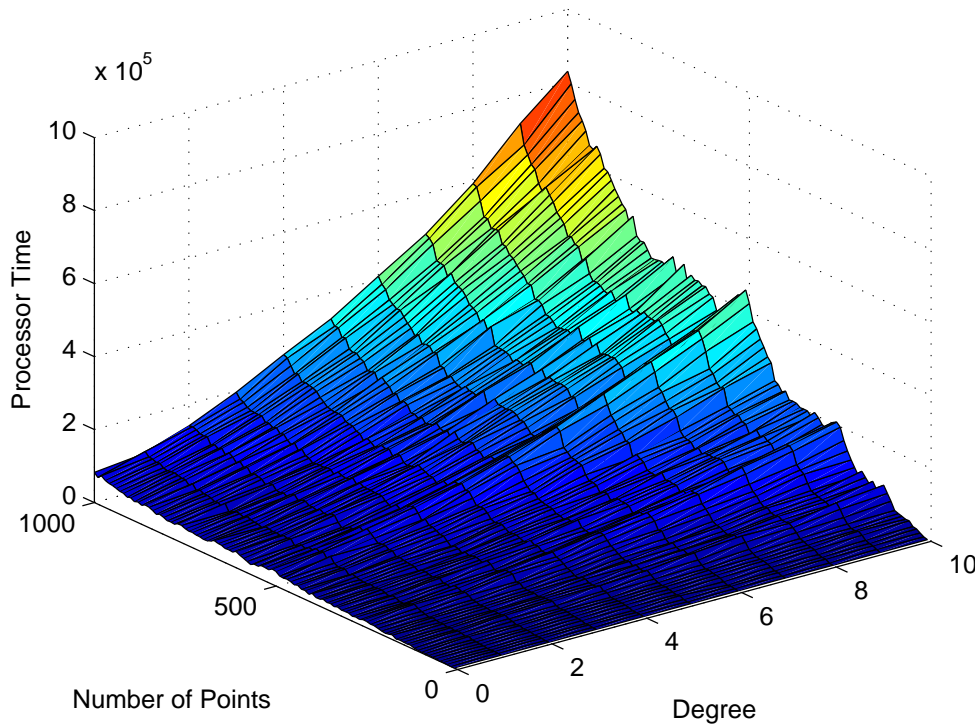
③: Far field analysis $\hat{=}$ Sum

$$\sum_{\substack{\sigma \in T_x \\ x_i \in \sigma}} \sum_{\substack{\mu \in T_y \\ (\sigma, \mu) \in P^{\text{far}}}} \mathbf{V}_\sigma \mathbf{X}_{\sigma, \mu} \mathbf{w}_\mu$$

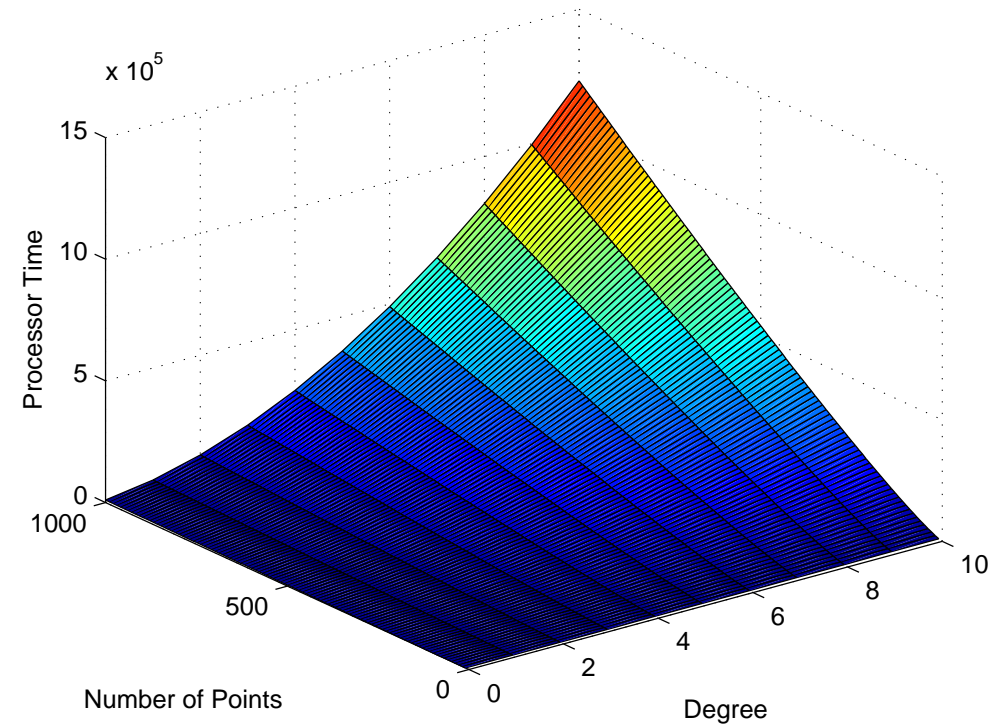
for each $\sigma \in T_x$ { Compute \mathbf{V}_σ ; $\mathbf{s} := 0$
 for each $\mu \in T_y, (\sigma, \mu) \in P^{\text{far}}$ { $\mathbf{s} \leftarrow \mathbf{s} + \mathbf{X}_{\sigma, \mu} \mathbf{w}_\mu$ }
 $\mathbf{f}_{|\sigma} \leftarrow \mathbf{f}_{|\sigma} + \mathbf{V}_\sigma \mathbf{s}$ }

Computational Effort $O(n \log_2 n (d + 1)^2)$, Memory Required $O(n(d + 1))$

Computational Time – Processor Time Far Field



Computational Time – Processor Time – Theoretical Far Field

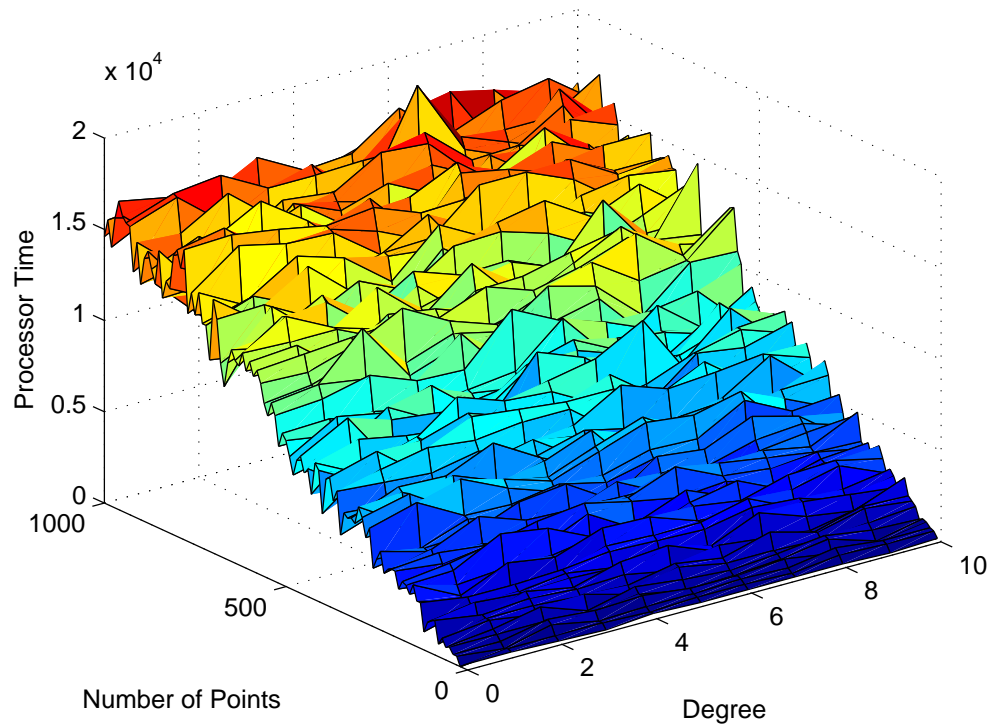


④: Near field Computation

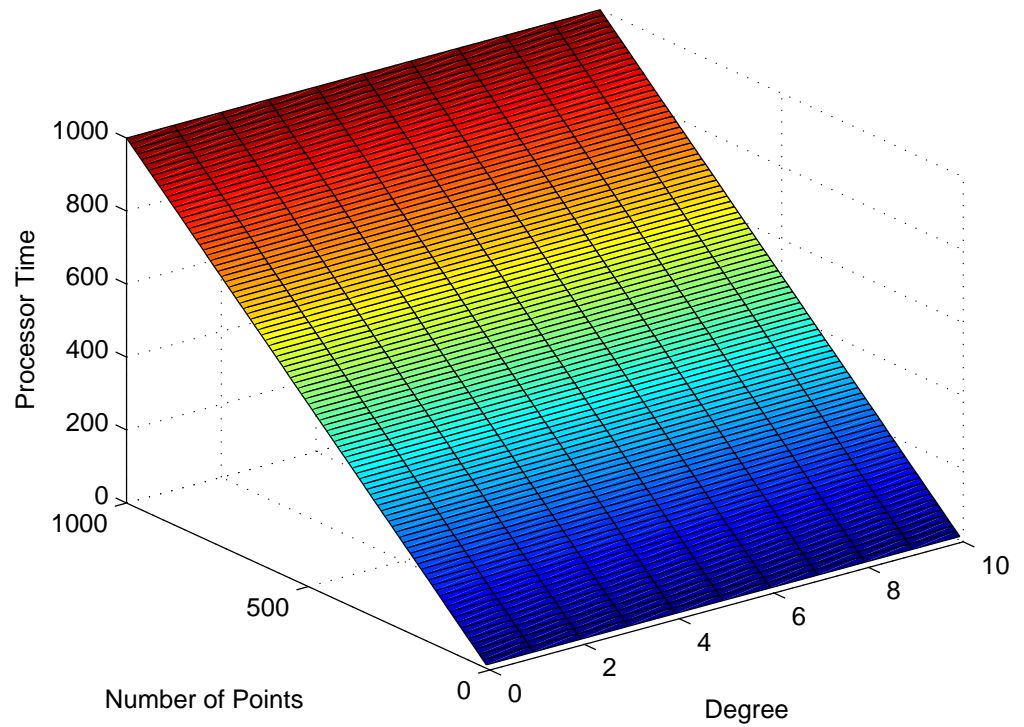
for each $(\sigma, \mu) \in P^{\text{near}}$ { for each $i: x_i \in \sigma$ { $\mathbf{f}_i \leftarrow \mathbf{f}_i + \sum_{j: y_j \in \mu} G(x_i, y_j)c_j$ } }

Computational Effort $O(n)$

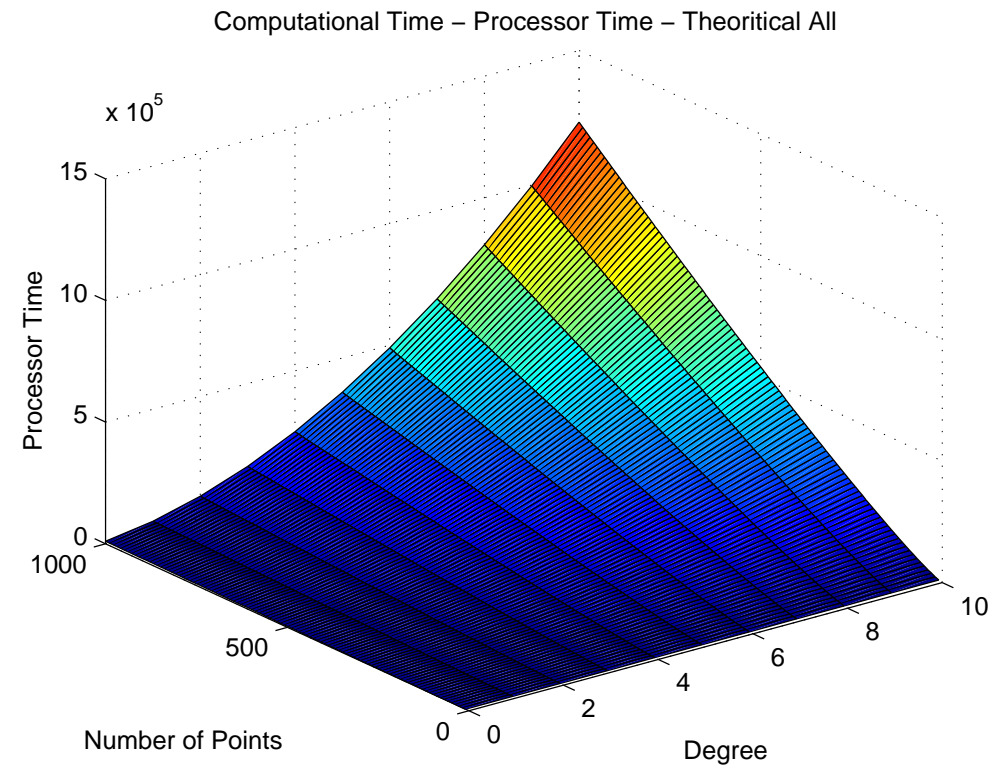
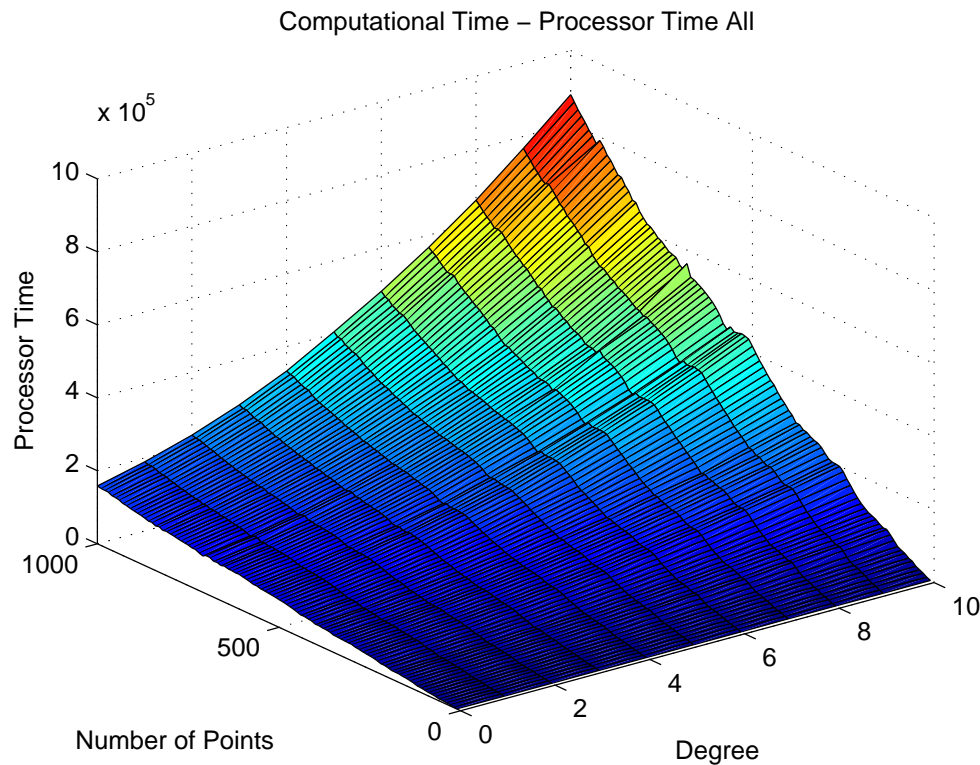
Computational Time – Processor Time Near Field



Computational Time – Processor Time – Theoretical Near Field



Total Computational Effort/Total Memory Required $O((d + 1)^2 n \log_2 n)$



Remark 11.4.5 (Convergence of cluster approximation).

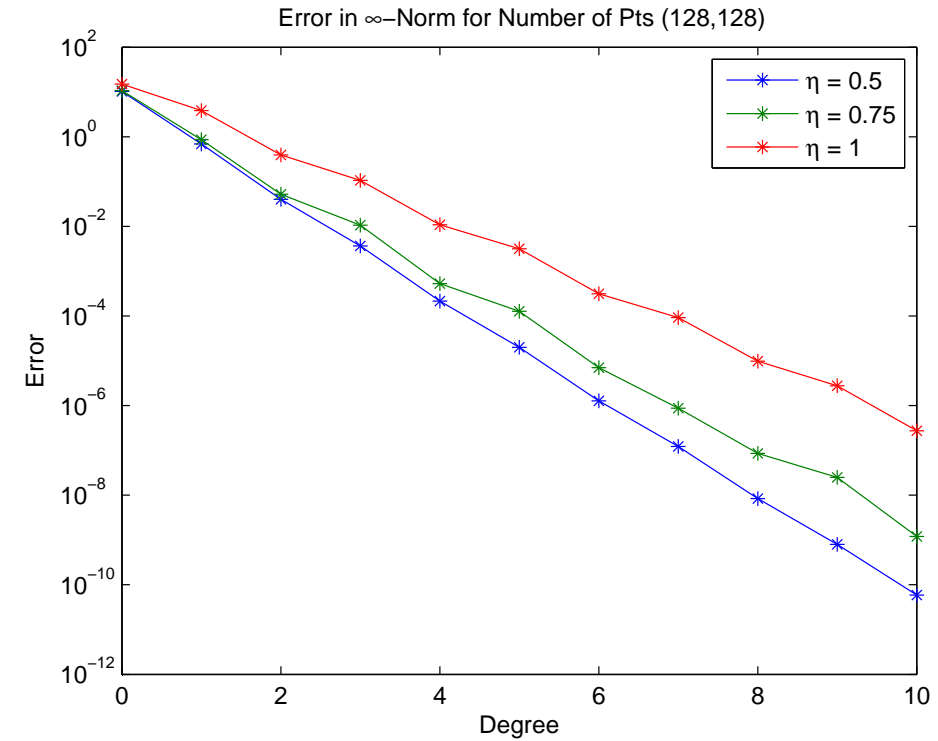
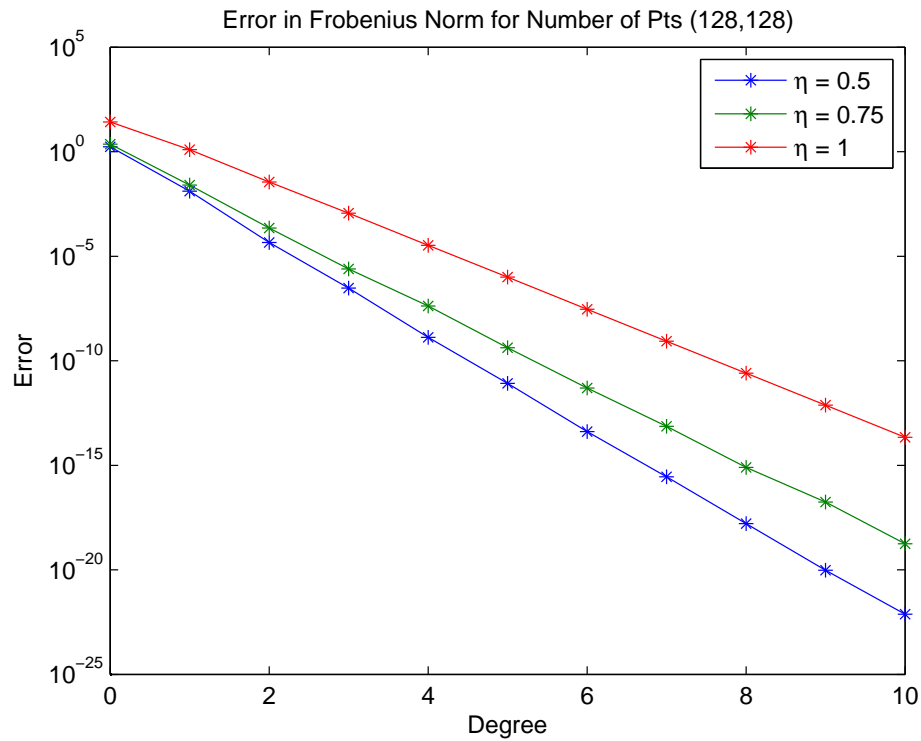
if $G(x, y)$ analytic (\rightarrow Def. 9.2.20) in $\{(x, y): x \neq y\}$ then

Cluster approximation inherits **exponential convergence** from Chebyshev interpolation



Example 11.4.6 (Convergence of clustering approximation with collocation matrix).

- $x = 0:1/128:1$; $G(x, y) = |x - y|^{-1}$ for $x \neq y$, $G(x, x) = 0$.



In file changeable.c

- put the admissibility condition in `admissible_ClusterPair`
- put kernel function in `kernel_function`
- put **NONE** as return values in `get_iOperation`

1) Run `main.cpp`, input as follows

Choice - 2, NumberofPts - 128 128 1,

Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1

Degree - 0 10 1, Admissibility Coefficient - **0.5 1.0 0.25**

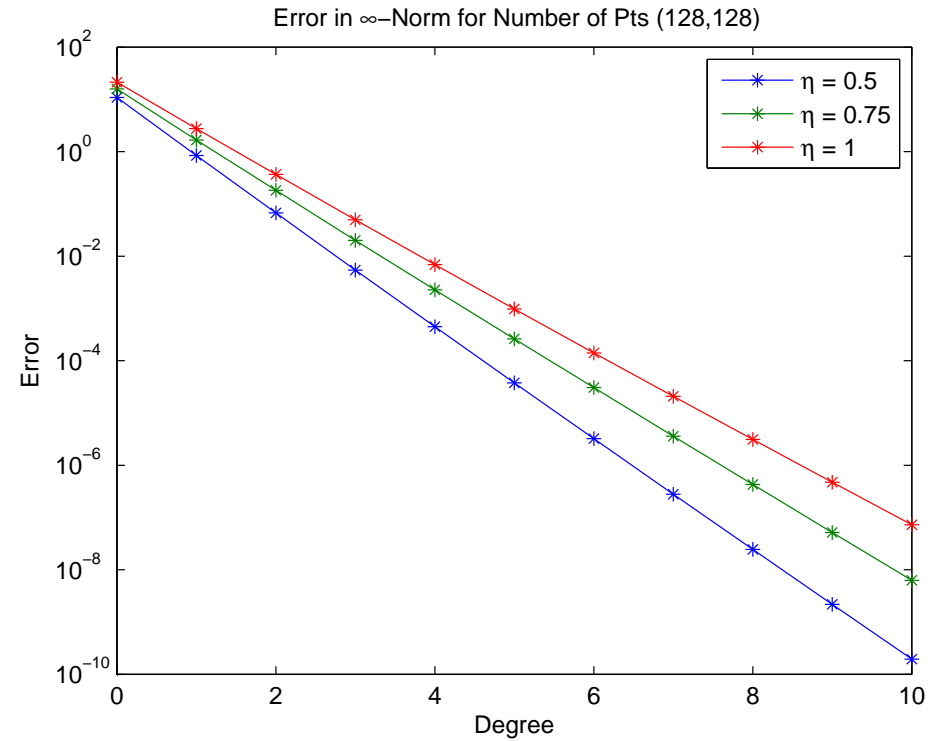
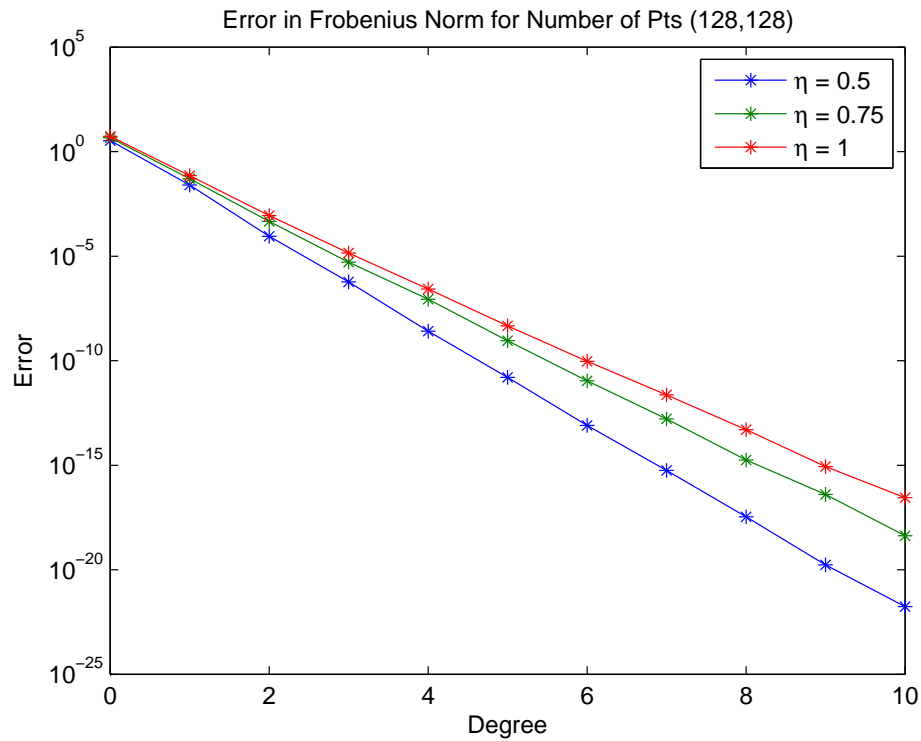
Output File from C++: `Error.txt`

OutputFile from MATLAB : `Error_in_Frobenius_Norm_128_128.eps`

2) Run `error_plot.m` from Matlab, arguments:

0 (argument here is for specifying the index of Number of Pts.)

• $x = \text{sqrt}(0:1/128:1)$; $G(x, y) = |x - y|^{-1}$ for $x \neq y$, $G(x, x) = 0$.



In file changeable.c

a) put the admissibility condition in `admissible_ClusterPair`

b) put kernel function in `kernel_function`

c) put **SQRT** as return values in `get_iOperation`

1) Run `main.cpp`, input as follows

Choice - 2, NumberofPts - 128 128 1,

Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1

Degree - 0 10 1, Admissibility Coefficient - **0.5 1.0 0.25**

Output File from C++: `Error.txt`

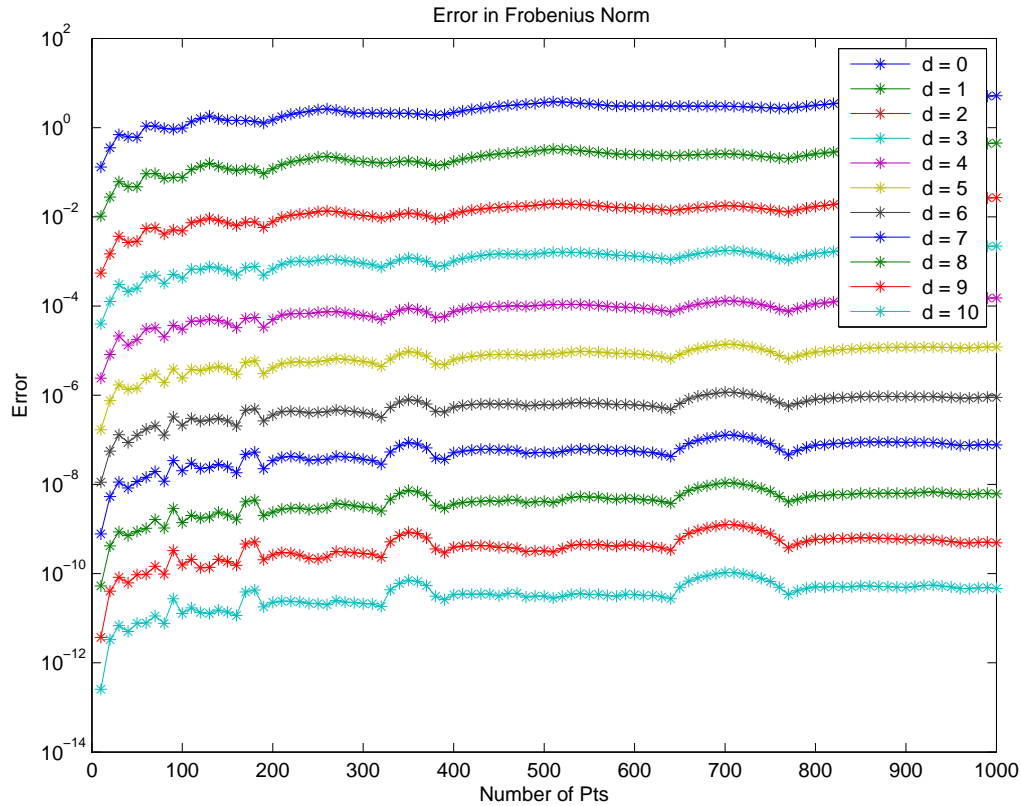
OutputFile from MATLAB : `Error_in_Frobenius_Norm_128_128.eps`

2) Run `error_plot.m` from Matlab, arguments:

0 (argument here is for specifying the index of Number of Pts.)



● Scaled Frobenius Error v/s Number of Points



- In file changeable.c
- a) put the admissibility condition in `admissible_ClusterPair`
 - b) put kernel function in `kernel_function`
 - c) put **NONE** as return values in `get_iOperation`
 - 1) Run `main.cpp`, input as follows
Choice - 2, NumberofPts - 10 1000 10,
Start of 1 - 0, End of 1 - 1, Start of 2 - 0, End of 2 - 1
Degree - 0 10 1, Admissibility Coefficient - 0.5 0.5 1.0
Output File from C++: `Error.txt`
 - 2) Run `error_plot_Degree.m` from Matlab
- OutputFile from MATLAB : `Error_in_Frobenius_Norm_Degree.eps`

Example 11.4.7 (Analysis of trigonometric polynomials). [19]

Given: $\{t_0, \dots, t_{n-1}\} \subset [0, 1[$, $\alpha_{-m+1}, \dots, \alpha_m \in \mathbb{C}$, $n = 2m$, $m \in \mathbb{N}$, compute

$$c_k := p(t_k), \quad j = 0, \dots, n-1 \quad \text{for} \quad p(t) := \underbrace{\sum_{j=-m+1}^m \alpha_j e^{-2\pi i j t}}_{\text{Trigonometric Polynomial}}.$$

Discrete Fourier transformation (DFT) \rightarrow Section 8.2:

$$f_l := \sum_{j=-m+1}^m \alpha_j e^{-\frac{2\pi i}{n} j l} \quad \stackrel{\text{Lemma 8.2.10}}{\Leftrightarrow} \quad \alpha_j = \frac{1}{n} \sum_{l=-m+1}^m f_l e^{\frac{2\pi i}{n} l j}$$

FFT: f_l , $l = -m+1, \dots, m$, calculated with computational effort $O(n \log_2 n)$

$$\begin{aligned} c_k &= \sum_{j=-m+1}^m \alpha_j e^{2\pi i j t_k} = \sum_{j=-m+1}^m \frac{1}{n} \left(\sum_{l=-m+1}^m f_l e^{\frac{2\pi i}{n} l j} \right) e^{-2\pi i j t_k} \\ &= \frac{1}{n} \sum_{l=-m+1}^m f_l \sum_{j=-m+1}^m e^{-2\pi i j (t_k - l/n)} \\ &= \frac{1}{n} \sum_{l=-m+1}^m f_l e^{-2\pi i (t_k - l/n)(-m+1)} \frac{1 - e^{-2\pi i n (t_k - l/n)}}{1 - e^{-2\pi i (t_k - l/n)}} \end{aligned}$$

$$= \frac{1}{n} \sum_{l=-m+1}^m f_l e^{-\pi i t_k} \sin(\pi n t_k) \frac{1}{\sin(\pi(t_k - l/n))} (-1)^l e^{-\pi i l/n} .$$

$$\mathbf{c} = \text{diag} \left(\frac{\sin(\pi n t_k)}{e^{\pi i t_k}} \right)_{k=-m+1}^m \mathbf{M} \text{diag} \left((-1)^l e^{-\pi i l/n} \right)_{l=-m+1, \dots, m} \mathbf{f} .$$

in accordance with collocation matrix (11.1.1)

$$\mathbf{M} := \left(\frac{1}{\sin(\pi(t_k - l/n))} \right)_{\substack{k=-m+1, \dots, m \\ l=-m+1, \dots, m}} \in \mathbb{R}^{2n, 2n} .$$

Approximative analysis of Clustering algorithms ! → USFFT



Remark 11.4.8. Clustering approximation example for **fast approximative** implementation of algorithms of numerical linear algebra (→ Chapter ??) → Trend in numerical linear algebra ?

| Problem | Exact/direct Method | Approximate/iterative Method |
|------------------------------|------------------------|--------------------------------------|
| Linear equation systems | Gaussian Elimination | CG-like iterative solvers → Sect. ?? |
| Eigenvalue problem | Transformation methods | Krylov-Subspace method → Sect. ?? |
| Collocations matrix × Vector | BLAS (SAXPY) | Clustering techniques |



12

Single Step Methods

12.1 Initial value problems (IVP) for ODEs

Acronym: **ODE = ordinary differential equation**

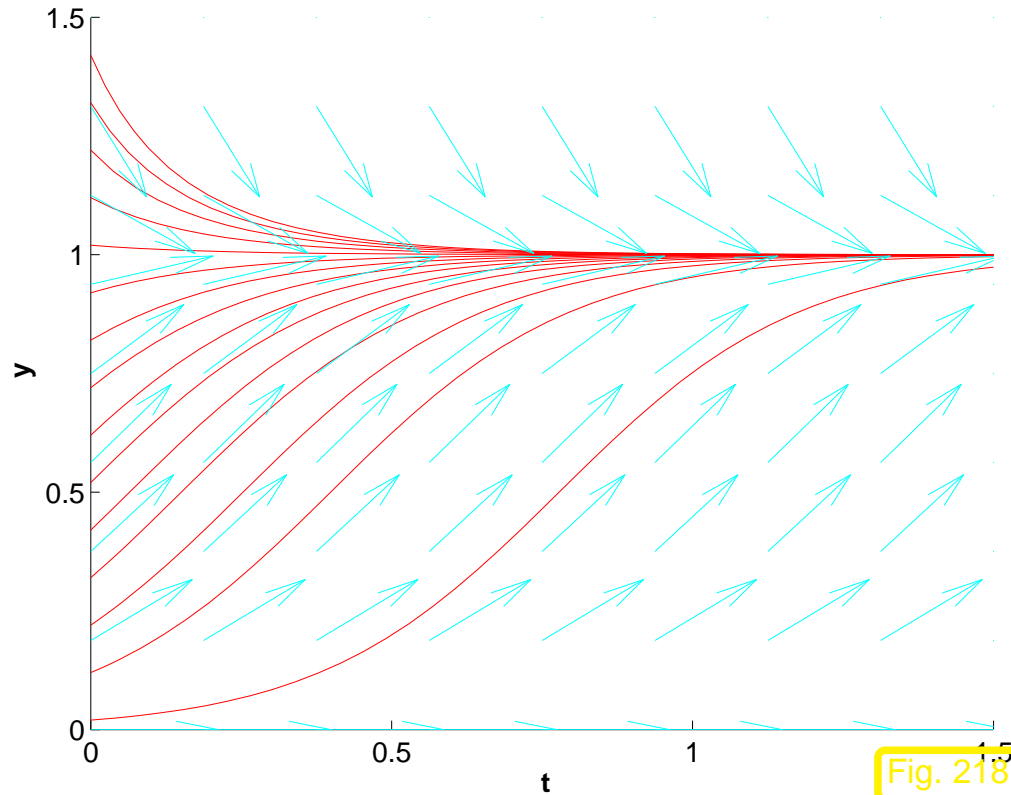
R. Hiptmair
rev 38355,
May 3,
2011

12.1.1 Examples

Example 12.1.1 (Growth with limited resources). [2, Sect. 1.1], [35, Ch. 60]

$y : [0, T] \mapsto \mathbb{R}$: bacterial population density as a function of time

$$\dot{y} = f(y) := (\alpha - \beta y) y \quad (12.1.2)$$



Solution for different $y(0)$ ($\alpha, \beta = 5$)

By separation of variables

→ solution of (12.1.2)

for $y(0) = y_0 > 0$

$$y(t) = \frac{\alpha y_0}{\beta y_0 + (\alpha - \beta y_0) \exp(-\alpha t)}, \quad (12.1.3)$$

for all $t \in \mathbb{R}$



Example 12.1.4 (Predator-prey model). [2, Sect. 1.1], [31, Sect. 1.1.1], [35, Ch. 60], [13, Ex. 11.3]

Model: autonomous Lotka-Volterra ODE:

$$\begin{aligned} \dot{u} &= (\alpha - \beta v)u \\ \dot{v} &= (\delta u - \gamma)v \end{aligned} \Leftrightarrow \dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) \quad \text{with} \quad \mathbf{y} = \begin{pmatrix} u \\ v \end{pmatrix}, \quad \mathbf{f}(\mathbf{y}) = \begin{pmatrix} (\alpha - \beta v)u \\ (\delta u - \gamma)v \end{pmatrix}. \quad (12.1.5)$$

population sizes:

$u(t)$ → no. of prey at time t ,

$v(t)$ → no. of predators at time t

vector field \mathbf{f} for Lotka-Volterra ODE ▷

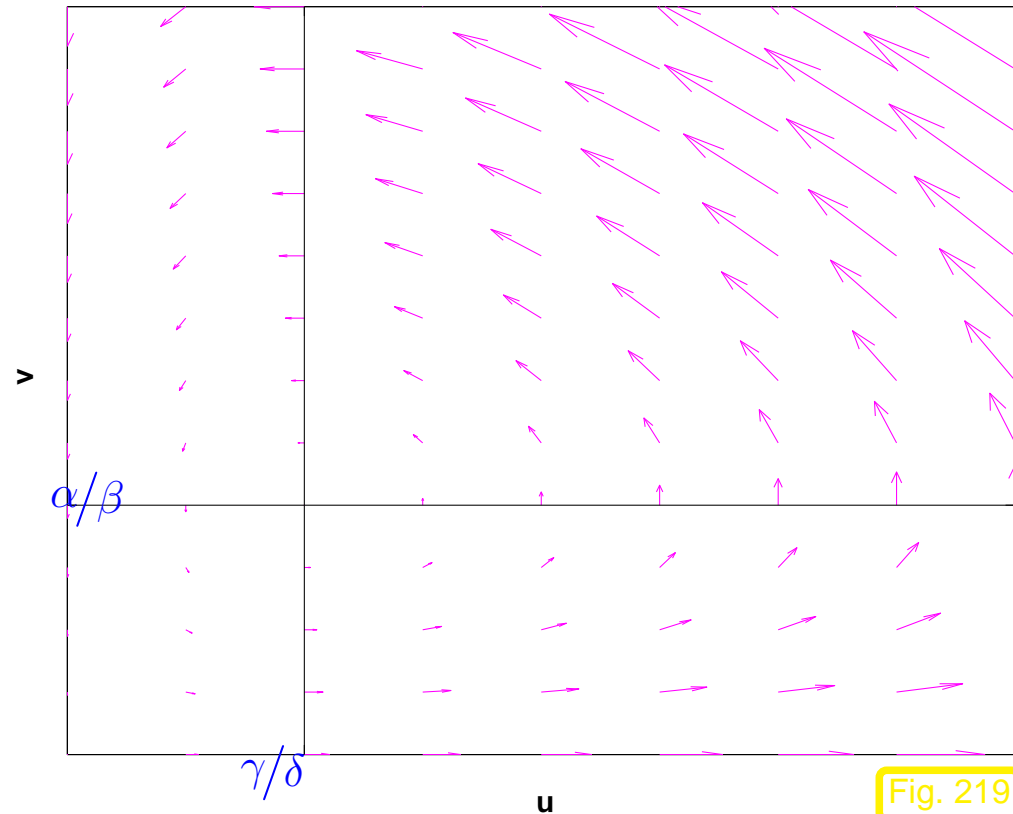


Fig. 219

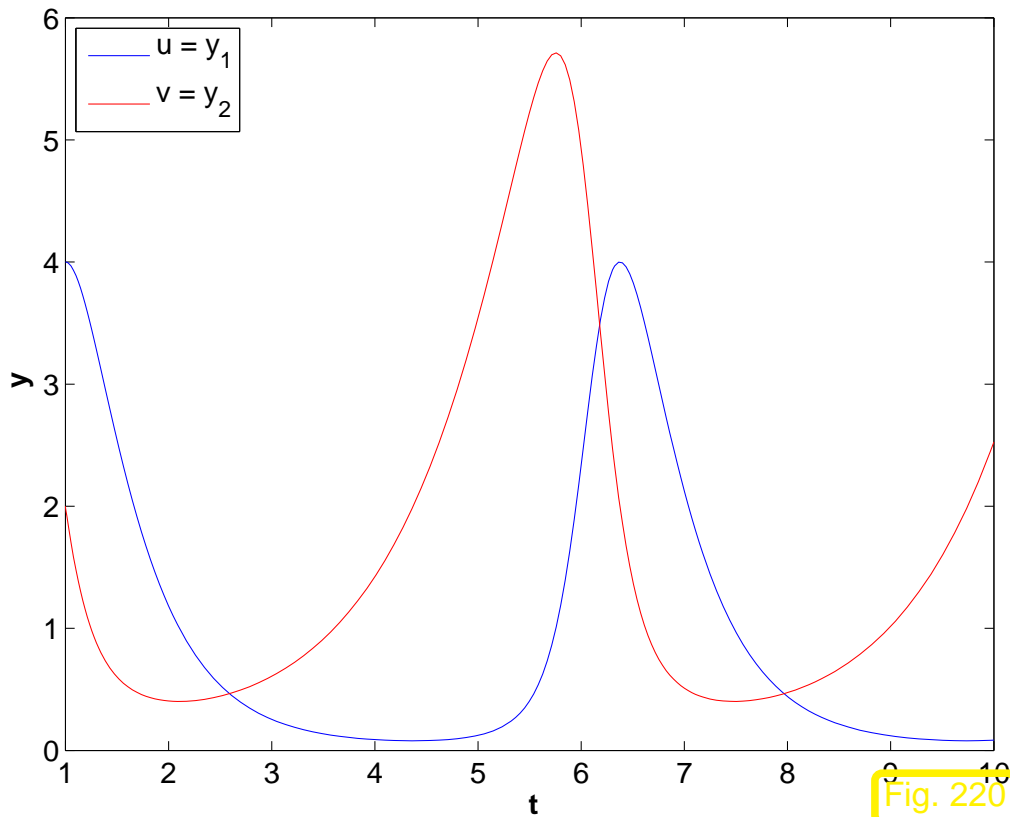


Fig. 220

Solution $\begin{pmatrix} u(t) \\ v(t) \end{pmatrix}$ for $\mathbf{y}_0 := \begin{pmatrix} u(0) \\ v(0) \end{pmatrix} = \begin{pmatrix} 4 \\ 2 \end{pmatrix}$

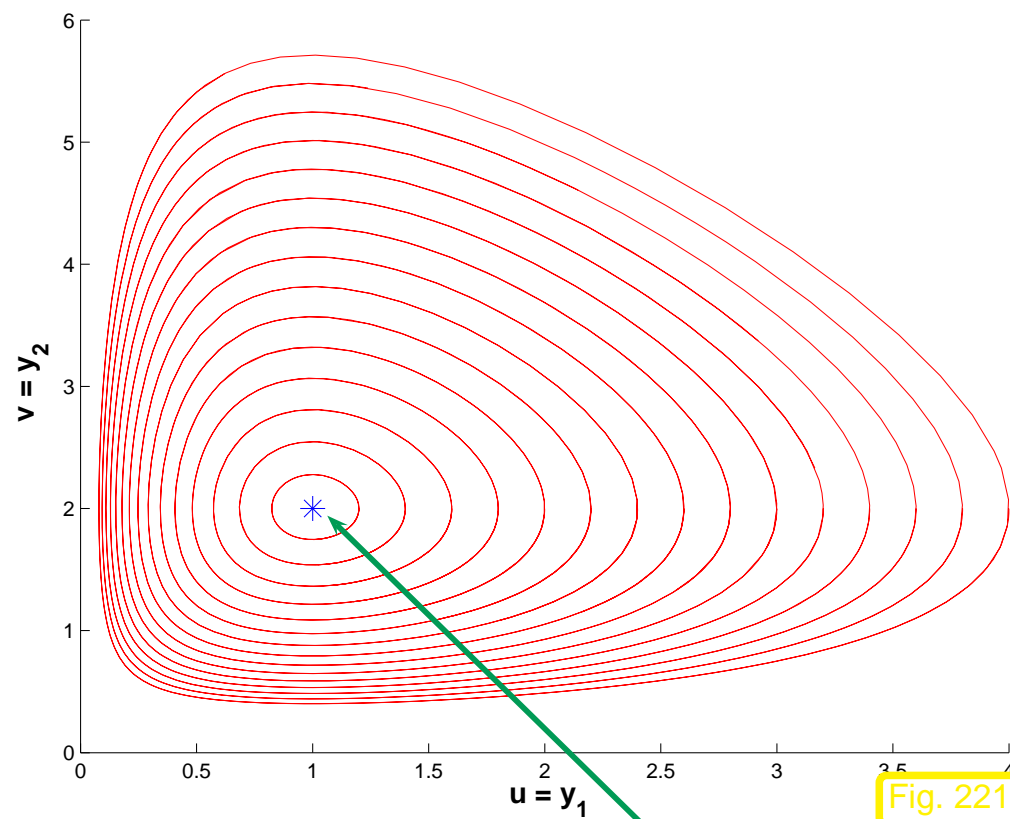


Fig. 221

Solution curves for (12.1.5)
stationary point



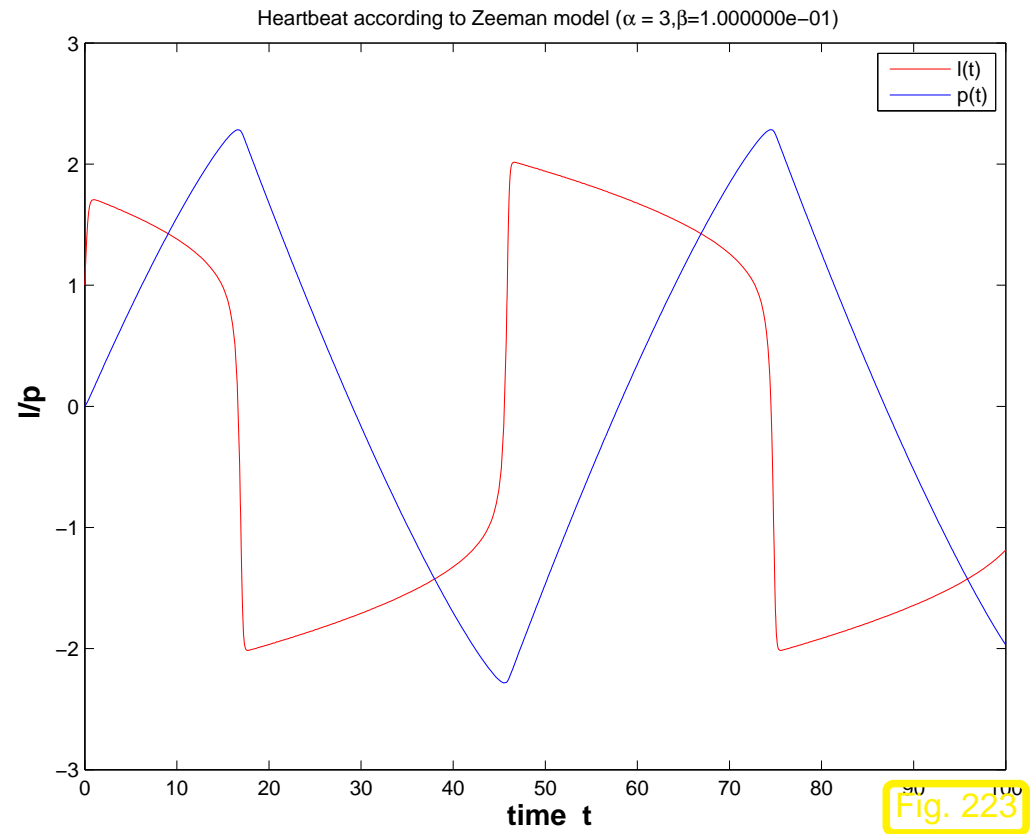
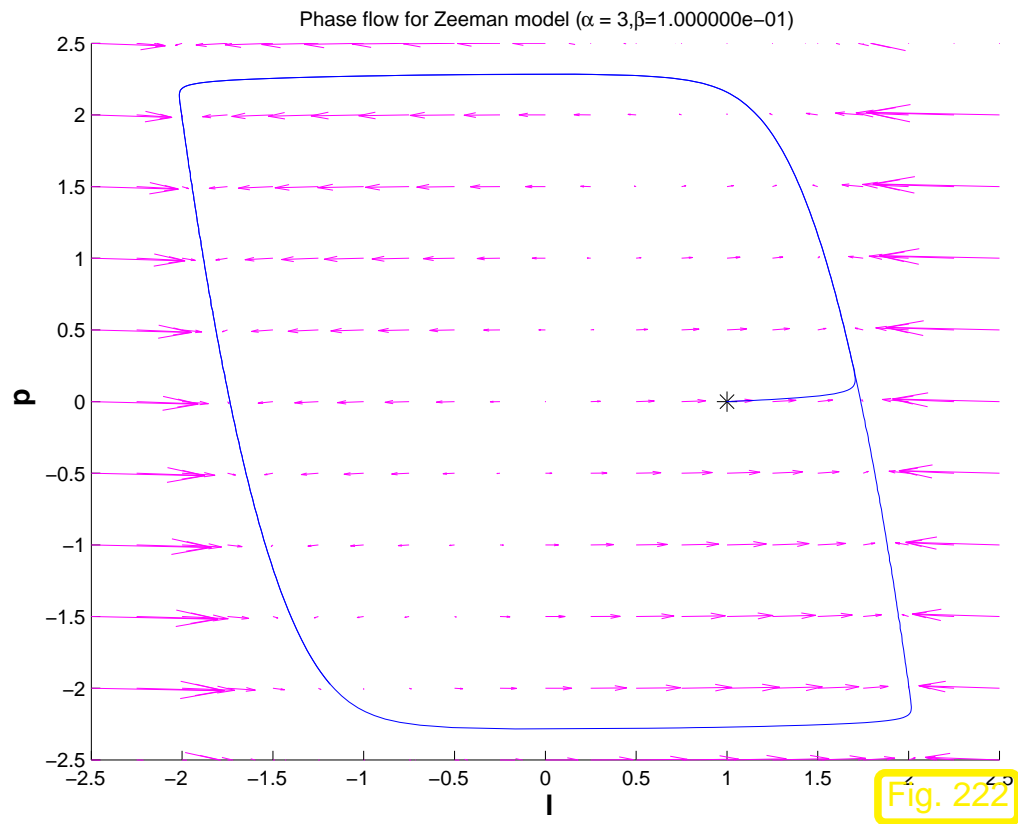
Example 12.1.6 (Heartbeat model). → [14, p. 655]

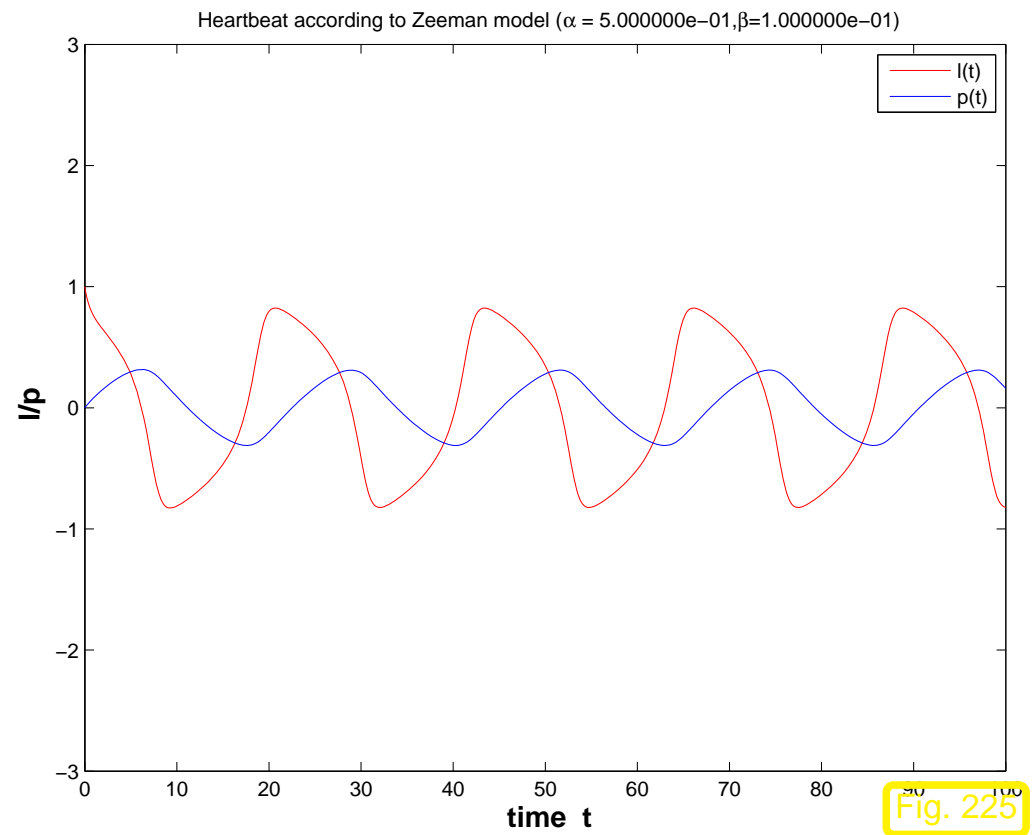
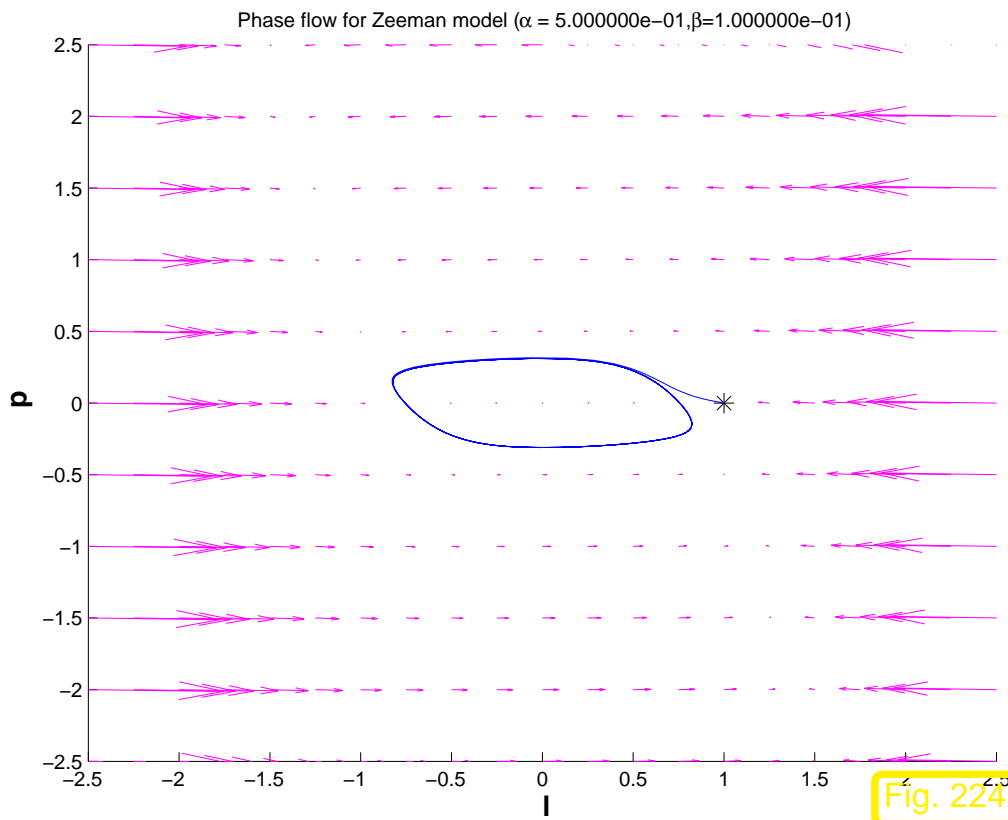
State of heart described by quantities: $l = l(t) \hat{=}$ length of muscle fiber
 $p = p(t) \hat{=}$ electro-chemical potential

Phenomenological model:

$$\begin{aligned} \dot{l} &= -(l^3 - \alpha l + p), \\ \dot{p} &= \beta l, \end{aligned} \tag{12.1.7}$$

with parameters: $\alpha \hat{=}$ pre-tension of muscle fiber
 $\beta \hat{=}$ (phenomenological) feedback parameter





Example 12.1.8 (Transient circuit simulation). [35, Ch. 64]

Transient nodal analysis, cf. Ex. 2.6.3:

Kirchhoff current law

$$i_R(t) - i_L(t) - i_C(t) = 0. \quad (12.1.9)$$

Transient constitutive relations:

$$i_R(t) = R^{-1}u_R(t), \quad (12.1.10)$$

$$i_C(t) = C \frac{du_C}{dt}(t), \quad (12.1.11)$$

$$u_L(t) = L \frac{di_L}{dt}(t). \quad (12.1.12)$$

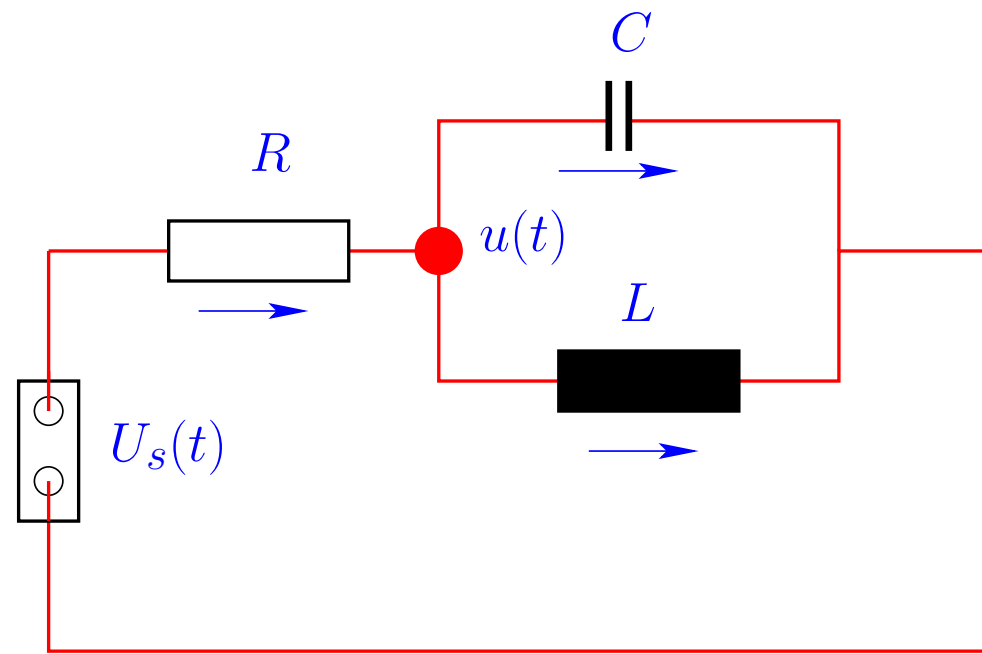


Fig. 226

Given: source voltage $U_s(t)$ ► autonomous **2nd-order** ordinary differential equation:

$$C\ddot{u} + R^{-1}\dot{u} + L^{-1}u = R^{-1}\dot{U}_s.$$



12.1.2 Theory [51, Sect. 11.1], [13, Sect. 11.3]

Abstract mathematical description

Initial value problem (IVP) for first-order ordinary differential equation (ODE):

(\rightarrow [63, Sect. 5.6], [13, Sect. 11.1])

$$\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}) \quad , \quad \mathbf{y}(t_0) = \mathbf{y}_0 \quad . \quad (12.1.13)$$

- $\mathbf{f} : I \times D \mapsto \mathbb{R}^d \hat{=} \text{right hand side (r.h.s.)}$ ($d \in \mathbb{N}$), given in procedural form

function $v = f(t, y)$.

- $I \subset \mathbb{R} \hat{=} \text{(time)interval}$ \leftrightarrow “time variable” t
- $D \subset \mathbb{R}^d \hat{=} \text{state space/phase space}$ \leftrightarrow “state variable” \mathbf{y} (*ger.:* Zustandsraum)
- $\Omega := I \times D \hat{=} \text{extended state space}$ (of tuples (t, \mathbf{y}))
- $t_0 \hat{=} \text{initial time}$, $\mathbf{y}_0 \hat{=} \text{initial state}$ \triangleright **initial conditions**

Remark 12.1.14 (Conversion into autonomous ODE).



Remark 12.1.15 (From higher order ODEs to first order systems). [13, Sect. 11.2]]



Basic assumption: right hand side $\mathbf{f} : I \times D \mapsto \mathbb{R}^d$ locally Lipschitz continuous in \mathbf{y}

Theorem 12.1.21 (Theorem of Peano & Picard-Lindelöf). [2, Satz II(7.6)], [63, Satz 6.5.1], [13, Thm. 11.10], [35, Thm. 73.1]

If $\mathbf{f} : \hat{\Omega} \mapsto \mathbb{R}^d$ is locally Lipschitz continuous (\rightarrow Def. 12.1.19) then for all initial conditions $(t_0, \mathbf{y}_0) \in \hat{\Omega}$ the IVP (12.1.13) 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 Rem. 12.1.14 and Thm. 12.1.21: we consider only

$$\text{autonomous IVP: } \boxed{\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) \quad , \quad \mathbf{y}(0) = \mathbf{y}_0} \quad , \quad (12.1.23)$$

with locally Lipschitz continuous (\rightarrow Def. 12.1.19) right hand side \mathbf{f} .

Assumption 12.1.24 (Global solutions).

All solutions of (12.1.23) are global: $J(\mathbf{y}_0) = \mathbb{R}$ for all $\mathbf{y}_0 \in D$.

Definition 12.1.25 (Evolution operator).

Under Assumption 12.1.24 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** for the autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$.

12.2 Euler methods

12.2.1 Explicit Euler method

Example 12.2.1 (Tangent field and solution curves).

Riccati differential equation

$$\dot{y} = y^2 + t^2 \quad \blacktriangleright \quad d = 1, \quad I, D = \mathbb{R}^+.$$

scalar ODE

(12.2.2)

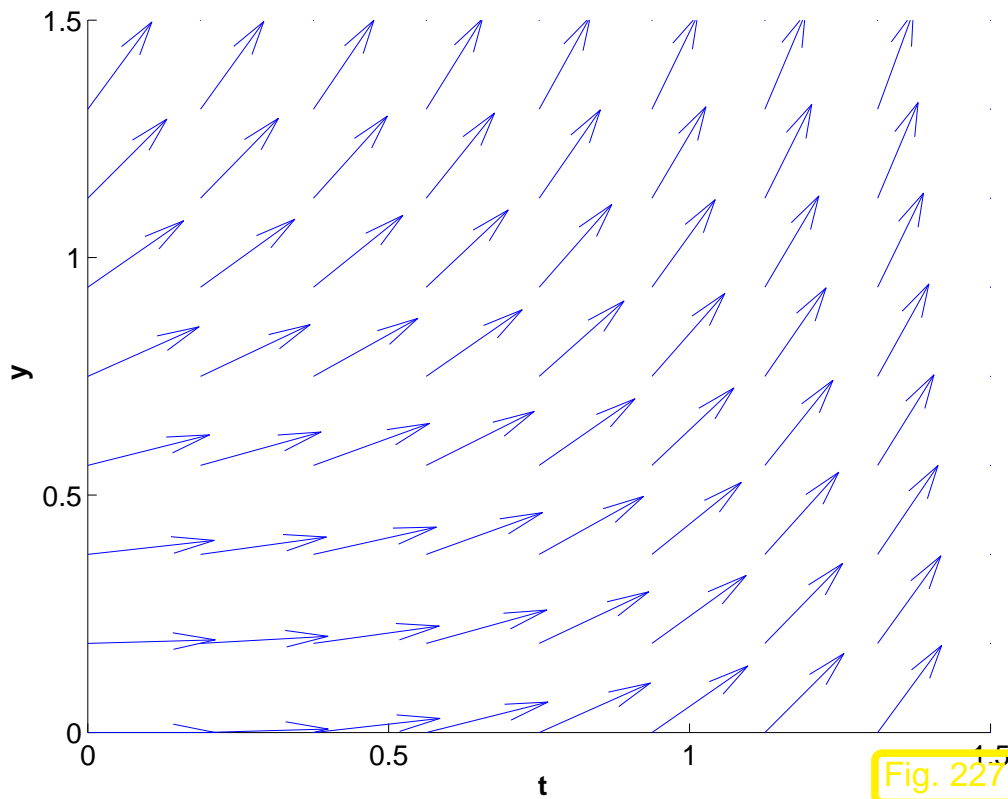


Fig. 227

tangent field

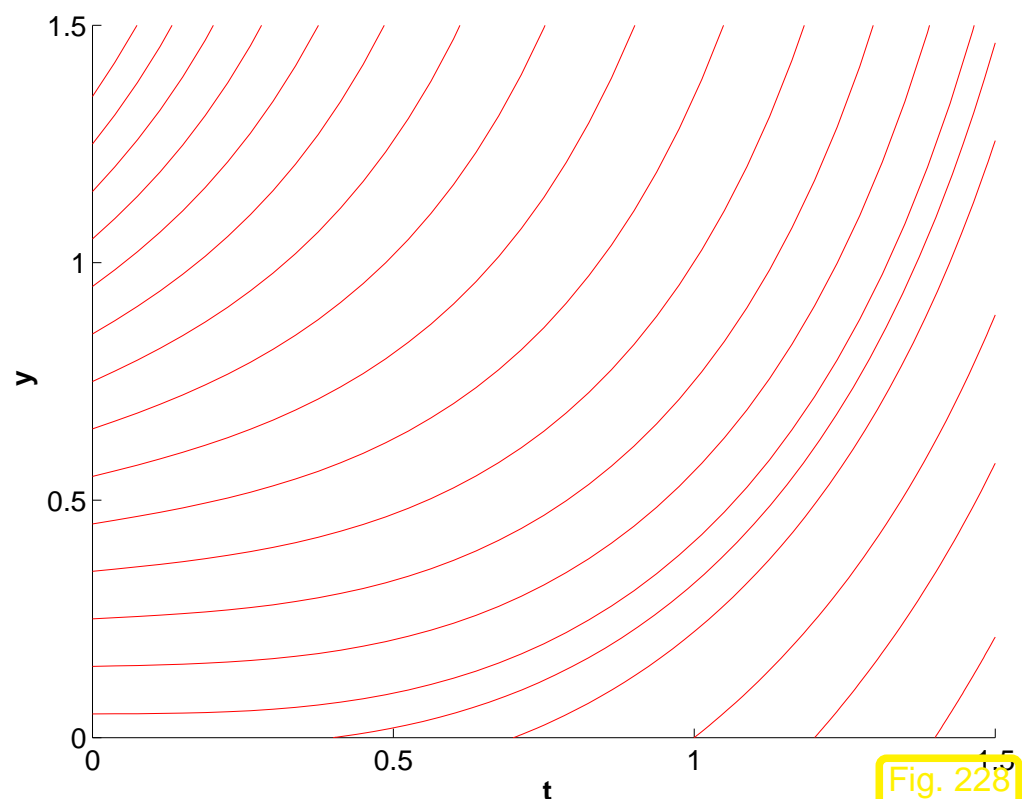
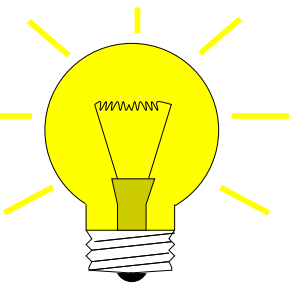


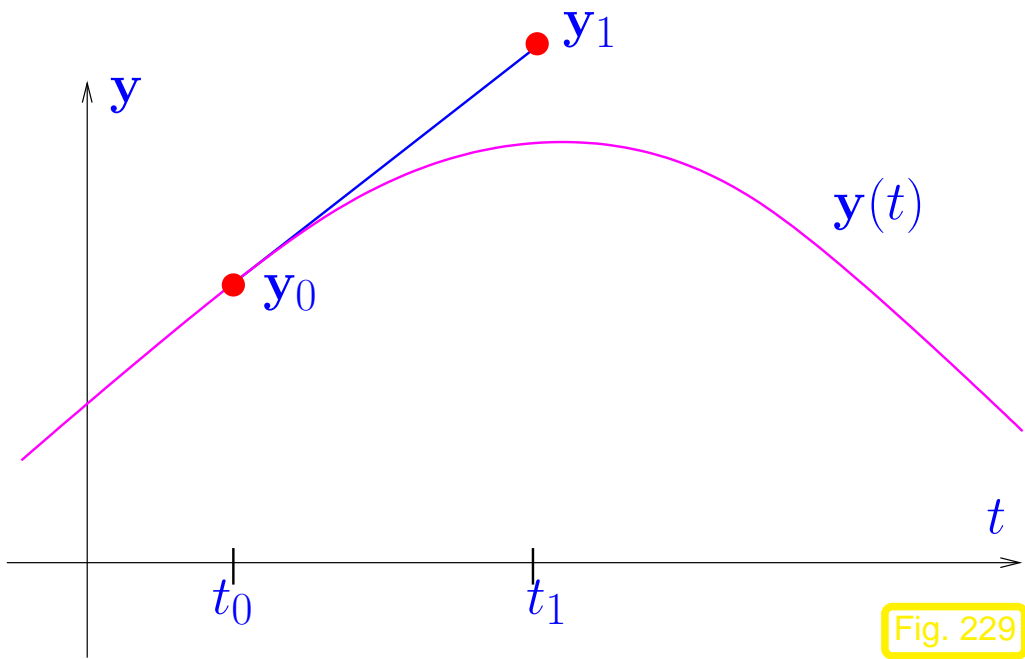
Fig. 228

solution curves



- Idea:
- ❶ **timestepping**: successive approximation of evolution on *small* intervals $[t_{k-1}, t_k]$, $k = 1, \dots, N$, $t_N := T$,
 - ❷ approximation of solution on $[t_{k-1}, t_k]$ by **tangent** curve to current initial condition.





explicit Euler method (Euler 1768)

◁ First step of explicit Euler method ($d = 1$):

Slope of tangent = $f(t_0, y_0)$

y_1 serves as initial value for next step!

Fig. 229

Example 12.2.3 (Visualization of explicit Euler method).

IVP for Riccati differential equation, see Ex. 12.2.1

$$\dot{y} = y^2 + t^2. \quad (12.2.2)$$

Here: $y_0 = \frac{1}{2}, t_0 = 0, T = 1,$
 — $\hat{=}$ “Euler polygon” for uniform timestep $h = 0.2$

$\mapsto \hat{=}$ tangent field of Riccati ODE

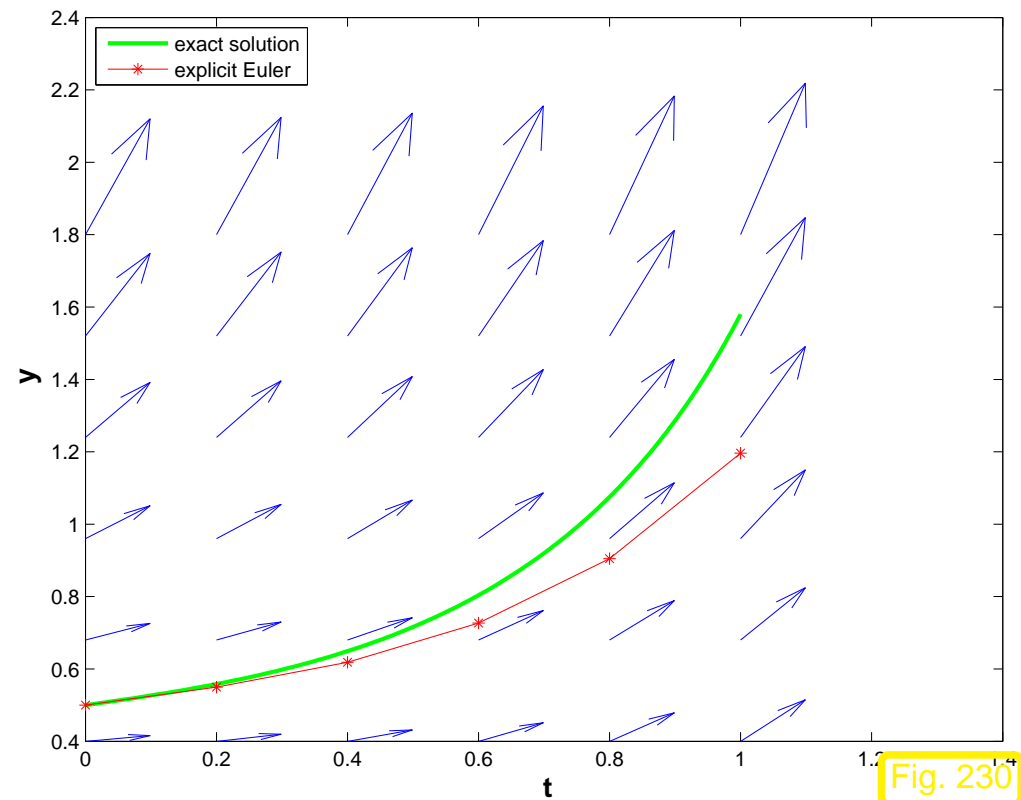


Fig. 230

Formula: explicit Euler method generates a sequence $(\mathbf{y}_k)_{k=0}^N$ by the recursion

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h_k \mathbf{f}(t_k, \mathbf{y}_k), \quad k = 0, \dots, N - 1, \quad (12.2.4)$$

with local (size of) **timestep** (stepsize) $h_k := t_{k+1} - t_k$.

Remark 12.2.5 (Explicit Euler method as **difference scheme**).

$$\mathbf{y}_{k+1} = \mathbf{y}_k + h_k \mathbf{f}(t_{k+1}, \mathbf{y}_{k+1}), \quad k = 0, \dots, N-1, \quad (12.2.8)$$

with local **timestep** (stepsize) $h_k := t_{k+1} - t_k$.

(12.2.8) = **implicit Euler method**

12.2.3 Abstract single step methods

Definition 12.2.12 (Single step method (for autonomous ODE)). \rightarrow [51, 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 $\mathcal{M} := \{t_0 < t_1 < \dots < t_N = T\}$ the recursion

$$\mathbf{y}_{k+1} := \Psi(t_{k+1} - t_k, \mathbf{y}_k), \quad k = 0, \dots, N-1, \quad (12.2.13)$$

defines a **single step method** (SSM, ger.: *Einschrittverfahren*) for the autonomous IVP $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$, $\mathbf{y}(0) = \mathbf{y}_0$.

Procedural view of discrete evolutions:

$$\Psi^h_{\mathbf{y}} \longleftrightarrow \begin{array}{l} \text{function } y1 = \text{esvstep}(h, y0) . \\ (\text{function } y1 = \text{esvstep}(@(\mathbf{y}) \text{ rhs}(\mathbf{y}), h, y0)) \end{array}$$

12.3 Convergence of single step methods [13, Sect. 11.5] [51, Sect. 11.3]

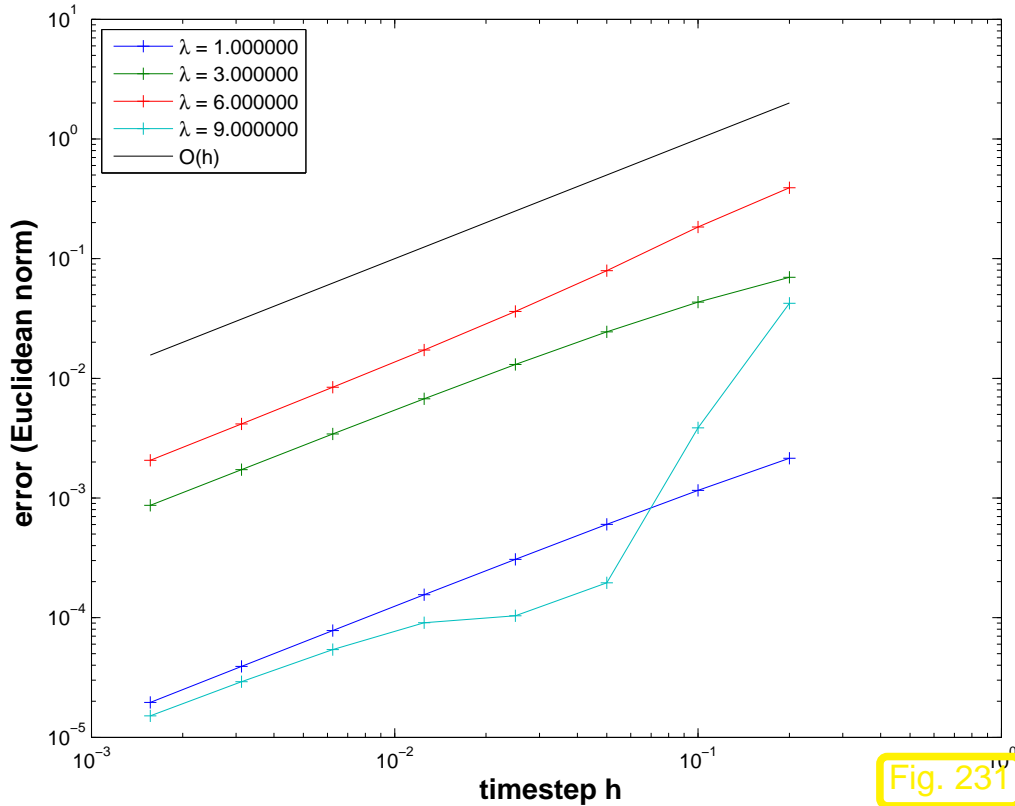
Example 12.3.2 (Speed of convergence of Euler methods).

- IVP for logistic ODE, see Ex. 12.1.1

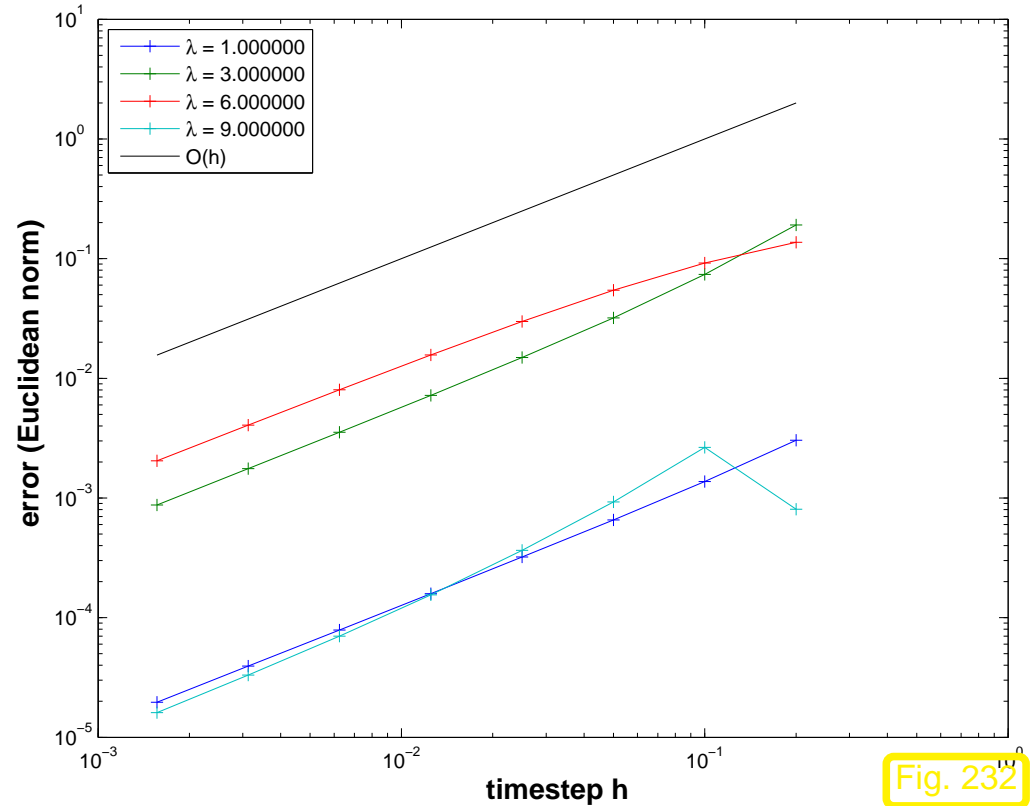
$$\dot{y} = \lambda y(1 - y) \quad , \quad y(0) = 0.01 .$$

- Explicit and implicit Euler methods (12.2.4)/(12.2.8) with uniform timestep $h = 1/N$, $N \in \{5, 10, 20, 40, 80, 160, 320, 640\}$.

- Monitored: Error at final time $E(h) := |y(1) - y_N|$



explicit Euler method



implicit Euler method

$O(h)$ algebraic convergence in both cases

The sequence $(\mathbf{y}_k)_k$ generated by a

single step method (\rightarrow Def. 12.2.12) of **order** (of consistency) $p \in \mathbb{N}$

for $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y})$ on a mesh $\mathcal{M} := \{t_0 < t_1 < \dots < t_N = T\}$ satisfies

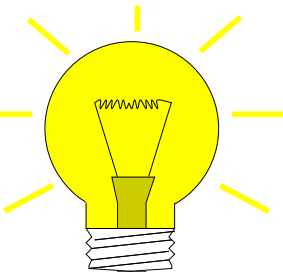
$$\max_k \|\mathbf{y}_k - \mathbf{y}(t_k)\| \leq Ch^p \quad \text{for } h := \max_{k=1, \dots, N} |t_k - t_{k-1}| \rightarrow 0, \quad (12.3.16)$$

with $C > 0$ independent of \mathcal{M} , provided that \mathbf{f} is *sufficiently smooth*, see [13, Thm. 11.25].

12.4 Runge-Kutta methods [13, Sect. 11.6], [35, Ch. 76], [51, Sect.

$$\text{IVP: } \begin{cases} \dot{\mathbf{y}}(t) = \mathbf{f}(t, \mathbf{y}(t)), \\ \mathbf{y}(t_0) = \mathbf{y}_0 \end{cases} \Rightarrow \mathbf{y}(t_1) = \mathbf{y}_0 + \int_{t_0}^{t_1} \mathbf{f}(\tau, \mathbf{y}(\tau)) \, d\tau$$

Idea: approximate integral by means of s -point quadrature formula (\rightarrow Sect. 10.1, defined on reference interval $[0, 1]$) with nodes c_1, \dots, c_s , weights b_1, \dots, b_s .



$$y(t_1) \approx y_1 = y_0 + h \sum_{i=1}^s b_i f(t_0 + c_i h, \boxed{y(t_0 + c_i h)}), \quad h := t_1 - t_0. \quad (12.4.2)$$

Obtain these values by **bootstrapping**

Example 12.4.3 (Construction of simple Runge-Kutta methods).

Quadrature formula = trapezoidal rule (10.2.7):

$$Q(f) = \frac{1}{2}(f(0) + f(1)) \quad \leftrightarrow \quad s = 2: \quad c_1 = 0, c_2 = 1, \quad b_1 = b_2 = \frac{1}{2}, \quad (12.4.4)$$

and $y(T)$ approximated by explicit Euler step (12.2.4)

$$\mathbf{k}_1 = \mathbf{f}(t_0, \mathbf{y}_0), \quad \mathbf{k}_2 = \mathbf{f}(t_0 + h, \mathbf{y}_0 + h\mathbf{k}_1), \quad \mathbf{y}_1 = \mathbf{y}_0 + \frac{h}{2}(\mathbf{k}_1 + \mathbf{k}_2). \quad (12.4.5)$$

(12.4.5) = **explicit trapezoidal rule** (for numerical integration of ODEs).

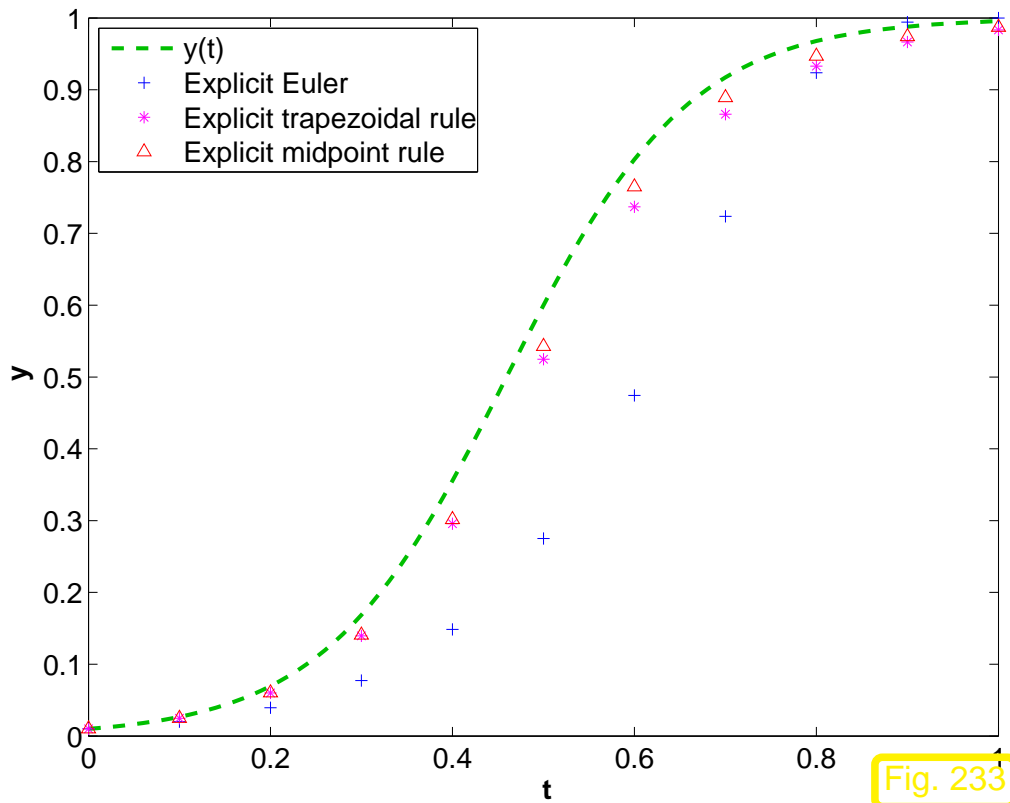
Quadrature formula \rightarrow simplest Gauss quadrature formula = midpoint rule (\rightarrow Ex. 10.2.5) & $y(\frac{1}{2}(t_1 - t_0))$ approximated by explicit Euler step (12.2.4)

$$\mathbf{k}_1 = \mathbf{f}(t_0, \mathbf{y}_0), \quad \mathbf{k}_2 = \mathbf{f}(t_0 + \frac{h}{2}, \mathbf{y}_0 + \frac{h}{2}\mathbf{k}_1), \quad \mathbf{y}_1 = \mathbf{y}_0 + h\mathbf{k}_2. \quad (12.4.6)$$

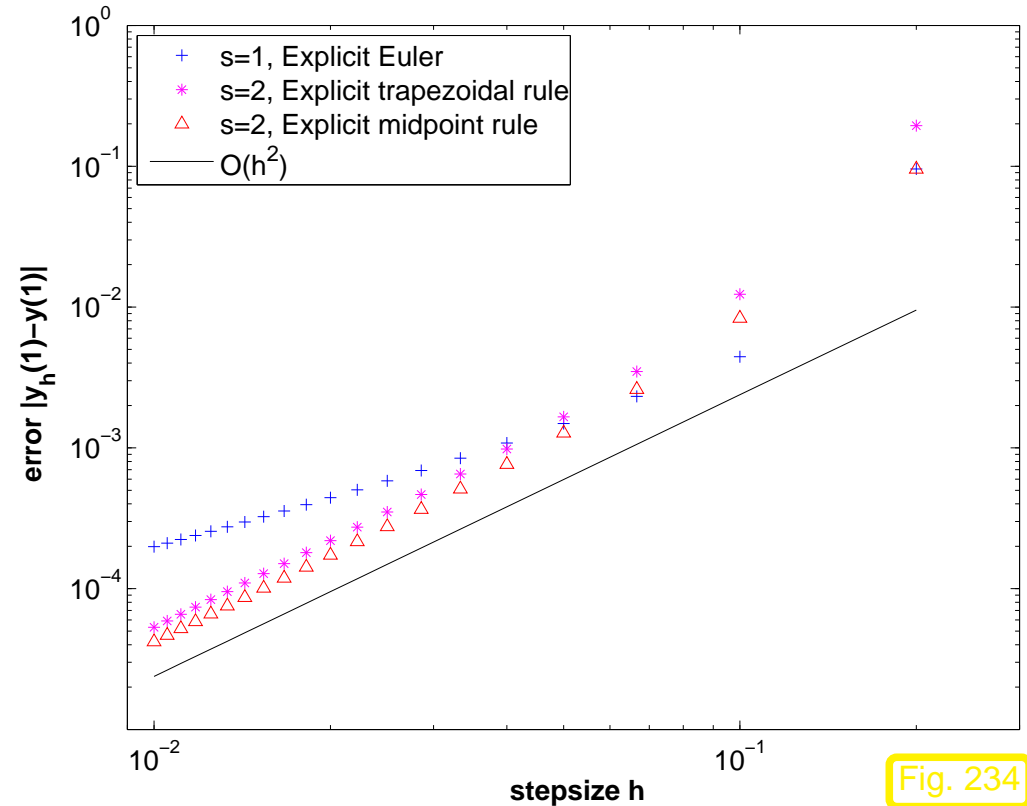
(12.4.6) = **explicit midpoint rule** (for numerical integration of ODEs) [13, Alg. 11.18].

Example 12.4.7 (Convergence of simple Runge-Kutta methods).

- IVP: $\dot{y} = 10y(1 - y)$ (logistic ODE (12.1.2)), $y(0) = 0.01$, $T = 1$,
- Explicit single step methods, uniform timestep h .



$y_h(j/10)$, $j = 1, \dots, 10$ for explicit RK-methods



Errors at final time $y_h(1) - y(1)$



Definition 12.4.8 (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 IVP (12.1.13) is defined by

$$\mathbf{k}_i := \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^{i-1} a_{ij} \mathbf{k}_j), \quad i = 1, \dots, s, \quad \mathbf{y}_1 := \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i.$$

The $\mathbf{k}_i \in \mathbb{R}^d$ are called *increments*.

Shorthand notation for (explicit) Runge-Kutta methods [13, (11.75)]

Butcher scheme

▷

$$\begin{array}{c|c} \mathbf{c} & \mathfrak{A} \\ \hline & \mathbf{b}^T \end{array} :=$$

$$\begin{array}{c|ccc} c_1 & 0 & \cdots & 0 \\ c_2 & a_{21} & \cdots & \vdots \\ \vdots & \vdots & & \ddots \\ c_s & a_{s1} & \cdots & a_{s,s-1} & 0 \\ \hline & b_1 & \cdots & b_s & \end{array} . \quad (12.4.9)$$

(Note: \mathfrak{A} is strictly lower triangular $s \times s$ -matrix)

Example 12.4.10 (Butcher scheme for some explicit RK-SSM). [13, Sect. 11.6.1]

- Explicit Euler method (12.2.4):

$$\begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array}$$



order = 1

- explicit trapezoidal rule (12.4.5):

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}$$



order = 2

- explicit midpoint rule (12.4.6):

$$\begin{array}{c|cc} 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \hline & 0 & 1 \end{array}$$



order = 2

- Classical 4th-order RK-SSM:

$$\begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ \hline & \frac{1}{6} & \frac{2}{6} & \frac{2}{6} & \frac{1}{6} \end{array}$$



order = 4

- Kutta's 3/8-rule:

$$\begin{array}{c|cccc}
 0 & 0 & 0 & 0 & 0 \\
 \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\
 \frac{2}{3} & -\frac{1}{3} & 1 & 0 & 0 \\
 1 & 1 & -1 & 1 & 0 \\
 \hline
 & \frac{1}{8} & \frac{3}{8} & \frac{3}{8} & \frac{1}{8}
 \end{array}
 \quad \triangleright \quad \text{order} = 4$$



Remark 12.4.11 (“Butcher barriers” for explicit RK-SSM).

| | | | | | | | | | |
|---------------------------|---|---|---|---|---|---|---|----|--------------|
| order p | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ≥ 9 |
| minimal no. s of stages | 1 | 2 | 3 | 4 | 6 | 7 | 9 | 11 | $\geq p + 3$ |



Remark 12.4.12 (Explicit ODE integrator in MATLAB).

Syntax:

```

[t, y] = ode45(odefun, tspan, y0);

```

odefun : **Handle** to a function of type $@(t, y) \leftrightarrow \text{r.h.s. } \mathbf{f}(t, \mathbf{y})$
 tspan : vector $(t_0, T)^T$, initial and final time for numerical integration
 y0 : (vector) passing initial state $\mathbf{y}_0 \in \mathbb{R}^d$

Return values:

t : temporal mesh $\{t_0 < t_1 < t_2 < \dots < t_{N-1} = t_N = T\}$
 y : sequence $(\mathbf{y}_k)_{k=0}^N$ (column vectors)

Code 12.4.13: code excerpts for MATLAB integrator ode45

```

1 function varargout = ode45(ode,tspan,y0,options,varargin)
2 % Processing of input parameters omitted
3 % :
4 % Initialize method parameters, c.f. Butcher scheme (12.4.9)
5 pow = 1/5;
6 A = [1/5, 3/10, 4/5, 8/9, 1, 1];
7 B = [
8     1/5      3/40      44/45      19372/6561      9017/3168      35/384
9     0        9/40      -56/15     -25360/2187     -355/33        0
10    0         0        32/9      64448/6561     46732/5247     500/1113
11    0         0         0        -212/729      49/176        125/192
12    0         0         0         0        -5103/18656   -2187/6784
13    0         0         0         0         0            11/84
14    0         0         0         0         0            0
15    ];
16 E = [71/57600; 0; -71/16695; 71/1920; -17253/339200; 22/525; -1/40];
17 % : (choice of stepsize and main loop omitted)

```

```

8  % ADVANCING ONE STEP.
9  hA = h * A;
10 hB = h * B;
11 f(:,2) = feval(odeFcn,t+hA(1),y+f*hB(:,1),odeArgs{:});
12 f(:,3) = feval(odeFcn,t+hA(2),y+f*hB(:,2),odeArgs{:});
13 f(:,4) = feval(odeFcn,t+hA(3),y+f*hB(:,3),odeArgs{:});
14 f(:,5) = feval(odeFcn,t+hA(4),y+f*hB(:,4),odeArgs{:});
15 f(:,6) = feval(odeFcn,t+hA(5),y+f*hB(:,5),odeArgs{:});
16
17 tnew = t + hA(6);
18 if done, tnew = tfinal; end % Hit end point exactly.
19 h = tnew - t; % Purify h.
20 ynew = y + f*hB(:,6);
21 % : (stepsize control, see Sect. 12.5 dropped

```



Example 12.4.14 (Numerical integration of logistic ODE in MATLAB).

MATLAB-CODE: usage of ode45

```

fn = @(t,y) 5*y*(1-y);
[t,y] = ode45(fn,[0 1.5],y0);
plot(t,y,'r-');

```

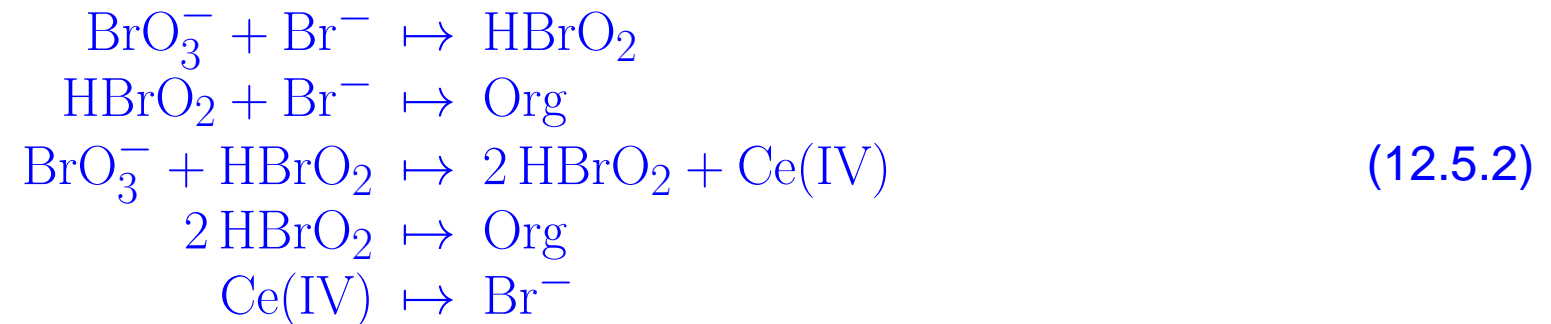
MATLAB-integrator: `ode45()`:

- Handle passing r.h.s.
- initial and final time
- initial state y_0



12.5 Stepsize control [13, Sect. 11.7], [51, Sect. 11.8.2]

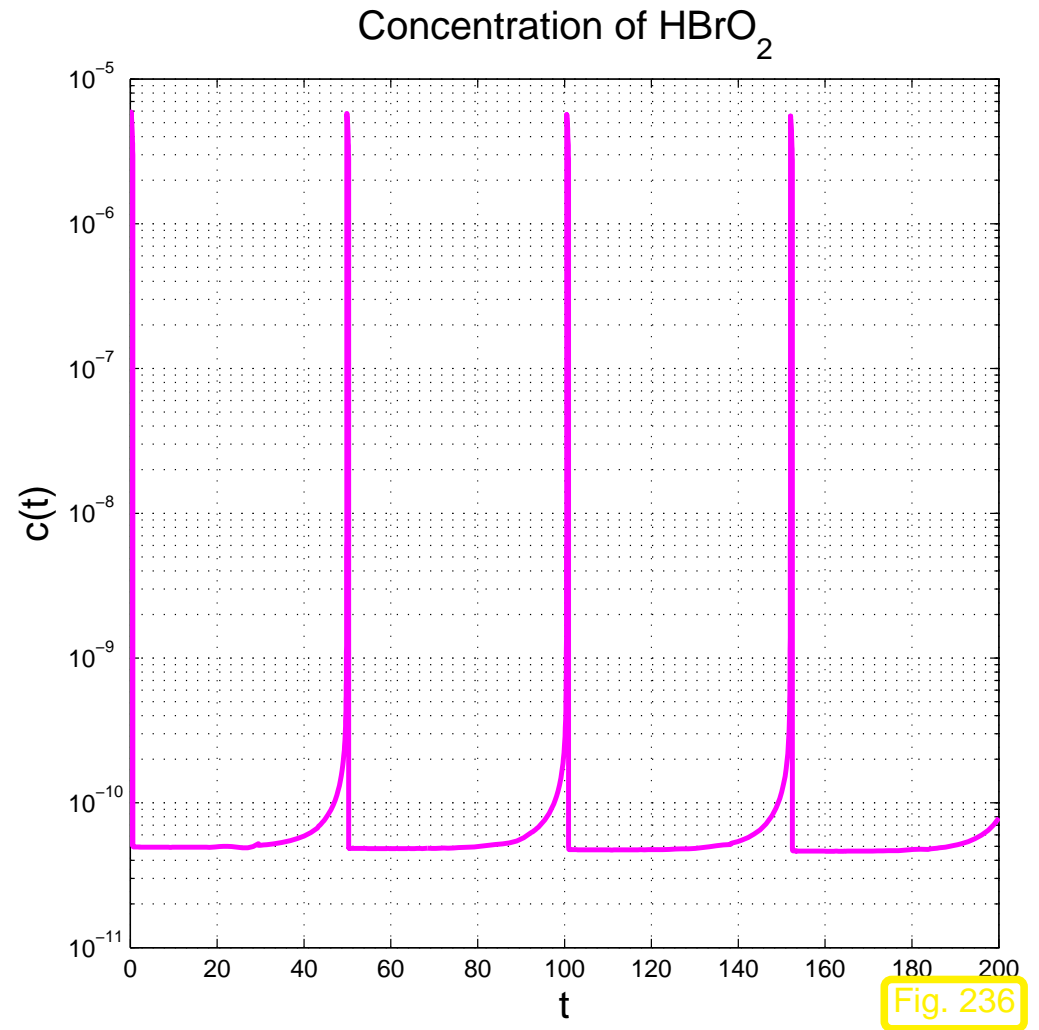
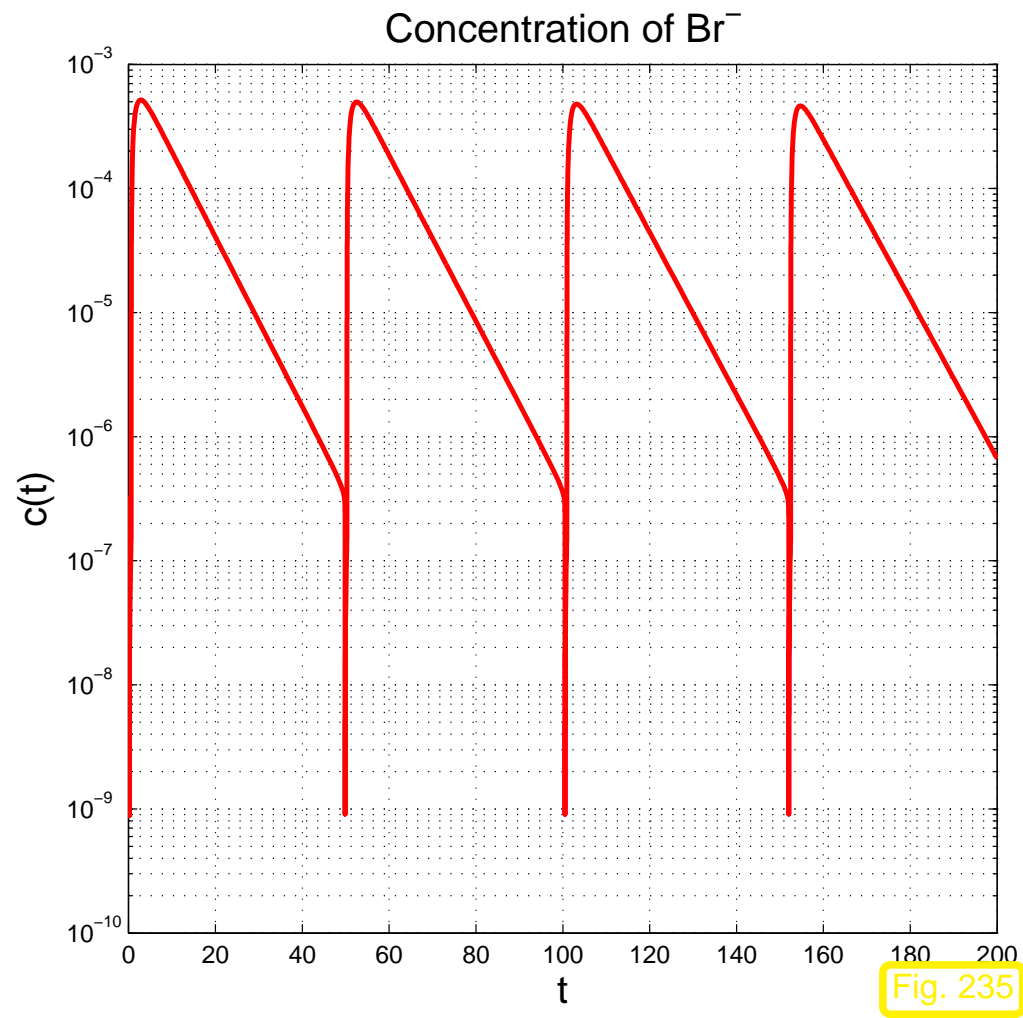
Example 12.5.1 (Oregonator reaction).



$$\begin{aligned} y_1 &:= c(\text{BrO}_3^-): & \dot{y}_1 &= -k_1 y_1 y_2 - k_3 y_1 y_3, \\ y_2 &:= c(\text{Br}^-): & \dot{y}_2 &= -k_1 y_1 y_2 - k_2 y_2 y_3 + k_5 y_5, \\ y_3 &:= c(\text{HBrO}_2): & \dot{y}_3 &= k_1 y_1 y_2 - k_2 y_2 y_3 + k_3 y_1 y_3 - 2k_4 y_3^2, \\ y_4 &:= c(\text{Org}): & \dot{y}_4 &= k_2 y_2 y_3 + k_4 y_3^2, \\ y_5 &:= c(\text{Ce(IV)}): & \dot{y}_5 &= k_3 y_1 y_3 - k_5 y_5, \end{aligned} \tag{12.5.3}$$

with (non-dimensionalized) reaction constants:

$$k_1 = 1.34, \quad k_2 = 1.6 \cdot 10^9, \quad k_3 = 8.0 \cdot 10^3, \quad k_4 = 4.0 \cdot 10^7, \quad k_5 = 1.0.$$

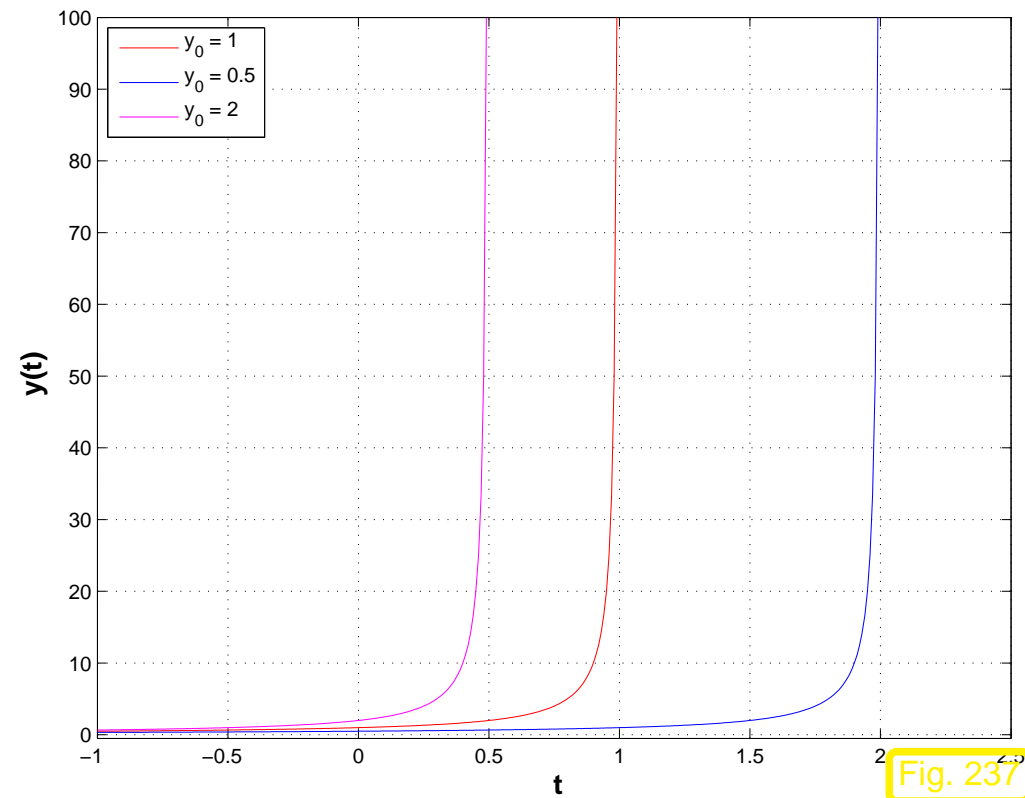


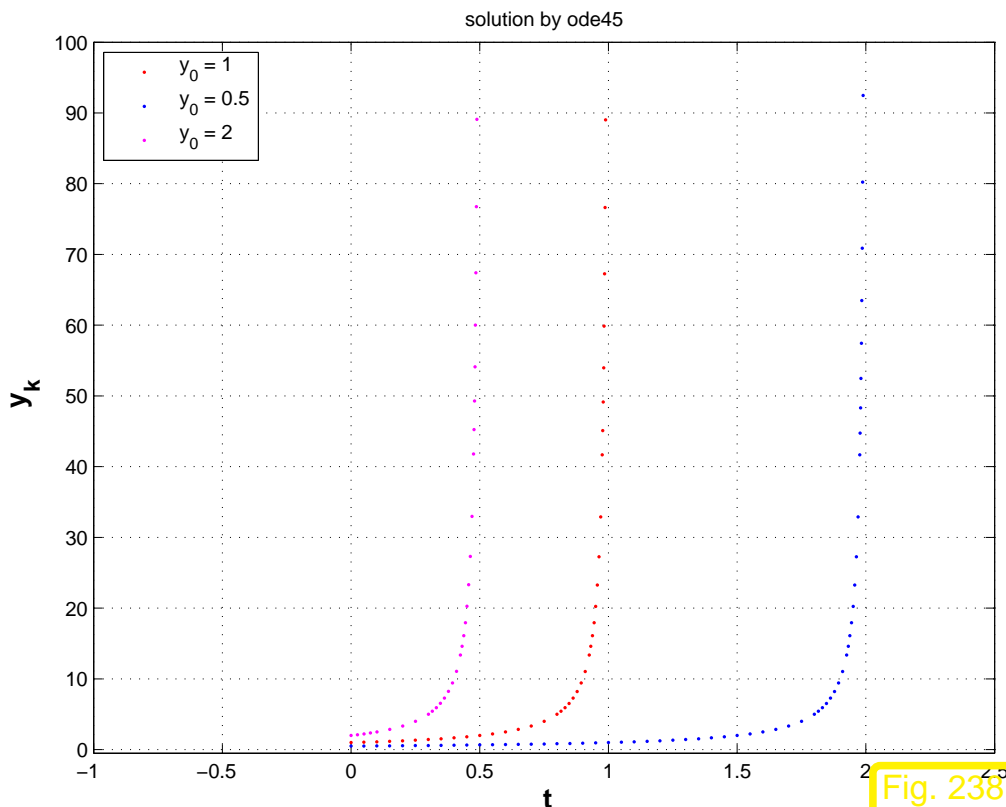
Example 12.5.4 (Blow-up).

Scalar autonomous IVP:

$$\dot{y} = y^2, \quad y(0) = y_0 > 0.$$

$$\blacktriangleright \quad y(t) = \frac{y_0}{1 - y_0 t}, \quad t < 1/y_0.$$





```

1 fun = @(t,y) y.^2;
2 [t1,y1] = ode45(fun,[0 2],1);
3 [t2,y2] = ode45(fun,[0 2],0.5);
4 [t3,y3] = ode45(fun,[0 2],2);

```

MATLAB warning messages:

```
Warning: Failure at t=9.999694e-01. Unable to meet integration
tolerances without reducing the step size below the smallest
value allowed (1.776357e-15) at time t.
```

```
> In ode45 at 371
   In simpleblowup at 22
```

```
Warning: Failure at t=1.999970e+00. Unable to meet integration
tolerances without reducing the step size below the smallest
value allowed (3.552714e-15) at time t.
```

```
> In ode45 at 371
   In simpleblowup at 23
```

```
Warning: Failure at t=4.999660e-01. Unable to meet integration
```


tolerances without reducing the step size below the smallest value allowed (8.881784e-16) at time t.

```
> In ode45 at 371
In simpleblowup at 24
```



be efficient

be accurate

Objective: N as small as possible & $\max_{k=1,\dots,N} \|\mathbf{y}(t_k) - \mathbf{y}_k\| < \text{TOL}$, TOL = tolerance
 or $\|\mathbf{y}(T) - \mathbf{y}_N\| < \text{TOL}$

Policy: Try to curb/balance **one-step error** by

- adjusting *current* stepsize h_k ,
- predicting suitable *next* timestep h_{k+1}

} **local-in-time
stepsize control**

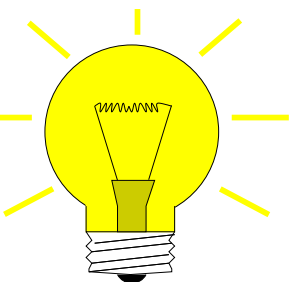
Tool: **Local-in-time** one-step error estimator (*a posteriori*, based on \mathbf{y}_k, h_{k-1})

Idea:

Estimation of one-step error, cf. Sect. 10.5

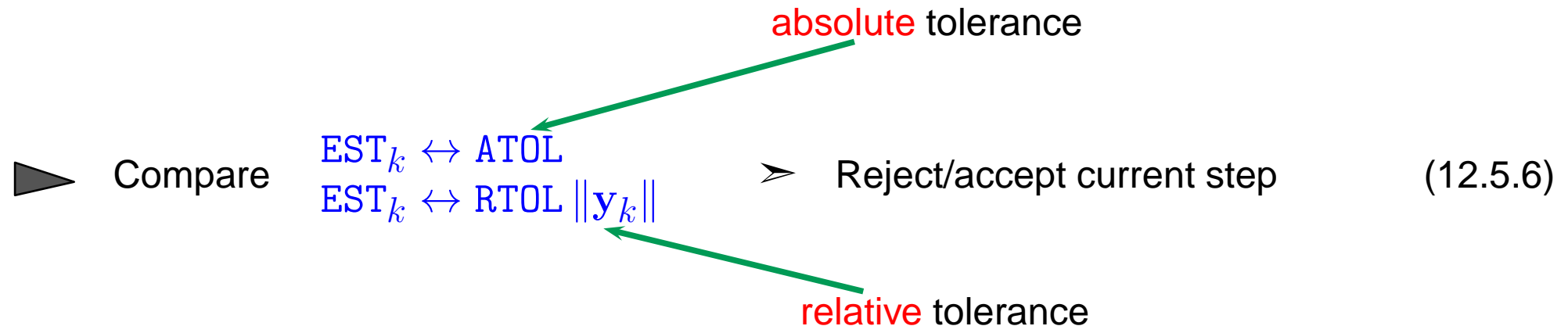
Compare two discrete evolutions $\Psi^h, \tilde{\Psi}^h$ of different order
 for *current timestep* h :

If $\text{Order}(\tilde{\Psi}) > \text{Order}(\Psi)$



$$\Rightarrow \underbrace{\Phi^h y(t_k) - \Psi^h y(t_k)}_{\text{one-step error}} \approx \text{EST}_k := \tilde{\Psi}^h y(t_k) - \Psi^h y(t_k). \quad (12.5.5)$$

Heuristics for concrete h



Code 12.5.7: simple local stepsize control for single step methods

```

1 function [t,y] =
   odeintadapt(Psilow,Psihigh,T,y0,h0,reltol,abstol,hmin)
2 t = 0; y = y0; h = h0; %
3 while ((t(end) < T) (h > hmin)) %
4   yh = Psihigh(h,y0); % high order discrete evolution  $\tilde{\Psi}^h$ 
5   yH = Psilow(h,y0); % low order discrete evolution  $\Psi^h$ 
6   est = norm(yH-yh); %  $\leftrightarrow \text{EST}_k$ 
7
8   if (est < max(reltol*norm(y0),abstol)) %

```

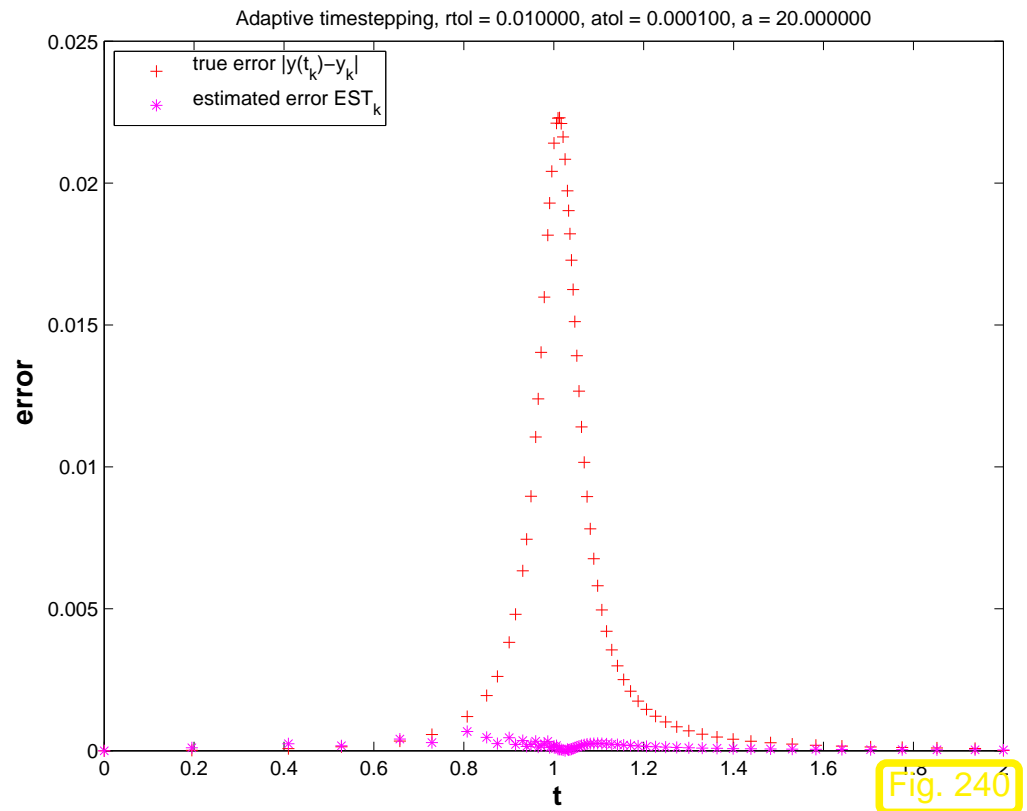
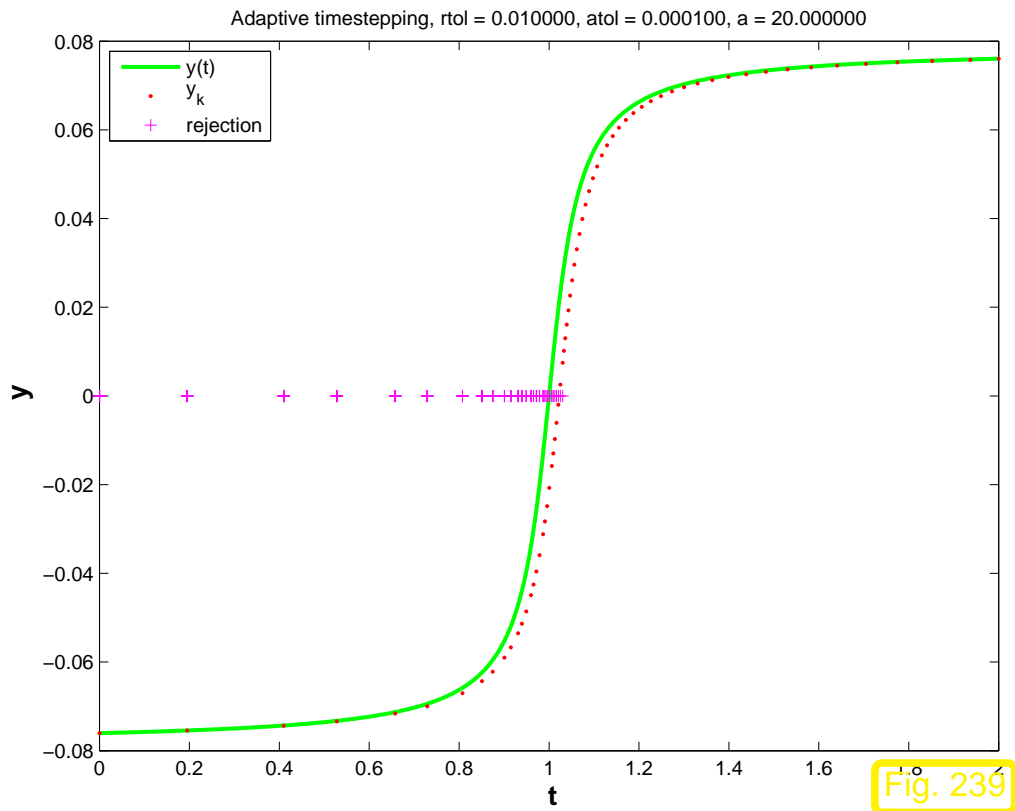
```

9   y0 = yh; y = [y,y0]; t = [t,t(end) + min(T-t(end),h)]; %
10  h = 1.1*h; % step accepted, try with increased stepsize
11  else , h = h/2; end % step rejected, try with half the stepsize
12 end

```

Example 12.5.8 (Simple adaptive stepsize control).

- 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 (12.2.4) and explicit trapezoidal rule (12.4.5)





Example 12.5.10 (Gain through adaptivity). → Ex. 12.5.8

Solving $d_t y = a \cos(y)^2$ with $a = 40.000000$ by simple adaptive timestepping

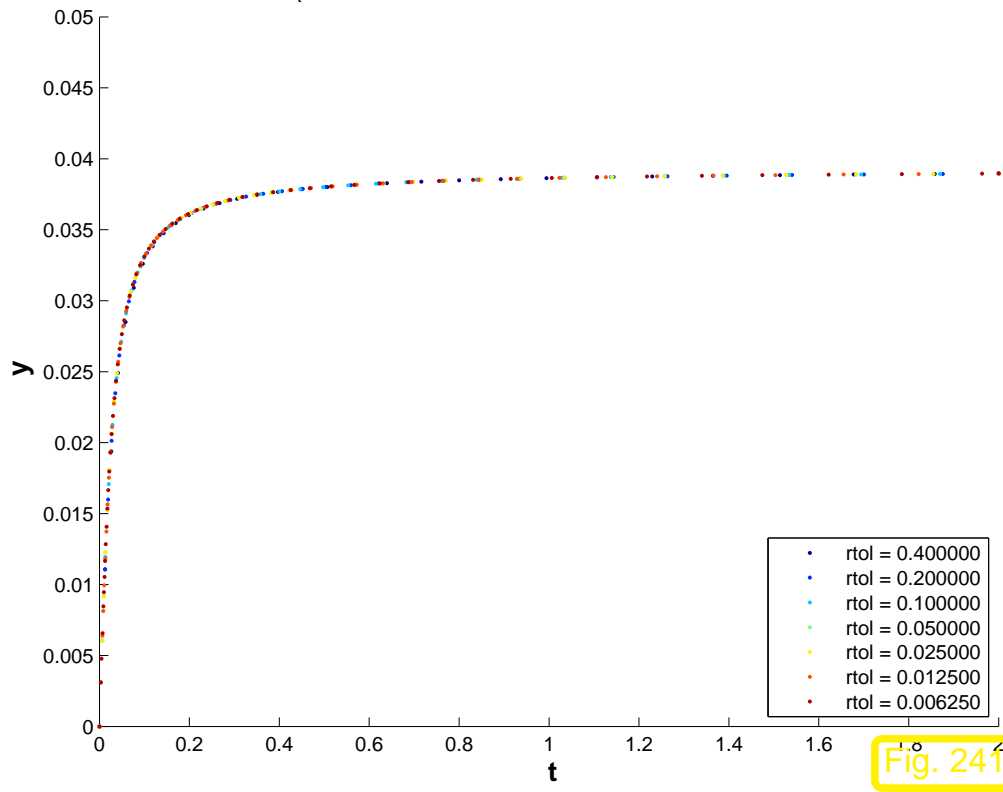


Fig. 241

Solutions $(y_k)_k$ for different values of $rtol$

Error vs. no. of timesteps for $d_t y = a \cos(y)^2$ with $a = 40.000000$

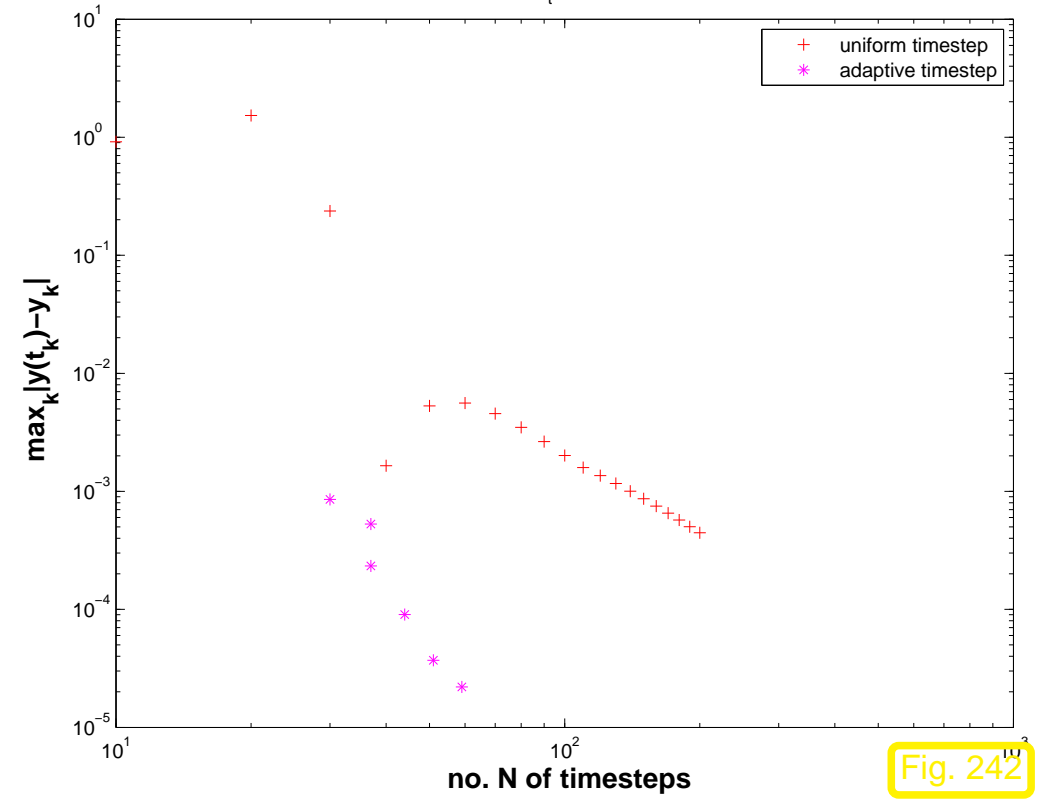


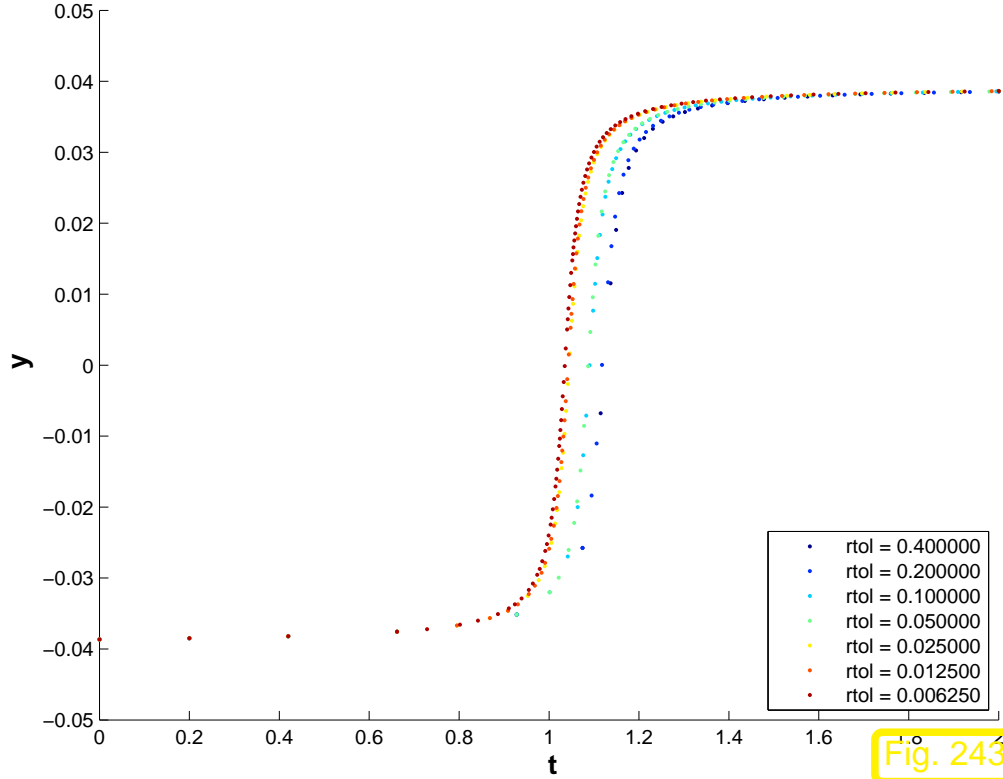
Fig. 242

Error vs. computational effort



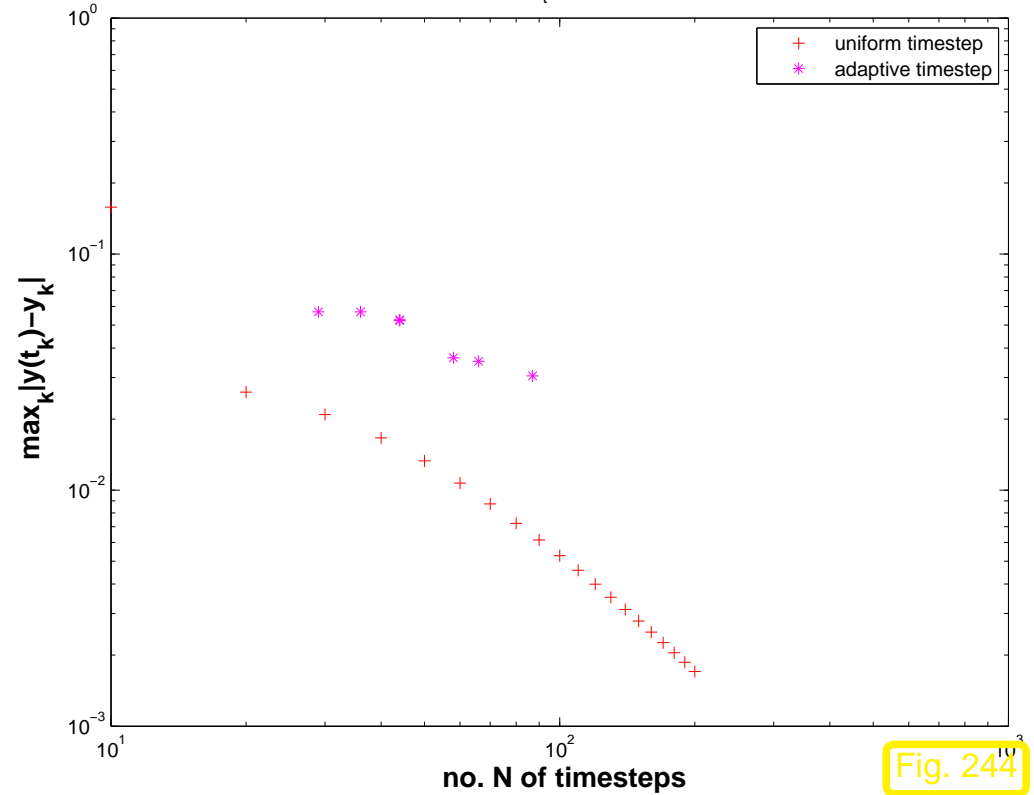
Example 12.5.12 (“Failure” of adaptive timestepping). → Ex. 12.5.10

Solving $d_t y = a \cos(y)^2$ with $a = 40.000000$ by simple adaptive timestepping



Solutions $(y_k)_k$ for different values of `rtol`

Error vs. no. of timesteps for $d_t y = a \cos(y)^2$ with $a = 40.000000$



Error vs. computational effort



Remark 12.5.20 (Stepsize control in MATLAB).



Specifying tolerances for MATLAB's integrators:

```
options = odeset('abstol',atol,'reltol',rtol,'stats','on');
[t,y] = ode45(@(t,x) f(t,x),tspan,y0,options);
(f = function handle, tspan  $\hat{=}$   $[t_0, T]$ ,  $y_0 \hat{=}$   $\mathbf{y}_0$ ,  $t \hat{=}$   $t_k$ ,  $y \hat{=}$   $\mathbf{y}_k$ )
```

△ R. Hiptmair
rev 38355,
May 3,
2011

Example 12.5.21 (Adaptive timestepping for mechanical problem).

Movement of a point mass in a conservative force field: $t \mapsto \mathbf{y}(t) \in \mathbb{R}^2 \hat{=}$ trajectory

Newton's law: $\ddot{\mathbf{y}} = F(\mathbf{y}) := -\frac{2\mathbf{y}}{\|\mathbf{y}\|_2^2}.$ (12.5.22)

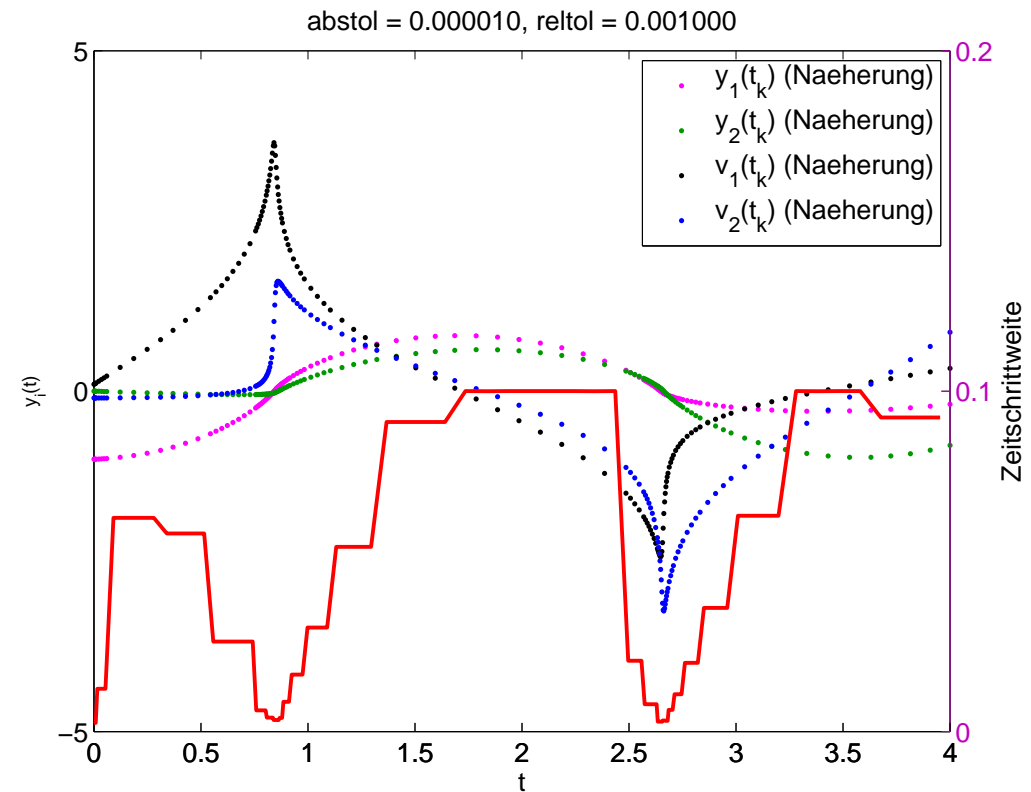
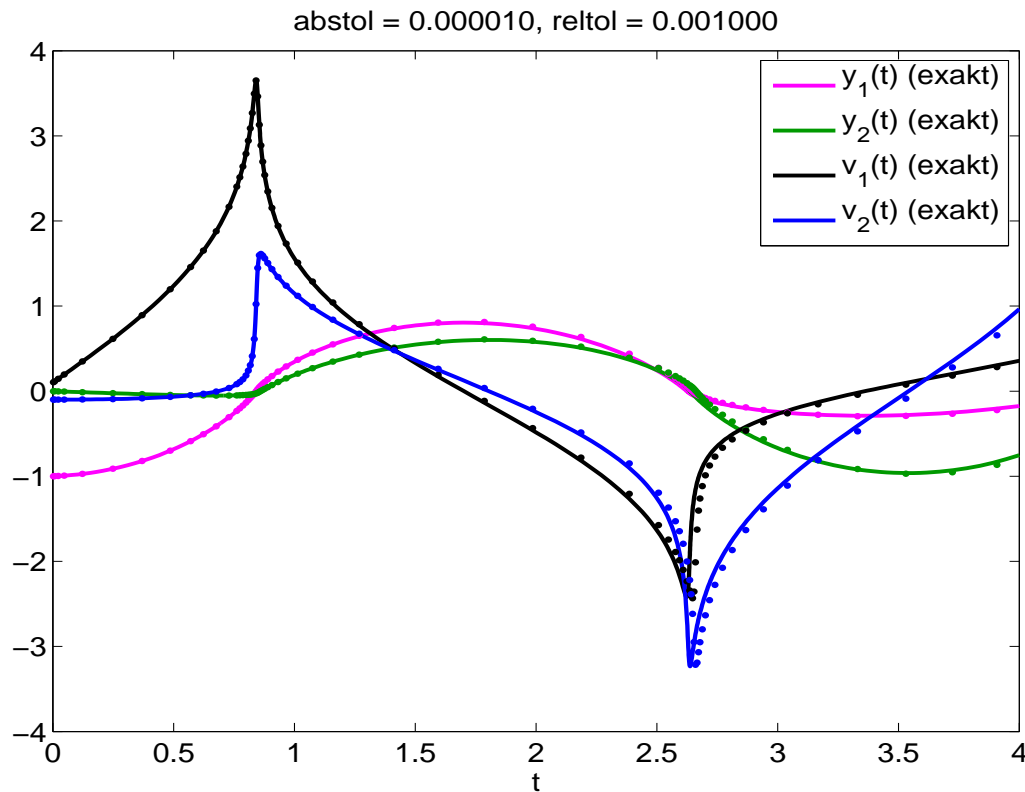
acceleration $\rightarrow \ddot{\mathbf{y}}$ $F(\mathbf{y})$ \leftarrow force

Equivalent 1st-order ODE, see Rem. 12.1.15: with velocity $\mathbf{v} := \dot{\mathbf{y}}$

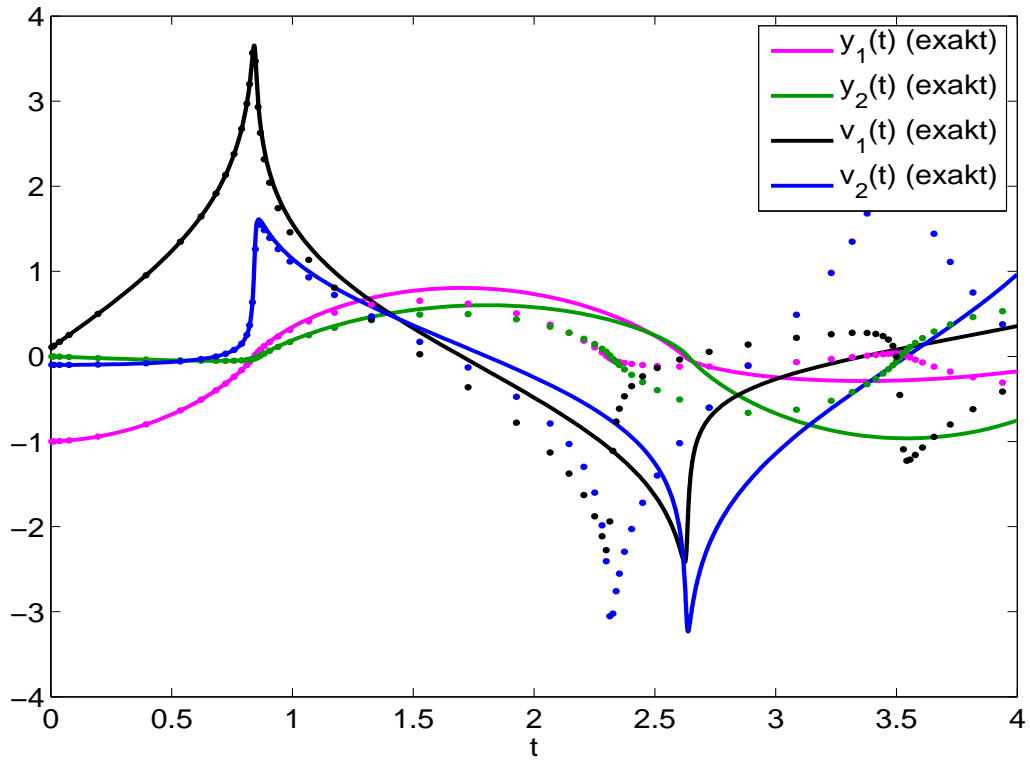
$$\begin{pmatrix} \dot{\mathbf{y}} \\ \dot{\mathbf{v}} \end{pmatrix} = \begin{pmatrix} \mathbf{v} \\ -\frac{2\mathbf{y}}{\|\mathbf{y}\|_2^2} \end{pmatrix}. \quad (12.5.23)$$

Adaptive integrator: `ode45(@(t,x) f(t,x), [0 4], [-1;0;0.1;-0.1], options):`

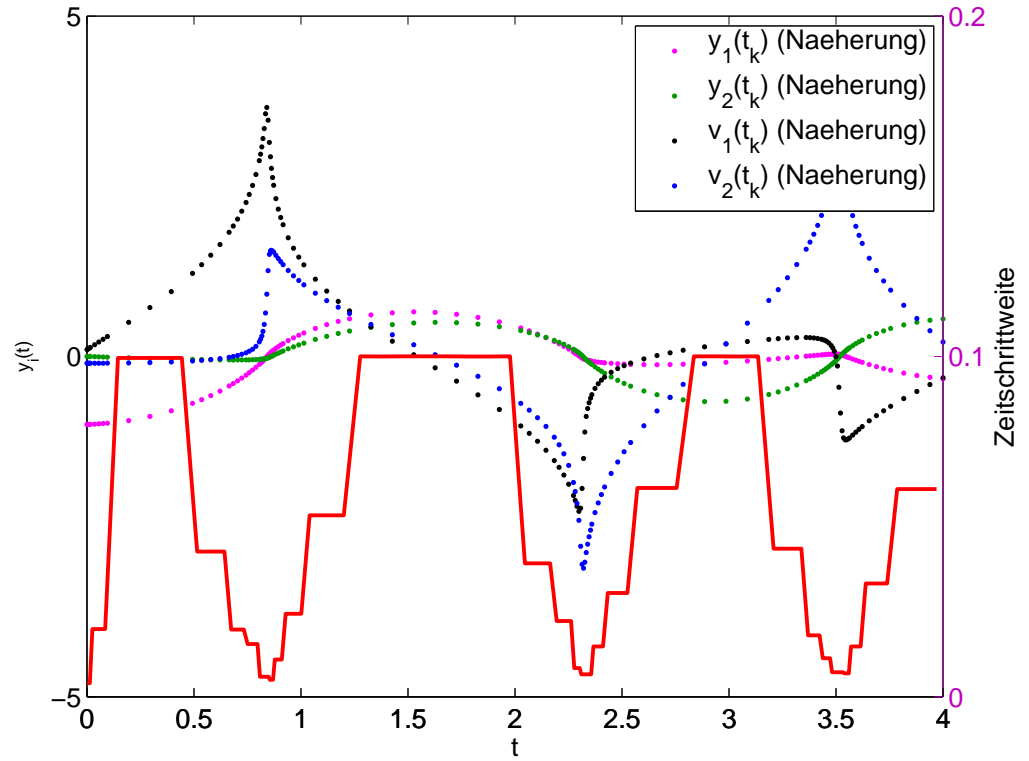
- ❶ `options = odeset('reltol',0.001,'abstol',1e-5);`
- ❷ `options = odeset('reltol',0.01,'abstol',1e-3);`

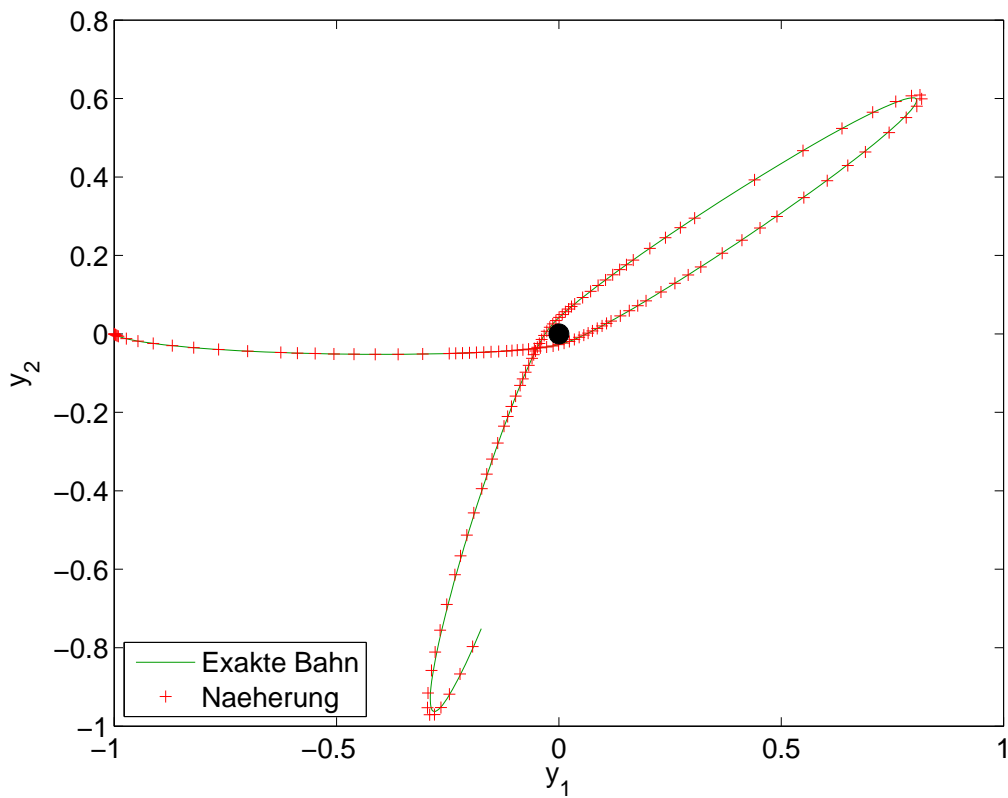


abstol = 0.001000, reltol = 0.010000

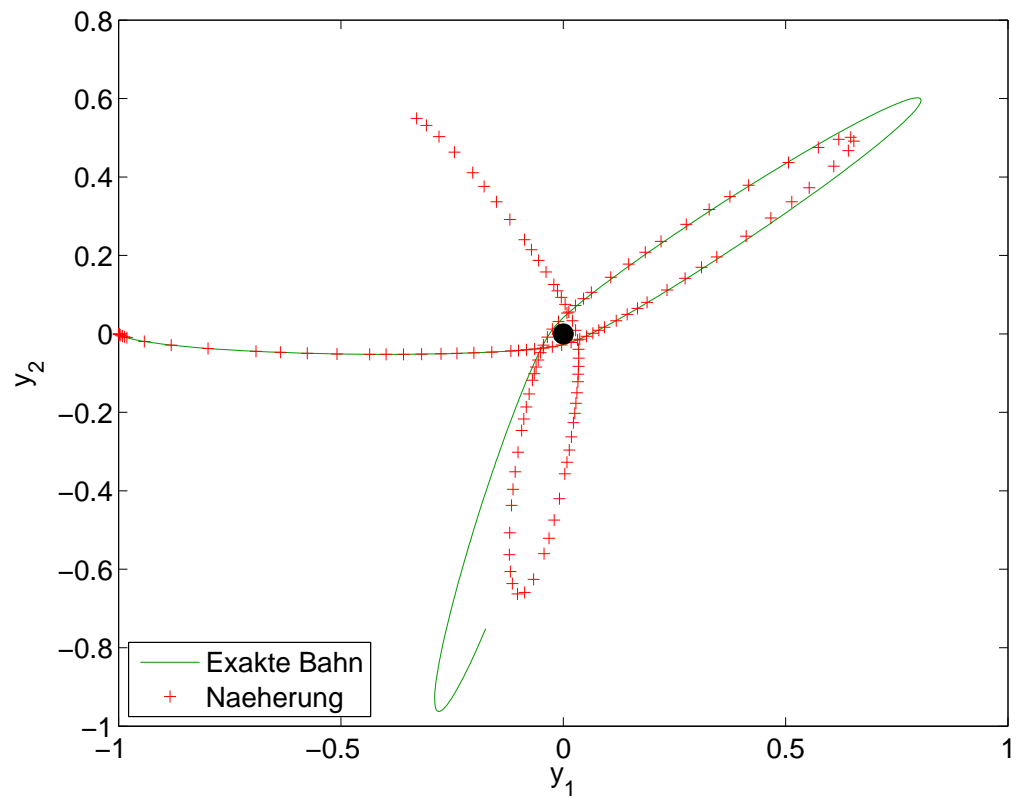


abstol = 0.001000, reltol = 0.010000





reltol=0.001, abstol=1e-5



reltol=0.01, abstol=1e-3



Absolute/relative tolerances do *not* allow to predict accuracy of solution!

13

Stiff Integrators [13, Sect. 11.9]

Example 13.0.1 (ode45 for stiff problem).

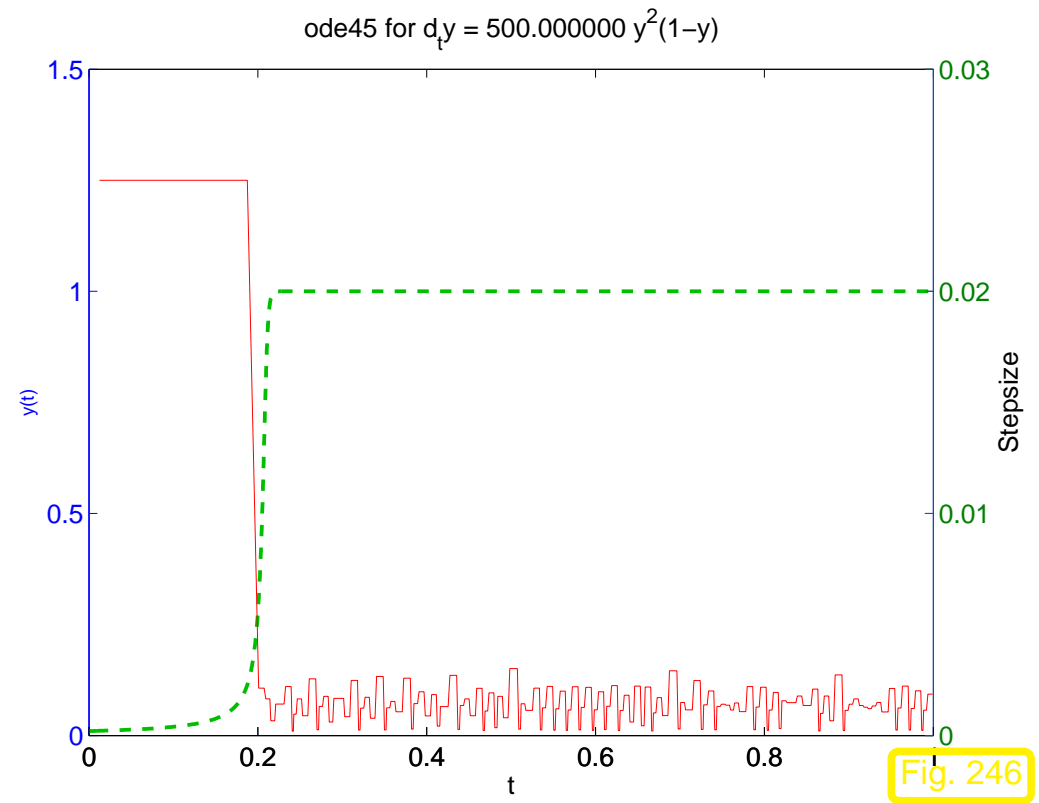
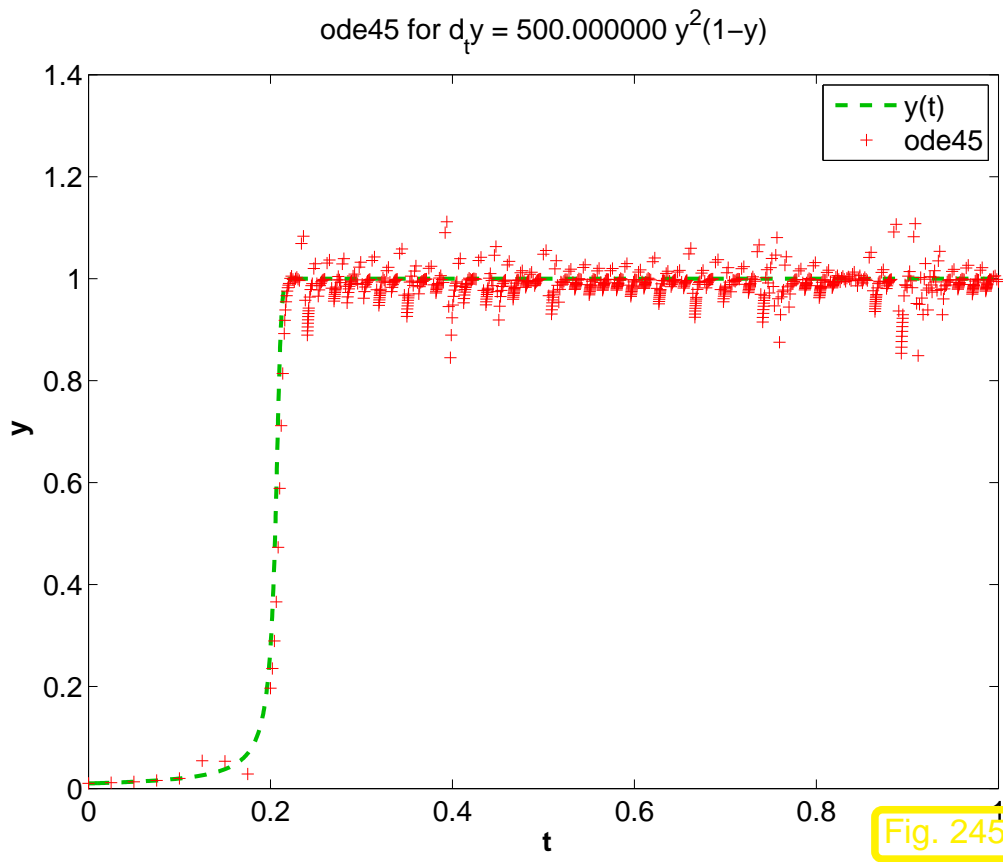
$$\text{IVP: } \dot{y} = \lambda y^2(1 - y), \quad \lambda := 500, \quad y(0) = \frac{1}{100}.$$

```
1 fun = @(t,x) 500*x^2*(1-x);  
2 options = odeset('reltol',0.1,'abstol',0.001,'stats','on');  
3 [t,y] = ode45(fun,[0 1],y0,options);
```

186 successful steps

55 failed attempts

1447 function evaluations



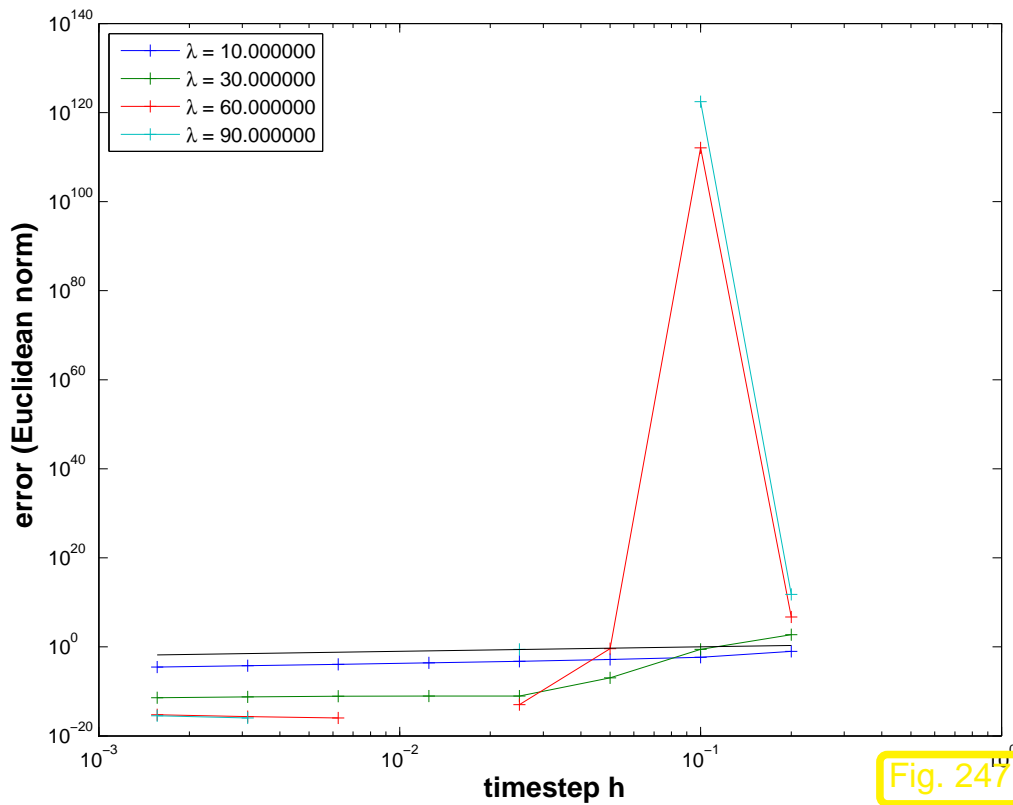
13.1 Model problem analysis [35, Ch. 77], [51, Sect. 11.3.3]

Example 13.1.1 (Blow-up of explicit Euler method).

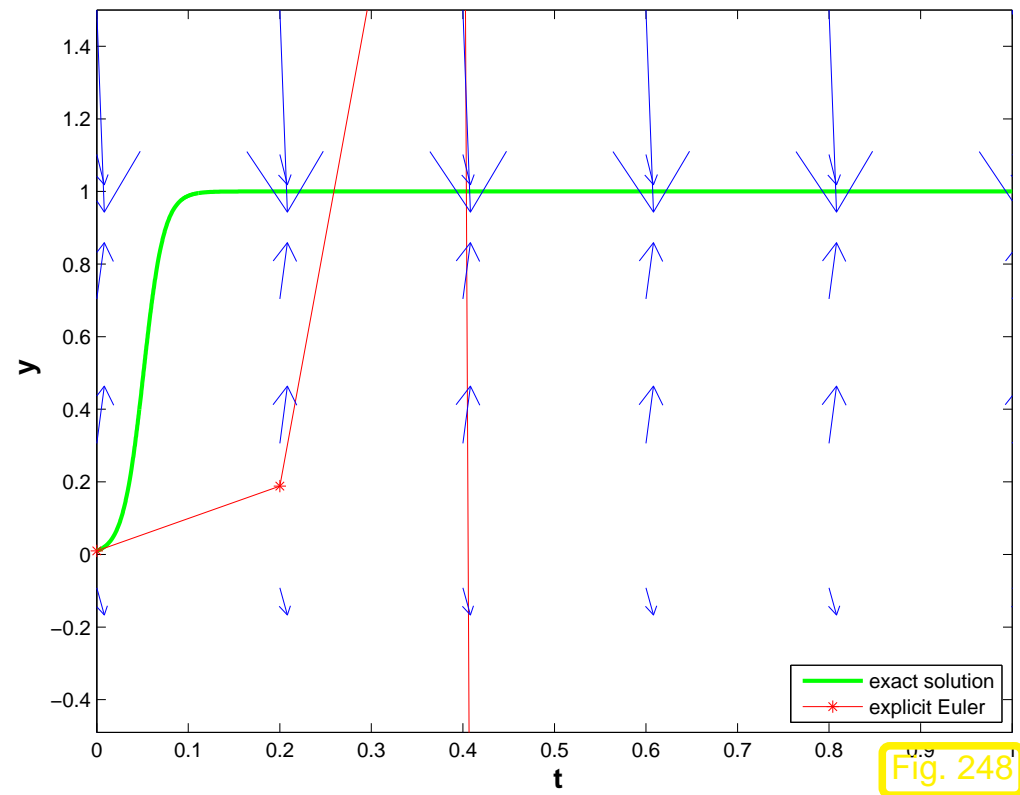
• IVP for logistic ODE, see Ex. 12.1.1

$$\dot{y} = f(y) := \lambda y(1 - y) \quad , \quad y(0) = 0.01 \quad .$$

- Explicit Euler method (12.2.4) with uniform timestep $h = 1/N$, $N \in \{5, 10, 20, 40, 80, 160, 320, 640\}$.



λ large: blow-up of y_k for large timestep h

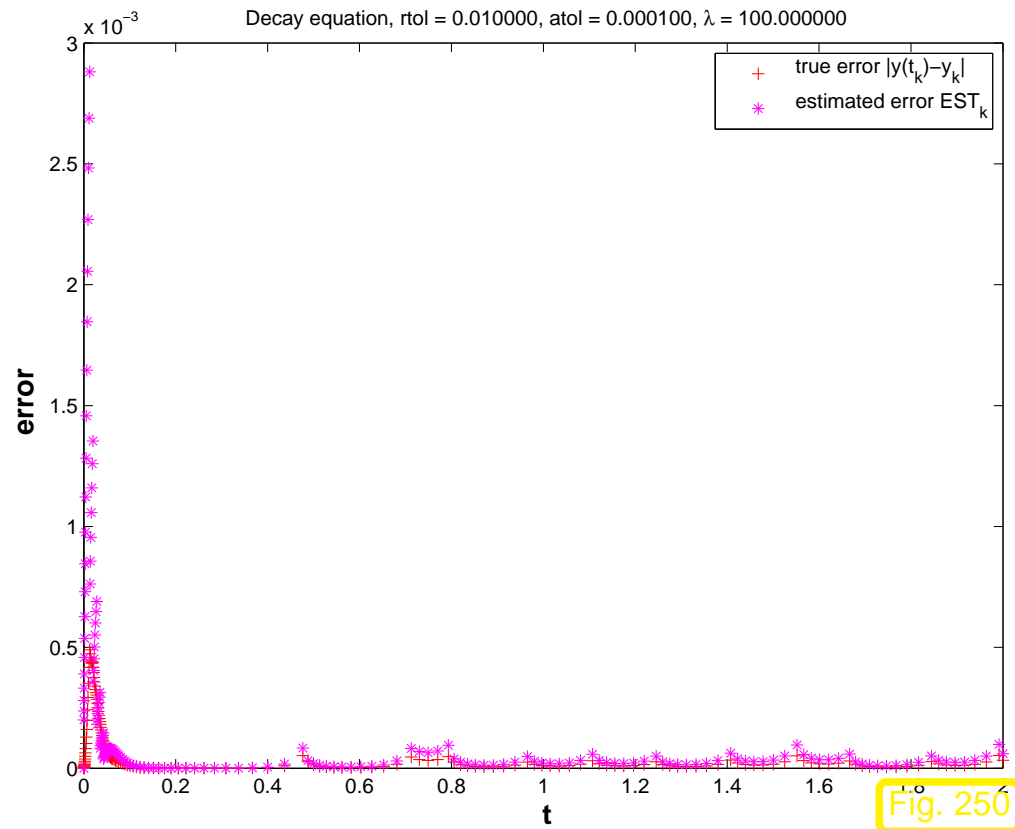
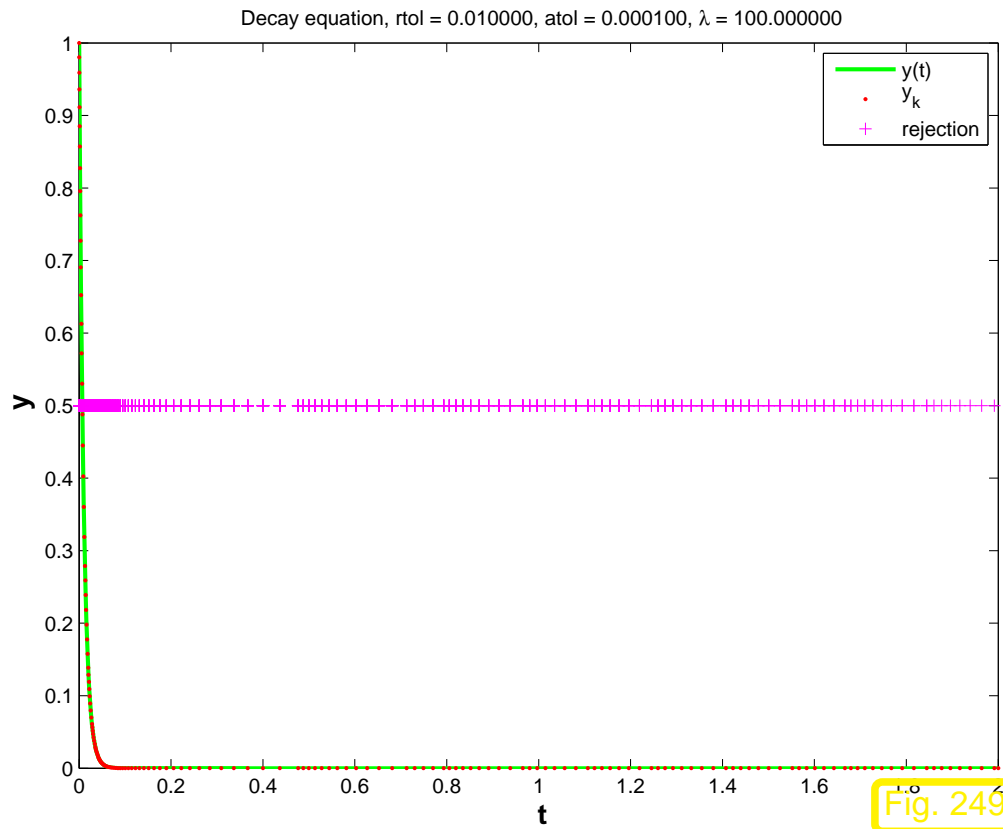


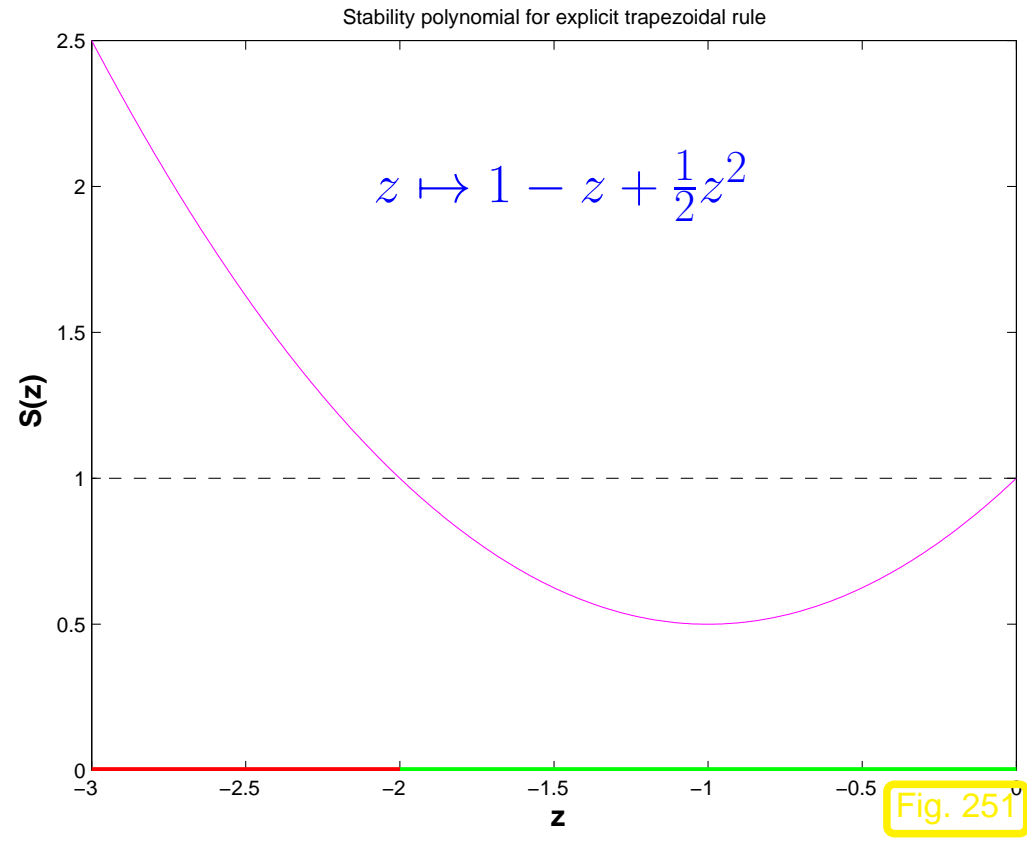
$\lambda = 90$: — $\hat{=}$ $y(t)$, —* $\hat{=}$ Euler polygon



Example 13.1.4 (Simple adaptive timestepping for fast decay).

- “Linear model problem IVP”: $\dot{y} = \lambda y, y(0) = 1, \lambda = -100$
- Simple adaptive timestepping method as in Ex. 12.5.8, see Code 12.5.6





Theorem 13.1.11 (Stability function of explicit Runge-Kutta methods). \rightarrow [35, Thm. 77.2], [51, Sect. 11.8.4]

The discrete evolution Ψ_λ^h of an explicit s -stage Runge-Kutta single step method (\rightarrow Def. 12.4.8) with Butcher scheme $\begin{array}{c|c} \mathbf{c} & \mathfrak{A} \\ \hline & \mathbf{b}^T \end{array}$ (see (12.4.9)) for the ODE $\dot{y} = \lambda y$ is a multiplication operator according to

$$\Psi_\lambda^h = \underbrace{1 + z\mathbf{b}^T (\mathbf{I} - z\mathfrak{A})^{-1} \mathbf{1}}_{\text{stability function } S(z)} = \det(\mathbf{I} - z\mathfrak{A} + z\mathbf{1b}^T), \quad z := \lambda h, \quad \mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^s.$$

13.2 Stiff problems [51, Sect. 11.10]

Example 13.2.1 (Transient simulation of RLC-circuit).



$$\ddot{u} + \alpha \dot{u} + \beta u = g(t),$$

$$\alpha := (RC)^{-1}, \beta = (LC)^{-1}, g(t) = \alpha \dot{U}_s.$$

Transformation to linear 1st-order ODE, see

Rem. 12.1.15, $v := \dot{u}$

$$\underbrace{\begin{pmatrix} \dot{u} \\ \dot{v} \end{pmatrix}}_{=:\dot{\mathbf{y}}} = \underbrace{\begin{pmatrix} 0 & 1 \\ -\beta & -\alpha \end{pmatrix} \begin{pmatrix} u \\ v \end{pmatrix} - \begin{pmatrix} 0 \\ g(t) \end{pmatrix}}_{=:\mathbf{f}(t,\mathbf{y})}.$$

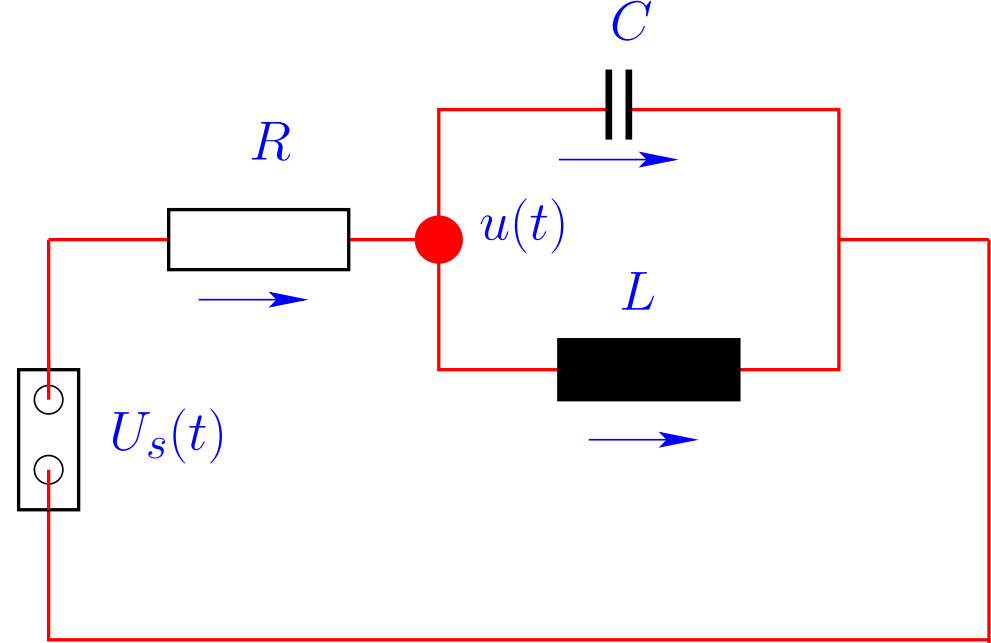
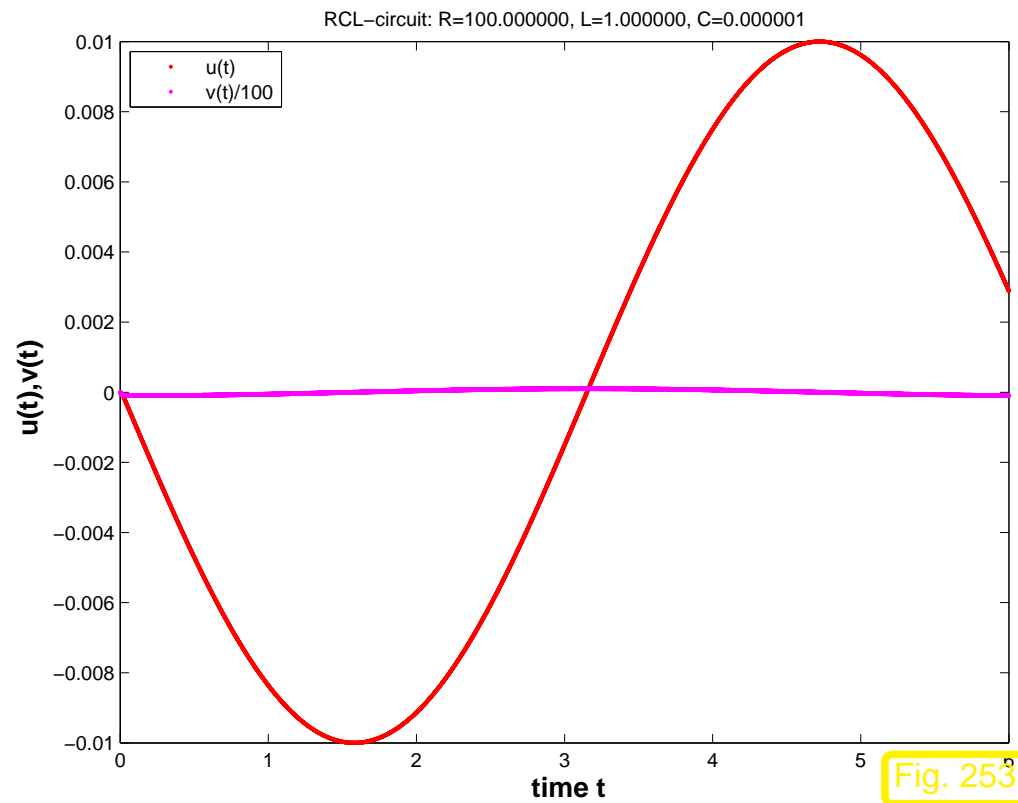


Fig. 252



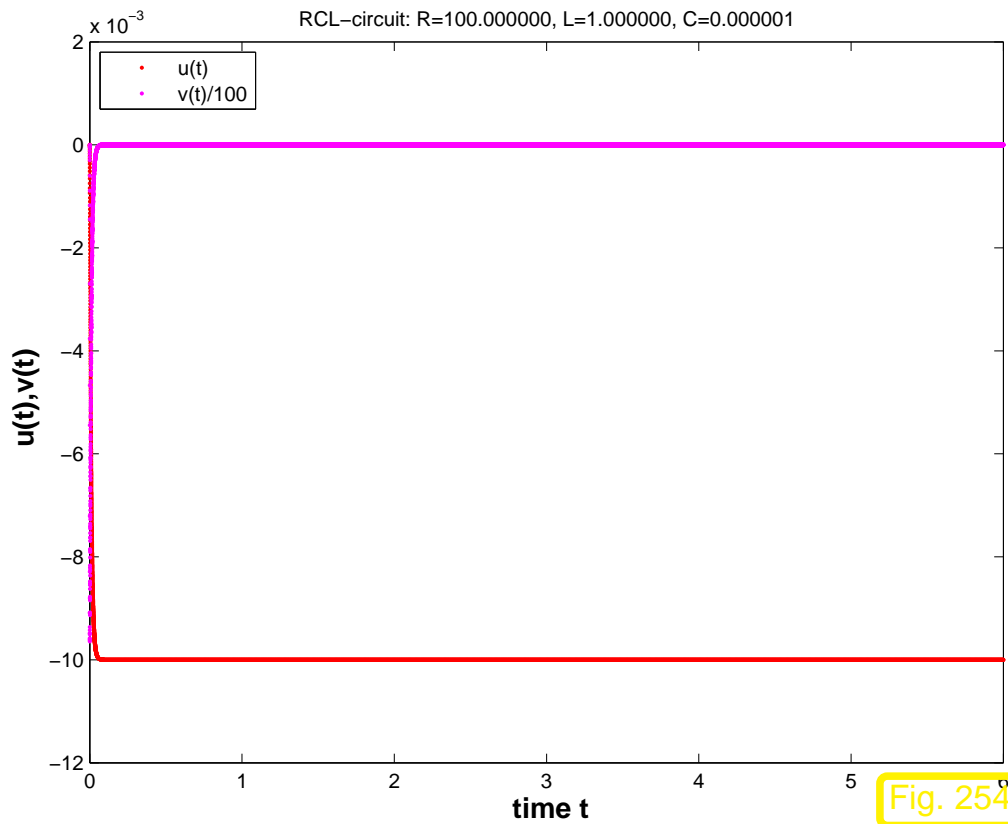
$R = 100\Omega$, $L = 1\text{H}$, $C = 1\mu\text{F}$, $U_s(t) = 1\text{V} \sin(t)$,
 $u(0) = v(0) = 0$ ("switch on")

ode45 statistics:

17897 successful steps

1090 failed attempts

113923 function evaluations



$R = 100\Omega$, $L = 1\text{H}$, $C = 1\mu\text{F}$, $U_s(t) = 1\text{V}$,
 $u(0) = v(0) = 0$ (“switch on”)

ode45 statistics:

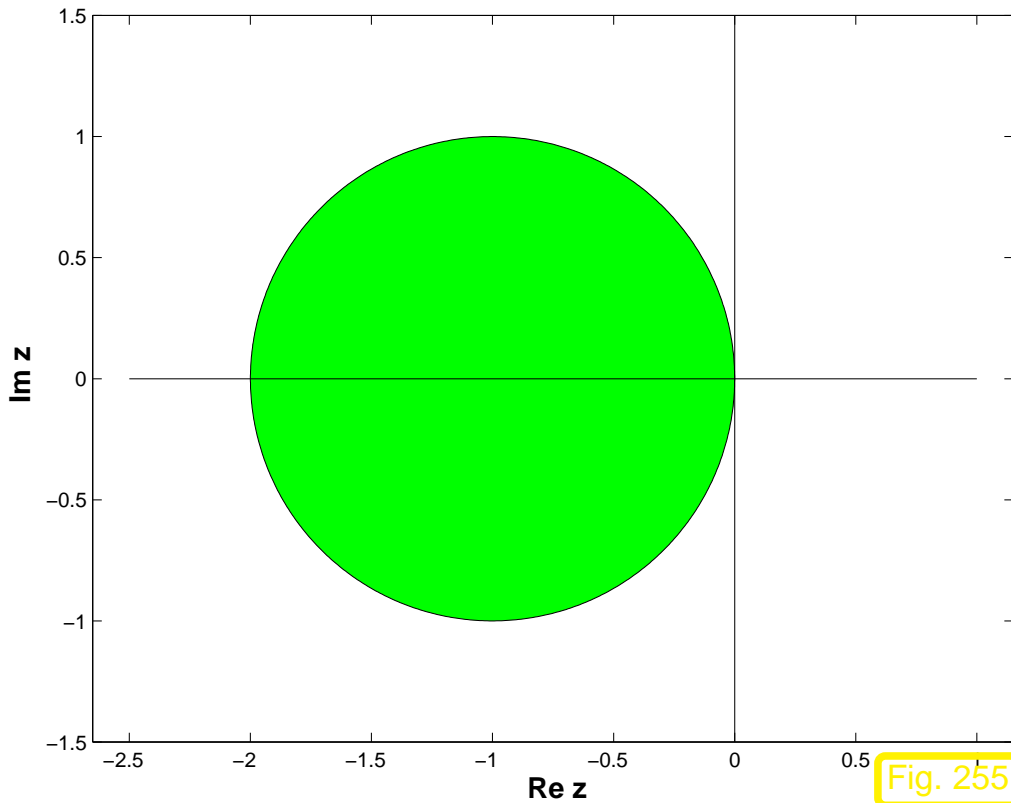
17901 successful steps

1210 failed attempts

114667 function evaluations



Example 13.2.7 (Explicit Euler method for damped oscillations).

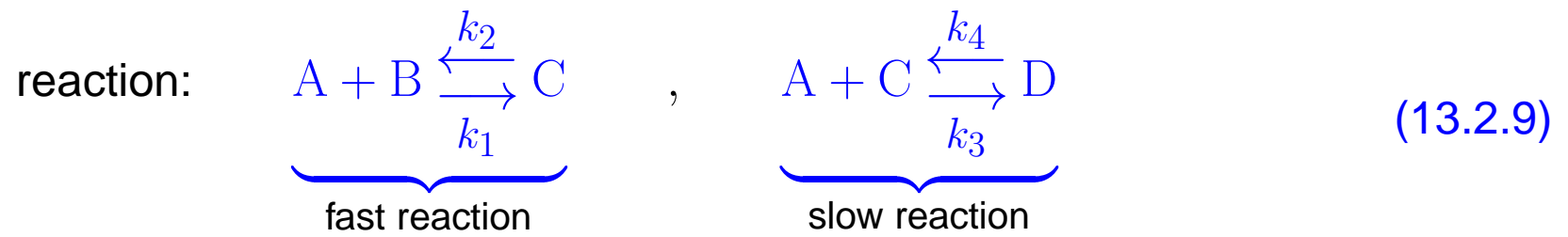


$$\triangleleft \{z \in \mathbb{C} : |1 + z| < 1\}$$



Can we predict this kind of difficulty ?

Example 13.2.8 (Kinetics of chemical reactions). \rightarrow [35, Ch. 62]



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: ODE involving **concentrations** $\mathbf{y}(t) = (c_A(t), c_B(t), c_C(t), c_D(t))^T$

$$\dot{\mathbf{y}} := \frac{d}{dt} \begin{pmatrix} c_A \\ c_B \\ c_C \\ c_D \end{pmatrix} = \mathbf{f}(\mathbf{y}) := \begin{pmatrix} -k_1 c_A c_B + k_2 c_C - k_3 c_A c_C + k_4 c_D \\ -k_1 c_A c_B + k_2 c_C \\ k_1 c_A c_B - k_2 c_C - k_3 c_A c_C + k_4 c_D \\ k_3 c_A c_C - k_4 c_D \end{pmatrix}. \quad (13.2.10)$$

Chemical reaction: concentrations

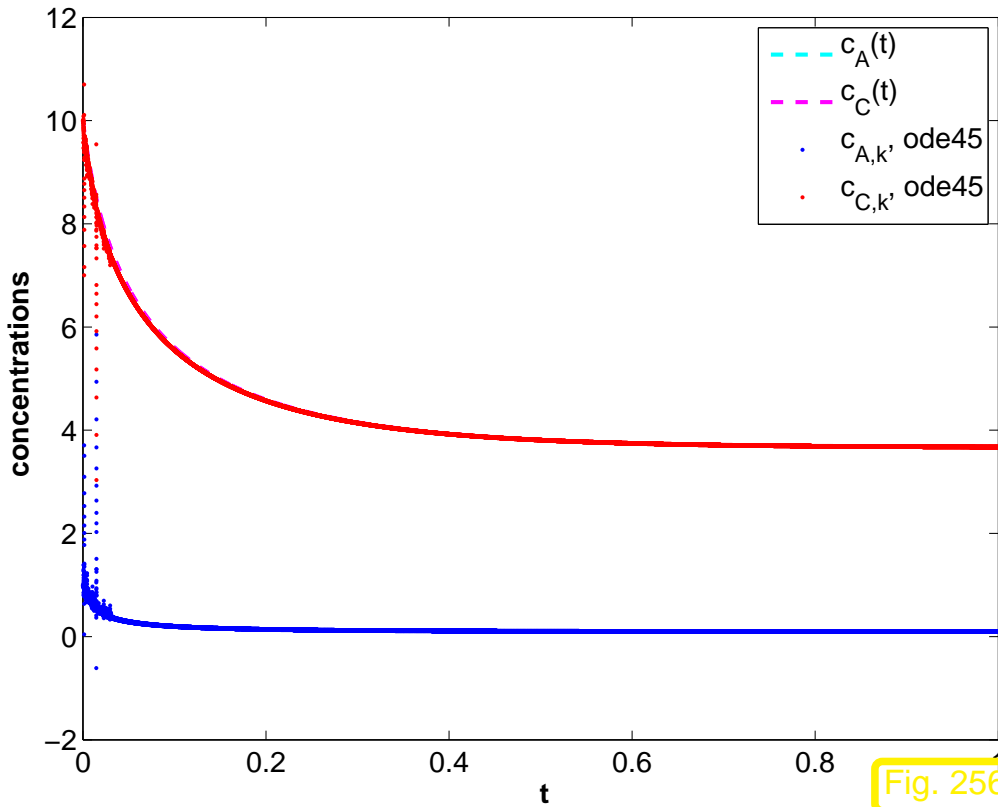


Fig. 256

Chemical reaction: stepsize

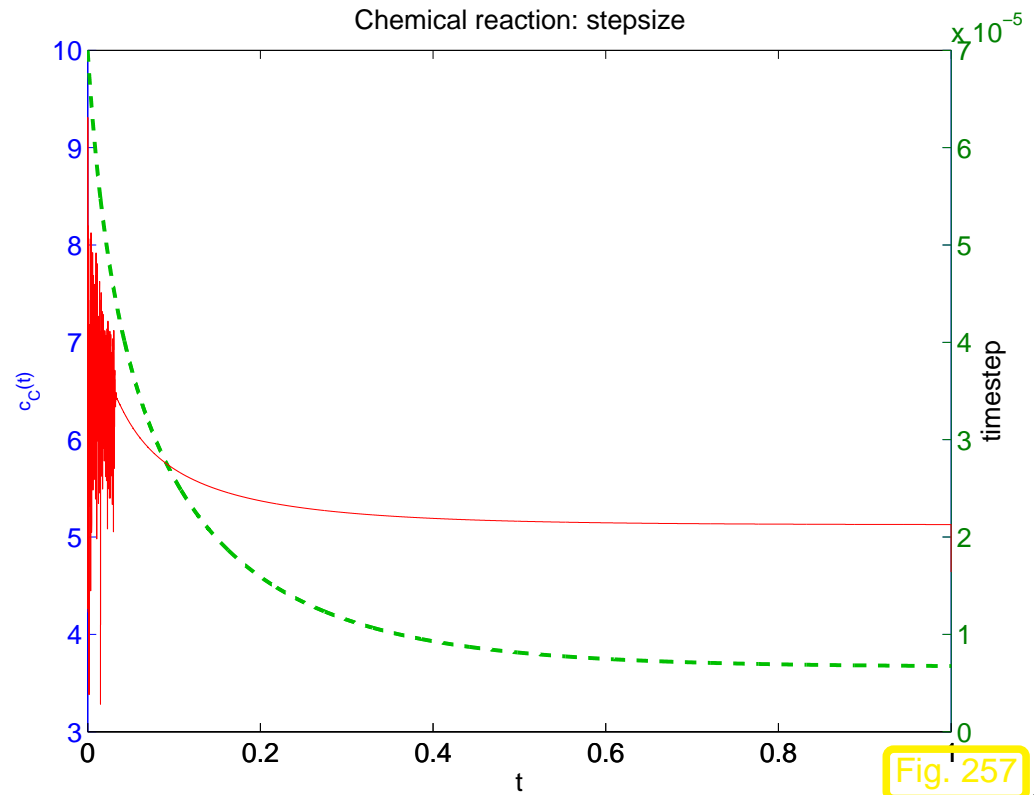


Fig. 257

Example 13.2.12 (Strongly attractive limit cycle).

Autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$

$$\mathbf{f}(\mathbf{y}) := \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \mathbf{y} + \lambda(1 - \|\mathbf{y}\|^2) \mathbf{y}, \quad (13.2.13)$$

on state space $D = \mathbb{R}^2 \setminus \{0\}$

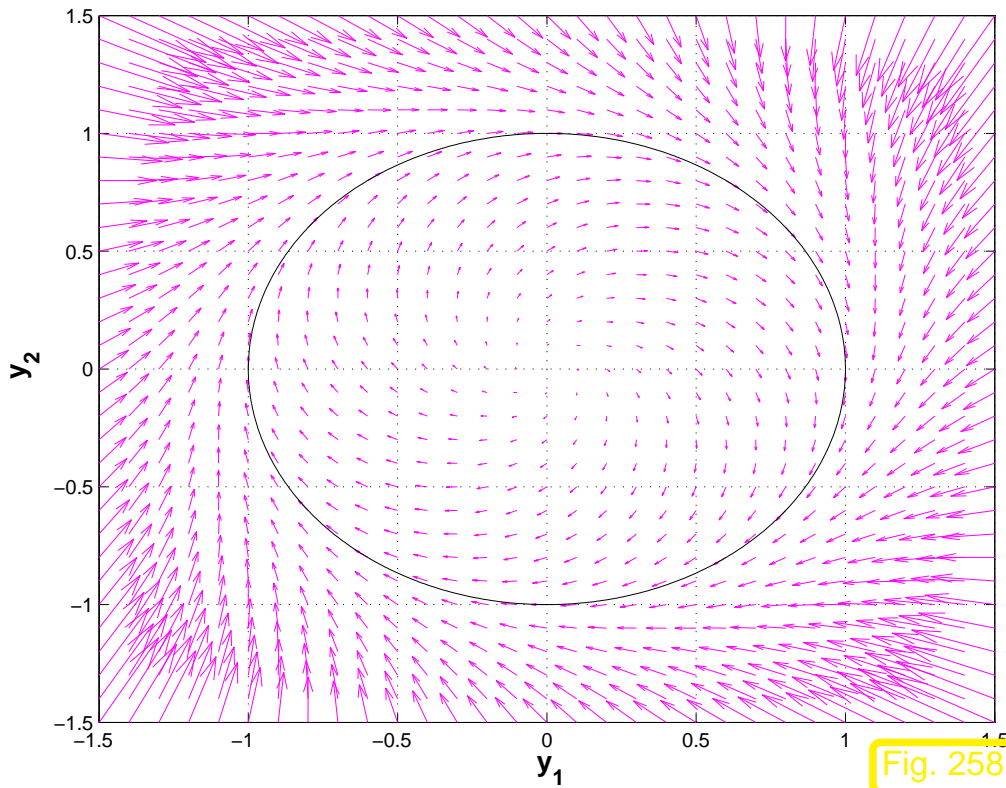


Fig. 258

vectorfield f ($\lambda = 1$)

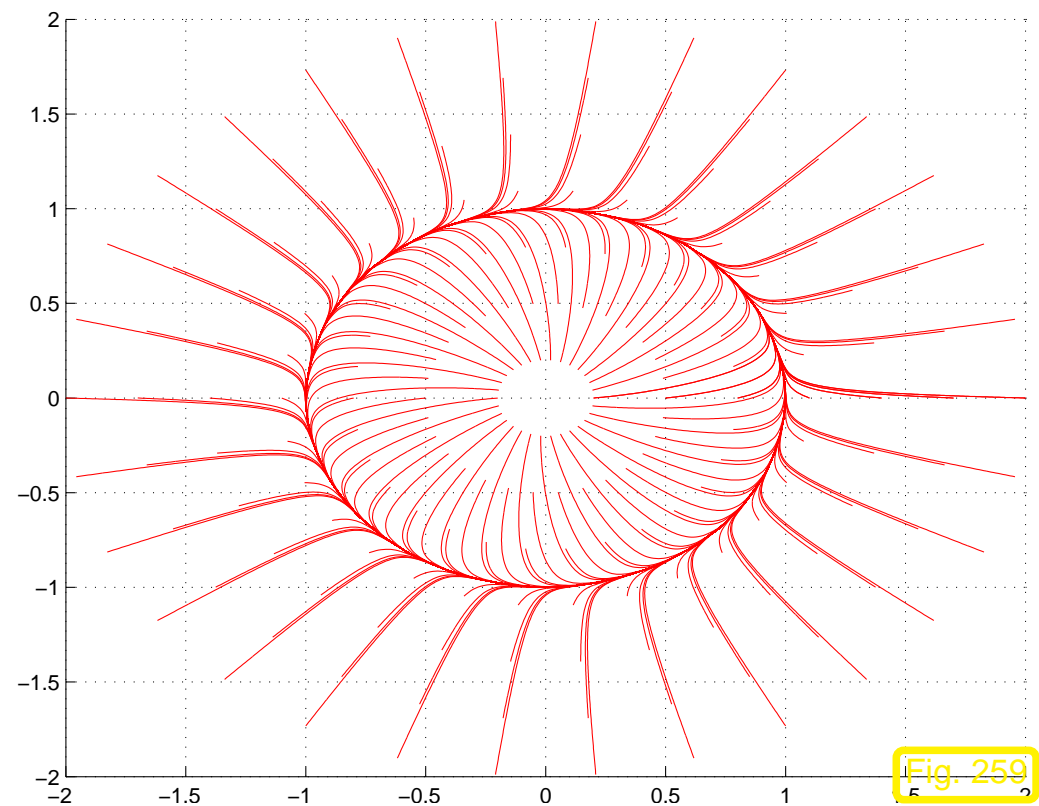
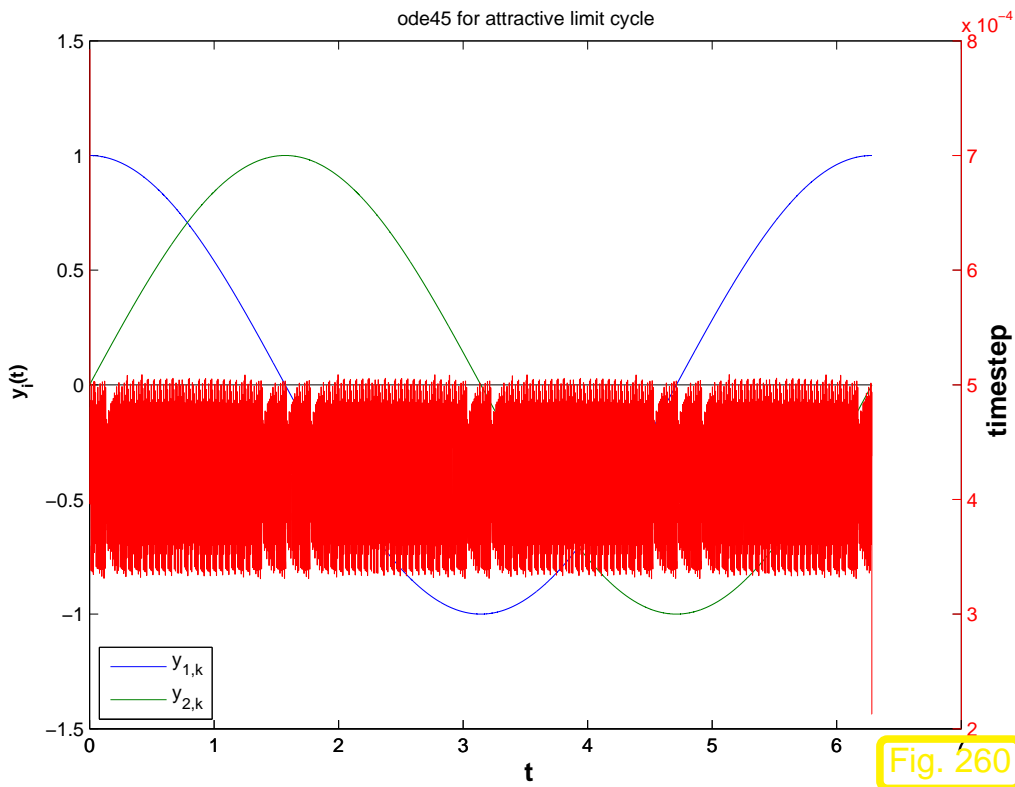
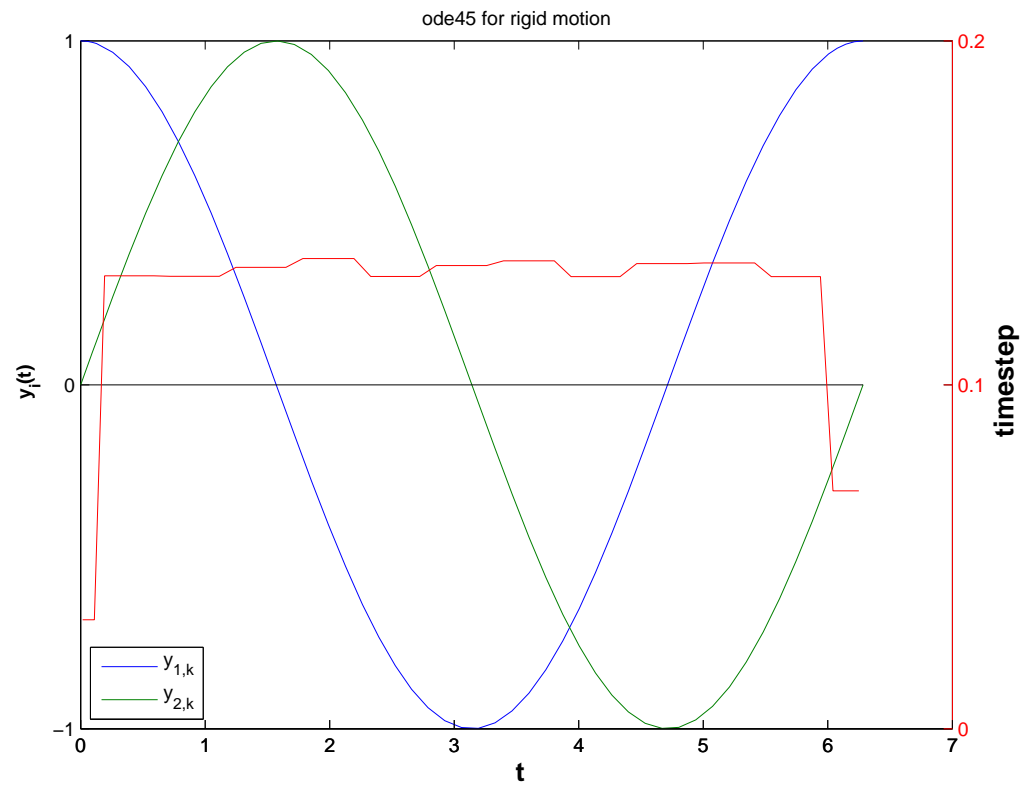


Fig. 259

solution trajectories ($\lambda = 10$)



many (3794) steps ($\lambda = 1000$)



accurate solution with few steps ($\lambda = 0$)



Notion 13.2.16 (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.*

Typical features of stiff IVPs:

- Presence of **fast transients** in the solution, see Ex. 13.1.1, 13.2.1,
- Occurrence of **strongly attractive** fixed points/limit cycles, see Ex. 13.2.13

13.3 Implicit Runge-Kutta methods [13, Sect. 11.6.2], [51, Sect. 11.

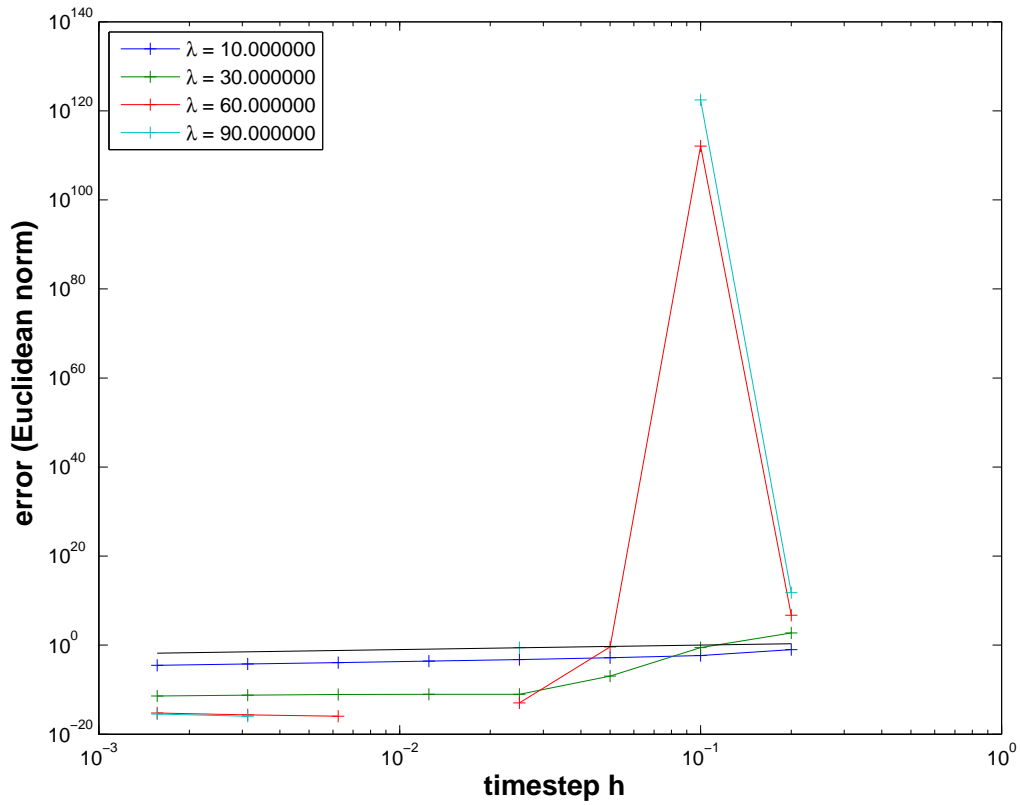
Example 13.3.1 (Implicit Euler timestepping for decay equation).



Example 13.3.4 (Euler methods for stiff logistic IVP).

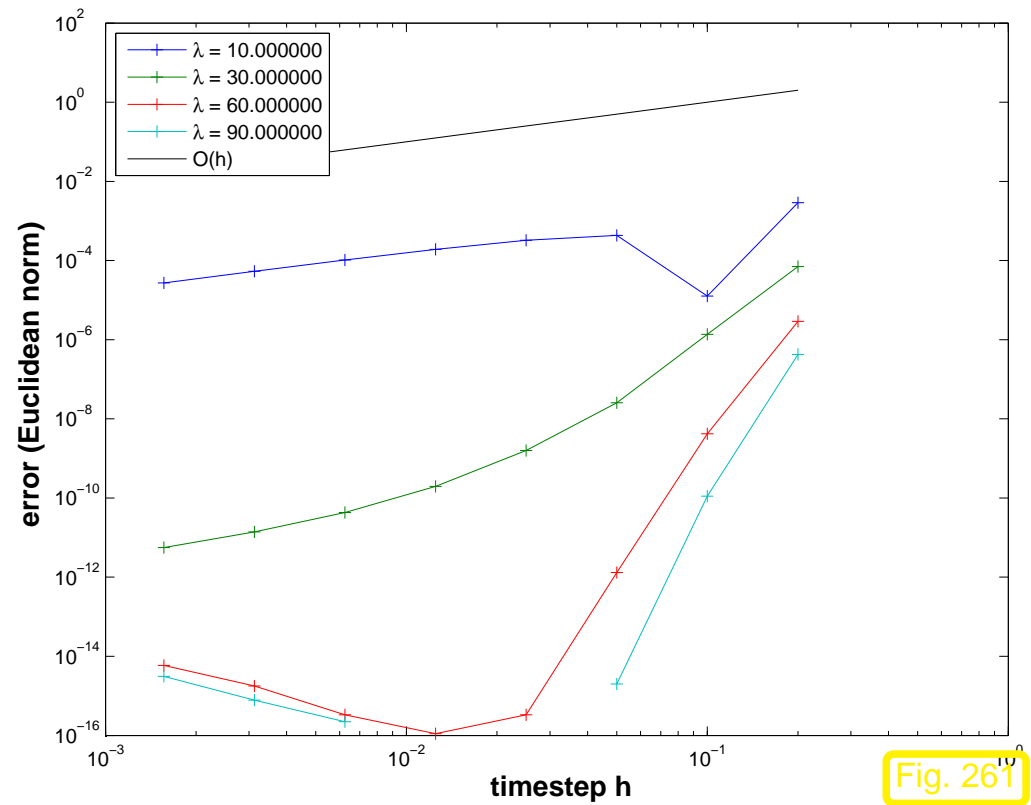
☞ Redo Ex. 13.1.1 for implicit Euler method:

Explicit Euler method (12.2.4)



λ large: blow-up of y_k for large timestep h

Implicit Euler method (12.2.8)



λ large: stable for all timesteps h !

Fig. 261



Definition 13.3.5 (General Runge-Kutta method). (cf. Def. 12.4.8)

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 (12.1.13) 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*.

Shorthand notation for Runge-Kutta methods

Butcher scheme



$$\frac{\mathbf{c} \mid \mathcal{A}}{\mathbf{b}^T} := \begin{array}{c|ccc} c_1 & a_{11} & \cdots & a_{1s} \\ \vdots & \vdots & & \vdots \\ c_s & a_{s1} & \cdots & a_{ss} \\ \hline & b_1 & \cdots & b_s \end{array}. \quad (13.3.6)$$

Theorem 13.3.7 (Stability function of Runge-Kutta methods).

The discrete evolution Ψ_λ^h of an s -stage Runge-Kutta single step method (\rightarrow Def. 13.3.5) with Butcher scheme $\begin{array}{c|c} \mathbf{c} & \mathbf{A} \\ \hline & \mathbf{b}^T \end{array}$ (see (13.3.6)) for the ODE $\dot{y} = \lambda y$ is a multiplication operator according to

$$\Psi_\lambda^h = \underbrace{1 + z\mathbf{b}^T (\mathbf{I} - z\mathbf{A})^{-1} \mathbf{1}}_{\text{stability function } S(z)} = \frac{\det(\mathbf{I} - z\mathbf{A} + z\mathbf{1b}^T)}{\det(\mathbf{I} - z\mathbf{A})}, \quad z := \lambda h, \quad \mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^s.$$

Definition 13.3.9 (L-stable Runge-Kutta method). \rightarrow [35, Ch. 77]

A Runge-Kutta method (\rightarrow Def. 13.3.5) is **L-stable/asymptotically stable**, if its stability function (\rightarrow Def. 13.3.7) satisfies

$$(i) \quad \operatorname{Re} z < 0 \quad \Rightarrow \quad |S(z)| < 1, \quad (13.3.10)$$

$$(ii) \quad \lim_{\operatorname{Re} z \rightarrow -\infty} S(z) = 0. \quad (13.3.11)$$

Example 13.3.15 (L-stable implicit Runge-Kutta methods).

$$\begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array}$$

Implicit Euler method

$$\begin{array}{c|cc} \frac{1}{3} & \frac{5}{12} & -\frac{1}{12} \\ \hline 1 & \frac{3}{4} & \frac{1}{4} \\ \hline & \frac{3}{4} & \frac{1}{4} \end{array}$$

Radau RK-SSM, order 3

$$\begin{array}{c|ccc} \frac{4-\sqrt{6}}{10} & \frac{88-7\sqrt{6}}{360} & \frac{296-169\sqrt{6}}{1800} & \frac{-2+3\sqrt{6}}{225} \\ \hline \frac{4+\sqrt{6}}{10} & \frac{296+169\sqrt{6}}{1800} & \frac{88+7\sqrt{6}}{360} & \frac{-2-3\sqrt{6}}{225} \\ \hline 1 & \frac{16-\sqrt{6}}{36} & \frac{16+\sqrt{6}}{36} & \frac{1}{9} \\ \hline & \frac{16-\sqrt{6}}{36} & \frac{16+\sqrt{6}}{36} & \frac{1}{9} \end{array}$$

Radau RK-SSM, order 5



13.4 Semi-implicit Runge-Kutta methods [35, Ch. 80]



Equations fixing increments $\mathbf{k}_i \in \mathbb{R}^d, i = 1, \dots, s$, for s -stage implicit RK-method

=

(Non-)linear system of equations with $s \cdot d$ unknowns

Example 13.4.1 (Linearization of increment equations).

- Initial value problem for logistic ODE, see Ex. 12.1.1

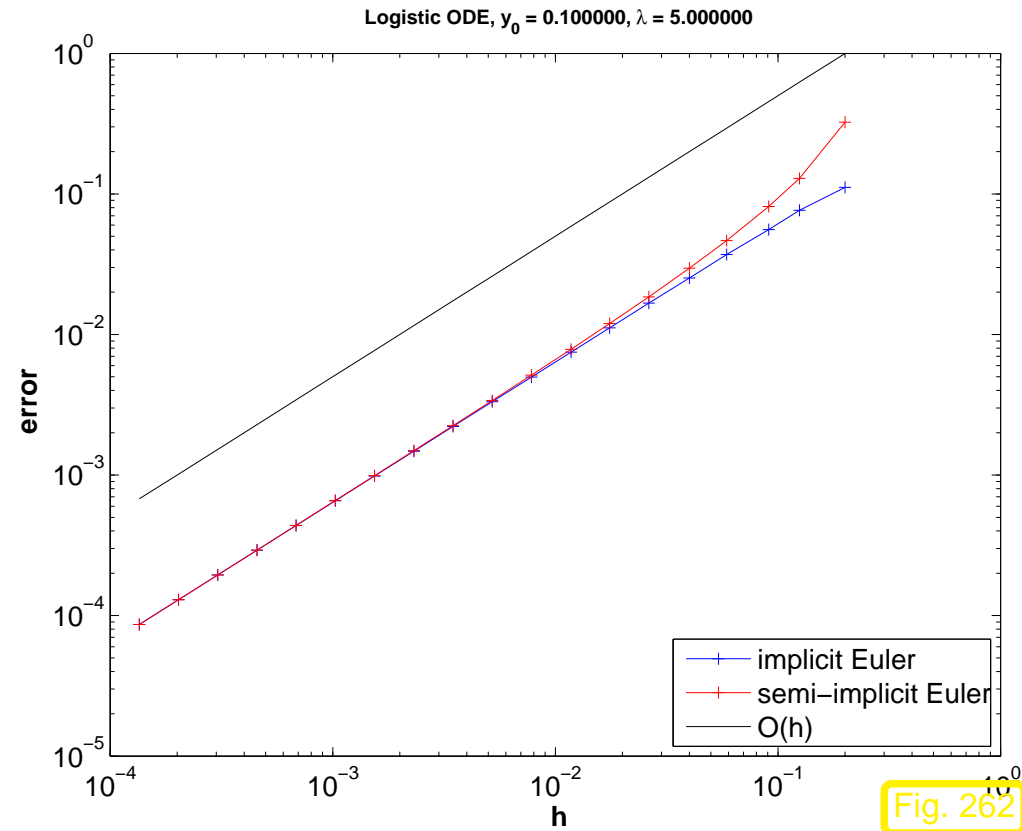
$$\dot{y} = \lambda y(1 - y) \quad , \quad y(0) = 0.1 \quad , \quad \lambda = 5 .$$

- Implicit Euler method (12.2.8) 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 $\text{err} = \max_{j=1, \dots, n} |y_j - y(t_j)|$



Idea: Use **linearized increment equations** for implicit RK-SSM

$$\mathbf{k}_i = \mathbf{f}(\mathbf{y}_0) + hD\mathbf{f}(\mathbf{y}_0) \left(\sum_{j=1}^s a_{ij} \mathbf{k}_j \right), \quad i = 1, \dots, s. \quad (13.4.2)$$

Linearization does nothing for linear ODEs ➤ stability function (→ Thm. 13.3.7) not affected!

▶ Class of **semi-implicit (linearly implicit) Runge-Kutta methods** (Rosenbrock-Wanner (ROW) methods):

$$(\mathbf{I} - ha_{ii}\mathbf{J})\mathbf{k}_i = \mathbf{f}(\mathbf{y}_0 + h \sum_{j=1}^{i-1} (a_{ij} + d_{ij})\mathbf{k}_j) - h\mathbf{J} \sum_{j=1}^{i-1} d_{ij}\mathbf{k}_j, \quad (13.4.3)$$

$$\mathbf{J} := D\mathbf{f}(\mathbf{y}_0 + h \sum_{j=1}^{i-1} (a_{ij} + d_{ij})\mathbf{k}_j), \quad (13.4.4)$$

$$\mathbf{y}_1 := \mathbf{y}_0 + \sum_{j=1}^s b_j \mathbf{k}_j. \quad (13.4.5)$$

R. Hiptmair
rev 38355,
December
23, 2010

Remark 13.4.6 (Adaptive integrator for stiff problems in MATLAB).

Handle of type $@(t, y) \mathbf{J}(t, y)$ to Jacobian $D\mathbf{f} : I \times D \mapsto \mathbb{R}^{d,d}$

```
opts = odeset('abstol',atol,'reltol',rtol,'Jacobian',J)
[t,y] = ode23s(odefun,tspan,y0,opts);
```

Stepsize control according to policy of Sect. 12.5:



14

Structure Preservation [31]

14.1 Dissipative Evolutions

14.2 Quadratic Invariants

14.3 Reversible Integrators

14.4 Symplectic Integrators

Course 401-0674-00: Numerical Methods for Partial Differential Equations

Many fundamental models in science & engineering boil down to

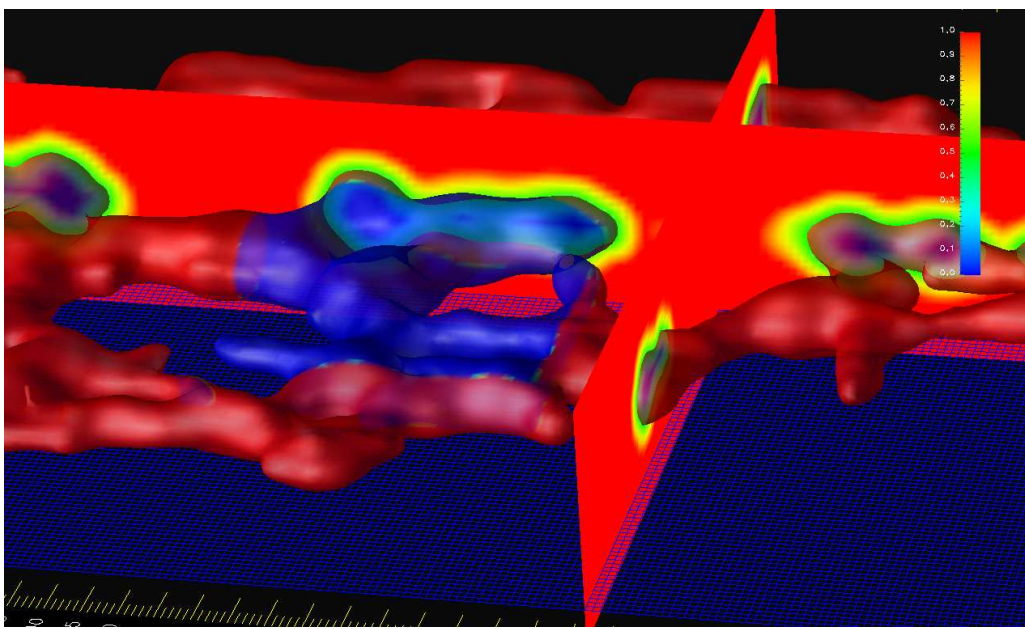
(initial) boundary value problems for **partial differential equations** (PDEs)

► Key role of numerical techniques for PDEs:

- Issue: Appropriate spatial (and temporal) **discretization** of PDE and boundary conditions
- Issue: **fast solution methods** for resulting *large* (non-)linear systems of equations

(initial) boundary value problems and techniques covered in the course:

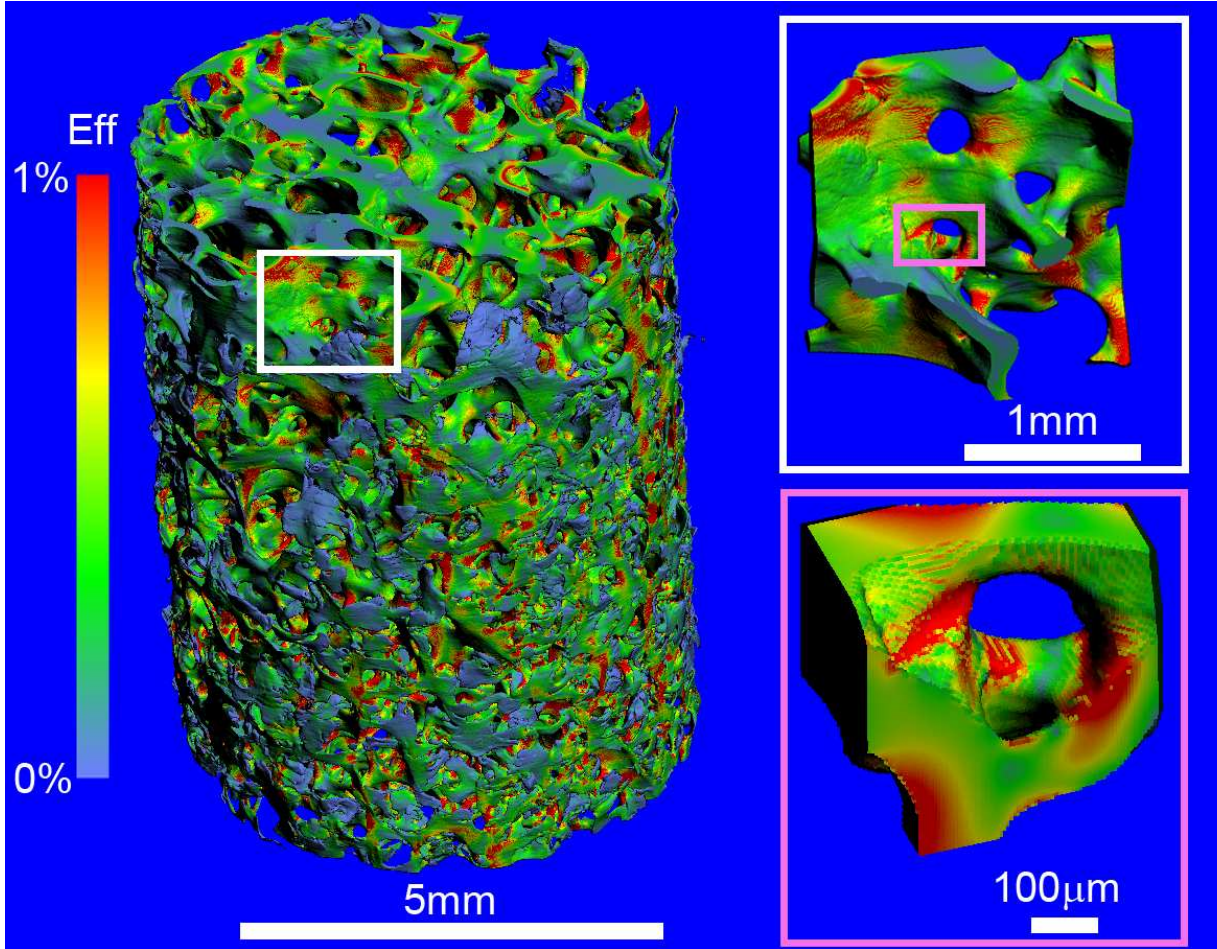
- ❶ **Stationary 2nd-order scalar elliptic boundary value problems**



Diffusion boundary value problem:

$$\begin{aligned}
 -\operatorname{div}(\mathbf{A}(\mathbf{x}) \operatorname{grad} u(\mathbf{x})) &= f(\mathbf{x}) \quad \text{in } \Omega \subset \mathbb{R}^d, \\
 u &= g \quad \text{on } \partial\Omega.
 \end{aligned}$$

◁ diffusion on the surface (membrane) of the endoplasmic reticulum (I. Sbalzarini, D-INFK, ETH Zürich)



◁ Elastic deformation of human bone
(P. Arbenz, D-INFK, ETH Zürich)

② Singularly perturbed elliptic boundary value problems

Stationary pollutant transport in water: find concentration $u = u(\boldsymbol{x})$ such that

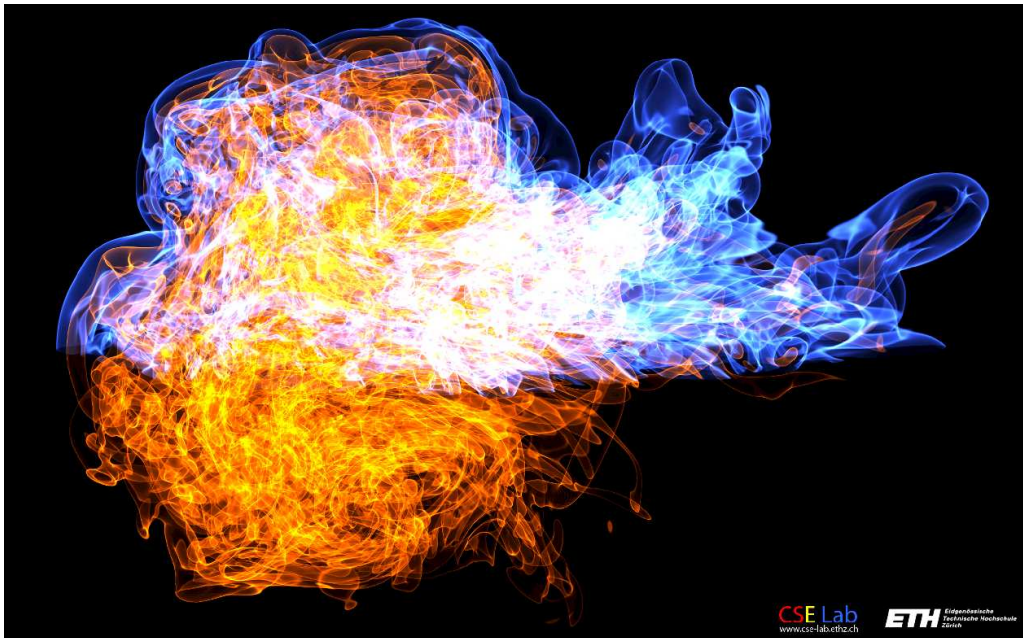
$$-\epsilon \Delta u + \mathbf{v}(\boldsymbol{x}) \cdot \mathbf{grad} u = 0 \quad \text{in } \Omega \quad , \quad u = g \quad \text{on } \partial\Omega .$$

③ 2nd-order parabolic evolution problems

Heat conduction: find temperature $u = u(\mathbf{x}, t)$

$$\frac{\partial}{\partial t} u(\mathbf{x}, t) - \operatorname{div}(\mathbf{A}(\mathbf{x}) \operatorname{grad} u(\mathbf{x}, t)) = 0 \quad \text{in } \Omega \times [0, T] \quad , \quad \begin{aligned} u(\cdot, t) &= g(t) \quad \text{on } \partial\Omega \quad , \\ u(\cdot, 0) &= u_0 \quad \text{in } \Omega \quad . \end{aligned}$$

④ Viscous fluid flow problems



Stokes equations:

$$\begin{aligned} -\Delta \mathbf{u} + \operatorname{grad} p &= \mathbf{f} \quad \text{in } \Omega \quad , \\ \operatorname{div} \mathbf{u} &= 0 \quad \text{in } \Omega \quad , \\ \mathbf{u} &= 0 \quad \text{on } \partial\Omega \quad . \end{aligned}$$

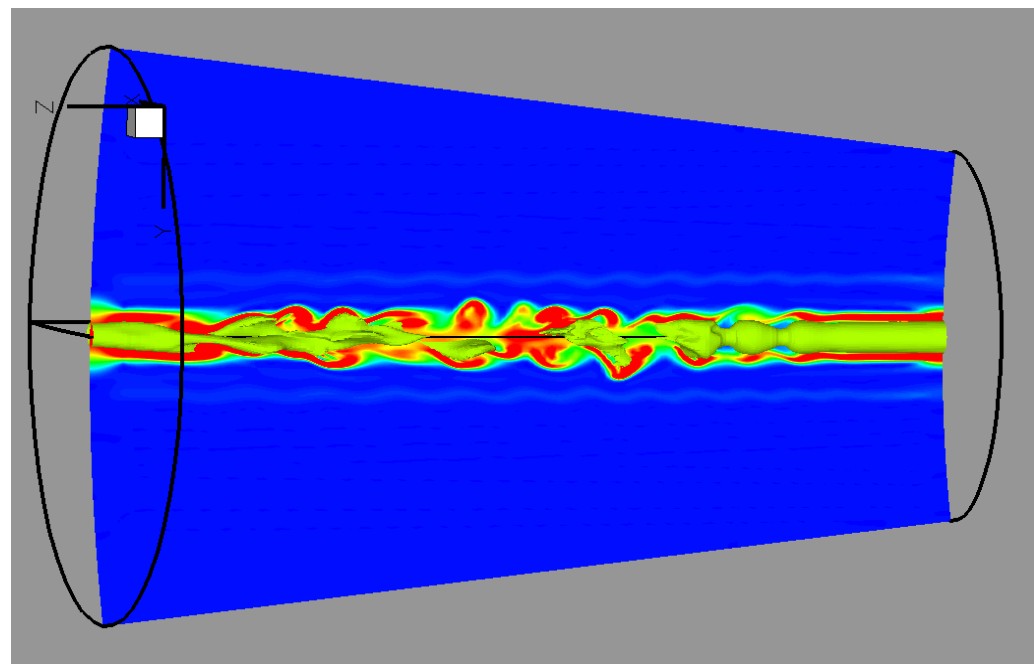
◁ Vortex ring in flow at $\operatorname{Re} = 7500$, (P. Koumoutsakos, D-INFK, ETH Zürich)

⑤ Conservation laws

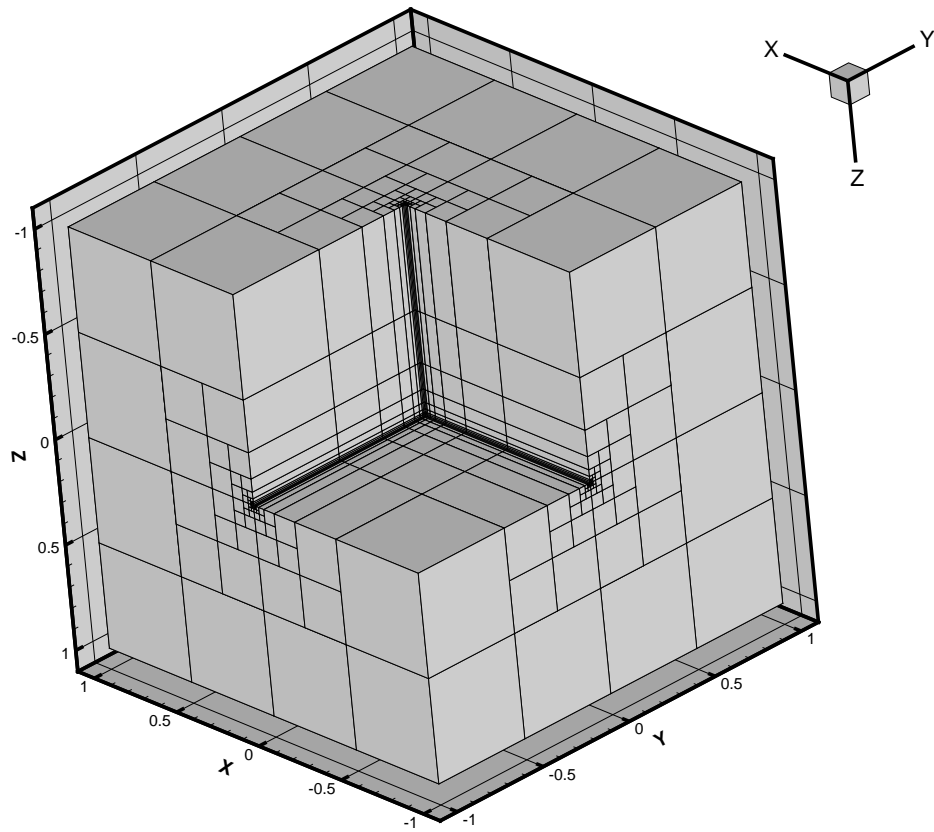
1D scalar conservation law with flux f :

$$\begin{aligned} \frac{\partial}{\partial t} u(x, t) + \frac{\partial}{\partial x} (f(u)) &= 0 && \text{in } \mathbb{R} \times \mathbb{R}^+, \\ u(x, 0) &= u_0(x) && \text{for } x \in \mathbb{R}. \end{aligned}$$

Inviscid fluid flow in 3D (SAM, D-MATH, ETH
Zürich) ▷



⑥ Adaptive finite element methods



◁ Adaptive FEM for diffusion problem:

Geometrically graded mesh at re-entrant corner (SAM, D-MATH, ETH Zürich)

7 Multilevel preconditioning

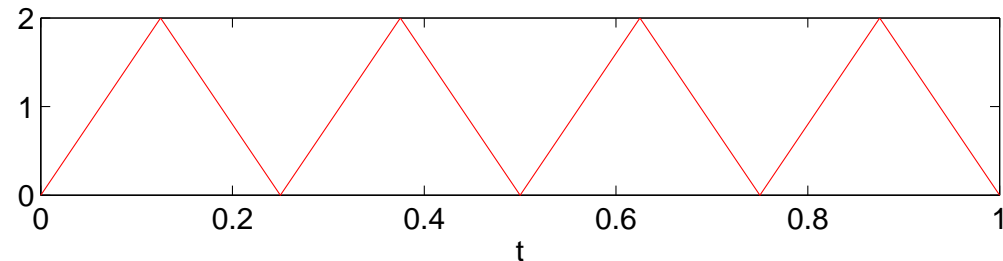
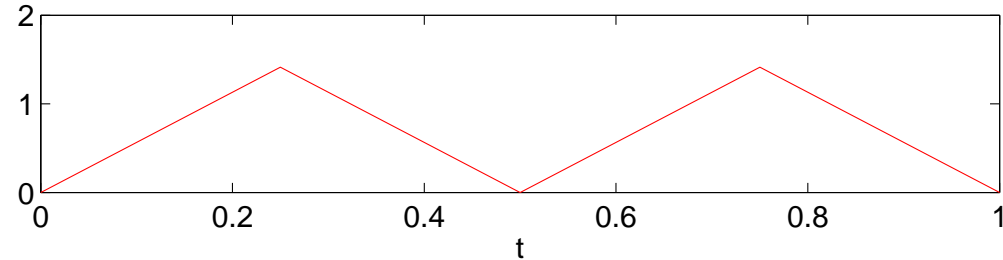
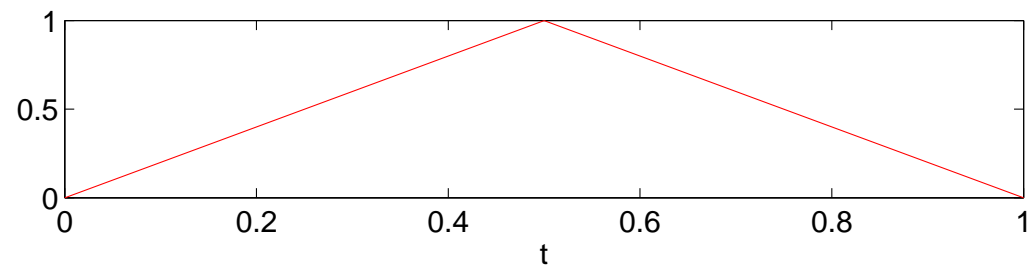
FEM, FD, FV



Huge sparse systems of equations

Efficient preconditioners required

1D hierarchical basis \triangleright



In FS12: Lecturers: Prof. Ralf Hiptmair (SAM, D-MATH)

Classes: Mon 15-17, HG E 7 and Fri 8-10, HG E 3

Tutorials: Mon 8-10, 13-15, Thu 13-15

[LINK to lecture notes](#)

Course: Parallel Computing for Scientific Simulations

This course will start in autumn term 2012 and deals with methods and techniques for numerical simulation on high performance (parallel) computers.

In HS12: Lecturers: Prof. Petros Koumoutsakos (D-MAVT)
Prof. Matthias Troyer (D-PHYS)

- Contents:
- Programming Models and Languages
 1. OpenCL
 2. CUDA
 3. Open MP
 4. MPI
 - Computers and Methods
 1. Hardware and Architectures
 2. Libraries
 3. Particles: N-Body solvers (→ numerical integration)
 4. Fields: PDEs

Bibliography

- [1] C. J. ALPERT, A. B. KAHNG, AND S.-Z. YAO, *Spectral partitioning with multiple eigenvectors*, Discrete Applied Mathematics, 90 (1999), pp. 3 – 26.
- [2] H. AMANN, *Gewöhnliche Differentialgleichungen*, Walter de Gruyter, Berlin, 1st ed., 1983.
- [3] S. AMAT, C. BERMUDEZ, S. BUSQUIER, AND S. PLAZA, *On a third-order Newton-type method free of bilinear operators*, Numerical Lin. Alg., (2010). DOI: 10.1002/nla.654.
- [4] Z.-J. BAI, J. DEMMEL, J. DONGARRA, A. RUHE, AND H. VAN DER VORST, *Templates for the Solution of Algebraic Eigenvalue Problems*, SIAM, Philadelphia, PA, 2000.
- [5] C. BISCHOF AND C. VAN LOAN, *The WY representation of Householder matrices*, SIAM J. Sci. Stat. Comput., 8 (1987).
- [6] S. BÖRM, L. GRASEDYCK, AND W. HACKBUSCH, *Introduction to hierarchical matrices with applications*, Engineering Analysis with Boundary Elements, 27 (2003), pp. 405–422.

- [7] F. BORNEMANN, *A model for understanding numerical stability*, IMA J. Numer. Anal., 27 (2006), pp. 219–231.
- [8] A. BRANDT AND A. LUBRECHT, *Multilevel matrix multiplication and fast solution of integral equations*, J. Comp. Phys., 90 (1990), pp. 348–370.
- [9] E. BRIGHAM, *The Fast Fourier Transform and Its Applications*, Prentice-Hall, Englewood Cliffs, NJ, 1988.
- [10] Q. CHEN AND I. BABUSKA, *Approximate optimal points for polynomial interpolation of real functions in an interval and in a triangle*, Comp. Meth. Appl. Mech. Engr., 128 (1995), pp. 405–417.
- [11] D. COPPERSMITH AND T. RIVLIN, *The growth of polynomials bounded at equally spaced points*, SIAM J. Math. Anal., 23 (1992), pp. 970–983.
- [12] D. COPPERSMITH AND S. WINOGRAD, *Matrix multiplication via arithmetic progression*, J. Symbolic Computing, 9 (1990), pp. 251–280.
- [13] W. DAHMEN AND A. REUSKEN, *Numerik für Ingenieure und Naturwissenschaftler*, Springer, Heidelberg, 2006.
- [14] M. DEAKIN, *Applied catastrophe theory in the social and biological sciences*, Bulletin of Mathematical Biology, 42 (1980), pp. 647–679.
- [15] P. DEUFLHARD, *Newton Methods for Nonlinear Problems*, vol. 35 of Springer Series in Computational Mathematics, Springer, Berlin, 2004.
- [16] P. DEUFLHARD AND F. BORNEMANN, *Scientific Computing with Ordinary Differential Equations*, vol. 42 of Texts in Applied Mathematics, Springer, New York, 2 ed., 2002.

- [17] P. DEUFLHARD AND A. HOHMANN, *Numerische Mathematik I*, DeGruyter, Berlin, 3 ed., 2002.
- [18] P. DUHAMEL AND M. VETTERLI, *Fast fourier transforms: a tutorial review and a state of the art*, *Signal Processing*, 19 (1990), pp. 259–299.
- [19] A. DUTT AND V. ROKHLIN, *Fast Fourier transforms for non-equispaced data II*, *Appl. Comput. Harmon. Anal.*, 2 (1995), pp. 85–100.
- [20] F. FRITSCH AND R. CARLSON, *Monotone piecewise cubic interpolation*, *SIAM J. Numer. Anal.*, 17 (1980), pp. 238–246.
- [21] M. GANDER, W. GANDER, G. GOLUB, AND D. GRUNTZ, *Scientific Computing: An introduction using MATLAB*, Springer, 2005. In Vorbereitung.
- [22] J. GILBERT, C. MOLER, AND R. SCHREIBER, *Sparse matrices in MATLAB: Design and implementation*, *SIAM Journal on Matrix Analysis and Applications*, 13 (1992), pp. 333–356.
- [23] G. GOLUB AND C. VAN LOAN, *Matrix computations*, John Hopkins University Press, Baltimore, London, 2nd ed., 1989.
- [24] C. GOTSMAN AND S. TOLEDO, *On the computation of null spaces of sparse rectangular matrices*, *SIAM J. Matrix Anal. Appl.*, 30 (2008), pp. 445–463.
- [25] C. GRAY, *An analysis of the Belousov-Zhabotinski reaction*, *Rose-Hulman Undergraduate Math Journal*, 3 (2002). <http://www.rose-hulman.edu/mathjournal/archives/2002/vol3-n1/paper1/v3n1-1pd.pdf>.
- [26] L. GREENGARD AND V. ROKHLIN, *A new version of the fast multipole method for the Laplace equation in three dimensions*, *Acta Numerica*, (1997), pp. 229–269.

- [27] M. GUTKNECHT, *Lineare Algebra*, lecture notes, SAM, ETH Zürich, 2009.
<http://www.sam.math.ethz.ch/~mhg/unt/LA/HS07/>.
- [28] W. HACKBUSCH, *Iterative Lösung großer linearer Gleichungssysteme*, B.G. Teubner–Verlag, Stuttgart, 1991.
- [29] W. HACKBUSCH, *Iterative solution of large sparse systems of equations*, vol. 95 of Applied Mathematical Sciences, Springer-Verlag, New York, 1994. Translated and revised from the 1991 German original.
- [30] W. HACKBUSCH AND S. BÖRM, *Data-sparse approximation by adaptive \mathcal{H}^2 -matrices*, Computing, 69 (2002), pp. 1–35.
- [31] E. HAIRER, C. LUBICH, AND G. WANNER, *Geometric numerical integration*, vol. 31 of Springer Series in Computational Mathematics, Springer, Heidelberg, 2 ed., 2006.
- [32] E. HAIRER, S. NORSETT, AND G. WANNER, *Solving Ordinary Differential Equations I. Nonstiff Problems*, Springer-Verlag, Berlin, Heidelberg, New York, 2 ed., 1993.
- [33] E. HAIRER AND G. WANNER, *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic Problems*, vol. 14 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, 2011.
- [34] C. HALL AND W. MEYER, *Optimal error bounds for cubic spline interpolation*, J. Approx. Theory, 16 (1976), pp. 105–122.
- [35] M. HANKE-BOURGEOIS, *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*, Mathematische Leitfäden, B.G. Teubner, Stuttgart, 2002.

- [36] N. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, PA, 2 ed., 2002.
- [37] I. IPSEN AND C. MEYER, *The idea behind Krylov methods*, Technical Report 97-3, Math. Dep., North Carolina State University, Raleigh, NC, January 1997.
- [38] S. JOHNSON, *Notes on the convergence of trapezoidal-rule quadrature*. MIT online course notes, <http://math.mit.edu/~stevenj/trapezoidal.pdf>, 2008.
- [39] D. KALMAN, *A singularly valuable decomposition: The SVD of a matrix*, The College Mathematics Journal, 27 (1996), pp. 2–23.
- [40] M. KOWARSCHIK AND W. C, *An overview of cache optimization techniques and cache-aware numerical algorithms*, in Algorithms for Memory Hierarchies, vol. 2625 of Lecture Notes in Computer Science, Springer, Heidelberg, 2003, pp. 213–232.
- [41] A. LALIENA AND F.-J. SAYAS, *Theoretical aspects of the application of convolution quadrature to scattering of acoustic waves*, Numer. Math., 112 (2009), pp. 637–678.
- [42] A. LENGVILLE AND C. MEYER, *Google's PageRank and Beyond: The Science of Search Engine Rankings*, Princeton University Press, Princeton, NJ, 2006.
- [43] D. MCALLISTER AND J. ROULIER, *An algorithm for computing a shape-preserving osculatory quadratic spline*, ACM Trans. Math. Software, 7 (1981), pp. 331–347.
- [44] M. MESSNER, M. SCHANZ, AND E. DARVE, *Fast directional multilevel summation for oscillatory kernels based on chebyshev interpolation*, Journal of Computational Physics, (2011), pp. –.
- [45] C. MOLER, *Numerical Computing with MATLAB*, SIAM, Philadelphia, PA, 2004.

- [46] K. NEYMEYR, *A geometric theory for preconditioned inverse iteration applied to a subspace*, Tech. Rep. 130, SFB 382, Universität Tübingen, Tübingen, Germany, November 1999. Submitted to Math. Comp.
- [47] —, *A geometric theory for preconditioned inverse iteration: III. Sharp convergence estimates*, Tech. Rep. 130, SFB 382, Universität Tübingen, Tübingen, Germany, November 1999.
- [48] K. NIPP AND D. STOFFER, *Lineare Algebra*, vdf Hochschulverlag, Zürich, 5 ed., 2002.
- [49] M. OVERTON, *Numerical Computing with IEEE Floating Point Arithmetic*, SIAM, Philadelphia, PA, 2001.
- [50] A. D. H.-D. QI, L.-Q. QI, AND H.-X. YIN, *Convergence of Newton's method for convex best interpolation*, Numer. Math., 87 (2001), pp. 435–456.
- [51] A. QUARTERONI, R. SACCO, AND F. SALERI, *Numerical mathematics*, vol. 37 of Texts in Applied Mathematics, Springer, New York, 2000.
- [52] C. RADER, *Discrete Fourier transforms when the number of data samples is prime*, Proceedings of the IEEE, 56 (1968), pp. 1107–1108.
- [53] R. RANNACHER, *Einführung in die numerische mathematik*. Vorlesungsskriptum Universität Heidelberg, 2000. <http://gaia.iwr.uni-heidelberg.de/>.
- [54] V. ROKHLIN, *Rapid solution of integral equations of classical potential theory*, J. Comp. Phys., 60 (1985), pp. 187–207.
- [55] Y. SAAD, *Iterative methods for sparse linear systems*, Society for Industrial and Applied Mathematics, Philadelphia, PA, second ed., 2003.

- [56] A. SANKAR, D. SPIELMAN, AND S.-H. TENG, *Smoothed analysis of the condition numbers and growth factors of matrices*, SIAM J. Matrix Anal. Appl., 28 (2006), pp. 446–476.
- [57] T. SAUER, *Numerical analysis*, Addison Wesley, Boston, 2006.
- [58] J.-B. SHI AND J. MALIK, *Normalized cuts and image segmentation*, IEEE Trans. Pattern Analysis and Machine Intelligence, 22 (2000), pp. 888–905.
- [59] D. SPIELMAN AND S.-H. 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.
- [60] M. STEWART, *A superfast toeplitz solver with improved numerical stability*, SIAM J. Matrix Analysis Appl., 25 (2003), pp. 669–693.
- [61] J. STOER, *Einführung in die Numerische Mathematik*, Heidelberger Taschenbücher, Springer, 4 ed., 1983.
- [62] V. STRASSEN, *Gaussian elimination is not optimal*, Numer. Math., 13 (1969), pp. 354–356.
- [63] M. STRUWE, *Analysis für Informatiker*. Lecture notes, ETH Zürich, 2009. <https://moodle-app1.net.ethz.ch/lms/mod/resource/index.php?id=145>.
- [64] F. TISSEUR AND K. MEERBERGEN, *The quadratic eigenvalue problem*, SIAM Review, 43 (2001), pp. 235–286.
- [65] L. TREFETHEN AND D. BAU, *Numerical Linear Algebra*, SIAM, Philadelphia, PA, 1997.

- [66] P. VERTESI, *On the optimal lebesgue constants for polynomial interpolation*, Acta Math. Hungaria, 47 (1986), pp. 165–178.
- [67] —, *Optimal lebesgue constant for lagrange interpolation*, SIAM J. Numer. Anal., 27 (1990), pp. 1322–1331.

Index

LU-decomposition

existence, 135

L^2 -inner product, 1066

h-convergence, 1006

MATLAB `arnoldi`, 710

MATLAB `bicgstab`, 570

MATLAB `broyden`, 498

MATLAB `cg`, 535

MATLAB `cholupdate`, 280

MATLAB `chol`, 237

MATLAB `costrans`, 927

MATLAB `ct_rect`, 1110

MATLAB `divide`, 1113

MATLAB `eigs`, 718

MATLAB `fftsolve`, 925

MATLAB `fft`, 857

MATLAB `fzero`, 446

MATLAB gallery, 693, 705, 714

MATLAB `gaussquad`, 1077

MATLAB `gmres`, 568

MATLAB `gn`, 816

MATLAB `icostrans`, 928

MATLAB `ifft`, 857

MATLAB `ipoleval`, 319

MATLAB `legendre`, 1076

MATLAB `lsqsvd`, 796

MATLAB `lsqttotal`, 803

MATLAB `lurec`, 119

MATLAB `lu`, 136

MATLAB `minres`, 567

MATLAB newton, 472
MATLAB ode15s, 1259
MATLAB ode23s, 1259
MATLAB ode23, 1215
MATLAB ode45, 1215
MATLAB odeset, 1215, 1259
MATLAB partition, 1113
MATLAB pcg, 535, 564
MATLAB pchip, 364
MATLAB polyfit, 767
MATLAB polyval, 319
MATLAB qmr, 570
MATLAB qrupdate, 277
MATLAB qr, 257
MATLAB quad1, 1092
MATLAB quad, 1092
MATLAB roudchol, 280
MATLAB rqui, 650
MATLAB secant, 459
MATLAB sinetrans, 921
MATLAB sinft2d, 924
MATLAB smw, 273
MATLAB spline, 375
MATLAB spowitrp, 693
MATLAB svds, 728
MATLAB svd, 728
MATLAB tic, toc, 925
MATLAB toeplitz, 935
MATLAB Adding eps to 1, 146
MATLAB Characteristic parameters of IEEE floating point numbers, 144
MATLAB Finding out eps in MATLAB, 146
MATLAB GE for “Wilkinson system”, 170
MATLAB Givens rotation, 250
MATLAB LU-factorization to solve LSE, 120
MATLAB Numerical differentiation, 325, 326
MATLAB gallery('circul' , v), 847
MATLAB newtoncotes, 1033
MATLAB ode45, 1187
MATLAB odeset, 1216
MATLAB pconvfft, 858
MATLAB pconv, 858
MATLAB polyfit, 948
MATLAB qrinsert, 285

MATLABrand, 601
MATLABrank, 731
MATLABsinetrans, 919
MATLABeig, 589
MATLABinput errors and rounding errors, 144
MATLABlinsolve, 239
MATLABqrllsqsolve, 791
MATLABrecursive Gaussian elimination, 104
MATLABrecursive LU-decomposition, 119
MATLABbisect , 445
MATLABblockgs, 104
MATLABchol, 221
MATLABdampnewton, 489
MATLABgausselim, 102
MATLABpevaltime, 318
MATLABplanerot, 250
MATLABpolyfit, 319
MATLABsa1, 204
MATLABsa2, 204
MATLABsa3, 206
MATLABsparse, 184
MATLABspdiags, 184
MATLABspeye, 184
MATLABspones, 184
MATLABspy, 198, 221
MATLABsymamd, 221
MATLABsymrcm, 221
3-term recursion
 for Chebychev polynomials, 959
 for Legendre polynomials, 1074
5-points-star-operator, 921
a posteriori error bound, 423
a posteriori termination, 419
a priori termination, 419
A-inner product, 503
A-orthogonal, 528
absolute tolerance, 419, 472, 1197
Acceptability Criteria, 1108
adaptive multigrid quadrature, 1086
adaptive quadrature, 1085
AGM, 416
Aitken-Neville scheme, 315
algebra, 52

algebraic convergence, 946
algebraic dependence, 61
analytic function, 974
Approximation
 Low rank, 1118
arrow matrix, 198, 200
asymptotic complexity, 54
asymptotic error behavior, 944
asymptotic rate of linear convergence, 438
Axiom of roundoff analysis, 145
AXPY operation, 534
axpy operation, 68
back substitution, 92
backward error analysis, 154
backward substitution, 120
Banach's fixed point theorem, 434
Bandbreite
 Zeilen-, 211
banded matrix, 209
bandwidth, 209
 lower, 209
 minimizing, 219
 upper, 209
barycentric interpolation formula, 313
basic linear algebra subroutines, 63
basis
 cosine, 926
 orthonormal, 260, 587
 sine, 917
 trigonometric, 852
Belousov-Zhabotinsky reaction, 1189
bending energy, 377
Besetzungsmuster, 221
best low rank approximation, 757
bicg, 569
BiCGStab, 570
bisection, 444
BLAS, 63
 axpy, 68
block LU-decomposition, 122
block matrix multiplication, 52
blow-up, 1192
blurring operator, 882

Bounding Box, 1109

Broyden

Quasi-Newton Method, 493

Broyden-Verfahren

convergence monitor, 496

Butcher scheme, 1184, 1252

cache miss, 40

cache thrashing, 40

cancellation, 326, 328

for Householder transformation, 249

capacitance, 177

capacitor, 177

cardinal

spline, 383

cardinal basis, 296

cardinal basis function, 382

cardinal interpolant, 382

Cauchy product

of power series, 835

causal filter, 826

cell

of a mesh, 1003

CG

convergence, 548

preconditioned, 555

termination criterion, 535

CG = conjugate gradient method, 522

CG algorithm, 532

chain rule, 474

characteristic polynomial, 583

Chebyshev expansion, 977

Chebyshev nodes, 963, 966, 968

Chebyshev polynomials, 543

3-term recursion, 959

Chebyshev-interpolation, 957

chemical reaction kinetics, 1240

Cholesky decomposition, 235

costs, 236

Cholesky factorization, 277

circuit simulation

transient, 1145

circulant matrix, 841

Classical Runge-Kutta method

- Butcher scheme, 1185, 1186
- Clenshaw algorithm, 978
- cluster analysis, 742
- Clustering Approximation
 - Task, 1093
- coil, 177
- Collocation Matrix, 1093, 1094
- column major matrix format, 35
- column sum norm, 150
- column transformation, 51
- complexity, 54
 - asymptotic, 54
 - linear, 58
 - of SVD, 729
- composite quadrature formulas, 1037
- compressed row storage, 182
- computational cost
 - Gaussian elimination, 97
- computational costs
 - LU-decomposition, 118
 - QR-decomposition, 262
- computational effort, 54, 466
- eigenvalue computation, 593
- condition number
 - of a matrix, 168
 - spectral, 521
- conjugate gradient method, 522
- consistency
 - of iterative methods, 406
 - fixed point iteration, 426
- constant
 - Lebesgue, 968
- constitutive relations, 177, 290
- constrained least squares, 804
- convergence
 - algebraic, 946
 - asymptotic, 461
 - exponential, 946, 954, 970
 - global, 407
 - iterative method, 406
 - linear, 409
 - linear in Gauss-Newton method, 820
 - local, 407
 - quadratic, 416

rate, 409

convergence monitor

- of Broyden method, 496

convolution

- discrete, 826, 834
- discrete periodic, 837
- of sequences, 835

cosine

- basis, 926
- transform, 926

cosine matrix, 926

cosine transform, 926

costs

- Cholesky decomposition, 236

Crout's algorithm, 116

CRS, 182

CRS format

- diagonal, 183

cubic Hermite interpolation, 310, 359

cubic spline interpolation

- error estimates, 1019

cyclic permutation, 203

damped Newton method, 485

damping factor, 487

data fitting, 764

- linear, 765
- polynomial, 767

data interpolation, 290

deblurring, 879

deblurring = Entrauschen, 879

definite, 148

dense matrix, 175

descent methods, 503

DFT, 845, 856

- two-dimensional, 877

Diagonal dominance, 227

diagonal matrix, 112

diagonalization

- of a matrix, 587

diagonalization of local translation invariant linear operators, 921

difference quotient, 325

- backward, 1161
- forward, 1160

difference scheme, 1160
 differential in non-linear least squares, 474
 direct power method, 615
 discrete convolution, 826, 834
 discrete Fourier transform, 845, 856
 discrete periodic convolution, 837
 divided differences, 332
 dot product, 41
 double precision, 142
 economical singular value decomposition, 729
 efficiency, 466
 eigenspace, 583
 eigenvalue, 583
 generalized, 588
 eigenvalue problem
 generalized, 588
 eigenvector, 583
 generalized, 588
 electric circuit, 176, 401
 resonant frequencies, 574
 elementary arithmetic operations, 139, 145
 elimination matrix, 110
 energy norm, 503
 envelope
 matrix, 211
 Equation
 non-linear, 404
 equidistant mesh, 1003
 ergodicity, 611
 error behavior
 asymptotic, 944
 error estimator
 a posteriori, 423
 Euler method
 explicit, 1158
 implicit, 1161
 semi implicit, 1257
 Euler polygon, 1159
 Euler's iteration, 457
 expansion
 asymptotic, 320
 explicit Euler method, 1158
 Butcher scheme, 1185

explicit midpoint rule
 Butcher scheme, 1185
 for ODEs, 1181
 explicit Runge-Kutta method, 1183
 explicit trapezoidal rule
 Butcher scheme, 1185
 exponential convergence, 970
 extended normal equations, 782
 extended state space
 of an ODE, 1148
 extrapolation, 320
 Far field, 1104
 fast Fourier transform, 906
 FFT, 906
 fill-in, 198
 filter
 high pass, 867
 low pass, 867
 finite filter, 826
 fixed point, 426
 fixed point form, 427
 fixed point iteration, 425
 fixed point iteration
 consistency, 426
 floating point number, 141
 floating point numbers, 140
 forward elimination, 91
 forward substitution, 120
 Fourier
 matrix, 853
 Fourier coefficient, 899
 Fourier series, 893
 Fourier transform, 893
 discrete, 845, 856
 isometry property, 903
 fractional order of convergence, 460
 frequency filtering, 859
 function representation, 292
 Funktion
 shandles, 472
 fzero, 446
 Gauss Quadrature, 1063

Gauss-Newton method, 814
 Gauss-Seidel preconditioner, 558
 Gaussian elimination, 89
 block version, 105
 by rank-1 modifications, 103
 for non-square matrices, 100
 instability, 170
 with partial pivoting, 127
 Gerschgorin circle theorem, 585
 Gibbs phenomenon, 998
 Givens rotation, 250, 274
 Givens-Rotation, 264, 284, 287
 global solution
 of an IVP, 1153
 GMRES, 568
 restarted, 568
 Golub-Welsch algorithm, 1077
 gradient, 475, 509
 Gram-Schmidt
 Orthonormalisierung, 709
 Gram-Schmidt orthogonalization, 530, 709
 Gram-Schmidt orthonormalization, 679, 787
 graph partitioning, 644
 grid, 1002
 grid cell, 1003
 grid function, 922
 grid interval, 1003
 Halley's iteration, 450, 457
 harmonic mean, 363
 hat function, 294
 heartbeat model, 1142
 Hermite interpolation
 cubic, 310
 Hessian, 231
 Hessian = Hesse-Matrix, 231
 Hessian matrix, 475
 high pass filter, 867
 homogeneous, 148
 Hooke's law, 661
 Horner scheme, 301
 Householder reflection, 247
 IEEE standard 754, 142
 ill conditioned, 172

image segmentation, 623
 implicit Euler method, 1161
 impulse response, 827
 of a filter, 826
 in place, 116, 118
 in situ, 104, 118
 in-situ, 132
 increment equations
 linearized, 1258
 increments
 Runge-Kutta, 1183, 1251
 inductance, 177
 inductor, 177
 inf, 143
 infinity, 143
 initial guess, 406, 425
 initial value problem
 stiff, 1247
 initial value problem (IVP), 1147
 initial value problem (IVP) = Anfangswertproblem, 1147
 inner product

A-, 503
 intermediate value theorem, 444
 interpolation
 barycentric formula, 313
 Chebychev, 957
 complete cubic spline, 374
 cubic Hermite, 359
 Hermite, 309
 Lagrange, 303
 natural cubic spline, 374
 periodic cubic spline, 375
 spline cubic, 369
 spline cubic, locality, 383
 spline shape preserving, 385
 trigonometric, 983
 interpolation operator, 298
 inverse interpolation, 463
 inverse iteration, 646
 preconditioned, 652
 inverse matrix, 87
 invertible matrix, 87, 88
 iteration

Halley's, 457
Euler's, 457
quadratical inverse interpolation, 457
iteration function, 406, 425
iterative method
 convergence, 406
IVP, 1147
Jacobi preconditioner, 558
Jacobian, 435, 470
Kernel Function, 1093, 1094
 separable, 1096
kinetics
 of chemical reaction, 1240
Kirchhoff (current) law, 177
knots
 spline, 367
Konvergenz
 Algebraische, Quadratur, 1082
Kronecker symbol, 31
Krylov space, 526
 for Ritz projection, 697
L-stable, 1254
Lagrange multiplier, 805
Lagrangian multiplier, 806
Landau-O, 55, 947
Lapack, 99
leading coefficient
 of polynomial, 300
Least squares
 with linear constraint, 804
least squares
 total, 801
least squares problem, 772, 773
 conditioning, 776
Lebesgue
 constant, 968
Lebesgue constant, 341
Levinson algorithm, 935
limit cycle, 1243
limiter, 362
line search, 508
linear complexity, 58
linear correlation, 738

linear data fitting, 765
linear electric circuit, 176
linear filter, 826
linear operator
 diagonalization, 921
linear ordinary differential equation, 581
linear regression, 62
linear system of equations, 86
 multiple right hand sides, 102
Lloyd-Max algorithm, 742
local a posteriori error estimation
 for adaptive quadrature, 1086
local convergence
 Newton method, 486
local linearization, 470
local mesh refinement
 for adaptive quadrature, 1086
locality
 of interpolation, 381
logistic differential equation, 1137
Lotka-Volterra ODE, 1140
low pass filter, 867
Low rank approximation, 1118
lower triangular matrix, 112
LU-decomposition
 blocked, 122
 computational costs, 118
 envelope aware, 215
 existence, 114
 in place, 118
LU-factorization
 envelope aware, 215
 of sparse matrices, 195
 with pivoting, 130
machine number, 141
 exponent, 141
machine numbers, 142
 distribution, 141
machine precision, 145
mantissa, 141
Markov chain, 599, 931
 stationary distribution, 602
mass matrix, 665

Matlab, 26

MATLAB \-operator, 99

Matrix

adjoint, 34

Collocation-, 1093

Hermitian, 587

Hermitian transposed, 34

normal, 587

obere Dreiecks, 281

obere Hessenberg, 275

skew-Hermitian, 587

transposed, 34

unitary, 587

matrix

banded, 209

condition number, 168

dense, 175

diagonal, 112

envelope, 211

Fourier, 853

Hermitian, 230

Hessian, 475

lower triangular, 112

normalized, 112

orthogonal, 243

positive definite, 230

positive semi-definite, 230

rank, 87

sine, 917

sparse, 175

storage formats, 35

structurally symmetric, 218

symmetric, 230

tridiagonal, 209

unitary, 243

upper triangular, 112

matrix algebra, 52

matrix block, 33

matrix factorization, 106, 110

matrix norm, 149

column sums, 150

row sums, 150

matrix product, 41

matrix storage

- envelope oriented, 218
- Matrixnorm, 149
 - Submultiplikativität, 150
- mesh, 1002
 - equidistant, 1003
 - in time, 1166
- mesh width, 1003
- Method
 - Quasi-Newton, 492
- midpoint rule, 1032, 1181
- Milleniums Algorithms, 1135
- Milne rule, 1034
- min-max theorem, 635
- minimal residual methods, 566
- model function, 447
- Modellfunktionsverfahren, 446
- modification technique
 - QR-factorization, 274
- monomial representation
 - of a polynomial, 300
- monomials, 300
- multi-point methods, 447, 458
- multiplicity
 - geometric, 583
 - of an interpolation node, 308
- NaN, 143
- Ncut, 626
- Near field, 1104
- nested
 - subspaces, 523
- Newton
 - basis, 330
 - damping, 487
 - damping factor, 487
 - monotonicity test, 488
 - simplified method, 476
- Newton correction, 470
 - simplified, 484
- Newton iteration, 470
 - numerical Differentiation, 477
 - termination criterion, 482
- Newton method
 - 1D, 447

damped, 485
local convergence, 486
local quadratic convergence, 477
region of convergence, 486
Newton's law of motion, 664
nodal analysis, 176, 401
nodal potentials, 177
node
 double, 309
 for interpolation, 303
 in electric circuit, 176
 multiple, 308
 multiplicity, 308
 of a mesh, 1003
 quadrature, 1026
nodes, 303
 Chebychev, 966
 Chebychev nodes, 968
 for interpolation, 303
non-linear data fitting, 810
non-normalized numbers, 143
norm, 148
 L^1 , 340
 L^2 , 340
 ∞ -, 149
 1-, 149
 energy-, 503
 Euclidean, 149
 of matrix, 149
 Sobolev semi-, 956
 supremum, 339
normal equations, 779
 extended, 782
 with constraint, 807
normalization, 615
normalized lower triangular matrix, 111
normalized triangular matrix, 112
not a number, 143
Nullstellenbestimmung
 Modellfunktionsverfahren, 446
numerical algorithm, 153
Numerical differentiation
 roundoff, 326

numerical Differentiation
 Newton iteration, 477
 numerical differentiation, 321, 325
 numerical quadrature, 1023
 numerical rank, 796
 Obere Hessenbergmatrix, 275
 ODE, 1147
 scalar, 1157
 Ohmic resistor, 177
 one-point methods, 447
 order
 of quadrature formula, 1048
 order of convergence, 413
 fractional, 460
 ordinary differential equation
 linear, 581
 ordinary differential equation (ODE), 1147
 oregonator, 1189
 orthogonal matrix, 243
 orthogonal polynomials, 1069
 orthonormal basis, 260, 587

overflow, 143
 page rank, 599
 stochastic simulation, 600
 partial pivoting, 130, 132
 choice of pivot, 132
 pattern
 of a matrix, 46
 PCA, 721
 PCG, 555
 Peano
 Theorem of, 1153
 penalization, 638
 penalty parameter, 639
 periodic sequence, 835
 permutation, 134
 permutation matrix, 134
 Permutationsmatrix, 282
 perturbation lemma, 167
 Petrov-Galerkin condition, 569
 phase space
 of an ODE, 1148

Picard-Lindelöf

Theorem of, 1153

piecewise quadratic interpolation, 353

PINVIT, 652

Pivot

choice of, 128

pivot, 92, 94, 95

pivot row, 92, 95

pivoting, 124

point spread function, 879

polynomial

characteristic, 583

Lagrange, 303

polynomial fitting, 767

polynomial interpolation

existence and uniqueness, 305

generalized, 308

polynomial space, 300

positive definite

criteria, 230

matrix, 230

potentials

nodal, 177

power spectrum

of a signal, 869

preconditioned CG method, 555

preconditioned inverse iteration, 652

preconditioner, 552

preconditioning, 550

predator-prey model, 1140

principal axis, 742

principal axis transformation, 515, 587

principal component, 738

principal component analysis, 721

principal minor, 122

problem

ill conditioned, 172

sensitivity, 172

well conditioned, 172

product rule, 475

pseudoinverse, 796

Punkt

stationär, 1142

pwer method

- direct, 615
- QR algorithm, 590
- QR-algorithm with shift, 590
- QR-decomposition
 - computational costs, 262
- QR-factorization, QR-decomposition, 253
- quadratic convergence, 441
- quadratic eigenvalue problem, 574
- quadratic functional, 504
- quadratic interpolation
 - piecewise, 353
- quadratic inverse interpolation, 464
- quadratical inverse interpolation, 457
- quadrature
 - adaptive, 1085
 - polynomial formulas, 1030
- quadrature formula, 1026
 - local, 1039
 - order, 1048
- quadrature node, 1026
- quadrature numerical, 1023
- quadrature weight, 1026
- Quasi-Newton Method, 493
- Quasi-Newton method, 492
- Radau RK-method
 - order 3, 1255
 - order 5, 1255
- radiative heat transfer, 838
- rank
 - column rank, 88
 - computation, 729
 - numerical, 796
 - of a matrix, 87
 - row rank, 88
- rank-1 modification, 103
- rank-1-modifications, 270
- rate
 - of algebraic convergence, 946
 - of convergence, 409
- Rayleigh quotient, 617, 635
- Rayleigh quotient iteration, 649
- regular matrix, 87

relative tolerance, 419, 472, 1197
rem:Fspec, 854
rem:polyerrrep, 952
residual quantity, 653
Riccati differential equation, 1157, 1159
right hand side
 of an ODE, 1148
Riemann sum, 899
right hand side vector, 86, 153
rigid body mode, 668
Ritz projection, 686, 696
Ritz value, 686
Ritz vector, 686
root of unity, 851
rounding, 145
rounding up, 145
roundoff
 for numerical differentiation, 326
row major matrix format, 35
row permutation, 287
row sum norm, 150
row transformation, 51, 90, 108, 109

Runge-Kutta
 increments, 1183, 1251
Runge-Kutta method, 1183, 1251
 L-stable, 1254
Runge-Kutta methods
 semi-implicit, 1256
 stability function, 1229, 1253

saddle point problem, 806
 matrix form, 807
scalar ODE, 1157
scaling
 of a matrix, 48
scheme
 Aitken-Neville, 315
 Horner, 301
Schur
 Komplement, 123
Schur complement, 123
Schur's lemma, 586
scientific notation, 140
secant condition, 492

secant method, 458, 492

segmentation
 of an image, 623

semi-implicit Euler method, 1257

seminorm, 956

sensitivity
 of a problem, 172

shape
 preservation, 353
 preserving spline interpolation, 385

Sherman-Morrison-Woodbury formula, 272

shifted inverse iteration, 647

similarity
 of matrices, 585

similarity function
 for image segmentation, 625

similarity transformations, 585

similarity transformation
 unitary, 590

Simpson rule, 1034

sine
 basis, 917

 matrix, 917

 transform, 918

Sine transform, 917

single precision, 142

single step method, 1166

singular value decomposition, 719, 726

sparse matrix, 175
 initialization, 187
 LU-factorization, 195
 multiplication, 191

sparse matrix storage formats, 181

spectral condition number, 521

spectral partitioning, 644

spectral radius, 584

spectrum, 583
 of a matrix, 514

spline, 367
 cardinal, 383
 complete cubic, 374
 cubic, 369
 cubic, locality, 383
 knots, 367

- natural cubic, 374
- periodic cubic, 375
- physical, 379
- shape preserving interpolation, 385
- square root
 - of a matrix, 551
- stability function
 - of explicit Runge-Kutta methods, 1229
 - of Runge-Kutta methods, 1253
- Stable
 - algorithm, 153
- stable
 - numerically, 153
- state space
 - of an ODE, 1148
- stationary distribution, 602
- steepest descent, 508
- stiff IVP, 1247
- stiffness matrix, 665
- stochastic matrix, 603
- stochastic simulation of page rank, 600
- Strassen's algorithm, 56
- structurally symmetric matrix, 218
- sub-matrix, 33
- sub-multiplicative, 150
- subspace correction, 523
- subspace iteration
 - for direct power method, 692
- subspaces
 - nested, 523
- SVD, 719, 726
- symmetry
 - structural, 218
- system matrix, 86, 153
- system of equations
 - linear, 86
- tangent field, 1157
- Taylor expansion, 439
- Taylor formula, 941
- Taylor's formula, 439
- Tensor product
 - Chebyshev interpolation polynomial, 1099, 1119, 14.4
 - interpolation polynomial, 1098

tensor product, 41
tent function, 294
Toeplitz matrices, 929
termination criterion, 418
 Newton iteration, 482
 reliable, 420
 residual based, 421
time-invariant filter, 826
timestep (size), 1159
timestep constraint, 1230
timestepping, 1158
Toeplitz solvers
 fast algorithms, 938
tolerance
 absolute, 445
tolerance, 421
 absolute, 1197
 absoute, 419, 472
 for adaptive timestepping for ODEs, 1195
 for termination, 419
 realtive, 1197
 relative, 419, 472
total least squares, 801
trajectory, 1141
transform
 cosine, 926
 fast Fourier, 906
 sine, 918
trapezoidal rule, 1033, 1180
 for ODEs, 1181
trend, 721
triangle inequality, 148
tridiagonal matrix, 209
trigonometric basis, 852
trigonometric function, 984
trigonometric interpolation, 983
trigonometric polynomial, 901
trigonometric polynomials, 984
trigonometric transformations, 916
truss structure
 vibrations, 660
trust region method, 822
underflow, 143

Uniform convergence

of Fourier series, 894

unit vector, 31

unitary matrix, 243

unitary similiary transformation, 590

upper Hessenberg matrix, 711

upper triangular matrix, 92, 111, 112

Vandermonde matrix, 306

variational calculus, 377

Weddle rule, 1034

weight

quadrature, 1026

well conditioned, 172

Young's modulus, 662

Zerlegung

LU, 120

QR, 281

zero padding, 843, 935, 994

List of Symbols

$(\mathbf{A})_{i,j} \hat{=}$ reference to entry a_{ij} of matrix \mathbf{A} , 33
 $(\mathbf{A})_{k:l,r:s} \hat{=}$ reference to submatrix of \mathbf{A} spanning rows k, \dots, l and columns r, \dots, s , 33
 $(\mathbf{x})_i \hat{=}$ i -th component of vector \mathbf{x} , 31
 $(x_k) * _n (y_k) \hat{=}$ discrete periodic convolution, 838
 $C_{\text{pw}}^0(I) \hat{=}$ space of piecewise continuous functions on interval I , 1004
 $C^1([a, b]) \hat{=}$ space of continuously differentiable functions $[a, b] \mapsto \mathbb{R}$, 355
 $D\Phi \hat{=}$ **Jacobian** of $\Phi : D \mapsto \mathbb{R}^n$ at $\mathbf{x} \in D$, 435
 $D_{\mathbf{y}}\mathbf{f} \hat{=}$ Derivative of \mathbf{f} w.r.t.. \mathbf{y} (Jacobian), 1152
 $J(t_0, \mathbf{y}_0) \hat{=}$ maximal domain of definition of a solution of an IVP, 1153
 $O \hat{=}$ zero matrix, 35

$O(n)$, 55
 $\mathcal{E} \hat{=}$ expected value of a random variable, 932
 $\mathcal{P}_n^T \hat{=}$ space of **trigonometric polynomials** of degree n , 984
 $\mathcal{R}_k(m, n)$, 758
 $\text{eps} \hat{=}$ machine precision, 145
 $\text{Eig}_{\mathbf{A}}(\lambda) \hat{=}$ eigenspace of \mathbf{A} for eigenvalue λ , 583
 $\text{Im}(\mathbf{A}) \hat{=}$ range/column space of matrix \mathbf{A} , 730
 $\text{Ker}(\mathbf{A}) \hat{=}$ nullspace of matrix \mathbf{A} , 730
 $\mathcal{K}_l(\mathbf{A}, \mathbf{z}) \hat{=}$ Krylov subspace, 526
 $\|\mathbf{Ax} - \mathbf{b}\|_2 \rightarrow \min \hat{=}$ minimize $\|\mathbf{Ax} - \mathbf{b}\|_2$, 772
 $\|\mathbf{A}\|_F^2$, 757

$\|\mathbf{x}\|_A \hat{=}$ energy norm induced by s.p.d. matrix \mathbf{A} , 503
 $\|f\|_{L^\infty(I)}$, 339
 $\|f\|_{L^1(I)}$, 340
 $\|f\|_{L^2(I)}^2$, 340
 \mathcal{P}_k , 300
 $\Psi^h \mathbf{y} \hat{=}$ discrete evolution for autonomous ODE, 1166
 $\mathcal{S}_{d,\mathcal{M}}$, 367
 \mathbf{A}^+ , 776
 $\mathbf{A}^\top \hat{=}$ transposed matrix, 34
 $\mathbf{I} \hat{=}$ identity matrix, 35
 $\mathbf{h} * \mathbf{x} \hat{=}$ discrete convolution of two vectors, 834
 $\mathbf{x} *_n \mathbf{y} \hat{=}$ discrete periodic convolution of vectors, 838
 $\bar{z} \hat{=}$ complex conjugation, 34
 $\mathbb{M} \hat{=}$ set of machine numbers, 138
 $\mathbb{S}^1 \hat{=}$ unit circle in \mathbb{C} , 987
 $\delta_{ij} \hat{=}$ Kronecker symbol, 31, 303
 $\mathbf{i} \hat{=}$ imaginary unit, “ $\mathbf{i} := \sqrt{-1}$ ”, 178
 $\kappa(\mathbf{A}) \hat{=}$ spectral condition number, 521

$\lambda_{\max} \hat{=}$ largest eigenvalue (in modulus), 521
 $\lambda_{\min} \hat{=}$ smallest eigenvalue (in modulus), 521
 $\mathbf{1} = (1, \dots, 1)^T$, 1229, 1253
 $\text{Ncut}(\mathcal{X}) \hat{=}$ normalized cut of subset of weighted graph, 626
 $\text{argmin} \hat{=}$ (global) minimizer of a functional, 507
 $\text{cond}(\mathbf{A})$, 168
 $\text{cut}(\mathcal{X}) \hat{=}$ cut of subset of weighted graph, 626
 $\text{env}(\mathbf{A})$, 211
 nnz , 175
 $\text{rank}(\mathbf{A}) \hat{=}$ rank of matrix \mathbf{A} , 87
 rd , 145
 $\text{weight}(\mathcal{X}) \hat{=}$ connectivity of subset of weighted graph, 626
 $\overline{m}(\mathbf{A})$, 209
 $\rho(\mathbf{A}) \hat{=}$ spectral radius of $\mathbf{A} \in \mathbb{K}^{n,n}$, 584
 $\rho_{\mathbf{A}}(\mathbf{u}) \hat{=}$ Rayleigh quotient, 617
 $\mathbf{f} \hat{=}$ right hand side of an ODE, 1148
 $\sigma(\mathbf{A}) \hat{=}$ spectrum of matrix \mathbf{A} , 583
 $\sigma(\mathbf{M}) \hat{=}$ spectrum of matrix \mathbf{M} , 514
 $\tilde{\star}$, 144
 $\underline{m}(\mathbf{A})$, 209

$m(\mathbf{A})$, 209

$y[t_i, \dots, t_{i+k}] \hat{=} \text{divided difference}$, 332

$\| \mathbf{x} \|_1$, 149

$\| \mathbf{x} \|_2$, 149

$\| \mathbf{x} \|_\infty$, 149

$\dot{\hat{=}}$ Derivative w.r.t. time t , 1138

TOL tolerance, 1195

List of Definitions

- Analytic function, 974
- Arrow matrix, 200
- Chebyshev polynomial, 958
- circulant matrix, 841
- Cluster
 - Tree, 1109
- concave
 - data, 347
 - function, 348
- Condition (number) of a matrix, 168
- Consistency of fixed point iterations, 426
- Consistency of iterative methods, 406
- Contractive mapping, 433
- Convergence, 406
 - global, 407
 - local, 407
- convex
 - data, 347
 - function, 348
- data
 - concave, 347
 - convex, 347
- Diagonally dominant matrix, 227
- Discrete convolution, 834
- discrete Fourier transform, 856
- discrete periodic convolution, 837
- eigenvalues and eigenvectors, 583
- energy norm, 503

equivalence of norms, 410
Evolution operator, 1155
Explicit Runge-Kutta method, 1183
Fill-In, 198
Frobenius norm, 757
function
 concave, 348
 convex, 348
Hessian matrix, 475
Inverse of a matrix, 87
Krylov space, 526
L-stable Runge-Kutta method, 1254
Lebesgue constant, 341
Legendre polynomials, 1070
Linear convergence, 409
Lipschitz continuous function, 1151, 1152
machine numbers, 141
matrix
 generalized condition number, 777
 s.p.d, 230
 symmetric positive definite, 230
Matrix envelope, 211
matrix norm, 149
monotonic data, 346
norm, 148
 Frobenius norm, 757
Normalized cut, 626
numerical algorithm, 153
Order of convergence, 413
orthogonal matrix, 243
Permutation matrix, 134
polynomial
 Chebychev, 958
 generalized Lagrange, 309
Polynomial
 Interpolation tensor product, 1098
pseudoinverse, 776
Rank of a matrix, 87
Rayleigh quotient, 617

residual, 161

Runge-Kutta method, 1251

Single step method, 1166

singular value decomposition (SVD), 726

sparse matrices, 175

sparse matrix, 175

splines, 367

stable algorithm, 153

Structurally symmetric matrix, 218

Tensor product

 interpolation polynomial, 1098

Toeplitz matrix, 933

Types of matrices, 112

unitary matrix, 243

Examples and Remarks

- LU -decomposition of sparse matrices, 196
- L^2 -error estimates for polynomial interpolation, 955
- h -adaptive numerical quadrature, 1089
- p -convergence of piecewise polynomial interpolation, 1011
- (Nearly) singular LSE in shifted inverse iteration, 648
- [Bad behavior of global polynomial interpolants, 349
- [Conditioning of row transformations, 241
- [From higher order ODEs to first order systems, 1150
- [Stable solution of LSE by means of QR-decomposition, 263
- ode45 for stiff problem, 1220
- “Annihilating” orthogonal transformations in 2D, 246
- “Butcher barriers” for explicit RK-SSM, 1186
- “Failure” of adaptive timestepping, 1209
- “Low” and “high” frequencies, 865
- “Software solution” of interpolation problem, 298
- “Squeezed” DFT of a periodically truncated signal, 887
- $\mathbf{B} = \mathbf{B}^H$ s.p.d. mit Cholesky-Zerlegung, 589
- L-stable implicit Runge-Kutta methods, 1255
- fft
 - Efficiency, 904
- q_1 -norm from eigenvalues, 520
- 3-Term recursion for Legendre polynomials, 1074

3-term recursion for Chebychev polynomials, 959

A posteriori error bound for linearly convergent iteration, 423

A posteriori termination criterion for linearly convergent iterations, 422

A posteriori termination criterion for plain CG, 535

Accessing rows and columns of sparse matrices, 184

Adapted Newton method, 453

Adaptive integrators for stiff problems in MATLAB, 1259

Adaptive quadrature in MATLAB, 1092

Adaptive timestepping for mechanical problem, 1215

Adding *EPS* to 1, 146

Affine invariance of Newton method, 473

Algorithm for cluster analysis, 742

Analytic solution of homogeneous linear ordinary differential equations, 581

Approximate computation of Fourier coefficients, 1059

Approximation

convergence Cluster- with collocation matrix, 1130

approximation

uses of, 942

Approximation and quadrature, 1030

Approximation by polynomials, 943

Arnoldi process Ritz projection, 712

Asymptotic perspective in convergence analysis, 1168

auxiliary construction for shape preserving quadrature interpolation, 389

Banach's fixed point theorem, 434

bandwidth, 209

BLAS calling conventions, 70

Block Gaussian elimination, 105

Block LU-factorization, 122

Block matrix product, 52

Blow-up, 1192

Blow-up of explicit Euler method, 1222

Blow-up solutions of vibration equations, 658

Bound for asymptotic rate of linear convergence, 438

Broyden method for a large non-linear system, 499

Broydens Quasi-Newton method: convergence, 494

Butcher scheme for some explicit RK-SSM, 1185

Cancellation in decimal floating point arithmetic, 328

CG convergence and spectrum, 548

Changing entries/rows/columns of a matrix, 270

Characteristic parameters of IEEE floating point numbers, 143

Chebyshev interpolation error, 968

Chebyshev interpolation of analytic function, 975

Chebyshev interpolation of analytic functions, 973

Chebyshev polynomials on arbitrary interval, 965

Chebyshev representation of built-in functions, 983

Chebyshev vs uniform nodes, 966

Choice of quadrature weights, 1060

Choice of unitary/orthogonal transformation, 256

Class PolyEval, 334

Classification from measured data, 720

Cluster Tree, 1112

Complexity of Householder QR-factorization, 257

Computational effort for eigenvalue computations, 593

Computing Gauss nodes and weights, 1077

condition
 extended system, 784

Conditioning and relative error, 169

Conditioning of normal equations, 780

Conditioning of the extended normal equations, 783

Conditioning of the least squares problem, 776

Constitutive relations from measurements, 290

Construction of simple Runge-Kutta methods, 118

Convergence der cluster approximation, 1129

Convergence of CG as iterative solver, 539

Convergence of clustering approximation with collocation matrix, 1130

Convergence of equidistant trapezoidal rule, 1052

- Convergence of Fourier sums, 897
- Convergence of gradient method, 516
- Convergence of Hermite interpolation, 1015
- Convergence of Hermite interpolation with exact slopes, 1012
- Convergence of Krylov subspace methods for non-symmetric system matrix, 571
- Convergence of Newton's method in 2D, 478
- Convergence of PINVIT, 655
- Convergence of simple Runge-Kutta methods, 1181
- Convergence of subspace variant of direct power method, 693
- Convergence rates for CG method, 546
- Convergence theory for PCG, 557
- Conversion into autonomous ODE, 1149
- Convolution of sequences, 835
- Cosine transforms for compression, 928
- CRS format, 182
- cubic Hermite interpolation, 310
- Damped Newton method, 490
- Data points confined to a subspace, 734
- Deblurring by DFT, 879
- Decay conditions for bi-infinite signals, 894
- Decimal floating point numbers, 140
- Details of Householder reflections, 248
- Detecting linear convergence, 411
- Detecting order of convergence, 415
- Detecting periodicity in data, 862
- Different implementations of matrix multiplication in MATLAB, 64
- Differentiation repetition, 474
- Direct power method, 618
- Divided differences and derivatives, 337
- Domain of definition of solutions of IVPs, 1153
- Efficiency of iterative methods, 468
- Efficiency of fft , 904
- Efficient associative matrix multiplication, 58
- Efficient evaluation of trigonometric interpolation polynomials, 993
- Efficient Initialization of sparse matrices in MATLAB, 187

- Eigenvalue computation with Arnoldi process, 716
- Eigenvectors of circulant matrices, 845
- Envelope of a matrix, 211
- Envelope oriented matrix storage, 218
- Error estimates for polynomial quadrature, 1036
- Error of Gauss quadrature, 1079
- Error of polynomial interpolation, 954
- Euler methods for stiff logistic IVP, 1249
- Explicit Euler method as difference scheme, 1160
- Explicit Euler method for damped oscillations, 1238
- Explicit integrator in MATLAB, 1187
- Explicit representation of error of polynomial interpolation, 952
- Explicit trapezoidal rule for decay equation, 1227
- Exploring convergence, 947
- Extended normal equations, 782
- Extremal properties of natural cubic spline interpolant, 376
- Failure of damped Newton method, 490
- Failure of Krylov iterative solvers, 570
- Fast evaluation of Chebychev expansion, 977
- Fast matrix multiplication, 56
- Fast Toeplitz solvers, 938
- Feasibility of implicit Euler timestepping, 1162
- FFT algorithm by matrix factorization, 909
- FFT based on general factorization, 912
- FFT for prime vector length, 914
- Filtering in Fourier domain, 900
- Finite linear time-invariant causal channel, 826
- Fit of hyperplanes, 797
- Fixed points in 1D, 431
- Fractional order of convergence of secant method, 461
- Frequency filtering by DFT, 868
- Frequency identification with DFT, 861
- Function representation, 292
- Gain through adaptivity, 1204
- Gaining efficiency through usage of BLAS, 71
- Gaussian elimination, 91
- Gaussian elimination and LU-factorization, 109
- Gaussian elimination for non-square matrices, 100

- Gaussian elimination via rank-1 modifications, 103
- Gaussian elimination with pivoting for 3×3 -matrix, 127
- Generalized eigenvalue problems and Cholesky factorization, 589
- Generalized polynomial interpolation, 308
- Gibbs phenomenon, 998
- Global separable approximation by non-smooth kernel function, 1101
- Global separable approximation by smooth kernel function, 1100
- Gradient method in 2D, 513
- Gravitational forces in galaxy, 1095
- Group property of autonomous evolutions, 1155
- Growth with limited resources, 1137
- Halley's iteration, 450
- Heartbeat model, 1142
- Heating generation in electrical circuits, 1025
- Hermite interpolation
 - theorem, 953
- Horner scheme, 301
- IEEE standard
 - special cases, 142
- IEEE standard 754 for machine numbers, 142
- Image compression, 761
- Image segmentation, 623
- Impact of choice of norm, 410
- Impact of data access patterns on runtime, 36
- Impact of roundoff errors on CG, 537
- Impact of roundoff on Lanczos process, 705
- Implicit Euler timestepping for decay equation, 1248
- Importance of numerical quadrature, 1024
- In-situ LU-decomposition, 118
- Initial guess for power iteration, 621
- Input errors and rounding errors, 144
- Instability of multiplication with inverse, 164
- Instability of normal equations, 781
- Interaction calculations for many body systems, 1094
- interpolation
 - cubic spline- locality, 382
 - piecewise cubic monotonicity preserving, 364

- shape preserving quadratic spline, 393
- Interpolation as linear mapping, 296
- interpolation error, 945
- Interpolation error: trigonometric interpolation, 996
- Intersection of lines in 2D, 173
- Justification of Ritz projection by min-max theorem, 686
- Keeping track of unitary transformations, 260
- Kinetics of chemical reactions, 1240
- Krylov methods for complex s.p.d. system matrices, 504
- Krylov subspace methods for generalized EVP, 718
- Lanczos process for eigenvalue computation, 704
- Least squares data fitting, 764
- Lebesgue constant for equidistant nodes, 341
- Linear data fitting, 765
- Linear filtering of periodic signals, 835
- linear regression, 770
- Linear regression for stationary Markov chains, 931
- Lineare zeitinvariante Systeme, 929
- Linearization of increment equations, 1257
- Linearly convergent iteration, 412
- Local convergence of Newton's method, 486
- local convergence of secant method, 462
- Loss of sparsity when forming normal equations, 782
- Machine precision for MATLAB, 146
- Magnetization curves, 346
- Many sequential solutions of LSE, 121
- Matrix algebra, 52
- Matrix storage formats, 35
- Meaningful " O -bounds" for complexity, 55
- Midpoint rule, 1032
- Min-max theorem, 635
- Multidimensional fixed point iteration, 438
- Multiplication of polynomials, 832
- Multiplication of sparse matrices, 191
- Necessary condition for L-stability, 1254

- Necessity of iterative approximation, 405
- Newton method and minimization of quadratic functional, 814
- Newton method in 1D, 448
- Newton method, modified , 457
- Newton-Cotes formulas, 1032
- Nodal analysis of (linear) electric circuit, 176
- Non-linear data fitting, 810
- Non-linear data fitting (II), 817
- Non-linear electric circuit, 401
- Non-linear interpolation, 365
- Normal equations vs. orthogonal transformations method, 797
- Notation for single step methods, 1167
- numerical differentiation through extrapolation, 321
- Numerical integration of logistic ODE in MATLAB, 1188
- Numerical summation of Fourier series, 895
- Occurrence of clusters in partition rectangles, 1124
- Options for fixed point iterations, 427
- Orders of simple polynomials quadrature formulas, 1049
- Oregonator reaction, 1189
- Origin of the term “Spline”, 379
- oscillating interpolation polynomial, 338
- Overdetermined linear systems, 771
- Page rank algorithm, 599
- PCA for data classification, 735
- PCA of stock prices, 748
- piecewise cubic Hermite interpolation, 359
- Piecewise cubic interpolation schemes, 376
- Piecewise linear interpolation, 293
- Piecewise polynomial interpolation, 1006
- Piecewise quadratic interpolation, 353
- Pivoting and numerical stability, 124
- Pivoting destroys sparsity, 207
- Polybomial interpolation vs. polynomial fitting, 767
- Polynomial
 - trigonometric analysis, 1133
- Polynomial evaluation: timing, 316
- Polynomial fitting, 767
- Polynomials in Matlab, 301

- Power iteration, 613
- Power iteration with Ritz projection, 688
- Predator-prey model, 1140
- Principal component analysis, 721
- Pseudoinverse, 776
- Pseudoinverse and SVD, 796
- QR
 - Orthogonalisierung, 259
- qr based orthogonalization, 682
- QR-Algorithm, 590
- QR-based solution of tridiagonal LSE, 264
- Quadratic convergence], 416
- Quadratic functional in 2D, 505
- quadratic inverse interpolation, 465
- Quadratur
 - Gauss-Legendre Ordnung 4, 1064
- Quadrature errors for composite quadrature rules, 1042
- Quality measure for kernel approximation, 1098
- Radiative heat transfer, 838
- Rank defect in linear least squares problems, 775
- Rationale for adaptive quadrature, 1085
- Rationale for high-order single step methods, 1177
- Rationale for partial pivoting policy, 132
- Rayleigh quotient iteration, 650
- Reading off polynomials complexity, 61
- Recursive LU-factorization, 119
- Reducing fill-in by reordering, 221
- reduction to periodic convolution, 842
- Refined local stepsize control, 1211
- Region of convergence of Newton method, 486
- Relevance of asymptotic complexity, 62
- Removing a singularity by transformation, 1051
- Resistance to currents map, 268
- Resonances of linear electrical circuits, 574
- Restarted GMRES, 568
- Ritz projections onto Krylov space, 697
- Roundoff errors and difference quotients, 325
- Row and column transformations, 51
- Row swapping commutes with forward elimination, 136

Row-wise & column-wise view of matrix product, 43
 Runge's example, 948, 955
 Runtime comparison for computation of coefficient of trigonometric interpolation polynomials, 991
 Runtime of Gaussian elimination, 97
 Runtimes of `eig`, 593
 S.p.d. Hessians, 231
 S.p.d. matrices from nodal analysis, 224
 Scalings, 48
 secant method, 459
 Sensitivity of linear mappings, 240
 Shape preservation, 381
 Shifted inverse iteration, 647
 Silly MATLAB, 193
 Simple adaptive stepsize control, 1200
 Simple adaptive timestepping for fast decay, 1225
 Simple composite polynomial quadrature rules, 1039
 Simple preconditioners, 557
 Simplified Newton method, 476
 Small residuals by Gaussian elimination, 161
 Solving LSE in the case of rank-1-modification, 272
 Sound filtering by DFT, 869
 Sparse LU -factors, 198
 Spectrum of Fourier matrix, 854
 Speed of convergence of explicit Euler method, 1168
 spline
 interpolants, approx. complete cubic, 1019
 natural cubic, locality, 382
 shape preserving quadratic interpolation, 393
 Square root of a s.p.d. matrix, 551
 Stability by small random perturbations, 158
 Stability of Arnoldi process, 714
 Stepsize control detects instability, 1231
 Stepsize control in MATLAB, 1214
 Storing the Q -factor, 266
 Strongly attractive limit cycle, 1243
 Subspace power iteration with orthogonal projection, 674

- Subspace power methods, 694
- SVD and additive rank-1 decomposition, 727
- Tangent field and solution curves, 1157
- Taylor approximation, 941
- Tensor product Chebyshev interpolation for variable rectangle sizes, 1107
- Tensor product Chebyshev interpolation on rectangles, 1106
- Termination criterion for contractive fixed point iteration, 442
- Termination criterion for direct power iteration, 622
- Termination criterion in `pcg`, 564
- Termination of PCG, 562
- Testing for near singularity of a matrix, 267
- Transformation of quadrature rules, 1027
- Transient circuit simulation, 1145
- Transient simulation of RLC-circuit, 1232
- Trend analysis, 719
- Tridiagonal preconditioner, 559
- Trigonometric interpolation of analytic functions, 999
- Understanding the structure of product matrices, 43
- Uniqueness of SVD, 728
- Unitary similarity transformation to tridiagonal form 591
- Vibrations of a truss structure, 660
- Visualization of explicit Euler method, 1159
- Why using $\mathbb{K} = \mathbb{C}$?, 850
- Wilkinson's counterexample, 157, 169
- Zerlegung
 - Teil-LU, 122

German terms

L^2 -inner product = L^2 -Skalarprodukt, 1066

], 92

bandwidth = Bandbreite, 209

cancellation = Auslöschung, 328

capacitance = Kapazität, 177

capacitor = Kondensator, 177

circulant matrix = zirkulante Matrix, 841

coil = Spule, 177

column (of a matrix) = (Matrix)spalte, 32

column transformation = Spaltenumformung, 51

column vector = Spaltenvektor, 30

composite quadrature formulas = zusammengesetzte Quadraturformeln, 1037

computational effort = Rechenaufwand, 54

consistency = Konsistenz, 406

constitutive relation = Kennlinie, 290

constitutive relations = Bauelementgleichungen,
177

constrained least squares = Ausgleichsproblem
mit Nebenbedingungen, 804

convergence = Konvergenz, 406

convolution = Faltung, 826

damped Newton method = gedämpftes Newton-
Verfahren, 485

dense matrix = vollbesetzte Matrix, 175

descent methods = Abstiegsverfahren, 503

divided differences = dividierte Differenzen, 332

dot product = (Euklidisches) Skalarprodukt, 41

- eigenspace = Eigenraum, 583
- eigenvalue = Eigenwert, 583
- electric circuit = elektrischer Schaltkreis/Netzwerk, 176
- energy norm = Energienorm, 503
- envelope = Hülle, 211
- extended normal equations = erweiterte Normalengleichungen, 782
- fill-in = "fill-in", 198
- fixed point iteration = Fixpunktiteration, 425
- floating point number = Gleitpunktzahl, 141
- forward elimination = Vorwärtselimination, 91
- Fourier series = Fourierreihe, 893
- Gaussian elimination = Gausselimination, 89
- high pass filter = Hochpass, 867
- identity matrix = Einheitsmatrix, 34
- image segmentation = Bildsegmentierung, 623
- impulse response = Impulsantwort, 827
- in situ = am Ort, 118
- in situ = an Ort und Stelle, 104
- initial guess = Anfangsnäherung, 406
- inverse iteration = inverse Iteration, 623
- Kirchhoff (current) law = Kirchhoffsche Knotenregel, 177
- knot = Knoten, 367
- Krylov space = Krylovraum, 526
- least squares = Methode der kleinsten Quadrate, 764
- line search = Minimierung in eine Richtung, 508
- linear system of equations = lineares Gleichungssystem, 86
- low pass filter = Tiefpassfilter, 867
- lower triangular matrix = untere Dreiecksmatrix, 112
- LU-factorization = LR-Zerlegung, 106
- machine number = Maschinenzahl, 141
- mass matrix = Massenmatrix, 665
- matrix factorization = Matrixzerlegung, 106
- mesh width = Gitterweite, 1003
- mesh/grid = Gitter, 1002

multiplicity = Vielfachheit, 583

nodal analysis = Knotenanalyse, 176

normal equations = Normalgleichungen, 779

order of convergence = Konvergenzordnung, 413

ordinary differential equation = gewöhnliche Differentialgleichung, 1147

partial pivoting = Spaltenpivotsuche, 130

pattern (of a matrix) = Besetzungsmuster, 48

power method = Potenzmethode, 599

preconditioning = Vorkonditionierung, 550

predator-prey model = Räuber-Beute-Modell, 1140

principal axis transformation = Hauptachsentransformation, 587

principal component analysis = Hauptkomponentenanalyse, 721

quadratic functional = quadratisches Funktional, 504

quadrature node = Quadraturknoten, 1026

quadrature weight = Quadraturgewicht, 1026

resistor = Widerstand, 177

right hand side vector = rechte-Seite-Vektor, 86

rounding = Rundung, 145

row (of a matrix) = (Matrix)zeile, 32

row transformation = Zeilenumformung, 51

row vector = Zeilenvektor, 30

saddle point problem = Sattelpunktproblem, 806

scaling = Skalierung, 48

singular value decomposition = Singulärwertzerlegung, 719

sparse matrix = dünnbesezte Matrix, 175

speed of convergence = Konvergenzgeschwindigkeit, 409

steepest descent = steilster Abstieg, 508

stiffness matrix = Steifigkeitsmatrix, 665

storage format = Speicherformat, 35

subspace correction = Unterraumkorrektur, 523

tensor product = Tensorprodukt, 41

tent/hat function = Hutfunktion, 294

termination criterion = Abbruchbedingung, 418

timestepping = Zeitdiskretisierung, 1158

total least squares = totales Ausgleichsproblem,
801

transpose = transponieren/Transponierte, 30

truss = Stabwerk, 660

unit vector = Einheitsvektor, 31

upper triangular matrix = obere Dreiecksmatrix,
112

variational calculus = Variationsrechnung, 377

zero padding = Ergänzen durch Null, 994