

Numerische Mathematik

für Studiengang Rechnergestützte Wissenschaften

Profs. Ralf Hiptmair und Rolf Jeltsch

Draft version 15. Januar 2007, Subversion rev #1769

(C) Seminar für Angewandte Mathematik, ETH Zürich

http://www.sam.math.ethz.ch/~hiptmair/tmp/NCSE_06.pdf

Inhaltsverzeichnis

0.1	Danksagung	11
1	Computerarithmetik und Konsequenzen	13
1.1	Beispiele	13
1.2	Zahldarstellung	23
1.3	Gleitpunktarithmetik und Rundungsfehler	27
1.5	Über- und Unterlauf	32
1.6	Kondition	33
1.7	Auslöschung	37
1.8	Stabilität	47

0.0
p. 1

2	Nichtlineare Gleichungen	53
2.1	Iterationsverfahren	54
2.2	Fixpunktiteration	68
2.3	Nullstellenbestimmung von Funktionen	79
2.3.1	Bisektionsverfahren	79
2.3.2	Modellfunktionsverfahren	81
2.3.2.1	Newton-Verfahren in 1D	81
2.3.2.2	Spezielle Einpunktverfahren	83
2.3.2.3	Mehrpunktverfahren	90
2.4	Effizienz	96
2.5	Newton-Verfahren für Gleichungssysteme	99
2.5.1	Die Newton-Iteration	99
2.6.1	Konvergenzanalyse des Newton-Verfahrens	103
2.6.2	Gedämpftes Newton-Verfahren	108
2.6.3	Quasi-Newton-Verfahren	111
2.7	Nichtlineare Ausgleichsrechnung	119
3	Numerische lineare Algebra	127
3.1	Grundbegriffe und -operationen	128
3.1.1	Operationen	129
3.1.2	Matrix-Speicherformate	131
3.1.3	Erinnerung: Normen	137
3.2	Numerische Lösung linearer Gleichungssysteme	139
3.2.1	Theorie und Kondition	140
3.2.2	Die Gaußelimination	146
3.2.3	Die LU -Zerlegung	150
3.2.4	Pivotsuche	157
3.2.5	Symmetrisch positiv definite Matrizen	168
3.2.6	Dünnbesetzte Gleichungssysteme	175
3.2.7	Die QR-Zerlegung	188
3.2.8	Modifikationstechniken	199
3.2.8.1	Rang-1-Modifikationen	199
3.2.8.2	Hinzufügen einer Spalte	207
3.2.8.3	Hinzufügen einer Zeile	211
3.3	Numerische Berechnung von Eigenwerten und Eigenvektoren	213
3.3.1	Theorie und Kondition	214
3.3.2	Transformationsmethoden	218
3.4.1	Potenzmethoden	225
3.4.2	Vorkonditionierte inverse Iteration	235
3.4.3	Krylov-Unterraumverfahren	243
3.4.4	Singulärwertzerlegungen	260
3.5	Numerik linearer Ausgleichsprobleme	272
3.5.1	Orthogonaltransformationsmethode	275

0.0
p. 2

0.0
p. 3

0.0
p. 4

3.5.2	Normalengleichungen	282
3.5.3	Totales Ausgleichsproblem	285
3.5.4	Ausgleichsrechnung mit linearen Nebenbedingungen	286
3.6	Krylov-Verfahren für lineare Gleichungssysteme	289
3.6.1	Das Verfahren der konjugierten Gradienten (CG)	290
3.6.1.1	Prinzip des CG-Verfahrens	291
3.6.1.2	Implementierung des CG-Verfahrens	294
3.6.1.3	Konvergenzgeschwindigkeit	303
3.6.2	Vorkonditionierung	309
3.6.3	Weitere Krylov-Unterraumverfahren	315
3.6.3.1	Residuenminimierende Verfahren	315
3.6.3.2	Verfahren mit kurzen Rekursionen	317
3.7	Spezielle Matrizen	320
3.7.1	Diskrete Fouriertransformationen	321
3.7.1.1	Schnelle Fouriertransformation	326
3.8.0.2	Sinustransformation	336
3.8.0.3	Kosinustransformation	346
3.8.1	Zirkulante Matrizen	348
3.9.1	Toeplitz-Matrizen	357
3.9.1.1	Toeplitz-Matrix-Arithmetik	360
3.9.1.2	Der Levinson-Algorithmus	362
4	Interpolation und Approximation	366
4.1	Polynomiale Techniken	367
4.1.1	Polynominterpolation	368
4.1.1.1	Theorie und Kondition	370
4.1.1.2	Algorithmen	377
4.1.2	Interpolationsfehlerabschätzungen	384
4.1.3	Tschebyscheff-Interpolation	391
4.1.4	Trigonometrische Interpolation	402
4.1.4.1	Trigonometrische Polynome	402
4.1.4.2	Trigonometrische Interpolation 1-periodischer Funktionen:	407
4.1.4.3	Trigonometrische Interpolation analytischer Funktionen	415
4.1.4.4	Trigonometrische Interpolation und Tschebyscheff-Interpolation	421
4.1.5	Approximation durch Polynome	425
4.1.5.1	Bestapproximation	425
4.1.5.2	Polynomiale Least-Squares Approximation	427
4.1.5.3	Tschebyscheff-Approximation	434
4.1.6	Clusteringapproximation	441
4.1.6.1	Separierte Kernapproximation	444
4.1.6.2	Clustertechnik	456
4.1.6.3	Multipol-Matrixmultiplikation	461
4.2	Dreitermrekursionen	472

0.0
p. 5

4.3	Stückweise Polynome	484
4.3.1	Stückweise Polynominterpolation	486
4.3.1.1	Stückweise lineare Interpolation	487
4.3.1.2	Stückweise polynomiale Interpolation von Funktionen	489
4.3.1.3	Kubische Hermite-Interpolation	494
4.3.2	Splines	499
4.3.2.1	Splineinterpolation	500
4.3.2.2	Formerhaltende Splineinterpolation	507
4.3.3	Bezier-Techniken	514
4.4	Numerische Quadratur	523
4.4.1	Polynomiale Quadraturformeln	524
4.4.2	Gauss-Quadratur	527
4.4.3	Zusammengesetzte Quadraturformeln	536
4.6.1	Adaptive Quadratur	544
4.6.2	Numerische Berechnung oszillatorischer Integrale	547
4.7	Multiskalenbasen	547
5	Numerik Gewöhnlicher Differentialgleichungen	564
5.1	Theorie gewöhnlicher Differentialgleichungen	566
5.2	Kondition von Anfangswertproblemen	572
5.3	Einschrittverfahren	575
5.3.1	Kollokation	576
5.3.2	Runge-Kutta-Verfahren	581
5.4	Konvergenz	588
5.4.1	Schrittweitensteuerung für Einschrittverfahren	595
5.6	Stabilität	606
5.6.1	Modellproblemanalyse	608
5.7.1	Steifheit	614
5.8.1	Einschrittverfahren für steife AWP	618
5.8.2	Implementierung impliziter Runge-Kutta-Einschrittverfahren	627
5.9	Differentiell-Algebraische Anfangswertprobleme	629
5.10	Strukturerhaltung	633
5.10.1	Nichtexpansivität	633
5.10.2	Quadratische erste Integrale	635
5.10.3	Symplektizität	638
5.10.4	Reversibilität	649
5.11	Splittingverfahren	658
5.12	Verfahren für oszillatorische Differentialgleichungen	664
	Verzeichnisse	679
	Stichwortverzeichnis	679
	Verzeichnis der Beispiele und Bemerkungen	696
	Verzeichnis der Definitionen und Konzepte	702
	Verzeichnis der MATLAB-CODE-Fragmente	705
	Symbolverzeichnis	708

0.0
p. 7

0.0
p. 6

0.0
p. 8

Reporting errors

Fehler in den Vorlesungsunterlagen bitte per [Wiki](#) melden!

<http://elbanet.ethz.ch/wikifarm/rhiptmair/index.php?n=Main.NCSECourse>

Bitte folgende Angaben in den Wiki eintragen:

- Abschnitt, in dem der Fehler gefunden wurde.
- Genauer Ort des Fehlers (z.B. „nach Gleichung (4)“)
- Kurzbeschreibung des Fehlers

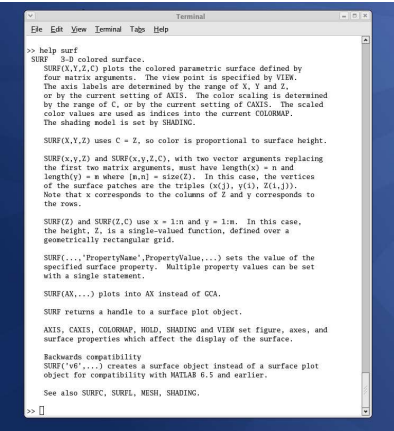
MATLAB

Beispielrechnungen und Übungen zur Vorlesung stützen sich auf die für numerische Berechnungen konzipierte Hochsprache **MATLAB** („Matrix Laboratory“) von der Firma MATHWORKS.

Matlab ist eine Interpretersprache und verfügt über folgende Möglichkeiten:

- Mächtige Funktionen für Vektor- und Matrixarithmetik und -manipulationen
- Kontrollstrukturen, wie `if ... else ... end`, `for/while ... end`
- Modularisierung durch Unterprogramme
- Bibliotheken für komplexe numerische Aufgaben
- Umfangreiche Visualisierungsfunktionen

Einarbeitung in MATLAB: Am besten durch „Learning by doing“ unter Zuhilfenahme der ausführlichen Online-Hilfsmenüs.



MATLAB-LINKS:

- Matlab Online Documentation from the Mathworks
- MATLAB guide
- MATLAB Primer

0.1 Danksagung

- 0.0 Wertvolle Beiträge zur Entstehung und Verbesserung der vorliegenden Präsentationsfolien haben
p. 9 geleistet:
- Herr Martin Knoller Stocker (Student RW): Erstellung des Index und der Verzeichnisse im SS05
 - Herr Christian Heitzmann (Student RW): Umfangreiche Fehlerlisten im WS 04/05
 - Frau Gisela Widmer (Doktorandin RW): Korrekturlesen als Assistentin der Vorlesung im WS 04/05, WS 05/06, WS 06/07
 - Herr Henning Avenhaus und Herr Christoph Scherrer (Studenten RW): Aufarbeitung der vorlesungsbegleitenden MATLAB-Skripten

0.1
p. 11

0.0
p. 10

0.1
p. 12

Computerarithmetik und Konsequenzen

Viele Beispiele aus diesem Kapitel sind den ersten Kapiteln von [16] und [28] entnommen. Zum weiterführenden Studium sind Kapitel 2 von [12] und das Buch [29] empfohlen.

1.1 Beispiele

[File: section-beispiele.tex, SVN: chapter-computerarithmetik-und-konsequenzen.tex 1035 2006-10-09 11:41:18Z kalai]

Rechner rechnen falsch !

Numerik \neq Analysis

Beispiel 1 (Berechnung der Eulerschen Zahl).

Bekannt aus Analysis I:

$$e = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n$$

MATLAB-CODE: Näherung für e

```

Approximation of e
for i=1:15
    n = 10^i; r = (1+1/n)^n;
    fprintf('10^%2d %20.15f ...
           %20.15f\n', i, r, r-exp(1));
end
    
```

n	Näherung e_n	Fehler $e_n - e$
10^1	2.593742460100002	-0.124539368359044
10^2	2.704813829421529	-0.013467999037517
10^3	2.716923932235520	-0.001357896223525
10^4	2.718145926824356	-0.000135901634689
10^5	2.718268237197528	-0.000013591261517
10^6	2.718280469156428	-0.000001359302618
10^7	2.718281693980372	-0.000000134478673
10^8	2.718281786395798	-0.000000042063248
10^9	2.718282030814509	0.000000202355464
10^{10}	2.718282053234788	0.000000224775742
10^{11}	2.718282053357110	0.000000224898065
10^{12}	2.718523496037238	0.000241667578192
10^{13}	2.716110034086901	-0.002171794372145
10^{14}	2.716110034087023	-0.002171794372023
10^{15}	3.035035206549262	0.316753378090216

Beispiel 2 („Quadratur des Kreises“).

Approximation durch reguläres n -Eck, $n \in \mathbb{N}$:

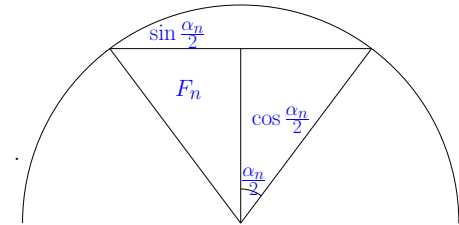
Fläche des n -Ecks:

$$A_n = n \cos \frac{\alpha_n}{2} \sin \frac{\alpha_n}{2} = \frac{n}{2} \sin \alpha_n = \frac{n}{2} \sin \left(\frac{2\pi}{n} \right).$$

REKURSION
Rekursion für A_n abgeleitet von

$$\sin \frac{\alpha_n}{2} = \sqrt{\frac{1 - \cos \alpha_n}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}},$$

Startwert: $A_6 = \frac{3}{2}\sqrt{3}.$



1.1
p. 13

MATLAB-CODE: Instabile Rekursion für Pi

```

s = sqrt(3)/2; A=3*s; n = 6; % initialization
z = [n A A-pi s]; % matrix for storing results
while (s > 1e-10) % terminate when s is small
    s = sqrt((1-sqrt(1-s*s))/2); % Recursion
    n=2*n; A = n/2*s;
    z = [z; n A A-pi s]; % store results
end
for i=1:length(z)
    fprintf('%10d %20.15f %20.15f %20.15f\n', ...
           z(i,1), z(i,2), z(i,3), z(i,4));
end
    
```

1.1
p. 15

1.1
p. 14

1.1
p. 16

n	A_n	$A_n - \pi$	$\sin \alpha_n$
6	2.598076211353316	-0.543516442236477	0.866025403784439
12	3.000000000000000	-0.141592653589794	0.500000000000000
24	3.105828541230250	-0.035764112359543	0.258819045102521
48	3.132628613281237	-0.008964040308556	0.130526192220052
96	3.139350203046872	-0.002242450542921	0.065403129230143
192	3.141031950890530	-0.000560702699263	0.032719082821776
384	3.141452472285344	-0.000140181304449	0.016361731626486
768	3.141557607911622	-0.000035045678171	0.008181139603937
1536	3.141583892148936	-0.000008761440857	0.004090604026236
3072	3.141590463236762	-0.000002190353031	0.002045306291170
6144	3.141592106043048	-0.000000547546745	0.001022653680353
12288	3.141592516588155	-0.000000137001638	0.000511326906997
24576	3.141592618640789	-0.000000034949004	0.000255663461803
49152	3.141592645321216	-0.000000008268577	0.000127831731987
98304	3.141592645321216	-0.000000008268577	0.000063915865994
196608	3.141592645321216	-0.000000008268577	0.000031957932997
393216	3.141592645321216	-0.000000008268577	0.000015978966498
786432	3.141593669849427	0.000001016259634	0.000007989485855
1572864	3.141592303811738	-0.000000349778055	0.000003994741190
3145728	3.141608696224804	0.000016042635011	0.000001997381017
6291456	3.141586839655041	-0.000005813934752	0.000000998683561
12582912	3.141674265021758	0.000081611431964	0.000000499355676
25165824	3.141674265021758	0.000081611431964	0.000000249677838
50331648	3.143072740170040	0.001480086580246	0.000000124894489
100663296	3.159806164941135	0.018213511351342	0.000000062779708
201326592	3.181980515339464	0.040387861749671	0.000000031610136
402653184	3.354101966249685	0.212509312659892	0.000000016660005
805306368	4.242640687119286	1.101048033529493	0.000000010536712
1610612736	6.000000000000000	2.858407346410207	0.000000007450581

Beispiel 3 (Auswertung der Exponentialfunktion).

Global konvergente Exponentialreihe :

$$\exp(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!} = 1 + x + \frac{x^2}{2} + \frac{x^3}{6} + \frac{x^4}{24} + \dots$$

```

MATLAB: Summation Exp.-Reihe
function y = expeval(x,tol)
% Evaluation of exponential series
s=1; term=1; k=1; % Initialization
while (abs(term)>tol*abs(s))
    term = term*x/k;
    s = s + term;
    k = k+1;
end
y = s;

```

Numerisches Experiment: Untersuchung des relativen Fehlers im Ergebnis für $\text{tol} = 10^{-8}$

Definition 1.1.1 (Relativer Fehler). $f(x) \in \mathbb{C} \hat{=}$ (theoretisch) exaktes Resultat einer Berechnung, $\tilde{f}(x) \in \mathbb{C} \hat{=}$ Näherung. Falls $f(x) \neq 0$ ➤ **relativer Fehler** (engl.: relative error) gegeben durch

$$\epsilon_r := \frac{|f(x) - \tilde{f}(x)|}{|f(x)|}$$

Warum betrachten wir den relative Fehler?

Relativer Fehler eines Resultats ↔ Anzahl der gültigen Stellen

x	Approximation $\widetilde{\exp}(x)$	$\exp(x)$	$\frac{ \exp(x) - \widetilde{\exp}(x) }{\exp(x)}$
-20	5.6218844674e-09	2.0611536224e-09	1.727542676201181
-18	1.5385415977e-08	1.5229979745e-08	0.010205938187564
-16	1.1254180496e-07	1.1253517472e-07	0.000058917020257
-14	8.3152907681e-07	8.3152871910e-07	0.00000430176956
-12	6.1442133148e-06	6.1442123533e-06	0.000000156480737
-10	4.5399929556e-05	4.5399929762e-05	0.00000004544414
-8	3.3546262817e-04	3.3546262790e-04	0.00000000788902
-6	2.4787521758e-03	2.4787521767e-03	0.00000000333306
-4	1.8315638879e-02	1.8315638889e-02	0.00000000530694
-2	1.3533528320e-01	1.3533528324e-01	0.00000000273603
0	1.0000000000e+00	1.0000000000e+00	0.000000000000000
2	7.3890560954e+00	7.3890560989e+00	0.00000000479969
4	5.4598149928e+01	5.4598150033e+01	0.000000001923058
6	4.0342879295e+02	4.0342879349e+02	0.000000001344248
8	2.9809579808e+03	2.9809579870e+03	0.000000002102584
10	2.2026465748e+04	2.2026465795e+04	0.000000002143800
12	1.6275479114e+05	1.6275479142e+05	0.000000001723845
14	1.2026042798e+06	1.2026042842e+06	0.000000003634135
16	8.8861105010e+06	8.8861105205e+06	0.000000002197990
18	6.5659968911e+07	6.5659969137e+07	0.000000003450972
20	4.8516519307e+08	4.8516519541e+08	0.000000004828738

1.1
p. 17

1.1
p. 19

1.1
p. 18

1.1
p. 20

Beispiel 4 (Summation der Harmonischen Reihe).

$$s := \sum_{k=1}^{100000} \frac{1}{k}$$

MATLAB-CODE: Summation vorwärts

```
s = 0;
for i = 1:1e6
    s = s+1/i;
end
```

← ? →

MATLAB-CODE: Summation rückwärts

```
s = 0;
for i = 1e6:-1:1
    s = s+1/i;
end
```

„+“ assoziativ !?

➤ $s = 14.39272672286498889$

➤ $s = 14.39272672286577226$

[„Exakte“ Auswertung durch MAPLE-Code: ➤ $s = 14.39272672286572363138112749 \dots$]

Beispiel 5 (Probleme bei Polynomauswertung).

p : Polynom, Grad 19

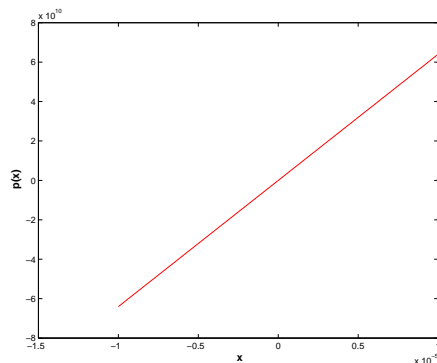
$$p(x) = x^{19} + \alpha_{18}x^{18} + \dots + \alpha_1x + \alpha_0$$

mit Nullstellen $0, 1, 2, \dots, 18$

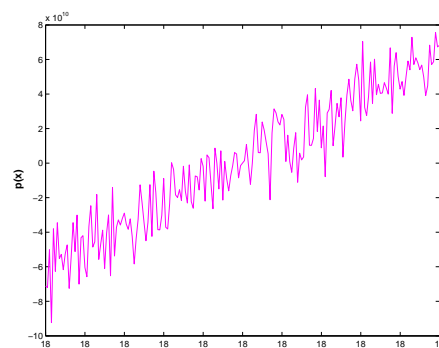
➤ Symmetrie $p(9-x) = p(9+x)$

MATLAB-CODE: Polynomauswertung

```
p = poly(0:18);
x = -1E-5:1E-7:1E-5;
y = polyval(p,x);
plot(x,y,'r-');
y = polyval(p,x+18);
plot(x+18,y,'m-');
```



Erwartet



„Numerisches Rauschen“

1.2 Zahldarstellung

[File: section-zahldarstellung.tex, SVN: section-zahldarstellung.tex 1371 2007-04-04 10:42:11Z hiptmair]

ZAHLDARSTELLUNG

Computer = endlicher Automat

rechnet mit endlich vielen Zahlen

Bemerkung 6 (Zahldarstellung durch Exponent und Mantisse: „wissenschaftliche Notation“).

Basis $B \in \mathbb{N} \setminus \{1\}$ fest: $\forall x \in \mathbb{R}: \exists_1 B^{-1} \leq d < 1, \exists_1 E \in \mathbb{Z}: x = d \cdot B^E$.
Mantisse Exponent

MSTELLEZZAHL

Beispiel 7 (m -stellige Dezimalzahlen).

3-stellige Dezimalzahlen $B = 10, 0.1 \leq d < 1$:

1.1 Gültig : $0.145 \cdot 10^5, -0.100 \cdot 10^7, 0.450 \cdot 10^{-20}$
p. 21 Ungültig : $1.345 \cdot 10^2, 0.001 \cdot 10^4$

Computer speichern Mantissen und Exponenten nur mit endlich vielen Stellen:

Definition 1.2.1 (Maschinenzahlen). Die Menge \mathbb{M} der **Maschinenzahlen** (engl. machine numbers) zur Basis $B \in \mathbb{N} \setminus \{1\}$, mit m -stelliger Mantisse d , $m \in \mathbb{N}$, und Exponentenbereich $\{e_{\min}, \dots, e_{\max}\}$, $e_{\min}, e_{\max} \in \mathbb{Z}$, ist gegeben durch

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

Niemals = 0!
Maschinenzahl $\in \mathbb{M}$: $x = \pm \underbrace{0.\square\square\square\dots\square}_{m \text{ Mantissenstellen}} \cdot B^{\underbrace{\square\dots\square}_{\text{Exponentenstellen}}}$

➤ Betragsgrösste Maschinenzahl : $x_{\max} = (1 - B^{-m}) \cdot B^{e_{\max}}$
Betragskleinste Maschinenzahl : $x_{\min} = B^{-1} \cdot B^{e_{\min}}$

Bemerkung 8. Endlichkeit von \mathbb{M} ➤ (Unvermeidliche) Eingabefehler (z.B. für $x = \pi$)

Beispiel 9 (IEEE Standard 754 für Gleitkommazahlen*). (engl. floating point numbers) ➔ [35]

1.2 Natürlich: $B = 2$ (Binärsystem)
p. 22

1.2
p. 23

1.2
p. 24

EINFACHE GENAUIGKEIT

Einfache Genauigkeit (single precision) : $m = 24^*$, $E \in \{-125, \dots, 128\}$ \rightarrow 4 Byte

DOBBELTE GENAUIGKEIT

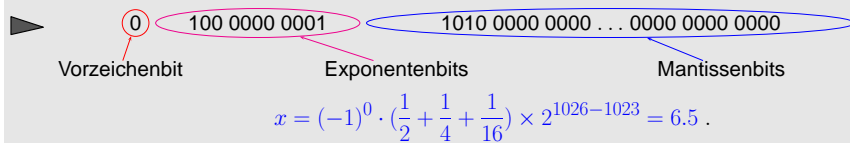
Doppelte Genauigkeit (double precision) : $m = 53^*$, $E \in \{-1021, \dots, 1024\}$ \rightarrow 8 Byte

*: Konsequent wird in dieser Vorlesung ein Dezimalpunkt verwendet.

*: Ein Vorzeichenbit ist inbegriffen.

Bemerkung 10 (Interne Binärdarstellung einer Gleitkommazahl doppelter Genauigkeit (MATLAB)).

```
>> format hex
>> 6.5
ans = 401a000000000000
```



Bemerkung 11 (Sonderfälle im IEEE Standard).

```
>> x = exp(1000), y = 3/x, z = x*sin(pi), w = x*log(1)
x = Inf
y = 0
z = Inf
w = NaN
```



$E = e_{\max}$, $M \neq 0 \hat{=}$ **NAN** = Not a number \rightarrow Ausnahme (engl. exception)
 $E = e_{\max}$, $M = 0 \hat{=}$ **Inf** = Infinity \rightarrow Überlauf (\rightarrow Abschnitt 11.5)
 $E = 0$ $\hat{=}$ Unnormalisierte Zahlen \rightarrow Unterlauf (\rightarrow Abschnitt 11.5)
 $E = 0$, $M = 0 \hat{=}$ Gleitkommazahl 0

Beispiel 12. Charakteristische Größen der IEEE Gleitpunktarithmetik (doppelte Genauigkeit)

Relevant für Rechnungen mit MATLAB

Parameter der Gleitkommaarithmetik in MATLAB

```
>> format hex; realmin, format long; realmin
ans = 0010000000000000
ans = 2.225073858507201e-308
```

```
>> format hex; realmax, format long; realmax
ans = 7fefffffffffffffff
ans = 1.797693134862316e+308
```

Konsequenzen der Gleitpunktarithmetik sind aktuelles Forschungsthema: [5], siehe auch die Monographie [29]

1.3 Gleitpunktarithmetik und Rundungsfehler

[File: section-gleitpunktarithmetik-und-rundungsfehler.tex, SVN: section-zahldarstellung.tex 1371 2007-04-04 10:42:11Z hiptmair]

GLEITPUNKTARITHMETIK

1.2
p. 25

1. Beobachtung: Maschinenzahlenmenge (\rightarrow Def. 11.2.1) \mathbb{M} ist endlich!

Beispiel 13. Abbruchkriterium von Iterationen unter Ausnutzung von Gleitpunktarithmetik:

WURZELITERATION

Wurzeliteration für $a > 0$:

$$x^{(k+1)} = \frac{1}{2} \left(x^{(k)} + \frac{a}{x^{(k)}} \right) \rightarrow \sqrt{a} \text{ for } k \rightarrow \infty.$$

Monotonie:

$$x^{(0)} > \sqrt{a} > \sqrt{a} < x^{(k)} < x^{(k-1)} \forall k \in \mathbb{N}.$$

MASCHINENUNABHÄNGIGES

Maschinenunabhängiges Abbruchkriterium

MATLAB-CODE: Wurzeliteration

```
function y=squareroot(a)
xo = 0.5*(1+a);
xn = (xo+a/xo)/2;
while (xn<xo)
    xo=xn; xn = (xo+a/xo)/2;
end
y = (xo+xn)/2;
```

2. Beobachtung: Maschinenzahlenmenge \mathbb{M} (\rightarrow Def. 11.2.1) nicht abgeschlossen unter den elementaren arithmetischen Operationen $+, -, \cdot, /$.

Unvermeidlich: **RUNDUNGSRUNDUNGSFEHLER**

Rundung:

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

(grössere Gleitkommazahl bei Nichteindeutigkeit, "Aufrunden")

Beispiel 14 (Eingabefehler und Rundungsfehler).

1.2
p. 26

1.3
p. 27

1.3
p. 28

MATLAB-CODE: Demonstration von Rundungsfehlern

```
>> format long;
>> a = 4/3; b = a-1; c = 3*b; e = 1-c
e = 2.220446049250313e-16
```

Notation: Realisierung von $\star \in \{+, -, \cdot, /\}$ in Gleitpunktarithmetik $\rightarrow \tilde{\star}$

Definition 1.3.1 (Maschinengenauigkeit). Die **Maschinengenauigkeit** (engl. machine precision) eines Maschinenzahlensystems nach Def. 1.2.1 ist gegeben durch

$$\text{eps} := \frac{1}{2}B^{-m+1}.$$

$$x_{\min} \leq x \leq x_{\max} \Rightarrow \text{Relativer Rundungsfehler} \frac{|\text{rd}(x) - x|}{|x|} \leq \text{eps}.$$

Beispiel 15 (Abfrage der Maschinengenauigkeit in MATLAB). (Prozessor Intel Pentium)

Abfrage der Maschinengenauigkeit in MATLAB

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

AXIOM **ELEMARITHOP**
„Axiom“ der Gleitpunktarithmetik: Für die elementaren arithmetischen Operationen $\star \in \{+, -, \cdot, /\}$ und elementare Funktionen $f \in \{\exp, \sin, \cos, \log, \dots\}$ gilt

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

mit $|\delta| < \text{eps}$, $\text{eps} = \text{Maschinengenauigkeit} \rightarrow \text{Def. 1.3.1}$

Relative Rundungsfehler elementarer Rechenschritte beschränkt durch Maschinengenauigkeit

Bemerkung 16 (Realisierung von $\tilde{+}, \tilde{-}, \tilde{\cdot}, \tilde{/}$).

$$\star \in \{+, -, \cdot, /\}: \quad x \tilde{\star} y := \text{rd}(x \star y) \quad (1.3.1)$$

Warum sollen uns diese winzigen Rundungsfehler kümmern ?

- Sie sollten (\rightarrow siehe Abschnitt 1.1):
- Akkumulation von Rundungsfehlern
 - Verstärkung von Rundungsfehlern

Bemerkung 17. Rundung \rightarrow Operationen $\tilde{+}, \tilde{-}: \mathbb{M} \times \mathbb{M} \mapsto \mathbb{M}$ nicht assoziativ ! (\rightarrow Beispiel 4)

Beispiel 18 (Nichtassoziativität der Maschinenaddition). 3-stellig Dezimalarithmetik, $\tilde{+}$ gemäss (1.3.1)

$$> 0.100 \cdot 10^4 \tilde{+} 1 = \text{rd}(0.1001 \cdot 10^4) = 0.100 \cdot 10^4 !$$

$$\blacktriangleright \quad 1000 \tilde{+} \underbrace{1 \tilde{+} \dots \tilde{+} 1}_{1000 \times} = 1000 \quad \leftrightarrow \quad \underbrace{1 \tilde{+} \dots \tilde{+} 1}_{1000 \times} \tilde{+} 1000 = 2000.$$

Bemerkung 19 (Numerische Nullabfrage).

Eine Konsequenz von Rundungsfehlern:



Abfragen wie `if (x == 0)` oft sinnlos/gefährlich, falls x Rechenergebnis

\blacktriangleright Besser: Abfrage `if (abs(x) < eps*s)` ...,
 $s = \text{positive Zahl, relativ zu der } x \text{ „klein“ sein soll.}$

1.3

p. 29

Übung 1.4. Gauge the relative error in the function evaluation $x \cdot \log(x)$ (can be plotted in MATLAB with `ezplot(@(x) x.*log(x), [0,5])`)

1.4

p. 31

1.5 Über- und Unterlauf

[File: section-ueber-und-unterlauf.tex, SVN: section-gleitpunktarithmetik-und-rundungsfehler.tex 1081 2006-10-27 14:13:07Z hiptmair]

ÜBERLAUF

Überlauf (engl. overflow) $\triangleq |\text{Resultat einer elementaren Operation}| > \max\{\mathbb{M}\}$

\triangleq IEEE Standard \Rightarrow Inf

UNTERLAUF

Unterlauf (engl. underflow) $\triangleq 0 < |\text{Resultat einer elementaren Operation}| < \min\{|\mathbb{M} \setminus \{0\}|\}$

\triangleq IEEE Standard \Rightarrow unnormalisierte Zahlen (!)

Bemerkung 20. Axiom der Gleitpunktarithmetik gilt nicht bei Unterlauf: MATLAB Beispiel

MATLAB-CODE: Unterlauf

```
>> format long; res=pi*realmin/123456789101112
res = 5.681754927174335e-322
>> res=res*123456789101112/realmin
res = 3.15248510554597
```

1.3

p. 30

1.5

p. 32

Vermeide Über- und Unterlauf wann immer möglich !

Beispiel 21 (Vermeidung von Überlauf).

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

Naive Auswertung: Überlauf für $x \gg 1, y \gg 1$

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

➤ Überlauf vermieden

1.6 Kondition

[File: section-kondition.tex, SVN: section-kondition.tex 1037 2006-10-09 13:23:05Z kalai]

KONDITION

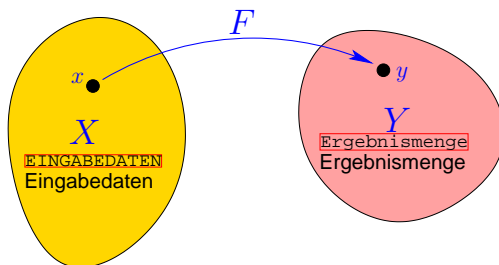
Die Kondition ist eine intrinsische Eigenschaft eines Problems

Hier:

KONDPBLEM
Problem = Mathematisch wohldefinierte Vorschrift, wie aus Eingabedaten (Input) ein Ergebnis (Output) erhalten werden kann.

X : Datenmenge
 Y : Ergebnismenge
 F : „Problemfunktion“ $X \mapsto Y$
 d_X : Metrik auf X
 d_Y : Metrik auf Y

Die Metriken (Abstandsmasse) auf X und Y sind „sinnvoll“ zu wählen.



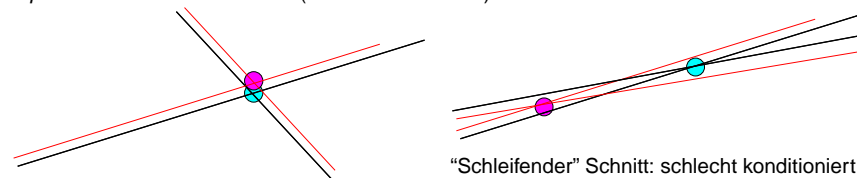
Definition 1.6.1 (Gut konditioniertes Problem). Ein Problem $F : X \mapsto Y$ heisst **gut konditioniert** zum Input $x \in X$ (engl. well conditioned), wenn „kleine Störungen in den Eingabedaten ebenso kleine Störungen im Ergebnis hervorrufen“, das heisst

$$\exists \kappa_{\text{abs}} \approx 1: d_Y(F(x), F(y)) \leq \kappa_{\text{abs}} d_X(x, y) \quad \forall y \in \mathcal{U}(x)$$

für eine Umgebung $\mathcal{U}(x)$ von x in X .

κ_{abs} = **absolute Konditionszahl**, misst Verstärkung von Fehlern in den Eingabedaten.

Beispiel 22. Schnitt von Geraden (mit Abstandsmetrik)



DIFKONDANAL

Differentielle Konditionsanalyse: betrachte Spezialfall: $X, Y \subset \mathbb{R}$, Betragsmetrik

Falls f differenzierbar: $\delta y := f(x + \delta x) - f(x) = f'(\xi)\delta x$, mit $\xi \in]x, x + \delta x[$

Falls f stetig differenzierbar: $\delta y = f'(x)\delta x + O(\delta x^2)$

Falls $x \neq 0, y := f(x) \neq 0, f$ stetig differenzierbar

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

1.6
p. 33

➤ $\kappa_{\text{rel}} := \left| f'(x) \frac{x}{f(x)} \right|$ = **relative Konditionszahl**, misst Verstärkung relativer Fehler

Beispiel 23 (Kondition von Nullstellen).

$f : I \subset \mathbb{R} \mapsto \mathbb{R}$ 2-fach stetig differenzierbar mit Nullstelle (engl. zero) $s \in I: f(s) = 0$.

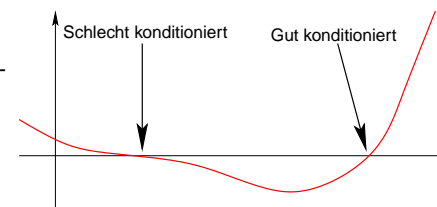
Additive Störung: $\tilde{f} := f + \epsilon \Rightarrow s = s(\epsilon), \tilde{f}(s(\epsilon)) = 0$

Implizites Differenzieren von $f(s(\epsilon)) = -\epsilon$

$$f'(s(\epsilon)) \frac{ds}{d\epsilon}(\epsilon) = -1 \Rightarrow \frac{ds}{d\epsilon}|_{\epsilon=0} = -\frac{1}{f'(s)}.$$

Je „flacher“ die Funktion desto schlechter die Kondition der Nullstelle bzgl. additiver Störungen

(betrifft **nicht** Funktionsauswertung !)



Kondition elementarer Operationen:

1.6
p. 34

1.6
p. 35

1.6
p. 36

- Multiplikation $f(x) = a \cdot x, a \in \mathbb{R}$

$$\kappa_{\text{rel}} = \left| f'(x) \frac{x}{f(x)} \right| = 1 !$$

➤ Keine Verstärkung relativer Fehler durch Multiplikation

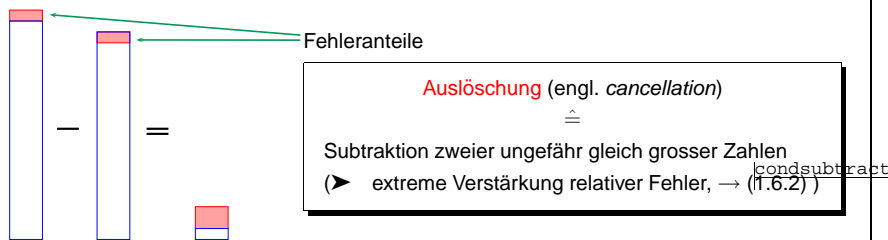
- Subtraktion $f(x) = x - a, a \in \mathbb{R}$. Relative Konditionszahl für $x \notin \{0, a\}$

$$\kappa_{\text{rel}} = \left| f'(x) \frac{x}{f(x)} \right| = \left| \frac{x}{x-a} \right| \quad \Rightarrow \quad \kappa_{\text{rel}} \gg 1, \text{ falls } x \approx a$$

1.7 Auslöschung

[File: section-ausloeschung.tex, SVN: section-ausloeschung.tex 1073 2006-10-24 19:22:33Z hiptmair]

AUSLOESCHUNG



Subtraktionen mit Auslöschung geschehen selbst rundungsfehlerfrei !

Beispiel 24 (Auslöschung). Rechnung mit 2-stelliger Dezimalarithmetik:

$$a = 1.3, \quad b = 1.2$$

$$a^2 - b^2 = (a + b)(a - b)$$

ALGA

Algorithmus A

$$x := \tilde{a} \cdot a = 1.7 \text{ (gerundet)}$$

$$y := \tilde{b} \cdot b = 1.4 \text{ (gerundet)}$$

$$x - y = 0.30 \text{ (exakt)}$$

Algorithmus B

$$x := \tilde{a} + \tilde{b} = 2.5 \text{ (exakt)}$$

$$y := \tilde{a} - \tilde{b} = 0.1 \text{ (exakt)}$$

$$x \cdot y = 0.25 \text{ (exakt)}$$

➤ Subtraktion im letzten Schritt ➤ grosse Verstärkung der (relativen) Rundungsfehler aus den ersten beiden Schritten

„Auslöschungsverdächtige“ Subtraktionen möglichst zu Beginn eines Algorithmus !
(Allgemein: Schlecht konditionierte Elementarschritte möglichst früh !)

Formale Analyse für $a, b, \in \mathbb{R}$: $f = a^2 - b^2$

Mit $\delta_1, \delta_2, \delta_3 \in \mathbb{R}, |\delta_i| < \text{eps}$ (\rightarrow Axiom der Gleitpunktarithmetik)

(1.6.2) condsub + Vernachlässigung von Termen höherer Ordnung in eps ($\hat{=}$ Linearisierung)

Algorithmus A:

$$x = a^2(1 + \delta_1), \quad y = b^2(1 + \delta_2)$$

$$\tilde{f} = (a^2(1 + \delta_1) - b^2(1 + \delta_2))(1 + \delta_3) = f + a^2\delta_1 + b^2\delta_2 + (a^2 - b^2)\delta_3 + O(\text{eps}^2)$$

$$\frac{|\tilde{f} - f|}{|f|} \leq \text{eps} \frac{a^2 + b^2 + |a^2 - b^2|}{|a^2 - b^2|} + O(\text{eps}^2) = \text{eps} \left(1 + \frac{|a^2 + b^2|}{|a^2 - b^2|} \right) + O(\text{eps}^2). \quad (1.7.1) \quad \text{ALGA}$$

Vernachlässigbar klein

1.7

p. 37

LANDAUO

Notation („Landau-O“): $f(h) = O(h^\alpha) \Leftrightarrow \exists C > 0, h_0 > 0: |f(h)| \leq C|h|^\alpha$ für alle $|h| < h_0$.

Algorithmus B:

$$x = (a + b)(1 + \delta_1), \quad y = (a - b)(1 + \delta_2)$$

$$\tilde{f} = (a + b)(a - b)(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) = f + (a^2 - b^2)(\delta_1 + \delta_2 + \delta_3) + O(\text{eps}^2)$$

$$\frac{|\tilde{f} - f|}{|f|} \leq |\delta_1 + \delta_2 + \delta_3| + O(\text{eps}^2) \leq 3\text{eps} + O(\text{eps}^2). \quad (1.7.2) \quad \text{ALGB}$$

➤ Relativer Ergebnisfehler immer $\approx \text{eps}$!

◇

Beispiel 25 (Fortsetzung von „Berechnung der Exponentialfunktion“, Beispiel 3).

lex:exp

1.7

p. 38

1.7

p. 40

Für $x > 0$:

$$\exp(-x) = \sum_{k=0}^{\infty} (-1)^k \frac{x^k}{k!} = 1 - x + \frac{x^2}{2} - \frac{x^3}{6} + \frac{x^4}{24} - \dots$$

Alternierende Reihe

► Ständige Auslöschung bei Summation

Partialsummen für $x = 20 \rightarrow$

partsum

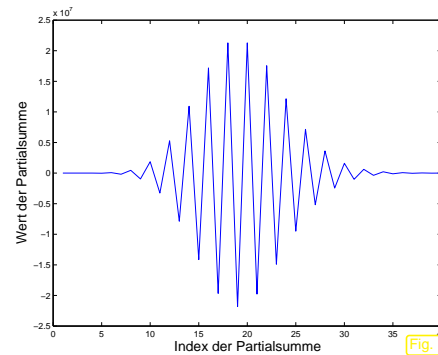


Fig. 1

Beispiel 26 (Fortsetzung „Quadratur des Kreises“, Beispiel 2).

ex:circarea

Rekursionsformel aus Beispiel 2:

ex:circarea

$$\sin \frac{\alpha_n}{2} = \sqrt{\frac{1 - \cos \alpha_n}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}}$$

Für $\alpha_n \ll 1$: $\sqrt{1 - \sin^2 \alpha_n} \approx 1$

Auslöschung!

Umformung in **auslöschungsfreie** Formel:

$$\begin{aligned} \sin \frac{\alpha_n}{2} &= \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2} \cdot \frac{1 + \sqrt{1 - \sin^2 \alpha_n}}{1 + \sqrt{1 - \sin^2 \alpha_n}}} \\ &= \sqrt{\frac{1 - (1 - \sin^2 \alpha_n)}{2(1 + \sqrt{1 - \sin^2 \alpha_n})}} = \frac{\sin \alpha_n}{\sqrt{2(1 + \sqrt{1 - \sin^2 \alpha_n})}} \end{aligned}$$

MATLAB Code:

(mit maschinenunabhängigem Abbruchkriterium)

```
MATLAB-CODE: Stabile Rekursion für Pi
s = sqrt(3)/2; Ao = 0; An=3*s; n = 6; % initialization
z = [n A An-pi s]; % matrix for storing results
while (An < Ao) % terminate when machine accuracy is reached
    s = s/sqrt(2*(1+sqrt((1+s)*(1-s)))); % Stabilized recursion
    n=2*n; Ao = An; An = n/2*s;
```

```
z = [z; n An An-pi s]; % store results
end
```

1.7
p. 41

n	A_n	$A_n - \pi$	$\sin \alpha_n$
6	2.598076211353316	-0.543516442236477	0.866025403784439
12	3.000000000000000	-0.141592653589793	0.500000000000000
24	3.105828541230249	-0.035764112359544	0.258819045102521
48	3.132628613281238	-0.008964040308555	0.130526192220052
96	3.139350203046867	-0.002242450542926	0.065403129230143
192	3.141031950890509	-0.000560702699284	0.032719082821776
384	3.141452472285462	-0.000140181304332	0.016361731626487
768	3.141557607911857	-0.000035045677936	0.008181139603937
1536	3.141583892148318	-0.000008761441475	0.004090604026235
3072	3.141590463228050	-0.000002190361744	0.002045306291164
6144	3.141592105999271	-0.000000547590522	0.001022653680338
12288	3.141592516692156	-0.000000136897637	0.000511326907014
24576	3.141592619365383	-0.000000034224410	0.000255663461862
49152	3.141592645033690	-0.000000008556103	0.000127831731976
98304	3.141592651450766	-0.000000002139027	0.000063915866118
196608	3.141592653055036	-0.000000000534757	0.000031957933076
393216	3.141592653456104	-0.000000000133690	0.000015978966540
786432	3.141592653556371	-0.000000000033422	0.000007989483270
1572864	3.141592653581438	-0.000000000008355	0.000003994741635
3145728	3.141592653587705	-0.000000000002089	0.000001997370818
6291456	3.141592653589271	-0.000000000000522	0.000000998685409
12582912	3.141592653589663	-0.000000000000130	0.000000499342704
25165824	3.141592653589761	-0.000000000000032	0.000000249671352
50331648	3.141592653589786	-0.000000000000008	0.000000124835676
100663296	3.141592653589791	-0.000000000000002	0.000000062417838
201326592	3.141592653589794	0.000000000000000	0.000000031208919
402653184	3.141592653589794	0.000000000000001	0.000000015604460
805306368	3.141592653589794	0.000000000000001	0.000000007802230
1610612736	3.141592653589794	0.000000000000001	0.000000003901115

1.7
p. 42

1.7
p. 43

1.7
p. 44

NUMDIFF
 Beispiel 27 (Numerische Differentiation).

Approximation der Ableitung durch Differenzenquotienten: $f'(x) \approx \frac{f(x+h) - f(x)}{h}$.

Je kleiner h desto besser, oder ?

MATLAB-CODE: Numerische Diff. von exp(x)

```
h = 0.1; x = 0.0;
for i = 1:16
    df = (exp(x+h)-exp(x))/h;
    fprintf('%d %16.14f\n',i,df-1);
    h = h*0.1;
end
```

Offensichtliche Auslöschung

$$\left. \begin{aligned} f'(x) - \frac{f(x+h) - f(x)}{h} &\rightarrow 0 \\ \text{Rundungsfehlereinfluss} &\rightarrow \infty \end{aligned} \right\} \text{ für } h \rightarrow 0.$$

Analyse für $f(x) = \exp(x)$:

$$\begin{aligned} df &= \frac{e^{x+h}(1+\delta_1) - e^x(1+\delta_2)}{h} \\ &= e^x \left(\frac{e^h - 1}{h} + \frac{\delta_1 e^h - \delta_2}{h} \right) \\ \Rightarrow |df| &\leq e^x \left(e^h + \text{eps} \frac{1+e^h}{h} \right) \end{aligned}$$

Korrekturfaktoren berücksichtigen Rundungsfehlereinfluss:
 (→ „Axiom der Gleitpunktarithmetik“, siehe 1.3)
 $|\delta_1|, |\delta_2| \leq \text{eps}$, eps = Maschinengenauigkeit
 $1 + O(h)$ $O(h^{-1})$ für $h \rightarrow 0$

► Relativer Fehler: $\left| \frac{e^x - df}{e^x} \right| \approx h + \frac{2\text{eps}}{h} \rightarrow \min \text{ für } h = \sqrt{2\text{eps}}.$

IEEE-Arithmetik doppelter Genauigkeit: $\sqrt{2\text{eps}} = 2.107342425544702 \cdot 10^{-8}$

$\log_{10}(h)$	Relativer Fehler
-1	0.05170918075648
-2	0.00501670841679
-3	0.00050016670838
-4	0.00005000166714
-5	0.00000500000696
-6	0.00000049996218
-7	0.00000004943368
-8	-0.00000000607747
-9	0.00000008274037
-10	0.00000008274037
-11	0.00000008274037
-12	0.00008890058234
-13	-0.00079927783736
-14	-0.00079927783736
-15	0.11022302462516
-16	-1.00000000000000

1.8 Stabilität

[File: section-stabilitaet.tex, SVN: section-stabilitaet.tex 1073 2006-10-24 19:22:33Z hiptmair]

STABILITÄT

Stabilität ist eine Eigenschaft eines Algorithmus für ein Problem

NUMALGO

Numerischer
Algorithmus

=

Spezifizierte Folge elementarer Operationen
(→ C(++)- oder FORTRAN-Programm)

Im folgenden: X, Y = normierte Vektorräume, z.B. $X = \mathbb{R}^n, Y = \mathbb{R}^m$

Definition 1.8.1 (Stabiler Algorithmus). Ein Algorithmus \tilde{F} zu einem Problem $F : X \mapsto Y$ heisst **numerisch stabil** (engl. *stable*), wenn sich für alle $x \in X$ sein (durch Rundungsfehlereinfluss gestörtes Ergebnis) $\tilde{F}(x)$ als exaktes Ergebnis zu leicht gestörten Eingabedaten interpretieren lässt, d.h.

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

1.7
p. 45

1.8
p. 47

Je besser die Kondition eines Problems, desto kritischer die Stabilität von Algorithmen !

Bemerkung: Def. 1.8.1 definiert Stabilität im Sinne der **Rückwärtsfehleranalyse** (engl. *backward error analysis*)

Theorem 1.8.2 (Stabilität elementarer Operationen).

Unter dem „Axiom der Gleitpunktarithmetik“ sind die elementaren Operationen $\tilde{+}, \tilde{-}, \tilde{\cdot}, \tilde{/}$ numerisch stabil.

Beweis. (für +)

$$x \tilde{+} y = (x + y)(1 + \delta) = \underbrace{x(1 + \delta)}_{=\tilde{x}} + \underbrace{y(1 + \delta)}_{=\tilde{y}}.$$

Beispiel 28. Stabilität von Algorithmus A aus Bsp. 24:

$$\tilde{f} = \underbrace{a^2(1 + \delta_1)(1 + \delta_3)}_{=\tilde{a}^2} - \underbrace{b^2(1 + \delta_2)(1 + \delta_3)}_{=\tilde{b}^2} \text{ mit } \begin{aligned} \tilde{a} &= a\sqrt{1 + \delta_a} \\ \tilde{b} &= b\sqrt{1 + \delta_b} \end{aligned}, |\delta_a|, |\delta_b| < 3\text{eps}.$$

► (Schlechter) Algorithmus ist trotzdem numerisch stabil ! Warum ?

1.8
p. 46

1.8
p. 48

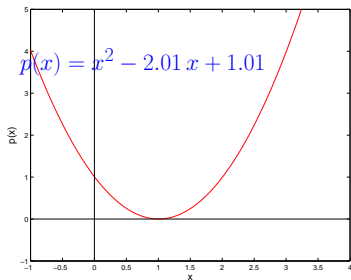
Betrachte (relative) Kondition von $(a, b) \mapsto f = a^2 - b^2$, z.B.

$$\frac{|\partial_a f(a, b)|}{|f(a, b)|} = \frac{2|a|^2}{|a^2 - b^2|} \gg 1 \quad \text{falls } a \approx b.$$

Algorithmus A, (1.7.1): grosser Rundungsfehlereinfluss nur im Fall schlechter Kondition

Beispiel 29 (Reelle Nullstellen eines quadratischen Polynoms).

$$p(x) = x^2 + px + q, \quad p, q \in \mathbb{R} \Rightarrow \text{Nullstellen } x_{\pm} = -\frac{p}{2} \pm \sqrt{(p/2)^2 - q}, \quad \text{falls } p^2 > 4q.$$



Fall: $p^2 \approx 4q$
 HILFE! Auslöschung unter der Wurzel
 Don't panic: **NULLSTELLENBEST**
 Nullstellenbestimmung schlecht konditioniert (*)
 !
 Formel auch mit Auslöschung stabil
 (*) Absolute Konditionszahl bzgl. q : $\kappa_{\text{abs}} = \frac{1}{\sqrt{p^2 - 4q}}$

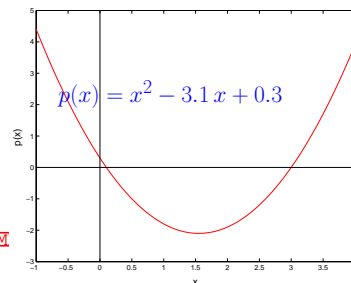
Fall: $p^2 \gg 4q$
 $[x_{\pm} \text{ gut konditioniert}] \Rightarrow$ Auslöschung gefährlich

$$\left| \frac{p}{2} \right| \approx \sqrt{(p/2)^2 - q}$$

Auslöschung bei x_- , falls $p < 0$,
 x_+ , falls $p > 0$.

Benutze $x_- x_+ = q$: **numerisch stabile Nullstellenformel**

$$\begin{cases} x_- = -\frac{p}{2} - \sqrt{(p/2)^2 - q}, & x_+ = \frac{q}{x_-}, \quad \text{falls } p > 0, \\ x_+ = -\frac{p}{2} + \sqrt{(p/2)^2 - q}, & x_- = \frac{q}{x_+}, \quad \text{falls } p < 0. \end{cases}$$



MATLAB-CODE: Nullstellen einer quadr. Fkt.

```
function [xm, xp] = quadzero(p, q)
D = 0.25*p^2 - q;
if (D < 0), error('No zeros'); end;
if (p > 0), xm = -0.5*p - sqrt(D); xp = q/xm;
else xp = -0.5*p + sqrt(D); xm = q/xp;
end
```

Stabile MATLAB-Implementierung

Resultate: Instabile Implementierung ↓ für $p = -1, q = a(1-a)$

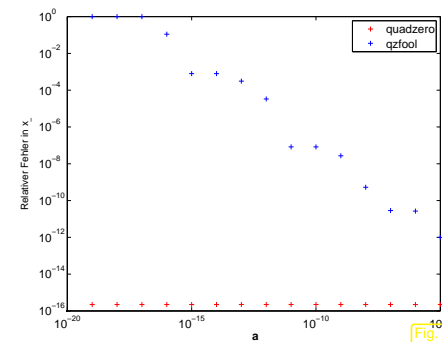
MATLAB-CODE: Nullstellen einer quadr. Fkt.

```
function [xm, xp] = qzfool(p, q)
D = 0.25*p^2 - q;
if (D < 0), error('No zeros'); end;
xp = -0.5*p + sqrt(D);
xm = -0.5*p - sqrt(D);
```

Instabile MATLAB-Implementierung

a	Relativer Fehler in x_-
10^{-5}	0.0000000000009995
10^{-6}	0.0000000000267561
10^{-7}	0.0000000000287549
10^{-8}	0.0000000005263567
10^{-9}	0.0000000272292189
10^{-10}	0.000000827403699
10^{-11}	0.000000827403699
10^{-12}	0.0000333894311084
10^{-13}	0.0003109451872646
10^{-14}	0.0007992778373605
10^{-15}	0.0007992778373608
10^{-16}	0.1102230246251547
10^{-17}	1.0000000000000000

1.8
p. 49



Relativer Fehler der Nullstellenbestimmung von

$$p(x) = x^2 - x + a(1-a) = (x-a)(x-1+a)$$

mittels `qzfool(-1.a*(1-a))`;
`quadzero(-1.a*(1-a))`;

(Relative Fehler $< \text{eps}$ ersetzt durch eps !)

1.8
p. 50

1.8
p. 51

1.8
p. 52

[File: chapter-nichtlineare-gleichungen.tex, SVN: chapter-nichtlineare-gleichungen.tex 1054 2006-10-18 11:43:04Z kalai]

Nichtlineare Gleichungen

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

Kann bedeuten: Wir verfügen über eine **Prozedur** function $y=F(x)$ zur Auswertung von F

Gesucht: Lösung der **NICHTLINGLEICHUNG** $F(\mathbf{x}) = 0$

Beachte: "Gleiche Anzahl von Gleichungen (engl. *equations*) und Unbekannten (engl. *unknowns*)"

Existenz und Eindeutigkeit von Lösungen ist für konkretes F jeweils zu klären !

2.1 Iterationsverfahren

[File: section-iterationsverfahren.tex, SVN: chapter-nichtlineare-gleichungen.tex 1054 2006-10-18 11:43:04Z kalai]

Iterationsverfahren zur (approximativen)
Lösung der nichtlinearen Gleichung $F(\mathbf{x}) = 0$ = Algorithmus, der eine Folge $(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$
von **Näherungslösungen** erzeugt.

- Iterierte $\mathbf{x}^{(k)}$ abhängig von F und (von einigen) $\mathbf{x}^{(n)}$, $n < k$, z.B.

$$\mathbf{x}^{(k)} = \underbrace{\mathbf{x}^{(k)}(F, \mathbf{x}^{(k-1)}, \dots, \mathbf{x}^{(k-m)})}_{\text{Iterationsvorschrift für } m\text{-Punkt-Verfahren}} \quad (2.1.1) \quad \text{eq:itmeth}$$

- $\mathbf{x}^{(0)}$ = **INITGUESS** (engl. *initial guess*)

Wir kennen bereits Iterationsverfahren: → Beispiele ex: ex: caracrit 2 & 13.

Definition 2.1.1 (Konvergenz von Iterationsverfahren).
Iterationsverfahren **konvergent** $\Leftrightarrow \mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ und $F(\mathbf{x}^*) = 0$.

Definition 2.1.2 (Konsistenz von Iterationsverfahren).

Iterationsverfahren **konsistent** mit $F(\mathbf{x}) = 0 \Leftrightarrow \{\mathbf{x}^{(k)} \rightarrow \mathbf{x}^* \Rightarrow F(\mathbf{x}^*) = 0\}$

Terminologie:

Fehler der Iterierten $\mathbf{x}^{(k)}$: $\mathbf{e}^{(k)} := \mathbf{x}^{(k)} - \mathbf{x}^*$

Definition 2.1.3 (Lokale und globale Konvergenz).

Ein Iterationsverfahren **konvergiert lokal** (engl. converges locally) gegen $\mathbf{x}^* \in \mathbb{R}^n$, falls eine Umgebung $U \subset D$ von \mathbf{x}^* existiert, so dass

$$\mathbf{x}^{(0)} \in U \Rightarrow \mathbf{x}^{(k)} \text{ wohldefiniert} \wedge \lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$$

für die erzeugten Iterationsfolgen $(\mathbf{x}^{(k)})_{k \in \mathbb{N}}$ gilt.

Falls $U = D$, so heisst die Iteration **global konvergent**.

Ziel: Finde Iterationsverfahren, die (lokal) gegen (eine) Lösung von $F(\mathbf{x}) = 0$ konvergieren.

2.1
p. 53

“Konvergenzgeschwindigkeit” ?

► Nötig: Norm $\| \cdot \|$ auf \mathbb{R}^n (\rightarrow Grundvorlesungen)
(Allgemeiner: Metrik auf \mathbb{R}^n)
Fest gewählte Norm auf \mathbb{R}^n wird im Folgenden vorausgesetzt.

2.1
p. 55

Definition 2.1.4 (Norm).

$X = \mathbb{K}$ -Vektorraum, $\mathbb{K} = \mathbb{C}, \mathbb{R}$. Eine Abbildung $\|\cdot\| : X \rightarrow \mathbb{R}_0^+$ ist eine **Norm** auf X , falls

- (i) $\forall \mathbf{x} \in X: \mathbf{x} \neq 0 \Leftrightarrow \|\mathbf{x}\| > 0$ (Definitheit),
- (ii) $\|\lambda \mathbf{x}\| = |\lambda| \|\mathbf{x}\| \quad \forall \mathbf{x} \in X, \lambda \in \mathbb{K}$ (Homogenität),
- (iii) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in X$ (Dreiecksungleichung).

Beispiele: (für Vektor $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{K}^n$)

Bezeichnung	Definition	MATLAB-Funktion
Euklidische Norm	$\ \mathbf{x}\ _2 := \sqrt{ x_1 ^2 + \dots + x_n ^2}$	norm(x)
1-Norm	$\ \mathbf{x}\ _1 := x_1 + \dots + x_n $	norm(x, 1)
Maximumnorm	$\ \mathbf{x}\ _\infty := \max\{ x_1 , \dots, x_n \}$	norm(x, inf)

Wahl der Norm "nicht so wichtig" im Endlichdimensionalen:

2.1
p. 54

2.1
p. 56

Theorem 2.1.5 (Normäquivalenz im Endlichdimensionalen).

Beliebige zwei Normen $\|\cdot\|, \|\cdot\|_0$ auf \mathbb{K}^n sind äquivalent

$$\exists C > 0: C^{-1} \|\mathbf{x}\| \leq \|\mathbf{x}\|_0 \leq C \|\mathbf{x}\| \quad \forall \mathbf{x} \in \mathbb{K}^n.$$

thm:normeq

Definition 2.1.6 (Lineare Konvergenz).

Eine Folge $\mathbf{x}^{(k)}, k = 0, 1, 2, \dots$, in \mathbb{R}^n konvergiert linear gegen $\mathbf{x}^* \in \mathbb{K}^n$, falls

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

def:cvglin

Terminologie: Kleinste obere Schranke für $L \rightarrow$ Konvergenzrate (engl. rate of convergence)

Erinnerung an zentralen Satz der Analysis:

Theorem 2.1.7 (Banachscher Fixpunktsatz).

Falls $D \subset \mathbb{K}^n$ abgeschlossen und für $\Phi: D \mapsto D$ gilt

$$\exists L < 1: \|\Phi(\mathbf{x}) - \Phi(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in D,$$

dann existiert genau ein Fixpunkt $\mathbf{x}^* \in D$ mit $\Phi(\mathbf{x}^*) = \mathbf{x}^*$, gegen den die Folge $\mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)})$ für beliebiges $\mathbf{x}^{(0)} \in D$ konvergiert

2.1
p. 57

Proof. Beweis verwendet 1-Punkt-Iteration $\mathbf{x}^{(k)} = \Phi(\mathbf{x}^{(k-1)}), \mathbf{x}^{(0)} \in D$:

$$\begin{aligned} \|\mathbf{x}^{(k+N)} - \mathbf{x}^{(k)}\| &\leq \sum_{j=k}^{k+N-1} \|\mathbf{x}^{(j+1)} - \mathbf{x}^{(j)}\| \leq \sum_{j=k}^{k+N-1} L^j \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \\ &\leq \frac{L^k}{1-L} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \xrightarrow{k \rightarrow \infty} 0. \end{aligned}$$

$(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0}$ Cauchy-Folge \rightarrow konvergent $\mathbf{x}^{(k)} \xrightarrow{k \rightarrow \infty} \mathbf{x}^*$.

Stetigkeit von $\Phi \rightarrow \Phi(\mathbf{x}^*) = \mathbf{x}^*$. Eindeutigkeit des Fixpunktes ist klar

□

Beispiel 30 (Linear Konvergente Iteration).

Iteration ($n = 1$):

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

```

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

```

2.1
p. 58

$x^{(15)}$ als Ersatz für \mathbf{x}^* bei der Berechnung der Konvergenzrate:

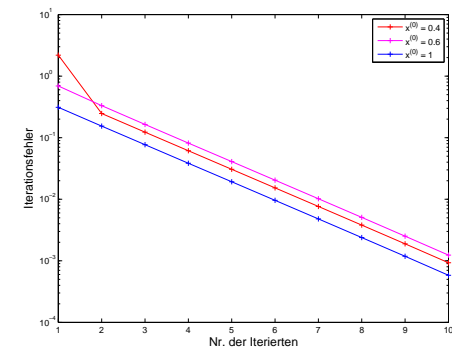
k	$x^{(0)} = 0.4$		$x^{(0)} = 0.6$		$x^{(0)} = 1$	
	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$	$x^{(k)}$	$\frac{ x^{(k)} - x^{(15)} }{ x^{(k-1)} - x^{(15)} }$
2	3.3887	0.1128	3.4727	0.4791	2.9873	0.4959
3	3.2645	0.4974	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

Konvergenzrate ≈ 0.5

Lineare Konvergenz \rightarrow Def. 2.1.6



Fehlerkurven = Geraden im lin-log Plot



2.1
p. 59

Abkürzung für Fehlernorm im k -ten Schritt $\epsilon_k := \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$. Im Falle linearer Konvergenz (\rightarrow Def. 2.1.6) nimmt man an ($0 < L < 1$)

$$\epsilon_{k+1} \approx L \epsilon_k \Rightarrow \log \epsilon_{k+1} \approx \log L + \log \epsilon_k \Rightarrow \log \epsilon_{k+1} \approx k \log L + \log \epsilon_0.$$

$\rightarrow \log L < 0$ gibt die Steigung der Geraden im lin-log-Plot.

2.1
p. 60

Im Fall von Konvergenz von p -ter Ordnung ($p > 1$) (\rightarrow Def. 2.1.8):

$$\epsilon_{k+1} \approx C \epsilon_k^p \Rightarrow \log \epsilon_{k+1} = \log C + p \log \epsilon_k \Rightarrow \log \epsilon_{k+1} = \log C + \sum_{l=0}^k p^l + p^{k+1} \log \epsilon_0$$

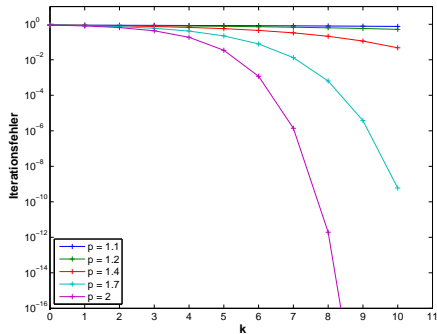
$$\Rightarrow \log \epsilon_{k+1} = \left(\frac{\log C}{p-1} + \log \epsilon_0 \right) p^{k+1}.$$

► Kurve ist negative Exponentialkurve im lin-log-Plot.

Definition 2.1.8 (Konvergenzordnung).

Eine **konvergente** Folge $\mathbf{x}^{(k)}$, $k = 0, 1, 2, \dots$, in \mathbb{R}^n konvergiert in p -ter Ordnung gegen $\mathbf{x}^* \in \mathbb{K}^n$, falls $p > 1$ und

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



◁ Qualitative Fehlerkurven für Konvergenz der Ordnung p (lin-log-Plot)

Beispiel 31 (Quadratische Konvergenz einer Iteration).

Wurzeliteration für $a > 0$, \rightarrow Bsp. 13, Monotonie & Beschränktheit \Rightarrow Konvergenz

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

AGM: $x^{(k)} > \sqrt{a}$ für $k \geq 1$ ► Konvergenzordnung 2 (**Quadratische Konvergenz**)
(AGM $\hat{=}$ Ungleichung zwischen arithmetischem und geometrischem Mittel: $\sqrt{ab} \leq \frac{1}{2}(a+b)$)

Kürze ab: $\epsilon_k := \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$

$$\epsilon_{k+1} \approx C \epsilon_k^p \Rightarrow \log \epsilon_{k+1} \approx \log C + p \log \epsilon_k \Rightarrow \frac{\log \epsilon_{k+1} - \log \epsilon_k}{\log \epsilon_k - \log \epsilon_{k-1}} \approx p.$$

Iterierte für $a = 2$:

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	2.0000000000000000	0.58578643762690485	
1	1.5000000000000000	0.08578643762690485	
2	1.4166666666666665	0.00245310429357137	1.850
3	1.4142156862745096	0.00000212390141452	1.984
4	1.4142135623746897	0.0000000000159472	2.000
5	1.4142135623730949	0.0000000000000022	0.630

Verdopplung der Anzahl der gültigen Stellen in jeder Iteration !

[Rundungsfehlereinfluss!]

2.1
p. 61

2.1
p. 63

Relativer Fehler: $x^{(k)} = x^*(1 + \delta_k) \Rightarrow x^{(k)} - x^* = \delta_k x^*$

$$|x^* \delta_{k+1}| = |x^{(k+1)} - x^*| \leq C |x^{(k)} - x^*|^2 = C |x^* \delta_k|^2$$

$$\Rightarrow |\delta_{k+1}| \leq C |x^*| \delta_k^2.$$

PERMRT

Abbruchkriterium (engl. *termination criterion*) ?

Einfach:

Abbruch konvergenter Iteration $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}_0}$ falls

$$\|F(\mathbf{x}^{(k)})\| \leq \tau = \text{vorgegebene Toleranz} > 0.$$

RESTERMCRIT
= Residuenbasiertes Abbruchkriterium

2.1
p. 62

2.1
p. 64

Schwieriger: Abbruch konvergenter Iteration $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}_0}$ mit Grenzwert \mathbf{x}^* falls

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \tau = \text{vorgegebene Toleranz} > 0.$$



\mathbf{x}^* leider unbekannt !

Strategien:

- ☞ Maschinenunabhängiges Abbruchkriterium $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ in \mathbb{M} (→ overkill ?)
- ☞ **A priori** Abbruchkriterium: fixiere Anzahl der Iterationen im Voraus (→ problemspezifisch !)
- ☞ **A posteriori** Abbruchkriterium: STOP, falls $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \tau'$ (→ zuverlässig ?)
- ☞ Ersatzlösung: **Residuenbasiertes** Abbruchkriterium:
STOP, falls $\|F(\mathbf{x}^{(k)})\| \leq \tau''$ (→ zuverlässig ?)

A posteriori Abbruchkriterium für linear konvergente Iterationen mit Konvergenzrate $0 < L < 1$:

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \stackrel{\Delta\text{-Ungl.}}{\leq} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| + \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| + L \|\mathbf{x}^{(k)} - \mathbf{x}^*\|.$$



$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \frac{L}{1-L} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|$$

(2.1.2)

Beispiel 32 (Fehlerschätzung für linear konvergente Iteration). → Bsp. 30

$$x^{(k+1)} = x^{(k)} + \frac{\cos x^{(k)} + 1}{\sin x^{(k)}} \Rightarrow x^{(k)} \rightarrow \pi \quad \text{für } x^{(0)} \text{ nahe bei } \pi.$$

Beobachtete Konvergenzrate: $L = 1/2$

Fehler und geschätzter Fehler für $x^{(0)} = 0.4$:

k	$ x^{(k)} - \pi $	$\frac{L}{1-L} x^{(k)} - x^{(k-1)} $	Schätzfehler
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

➤ A posteriori Schätzung sehr genau !

Bemerkung. L nicht genau bekannt ➤ $\tilde{L} > L$ für Fehlerschätzung garantiert Genauigkeit

2.2 Fixpunktiteration

[File: section-fixpunktiteration.tex, SVN: section-fixpunktiteration.tex 1103 2006-11-07 05:33:57Z hiptmainr]

Eine **Fixpunktiteration** ist definiert durch eine **Iterationsfunktion** $\Phi : \text{dom}(\Phi) \subset \mathbb{R}^n \mapsto \mathbb{R}^n$.

Iterationsfunktion $\Phi : \text{dom}(\Phi) \subset \mathbb{R}^n \mapsto \mathbb{R}^n$
Startwert $\mathbf{x}^{(0)} \in \text{dom}(\Phi) \Rightarrow$ Iterationsfolge $(\mathbf{x}^{(k)})_{k \in \mathbb{N}_0} : \mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)})$.

Iterationsfolge nicht unbedingt wohldefiniert: $\mathbf{x}^{(k)} \notin \text{dom}(\Phi)$ möglich !

Definition 2.2.1 (Spezialisierung von Def. 2.1.2). Eine **Fixpunktiteration** $\mathbf{x}^{(k+1)} = \Phi(\mathbf{x}^{(k)})$ ist **konsistent** (engl. consistent) mit $F(\mathbf{x}) = 0$, falls

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

Beachte: Φ stetig & Fixpunktiteration (lokal) konvergent gegen Nullstelle \mathbf{x}^* ➤ \mathbf{x}^* ist **Fixpunkt** (engl. fixed point) der Iterationsfunktion.

Konstruktion von Fixpunktiteration für $F(\mathbf{x}) = 0$:

$F(\mathbf{x}) = 0 \Leftrightarrow \Phi(\mathbf{x}) = \mathbf{x} \quad \blacktriangleright \quad \text{Fixpunktiteration} \quad \mathbf{x}^{(k+1)} := \Phi(\mathbf{x}^{(k)}) . \quad (2.2.1) \quad \text{fpt}$

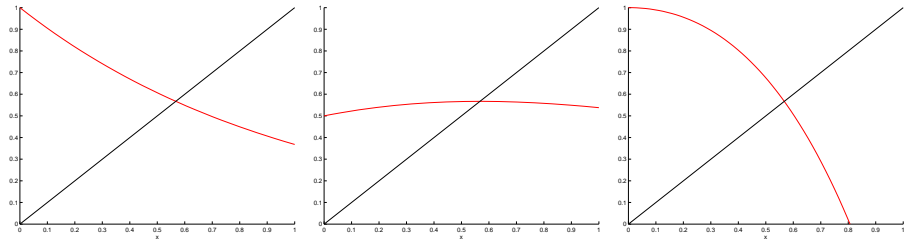
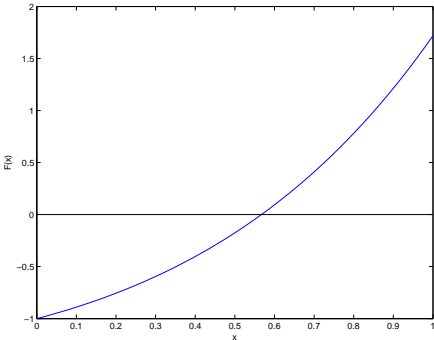
Transformation in Fixpunktform (nicht eindeutig):

Beispiel 33 (Fixpunktiteration).

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

Verschiedene Fixpunktformen:

$\Phi_1(x) = e^{-x} ,$
 $\Phi_2(x) = \frac{1+x}{1+e^x} ,$
 $\Phi_3(x) = x + 1 - xe^x .$



Funktion Φ_1

Funktion Φ_2

Funktion 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	0.000000000000003	0.288178118764323
4	0.007078662470882	0.000000000000000	0.723649245792953
5	0.004028858567431	0.000000000000000	0.410183132337935
6	0.002280343429460	0.000000000000000	1.186907542305364
7	0.001294757160282	0.000000000000000	0.146569797006362
8	0.000733837662863	0.000000000000000	0.310516641279937
9	0.000416343852458	0.000000000000000	0.357777386500765
10	0.000236077474313	0.000000000000000	0.974565695952037

Beobachtung: Lineare Konvergenz von $x_1^{(k)}$, quadratische Konvergenz von $x_2^{(k)}$,
Keine Konvergenz (chaotisches Verhalten) von $x_3^{(k)}$, $x_i^{(0)} = 0.5$.

KONV ANALYSE FIXP
Konvergenzanalyse, Fall $n = 1 \quad (\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R})$

2.2
p. 69

THEOREM 2.2.2 (Taylorformel). Falls $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$ $(m + 1)$ -mal stetig differenzierbar, $x \in \text{dom}(\Phi)$

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

Anwendung auf Fixpunktiteration mit Iterationsfunktion Φ :

Falls $\Phi(x^*) = x^*$, $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$ "hinreichend glatt"

$$x^{(k+1)} - x^* = \Phi(x^{(k)}) - \Phi(x^*) = \sum_{l=1}^m \frac{1}{l!} \Phi^{(l)}(x^*)(x^{(k)} - x^*)^l + O(|x^{(k)} - x^*|^{m+1}) . \quad (2.2.2) \quad \text{eq:FPT}$$

LEMMA 2.2.3. Falls $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$ $(m + 1)$ -mal stetig differenzierbar, $\Phi(x^*) = x^*$ für ein x^* im Innern von $\text{dom}(\Phi)$, und $\Phi^{(l)}(x^*) = 0$ für $l = 1, \dots, m, m \geq 1$, dann konvergiert die Fixpunktiteration (2.2.1) lokal mit Ordnung $m + 1$ (\rightarrow Def. 2.1.8) gegen x^* .

2.2
p. 70

2.2
p. 71

2.2
p. 72

Beweis. Für Umgebung \mathcal{U} von x^*

$$\stackrel{\text{lex:FPT}}{(2.2.2)} \Rightarrow \exists C > 0: |\Phi(y) - \Phi(x^*)| \leq C |y - x^*|^{m+1} \quad \forall y \in \mathcal{U}.$$

$$\delta^m C < 1/2: |x^{(0)} - x^*| < \delta \Rightarrow |x^{(k)} - x^*| < 2^{-k} \delta \Rightarrow \text{lokale Konvergenz} . \quad \square$$

Fortsetzung Bsp. 33: $\stackrel{\text{lex:fixp}}{}$

$$\Phi_2'(x) = \frac{1 - xe^x}{(1 + e^x)^2} = 0 \quad \text{wenn} \quad xe^x - 1 = 0 \quad \blacktriangleright \quad \text{Quadratische Konvergenz} .$$

TAYLORMEHRDIM

Lemma 2.2.4 (Taylorformel im Mehrdimensionalen). Falls $\Phi : \text{dom}(\Phi) \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ 2-fach stetig differenzierbar

$$\Phi(y) - \Phi(x) = D\Phi(x)(y - x) + O(\|y - x\|^2) \quad \forall y \in \mathcal{U}(x) .$$

Terminologie: $D\Phi(x) \in \mathbb{R}^{n,n}$ = Jacobimatrix von Φ an der Stelle x

Lemma 2.2.5. Falls $\Phi : \text{dom}(\Phi) \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ 2-fach stetig differenzierbar, $\Phi(x^*) = x^*$, $\|D\Phi(x^*)\| < 1$, dann konvergiert die Fixpunktiteration $\stackrel{\text{fixp}}{(2.2.1)}$ lokal linear.

MATNORM
Matrixnorm !

Definition 2.2.6. Der Norm $\|\cdot\|$ auf \mathbb{R}^n zugeordnete **MATRIXNORM** Matrixnorm

$$\mathbf{M} \in \mathbb{R}^{n,n}: \|\mathbf{M}\| := \sup_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|\mathbf{M}x\|}{\|x\|} .$$

Notation: $\|x\|_2 \rightarrow \|\mathbf{M}\|_2, \quad \|x\|_1 \rightarrow \|\mathbf{M}\|_1, \quad \|x\|_\infty \rightarrow \|\mathbf{M}\|_\infty$

Beispiel 34 (Matrixnorm zur Maximumnorm und 1-Norm).

$$\begin{aligned} \text{z.B. für } \mathbf{M} \in \mathbb{K}^{2,2}: \quad \|\mathbf{M}x\|_\infty &= \max\{|m_{11}x_1 + m_{12}x_2|, |m_{21}x_1 + m_{22}x_2|\} \\ &\leq \max\{|m_{11}| + |m_{12}|, |m_{21}| + |m_{22}|\} \|x\|_\infty , \\ \|\mathbf{M}x\|_1 &= |m_{11}x_1 + m_{12}x_2| + |m_{21}x_1 + m_{22}x_2| \\ &\leq \max\{|m_{11}| + |m_{21}|, |m_{12}| + |m_{22}|\} (|x_1| + |x_2|) . \end{aligned}$$

$$\blacktriangleright \text{Matrixnorm zu } \|\cdot\|_\infty = \text{Zeilensummennorm} \quad \|\mathbf{M}\|_\infty := \max_{i=1,\dots,n} \sum_{j=1}^n |m_{ij}| , \quad (2.2.3) \quad \text{eq: defrow}$$

$$\blacktriangleright \text{Matrixnorm zu } \|\cdot\|_1 = \text{Spaltensummennorm} \quad \|\mathbf{M}\|_1 := \max_{j=1,\dots,m} \sum_{i=1}^n |m_{ij}| . \quad (2.2.4) \quad \text{eq: defcol}$$

◇

Beweis von Lemma 2.2.5. $\stackrel{\text{lem:fixp}}{}$ Für Umgebung $\mathcal{U}(x^*)$:

$$\exists C > 0: \|\Phi(y) - \Phi(x) - D\Phi(x)(y - x)\| \leq C \|y - x\|^2 \quad \forall y \in \mathcal{U}(x^*) .$$

Wähle $\delta > 0$ so dass $L := \delta C + \|D\Phi(x^*)\| < 1, \{y \in \mathbb{R}^n: \|y - x^*\| < \delta\} \subset \mathcal{U}(x^*)$

$$\blacktriangleright \quad \|x - x^*\| < \delta \Rightarrow \|\Phi(x) - x^*\| \leq \underbrace{(\|D\Phi(x^*)\| + C \|x - x^*\|)}_{\leq L} \|x - x^*\| . \quad \square$$

Bemerkung 35.

$$\text{Falls } \|D\Phi(x^*)\| < 1, x^{(k)} \approx x^* \quad \blacktriangleright \quad \text{ASYMPKONVRTE} \quad \text{Asymptotische Konvergenzrate} \quad L = \|D\Phi(x^*)\|$$

2.2
p. 73

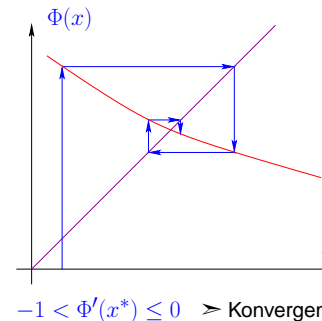
2.2
p. 75

Fortsetzung Bsp. 33: $\stackrel{\text{lex:fixp}}{}$ Da $x^*e^{x^*} - 1 = 0$

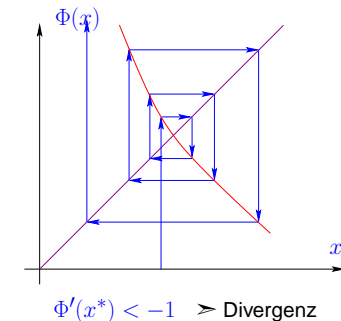
$$\Phi_1'(x) = -e^{-x} \Rightarrow \Phi_1'(x^*) = -x^* \approx -0.56 \quad \blacktriangleright \quad \text{Lokal lineare Konvergenz} .$$

$$\Phi_3'(x) = 1 - xe^x - e^x \Rightarrow \Phi_3'(x^*) = -\frac{1}{x^*} \approx -1.79 \quad \blacktriangleright \quad \text{Keine Konvergenz} .$$

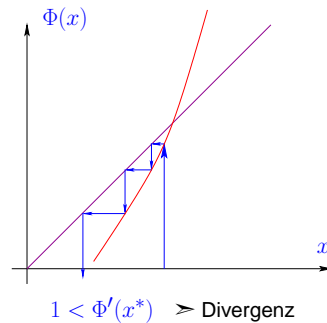
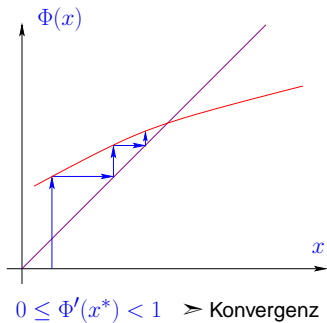
Definition 2.2.7. Für stetig differenzierbares $\Phi : \text{dom}(\Phi) \subset \mathbb{R} \mapsto \mathbb{R}$ heisst ein Fixpunkt $x^* \in \text{dom}(\Phi)$ **FIXPUNKT** ~~FIXPUNKT~~ **anziehend** (engl. attractive), falls $|\Phi'(x^*)| < 1$ und **abstossend**, falls $|\Phi'(x^*)| > 1$



2.2
p. 74



2.2
p. 76



Beispiel 36 (Mehrdimensionale Fixpunktiteration).

$$\begin{aligned} \text{Gleichungssystem} & \Rightarrow \text{Fixpunktform} \\ \begin{aligned} x_1 - c(\cos x_1 - \sin x_2) &= 0 \\ (x_1 - x_2) + c \sin x_2 &= 0 \end{aligned} & \Rightarrow \begin{aligned} c(\cos x_1 - \sin x_2) &= x_1 \\ c(\cos x_1 - 2 \sin x_2) &= x_2 \end{aligned} \\ \Rightarrow \Phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = c \begin{pmatrix} \cos x_1 - \sin x_2 \\ \cos x_1 - 2 \sin x_2 \end{pmatrix} & \Rightarrow D\Phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = -c \begin{pmatrix} \sin x_1 & \cos x_2 \\ \sin x_1 & 2 \cos x_2 \end{pmatrix} \end{aligned}$$

Wahl einer geeigneten Vektornorm: $\|\cdot\| = \text{Maximumnorm}$ $\|\cdot\|_\infty$ (\rightarrow Bsp. 65) eq:infmatrixnorm

Wenn $c < \frac{1}{3} \Rightarrow \|D\Phi(x)\|_\infty < 1 \quad \forall x \in \mathbb{R}^2 \Rightarrow$ Lineare Konvergenz der Fixpunktiteration

◇

Abbruchkriterien für kontrahierende Iterationen mit Kontraktionsrate $0 < L < 1$:

$$\begin{aligned} \|x^{(k+m)} - x^{(k)}\| &\stackrel{\Delta\text{-Ungl.}}{\leq} \sum_{j=k}^{k+m-1} \|x^{(j+1)} - x^{(j)}\| \leq \sum_{j=k}^{k+m-1} L^{j-k} \|x^{(k+1)} - x^{(k)}\| \\ &= \frac{1-L^m}{1-L} \|x^{(k+1)} - x^{(k)}\| \leq \frac{1-L^m}{1-L} L^{k-l} \|x^{(l+1)} - x^{(l)}\|. \end{aligned}$$

► für $m \rightarrow \infty$, mit $x^* := \lim_{k \rightarrow \infty} x^{(k)}$:

$$\|x^* - x^{(k)}\| \leq \frac{L^{k-l}}{1-L} \|x^{(l+1)} - x^{(l)}\|. \quad (2.2.5) \quad \text{eq:linvcvg}$$

Setze $l = 1$ in (2.2.5) eq:linvcvg

⚠ a priori Abbruchkriterium

$$\|x^* - x^{(k)}\| \leq \frac{L^k}{1-L} \|x^{(1)} - x^{(0)}\| \quad (2.2.6) \quad \text{eq:linvcvg}$$

Setze $l = k$ in (2.2.5) eq:linvcvg

⚠ a posteriori Abbruchkriterium

$$\|x^* - x^{(k)}\| \leq \frac{L}{1-L} \|x^{(k)} - x^{(k-1)}\| \quad (2.2.7) \quad \text{eq:fixapost}$$

2.3 Nullstellenbestimmung von Funktionen

[File: section-nullstellenbestimmung-von-funktionen.tex, SVN: section-nullstellenbestimmung-von-funktionen.tex 1183 2006-11-29 10:46:03Z hiptmair]

Spezialfall $n = 1$:

$F : I \subset \mathbb{R} \rightarrow \mathbb{R}$ stetig, I Intervall

Gesucht:

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

2.3.1 Bisektionsverfahren

Idee: Nutze Anordnung der reellen Zahlen aus (Intervallschachtelung) & Zwischenwertsatz

Input: $a, b \in I$ mit $F(a)F(b) < 0$ (unterschiedliche Vorzeichen)

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

MATLAB-Implementierung:

2.2
p. 77

```
MATLAB-CODE: Bisektionsverfahren
function x = bisection(F,a,b,tol)
% Searching zero by bisection
if (a>b), t=b; a=b; b=t; end;
fa = F(a); fb = F(b);
if (fa*fb>0)
    error('f(a), f(b) same sign'); end;
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
```

Handle auf MATLAB-Funktion, die F implementiert.

Vermeide Endlosschleife, falls $\text{tol} < \text{Abstand von Maschinenzahlen bei Nullstelle } x^*$.



- „narrensicher“
- Nur Auswertungen von F erforderlich



Nur lineare Konvergenz: $|x^{(k)} - x^*| \leq 2^{-k}|b - a| \Rightarrow \log_2 \left(\frac{|b - a|}{\text{tol}} \right)$ Schritte

Bemerkung 37. Bisektion implementiert in MATLAB-Funktion `fzero`.

2.3
p. 78

2.3
p. 79

△

2.3
p. 80

2.3.2 Modellfunktionsverfahren

MODELLFUNKTVERF

Iterationenverfahren zur Nullstellenbestimmung von F :



Idee: Gegeben approximative Nullstellen $x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}$

- ersetze F durch **MODELFUNCT** Modellfunktion \tilde{F}
(verwende Funktionswerte/Ableitungen in $x^{(k)}, x^{(k-1)}, \dots, x^{(k-m)}$)
- $x^{(k+1)} :=$ Nullstelle von \tilde{F}
(muss "einfach" \leftrightarrow analytisch zu bestimmen sein)

Unterscheidung:

EINPUNKTVERFAHREN

Einpunktverfahren: $x^{(k+1)} = x^{(k+1)}(F, x^{(k)}), k \in \mathbb{N}$ (z.B. Fixpunktiteration \rightarrow Abschn. 2.2)

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

2.3.2.1 Newton-Verfahren in 1D

NEWTONONED

Annahme: $F : I \mapsto \mathbb{R}$ stetig differenzierbar.

Modellfunktion = Tangente an F in $x^{(k)}$:

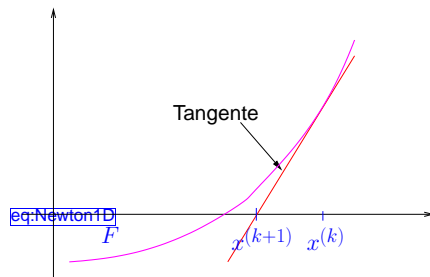
$$\tilde{F}(x) := F(x^{(k)}) + F'(x^{(k)})(x - x^{(k)})$$

► $x^{(k+1)}$ als Nullstelle der Tangente

► **NEWTONIT**
Newton-Iteration

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

(2.3.1)



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

Beispiel 38 (Newton-Verfahren in 1D). (\rightarrow Beispiel 33)

$$F(x) = xe^x - 1 \Rightarrow F'(x) = e^x(1+x) \Rightarrow x^{(k+1)} = x^{(k)} - \frac{x^{(k)}e^{x^{(k)}} - 1}{e^{x^{(k)}}(1+x^{(k)})} = \frac{(x^{(k)})^2 - e^{-x^{(k)}}}{1+x^{(k)}}$$

$$F(x) = x - e^{-x} \Rightarrow F'(x) = 1 + e^{-x} \Rightarrow x^{(k+1)} = x^{(k)} - \frac{x^{(k)} - e^{-x^{(k)}}}{1 + e^{-x^{(k)}}} = \frac{1 + x^{(k)}}{1 + e^{x^{(k)}}}$$

Beispiel 33 quadratische Konvergenz! (\rightarrow Def. 2.1.8)

◇

Newton-Iteration (2.5.1) \equiv Fixpunktiteration (\rightarrow Abschnitt 2.2) mit Iterationsfunktion

$$\Phi(x) = x - \frac{F(x)}{F'(x)} \Rightarrow \Phi'(x) = \frac{F(x)F''(x)}{(F'(x))^2} \Rightarrow \Phi'(x^*) = 0, \text{ wenn } F(x^*) = 0, F'(x^*) \neq 0.$$

Lemma 2.2.3 $\xrightarrow{\text{lem:fpcvg}}$ Vermutung (bestätigt durch Theorie, [12, Thm. 4.10], [10, Sect. 2.1]):

Newton-Verfahren konvergiert **LOKQUADKONV** lokal quadratisch (\rightarrow Def. 2.1.8), falls $F'(x^*) \neq 0$

2.3.2.2 Spezielle Einpunktverfahren

Konstruktionsidee (für weitere Einpunktverfahren): Nichtlineare lokale Approximation

Voraussetzung: Glattheit von F : $F \in C^m(I)$ für ein $m > 1$

Beispiel 39 (Hallesche Iteration).

Gegeben $x^{(k)} \in I$, Modellfunktion $h(x^{(k+1)}) = 0$ wobei

$$h(x) = \frac{a}{x+b} + c \text{ (Rationale Funktion)}, F^{(j)}(x^{(k)}) = h^{(j)}(x^{(k)}), j = 0, 1, 2.$$

2.3
p. 81

$$\frac{a}{x^{(k)}+b} + c = F(x^{(k)}), -\frac{a}{(x^{(k)}+b)^2} = F'(x^{(k)}), \frac{2a}{(x^{(k)}+b)^3} = F''(x^{(k)}).$$

$$x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})} \cdot \frac{1}{1 - \frac{1}{2} \frac{F(x^{(k)})F''(x^{(k)})}{F'(x^{(k)})^2}}.$$

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

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

2.3 Vergleich: Newton-Verfahren (2.3.1) für das gleiche Problem
p. 82

2.3
p. 83

2.3
p. 84

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

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

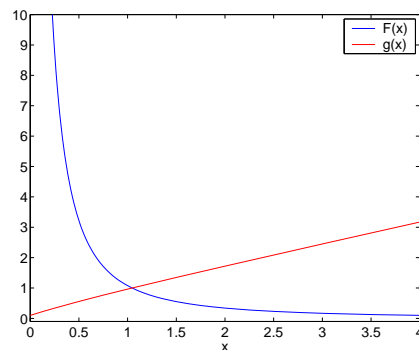
► Überlegenheit der Halley-Iteration, da F rationale Funktion.

! Newton-Verfahren langsamer, aber auch billiger (→ Abschn. 2.4) sec:effizienz

Beispiel 40 (Massgeschneidertes Newton-Verfahren).

Wie in Bsp. 39: lex:Halley

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



Beobachtung:

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

► $g(x) := \frac{1}{\sqrt{F(x)+1}}$ **fast linear** für $x \gg 1$

► Anstelle von $F(x) \stackrel{!}{=} 0$ löse $g(x) \stackrel{!}{=} 1$ mit Newton-Verfahren leg:Newton1D

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

Konvergenzhistorie für $x^{(0)} = 0$:

Weitere Idee: **Konsistente Modifikation** der Newton-Iterationsvorschrift:

► Fixpunktiteration: $\Phi(x) = x - \frac{F(x)}{F'(x)} H(x)$ mit „geeignetem“ $H: I \mapsto \mathbb{R}$.

Ziel: Bestimme H so, dass Verfahren p -ter Ordnung entsteht, Werkzeug: Lemma 2.2.3. lem:fpcv

Annahme: F „genügend“ glatt, $\exists x^* \in I: F(x^*) = 0, F'(x^*) \neq 0$.

$$\begin{aligned} \Phi &= x - uH, \quad \Phi' = 1 - u'H - uH', \quad \Phi'' = -u''H - 2u'H - uH'', \\ \text{mit } u &= \frac{F}{F'} \Rightarrow u' = 1 - \frac{FF''}{(F')^2}, \quad u'' = -\frac{F''}{F'} + 2\frac{F(F'')^2}{(F')^3} - \frac{FF'''}{(F')^2}. \end{aligned}$$

2.3
p. 85

2.3
p. 87

$$F(x^*) = 0 \quad u(x^*) = 0 \quad u'(x^*) = 1 \quad u''(x^*) = -\frac{F''(x^*)}{F'(x^*)}$$

► $\Phi'(x^*) = 1 - H(x^*) \quad \Phi''(x^*) = \frac{F''(x^*)}{F'(x^*)} H(x^*) - 2H'(x^*) \quad (2.3.2)$ MN31

lem:fpcv
2.2.3

p

$p = 2$ (quadratische Konvergenz): $H(x^*) = 1$,

$p = 3$ (**KUBKONV** kubische Konvergenz): $H(x^*) = 1 \wedge H'(x^*) = \frac{1}{2} \frac{F''(x^*)}{F'(x^*)}$.

$$H(x) = G(1 - u'(x)) \quad G$$

► Fixpunktiteration $x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})} G \left(\frac{F(x^{(k)})F''(x^{(k)})}{(F'(x^{(k)}))^2} \right) \quad (2.3.3)$ eq:modN

Lemma 2.3.1. Falls $F \in C^2(I)$, $F(x^*) = 0$, $F'(x^*) \neq 0$, $G \in C^2(U)$ in einer Umgebung U von 0, $G(0) = 1$, $G'(0) = \frac{1}{2}$, dann konvergiert die Fixpunktiteration leg:modNewt (2.3.3) lokal kubisch gegen x^* .

2.3
p. 86

2.3
p. 88

Beweis: Lemma 2.2.3, (2.3.2) und

$$H(x^*) = G(0) \quad , \quad H'(x^*) = -G'(0)u''(x^*) = G'(0)\frac{F''(x^*)}{F'(x^*)}.$$

Beispiel 41 (Beispiele für modifizierte Newton-Verfahren).

- $G(t) = \frac{1}{1 - \frac{1}{2}t} \Rightarrow$ Halleysche Methode (\rightarrow Beispiel 39)
- $G(t) = \frac{2}{1 + \sqrt{1 - 2t}} \Rightarrow$ Eulersche Methode
- $G(t) = 1 + \frac{1}{2}t \Rightarrow$ Quadratische inverse Interpolation

Numerisches Experiment:

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

k	$e^{(k)} := x^{(k)} - x^*$		
	Halley	Euler	Quad. Inv.
1	2.81548211105635	3.57571385244736	2.03843730027891
2	1.37597082614957	2.76924150041340	1.02137913293045
3	0.34002908011728	1.95675490333756	0.28835890388161
4	0.00951600547085	1.25252187565405	0.01497518178983
5	0.00000024995484	0.51609312477451	0.00000315361454
6		0.14709716035310	
7		0.00109463314926	
8		0.00000000107549	

2.3
p. 89

2.3.2.3 Mehrpunktverfahren



Idee:

Ersetze F durch **Interpolationspolynom**
Interpolationsverfahren

Einfachster Vertreter: **Sekantenverfahren**

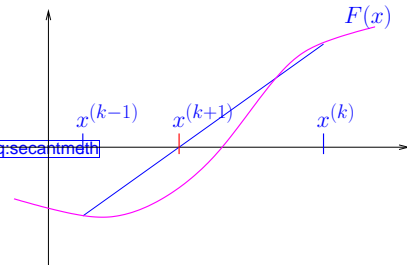
$x^{(k+1)}$ Nullstelle der Sekante

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

(2.3.4) [eq:secantmeth](#)

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

(2.3.5)



2.3
p. 90

Sekantenverfahren

(MATLAB-Implementierung)

- Nur eine Funktionsauswertung von F pro Schritt!
- **Ableitungsfrei!**

```
function x = secant(x0,x1,F,tol)
fo = feval(F,x0);
for i=1:MAXIT
    fn = feval(F,x1);
    s = fn*(x1-x0)/(fn-fo);
    x0 = x1; x1 = s;
    if (abs(s) < tol), x = x1; return; end
    fo = fn;
end
```

MATLAB-CODE Sekantenverfahren

Beispiel 42 (Sekantenverfahren). $F(x) = xe^x - 1$, $x^{(0)} = 0$, $x^{(1)} = 5$.

k	$x^{(k)}$	$F(x^{(k)})$	$e^{(k)} := x^{(k)} - x^*$	$\frac{\log e^{(k+1)} - \log e^{(k)} }{\log e^{(k)} - \log e^{(k-1)} }$
2	0.00673794699909	-0.99321649977589	-0.56040534341070	
3	0.01342122983571	-0.98639742654892	-0.55372206057408	24.43308649757745
4	0.98017620833821	1.61209684919288	0.41303291792843	2.70802321457994
5	0.38040476787948	-0.44351476841567	-0.18673852253030	1.48753625853887
6	0.50981028847430	-0.15117846201565	-0.05733300193548	1.51452723840131
7	0.57673091089295	0.02670169957932	0.00958762048317	1.70075240166256
8	0.56668541543431	-0.00126473620459	-0.00045787497547	1.59458505614449
9	0.56713970649585	-0.00000990312376	-0.0000358391394	1.62641838319117
10	0.56714329175406	0.00000000371452	0.0000000134427	
11	0.56714329040978	-0.00000000000001	-0.00000000000000	

2.3
p. 91

Konvergenzanalyse: Fehler $e^{(k)} := x^{(k)} - x^*$

$$e^{(k+1)} = \Phi(x^* + e^{(k)}, x^* + e^{(k-1)}) - x^* \quad \text{mit} \quad \Phi(x, y) := x - \frac{F(x)(x - y)}{F(x) - F(y)}. \quad (2.3.6) \quad \text{eq:secm}$$

Taylorentwicklung mit MAPLE:

```
> Phi := (x,y) -> x-F(x)*(x-y)/(F(x)-F(y));
> F(s) := 0;
> e2 = normal(mttaylor(Phi(s+e1,s+e0)-s,[e0,e1],4));
```

$$\text{▶ (Linearisierte Fehlerrekursion: } e^{(k+1)} \doteq \frac{1}{2} \frac{F''(x^*)}{F'(x^*)} e^{(k)} e^{(k-1)} = C e^{(k)} e^{(k-1)} \text{)}$$

Heuristik zur Bestimmung der Konvergenzordnung: $e^{(k+1)} = K^2(e^{(k-1)})^2$, $e^{(k)} = K(e^{(k-1)})^p$
(vgl. Def. 2.1.8) [def:cvgord](#)

$$\text{▶ } (e^{(k-1)})^{p^2-p-1} = K^{-1}C \Rightarrow p^2 - p - 1 = 0 \Rightarrow p = \frac{1}{2}(1 \pm \sqrt{5}).$$

Da $e^{(k)} \rightarrow 0$ für $k \rightarrow \infty$ \Rightarrow Konvergenzrate $p = \frac{1}{2}(1 + \sqrt{5}) \approx 1.62$ (siehe Bsp. 42!) [ex:sec](#)

2.3
p. 90

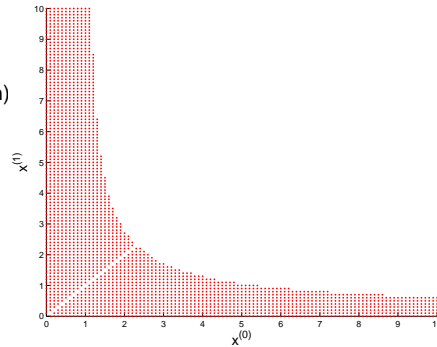
2.3
p. 92

Beispiel 43 (Lokale Konvergenz, Sekantenverfahren)

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

• : Konvergenz des Sekantenverfahrens
für ein Paar $(x^{(0)}, x^{(1)})$ von Startwerten. \triangleright

= lokale Konvergenz \rightarrow Def. 2.1.3 ^{def:cvgit}



INVINTERP

Variante von Mehrpunktverfahren: *Inverse Interpolation*

Annahme:

$F : I \subset \mathbb{R} \mapsto \mathbb{R}$ umkehrbar

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

- Interpoliere F^{-1} durch Polynom p vom Grad d festgelegt durch

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

- Neue Näherung $x^{(k+1)} := p(0)$

Fall $m = 1 \rightarrow$ Sekantenverfahren

Fall $m = 2$: Quadratische inverse Interpolation

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);`

$$x^{(k+1)} = \frac{F_0^2(F_1x_2 - F_2x_1) + F_1^2(F_2x_0 - F_0x_2) + F_2^2(F_0x_1 - F_1x_0)}{F_0^2(F_1 - F_2) + F_1^2(F_2 - F_0) + F_2^2(F_0 - F_1)}.$$

($F_0 := F(x^{(k-2)}), F_1 := F(x^{(k-1)}), F_2 := F(x^{(k)}), x_0 := x^{(k-2)}, x_1 := x^{(k-1)}, x_2 := x^{(k)}$)

Beispiel 44 (Quadratische 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	0.00000000000001	

2.3
p. 93

2.4 Effizienz

[File: section-effizienz.tex, SVN: section-effizienz.tex 1054 2006-10-18 11:43:04Z kalai]

EFFIZIENZ

Effizienz eines Iterationsverfahrens
(zur Lösung von $F(x) = 0$)

RECHNAUFW

Rechenaufwand um eine vorgegebene
Anzahl gültiger Stellen zu erreichen.

Abstrakt:

$W \triangleq$ Rechenaufwand pro Schritt

$$(\text{etwa } W \approx \frac{\#\{\text{Auswertungen von } F\}}{\text{Schritt}} + n \cdot \frac{\#\{\text{Auswertungen von } F'\}}{\text{Schritt}} + \dots)$$

Quantitative Analyse für Iterationsverfahren der Konvergenzordnung p , $p \geq 1$ (\rightarrow Def. 2.1.8), d.h., ^{def:cvgord}

$$\exists C > 0: \|e^{(k)}\| \leq C \|e^{(k-1)}\|^p \quad \forall k \geq 1.$$

2.3
p. 94

2.4
p. 95

2.4
p. 96

Anzahl Schritte k für $\|e^{(k)}\| \leq \rho \|e^{(0)}\|$, $\rho > 0$ vorgegeben?: (Annahme $C \|e^{(0)}\|^{p-1} < 1$)

$$p = 1: \quad \|e^{(k)}\| \leq C^k \|e^{(0)}\| \Rightarrow k \geq \frac{\log \rho}{\log C}, \quad (2.4.1) \quad \text{eq:IS}$$

$$p > 1: \quad \|e^{(k)}\| \leq C^{\frac{p^k - 1}{p-1}} \|e^{(0)}\|^{p^k} \Rightarrow p^k \geq 1 + \frac{\log \rho}{\log C/p - 1 + \log(\|e^{(0)}\|)} \\ \Rightarrow k \geq \log(1 + \frac{\log \rho}{\log L_0}) / \log p, \quad L_0 := C^{1/p-1} \|e^{(0)}\|. \quad (2.4.2) \quad \text{eq:IS}$$

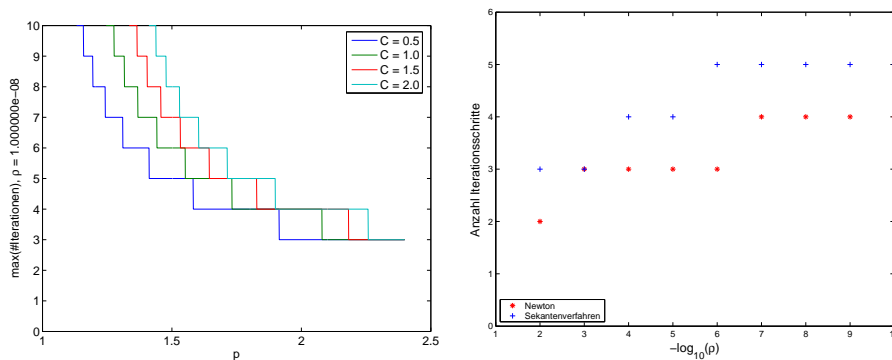
Beachte: $|\log \rho| \leftrightarrow$ Anzahl gültiger Stellen in $x^{(k)}$

EFFICIENCY

Effizienzmass (asymptotisch für $\rho \rightarrow 0 \rightarrow |\log \rho| \rightarrow \infty$):

$$\text{Effizienz} = \begin{cases} -\frac{\log C}{W} & , \text{ falls } p = 1, \\ \frac{\log p |\log \rho|}{W \log |\log \rho|} & , \text{ falls } p > 1. \end{cases} \quad (2.4.3) \quad \text{eq:effie}$$

Beispiel 45 (Effizienz von Iterationsverfahren).



Auswertung (2.4.2) für $\|e^{(0)}\| = 0.1$, $\rho = 10^{-8}$

Newton \leftrightarrow Sekantenverfahren, $C = 1$,
Anfangsfehler $\|e^{(0)}\| = 0.1$

$$\frac{W_{\text{Newton}}}{p_{\text{Newton}}} = 2, \quad \frac{W_{\text{Sekant}}}{p_{\text{Sekant}}} = 1.62 \quad \Rightarrow \quad \frac{\log p_{\text{Newton}}}{W_{\text{Newton}}} : \frac{\log p_{\text{Sekant}}}{W_{\text{Sekant}}} = 0.71.$$

Sekantenverfahren effizienter als Newton-Verfahren

2.5 Newton-Verfahren für Gleichungssysteme

[File: section-newton-verfahren.tex, SVN: section-newton-verfahren.tex 1234 2006-12-19 12:04:51Z hiptmair]

Nichtlineares Gleichungssystem: zu $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ finde $\mathbf{x}^* \in D$: $F(\mathbf{x}^*) = 0$

Annahme: $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ stetig differenzierbar

2.5.1 Die Newton-Iteration

Idee, siehe Abschn. 2.3.2.1: Lokale Linearisierung:

Gegeben $x^{(k)} \in D$ bestimme $x^{(k+1)}$ als Nullstelle der affin linearen Modellfunktion

$$F(\mathbf{x}) \approx \tilde{F}(\mathbf{x}) := F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}),$$

$$DF(\mathbf{x}) \in \mathbb{R}^{n,n} = \text{Jacobi-Matrix}, \quad DF(\mathbf{x}) = \left(\frac{\partial F_j}{\partial x_k}(\mathbf{x}) \right)_{j,k=1}^n.$$



Newton-Iteration: (\rightarrow (2.3.1) für 1D-Fall)

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - DF(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)}), \quad [\text{falls } DF(\mathbf{x}^{(k)}) \text{ regulär}] \quad (2.5.1) \quad \text{eq:Newt}$$

2.4

p. 97

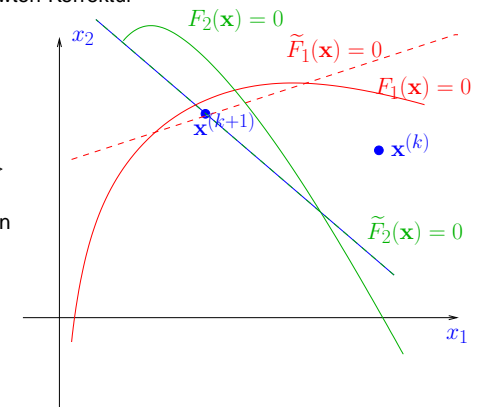
Terminologie: $-DF(\mathbf{x}^{(k)})^{-1} F(\mathbf{x}^{(k)}) = \text{Newton-Korrektur}$

2.5

p. 99

Veranschaulichung in 2D:

$\mathbf{x}^{(k+1)}$ als Schnitt von Geraden (= Nullstellen
mengen der Komponenten der Modellfunktion)



Beispiel 46 (Newton-Verfahren in 2D, [32]).

$$F(\mathbf{x}) = \begin{pmatrix} x_1^2 - x_2^4 \\ x_1 - x_2^3 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2 \quad \Rightarrow \quad F \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 0.$$

$$DF(\mathbf{x}) = \begin{pmatrix} \partial_{x_1} F_1(x) & \partial_{x_2} F_1(x) \\ \partial_{x_1} F_2(x) & \partial_{x_2} F_2(x) \end{pmatrix} = \begin{pmatrix} 2x_1 & -4x_2^3 \\ 1 & -3x_2^2 \end{pmatrix}$$

Realisierung der Newton-Iteration (2.5.1):

2.5

p. 98

2.5

p. 100

1. Mit $\mathbf{x}^{(k)} = (x_1, x_2)^T$ löse das Gleichungssystem $\begin{pmatrix} 2x_1 & -4x_2^3 \\ 1 & -3x_2^2 \end{pmatrix} s = -\begin{pmatrix} x_1^2 - x_2^4 \\ x_1 - x_2^2 \end{pmatrix}$.
2. Setze $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + s$

k	$\mathbf{x}^{(k)}$	$\ \mathbf{x}^* - \mathbf{x}^{(k)}\ _2$
0	$(0.7, 0.7)^T$	4.24e-01
1	$(0.87850000000000, 1.06428571428571)^T$	1.37e-01
2	$(1.01815943274188, 1.00914882463936)^T$	2.03e-02
3	$(1.00023355916300, 1.00015913936075)^T$	2.83e-04
4	$(1.00000000583852, 1.00000002726552)^T$	2.79e-08
5	$(0.99999999999998, 1.00000000000000)^T$	2.11e-15
6	$(1, 1)^T$	



Neuer Aspekt für $n \gg 1$ (im Vergleich zu $n = 1$, Abschn. 2.3.2.1)

Berechnung der Newton-Korrektur (evtl.) aufwändig, cf. (3.2.6)

MATLAB-Shell für Newton-Verfahren:

Löse lineares Gleichungssystem:

$A \setminus b = A^{-1}b \rightarrow$ Abschn. 3.2

F, DF : Funktionshandles

```
MATLAB-CODE: Newtonverfahren
function x = newton(x,F,DF,tol)
% MATLAB template for Newton method
for i=1:MAXIT
    s = DF(x) \ F(x);
    x = x-s;
    if (norm(s) < tol), return; end;
end
```

Bemerkung 47. Wichtige Eigenschaft des Newton-Verfahrens: **AFFININV** **Affin-Invarianz** [10, Sect. 1.2.2]

$$G(\mathbf{x}) := \mathbf{A}F(\mathbf{x}) \text{ mit } \mathbf{A} \in \mathbb{R}^{n,n} \text{ regulär} \Rightarrow F(\mathbf{x}^*) = 0 \Rightarrow G(\mathbf{x}^*) = 0.$$

Die Newton-Iteration für G ist unabhängig von \mathbf{A} !

Bemerkung 48. Vereinfachtes Newton-Verfahren: Verwende $DF(\mathbf{x}^{(k)})$ für mehrere Schritte (\rightarrow nur lokal lineare Konvergenz)

Bemerkung 49. Falls $DF(\mathbf{x})$ nicht verfügbar (z.B. wenn $F(\mathbf{x})$ nur als Prozedur gegeben):

$$\text{Numerisches Differenzieren: } \frac{\partial F_i}{\partial x_j}(\mathbf{x}) \approx \frac{F_i(\mathbf{x} + h\vec{e}_j) - F_i(\mathbf{x})}{h}.$$

Achtung: Auslöschung (\rightarrow Beispiel 27) $\frac{\partial F_i}{\partial x_j}(\mathbf{x}) \approx \frac{F_i(\mathbf{x} + h\vec{e}_j) - F_i(\mathbf{x})}{h} \approx \sqrt{\text{EPS}}$.

Übung 2.6. Given a matrix \mathbf{A} implement a numerical method for computing one eigenvalue/eigenvector of \mathbf{A} by applying Newton's method to

$$F\left(\begin{matrix} \mathbf{x} \\ \lambda \end{matrix}\right) := \left(\begin{matrix} \mathbf{A}\mathbf{x} - \lambda\mathbf{x} \\ \|\mathbf{x}\|^2 - 1 \end{matrix}\right).$$

2.6.1 Konvergenzanalyse des Newton-Verfahrens

Newton-Iteration (2.5.1) $\stackrel{\text{leg:Newton}}{=} \text{Fixpunktiteration } (\rightarrow \text{Abschnitt 2.2}) \text{ mit}$

$$\Phi(\mathbf{x}) = \mathbf{x} - DF(\mathbf{x})^{-1}F(\mathbf{x}).$$

$$\text{„Produktregel“: } D\Phi(\mathbf{x}) = \mathbf{I} - D(DF(\mathbf{x})^{-1})F(\mathbf{x}) - DF(\mathbf{x})^{-1}DF(\mathbf{x})$$

$$F(\mathbf{x}^*) = 0 \Rightarrow D\Phi(\mathbf{x}^*) = 0.$$

2.5 Lemma 2.2.3 $\stackrel{\text{lem:fpcvg}}{\Rightarrow}$ Vermutung (bestätigt durch Theorie, [12, Thm. 4.10], [10, Sect. 2.1]):

Newton-Verfahren konvergiert lokal quadratisch, falls $DF(\mathbf{x}^*)$ regulär

Beispiel für theoretische Aussage über Newton-Verfahren:

Theorem 2.6.1 (Thm. 4.10 in [12]). Sei $D \subset \mathbb{R}^n$ offen und konvex und $F : D \mapsto \mathbb{R}^n$ stetig differenzierbar mit invertierbarer Jacobi-Matrix $DF(x) \forall x \in D$.

$$\exists \omega \geq 0: \quad \|DF(x)^{-1}(DF(x + \sigma v) - DF(x))v\| \leq \sigma \omega \|v\|^2 \quad \forall \sigma \in [0, 1], x \in D, v \in \mathbb{R}^n,$$

$$\exists x^*, x^{(0)} \in D: \quad F(x^*) = 0 \quad \wedge \quad \rho := \|x^* - x^{(0)}\| < \frac{2}{\omega} \quad \wedge \quad B_\rho(x^*) \subset D.$$

Dann erfüllt die Newton-Folge aus (2.5.1):

$$1. x^{(k)} \in B_\rho(x^*) := \{y \in \mathbb{R}^n, \|y - x^*\| < \rho\} \text{ für alle } k \in \mathbb{N},$$

$$2. \lim_{k \rightarrow \infty} x^{(k)} = x^*,$$

$$3. \|x^{(k+1)} - x^*\| \leq \frac{\omega}{2} \|x^{(k)} - x^*\|^2.$$

Problem: Normalerweise sind weder ω noch x^* bekannt!

A priori Abschätzungen (wie in Thm. 2.6.1) von geringem praktischen Nutzen

ABBRUCHKRIEITNEWTN

Abbruchkriterien für Newton-Verfahren:

Asymptotisch wegen quadratischer Konvergenz:

$$\|x^{(k+1)} - x^*\| \ll \|x^{(k)} - x^*\| \Rightarrow \|x^{(k)} - x^*\| \approx \|x^{(k+1)} - x^{(k)}\|.$$

► Beende Iteration, wenn $\|x^{(k+1)} - x^{(k)}\| = \|DF(x^{(k)})^{-1}F(x^{(k)})\| < \tau \|x^{(k)}\|$, τ = Toleranz
→ Überflüssiger (letzter) Schritt!

Beachte: für $n \gg 1$ kann jeder Newtonschritt mit grossem Rechenaufwand verbunden sein.

Alternativ: $DF(x^{(k-1)}) \approx DF(x^{(k)})$ ► Beende Iteration, wenn $\|DF(x^{(k-1)})^{-1}F(x^{(k)})\| < \tau$

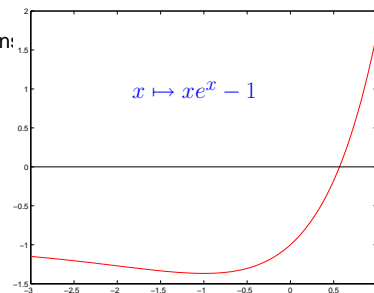
Terminologie: $DF(x^{(k-1)})^{-1}F(x^{(k)})$ heisst vereinfachte Newton-Korrektur
(Billig verfügbar, falls LU-Zerlegung (→ Abschn. 3.2.3) von $DF(x^{(k-1)})$ bereits berechnet)

Beispiel 50 (Lokale Konvergenz des Newton-Verfahrens:
Was heisst lokale Konvergenz? Beispiel:

$$F(x) = xe^x - 1 \Rightarrow F'(-1) = 0$$

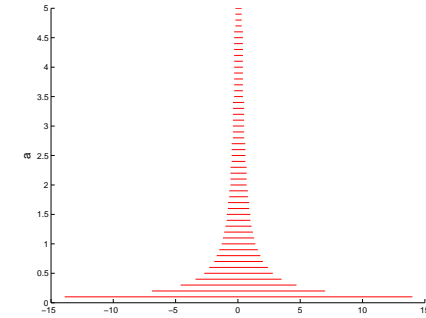
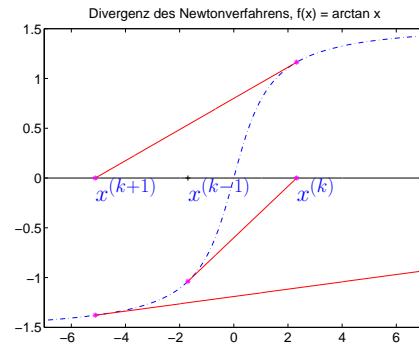
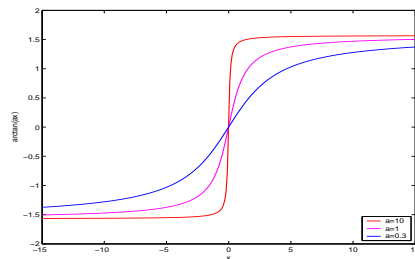
$$\begin{aligned} x^{(0)} < -1 &\Rightarrow x^{(k)} \rightarrow -\infty \\ x^{(0)} > -1 &\Rightarrow x^{(k)} \rightarrow x^* \end{aligned}$$

◇



Beispiel 51 (Konvergenzbereich des Newton-Verfahrens).

$$\begin{aligned} F(x) &= \arctan(ax), \quad a > 0, x \in \mathbb{R} \\ &\Rightarrow x^* = 0. \end{aligned}$$



roter Bereich = $\{x^{(0)} \in \mathbb{R}, x^{(k)} \rightarrow 0\}$

Zusammenfassung: Das Newton-Verfahren



konvergiert *asymptotisch* sehr schnell: Verdopplung Zahl gültiger Stellen in jedem Schritt



hat u.U. nur sehr kleinen (lokalen) Konvergenzbereich

2.6
p. 105

2.6.2 Gedämpftes Newton-Verfahren

Beispiel 51: Problem ist „Überschiessen“ der Newton-Korrektur



Idee:

DÄMPFUNG

Dämpfung der Korrektur:

$$\text{Mit } \lambda^{(k)} > 0: x^{(k+1)} := x^{(k)} - \lambda^{(k)} DF(x^{(k)})^{-1}F(x^{(k)}). \quad (2.6.1) \quad \text{eq:dam}$$

Terminologie: $\lambda^{(k)}$ = Dämpfungsfaktor

Wahl des Dämpfungsfaktors: Affinvarianter natürlicher Monotonietest [10, Ch. 3]

$$\text{„grösstmögliches“ } \lambda^{(k)} > 0: \|\bar{s}(\lambda^{(k)})\| \leq (1 - \frac{\lambda^{(k)}}{2}) \|s\| \quad (2.6.2) \quad \text{eq:mt}$$

mit $s := DF(x^{(k)})^{-1}F(x^{(k)})$ → aktuelle Newton-Korrektur,
 $\bar{s}(\lambda^{(k)}) := DF(x^{(k)})^{-1}F(x^{(k)}) + \lambda^{(k)}s$ → versuchsweise vereinfachte Newton-Korrektur.

Heuristik: Konvergenz \Leftrightarrow Abnahme der „Länge“ der Newton-Korrektur

MATLAB-CODE : Gedämpftes Newton-Verfahren

```
function [x,cvg] = dampnewton(x,F,DF,q,tol)
A = DF(x); s = A\F(x);
xn = x-s;
lambda = 1; cvg = 0;
while (norm(s) > tol)
```

2.6
p. 106

2.6
p. 107

2.6
p. 108

```

f = F(xn); st = A\f;
while (norm(st) > (1-lambda/2)*norm(s))
    lambda = q*lambda;
    if (lambda < LMIN), cvg = -1; return; end
    xn = x-lambda*s;
    f = F(xn); st = A\f;
end
x = xn;
A = DF(x); s = A\f;
lambda = min(lambda/q,1);
xn = x-lambda*s;
end
x = xn;

```

Strategie: Reduziere Dämpfungsparameter um Faktor $q \in]0, 1[$ (üblich $q = \frac{1}{2}$) solange bis affinvarianter natürlicher Monotonietest (2.6.2) erfüllt

Beispiel 52 (Gedämpftes Newton-Verfahren). (→ Beispiel 51)

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

- $x^{(0)} = 20$
- $q = \frac{1}{2}$
- $LMIN = 0.001$

Beobachtung: asymptotische quadratische Konvergenz

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	-0.00000000000001	-0.00000000000001

Beispiel 53 (Versagen des gedämpften Newton-Verfahrens).

- Wie in Bsp. 50: $F(x) = xe^x - 1$,
- Startwert für gedämpftes Newton-Verfahren $x^{(0)} = -1.5$

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

Abbruch wegen $\lambda < LMIN$!

Beobachtung: Newton-Korrektur hat „falsche Richtung“ ➤ Keine Konvergenz

2.6.3 Quasi-Newton-Verfahren

QUASINewTON

Was tun, wenn $DF(x)$ nicht verfügbar und numerisches Differenzieren (→ Bem. 49) zu aufwändig?

Idee: in einer Dimension ($n = 1$, → Sekantenverfahren (2.3.4), Abschn. 2.3.2.3)



$$F'(x^{(k)}) \approx \frac{F(x^{(k)}) - F(x^{(k-1)})}{x^{(k)} - x^{(k-1)}} \quad \text{„Differenzenquotient“} \quad (2.6.3)$$



Verallgemeinerung für $n > 1$?

Sekantenbedingung für Approximation $J_k \approx DF(x^{(k)})$, $x^{(k)} \triangleq$ Iterierte:

$$J_k(x^{(k)} - x^{(k-1)}) = F(x^{(k)}) - F(x^{(k-1)}) \quad (2.6.4)$$

Klar:

Viele J_k erfüllen (2.6.4)

➤ Nötig: Weitere Bedingungen an $J_k \in \mathbb{R}^{n,n}$

2.6
p. 109

Idee: J_k durch Modifikation von J_{k-1}

Nebenbedingung (Broyden): $J_k z = J_{k-1} z \quad \forall z: z \cdot (x^{(k)} - x^{(k-1)}) = 0$ (2.6.5)



$$J_k = J_{k-1} + \frac{F(x^{(k)}) - F(x^{(k-1)})}{\|x^{(k)} - x^{(k-1)}\|_2} \quad (2.6.6)$$

RANG-1-Modifikation, Abschn. 3.2.8.1

BRPYON

Broydens Quasi-Newton-Verfahren zur Lösung von $F(x) = 0$:

$$x^{(k+1)} := x^{(k)} + \Delta x^{(k)}, \Delta x^{(k)} := -J_k^{-1} F(x^{(k)}), J_{k+1} := J_k + \frac{F(x^{(k+1)}) - F(x^{(k)})}{\|\Delta x^{(k)}\|_2^2} \quad (2.6.7)$$

Startwert $J_0 = ?$ → Verwende etwa exakte Jacobimatrix $DF(x^{(0)})$

2.6
p. 110

2.6
p. 111

2.6
p. 112

Bemerkung 54 (Minimalität von Broydens Rang-1-Modifikation).

$$\mathbf{J} \in \mathbb{R}^{n,n} \text{ erfüllt (2.6.4)} \quad \mathbf{J}_k, \mathbf{x}^{(k)} \text{ aus (2.6.7)} \quad \Rightarrow \quad (\mathbf{I} - \mathbf{J}_k^{-1} \mathbf{J})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = -\mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)})$$

$$\|\mathbf{I} - \mathbf{J}_k^{-1} \mathbf{J}_{k+1}\|_2 = \left\| \frac{-\mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)}) \Delta \mathbf{x}^{(k)}}{\|\Delta \mathbf{x}^{(k)}\|_2^2} \right\|_2 = \left\| (\mathbf{I} - \mathbf{J}_k^{-1} \mathbf{J}) \frac{\Delta \mathbf{x}^{(k)} (\Delta \mathbf{x}^{(k)})^T}{\|\Delta \mathbf{x}^{(k)}\|_2^2} \right\|_2$$

$$\leq \|\mathbf{I} - \mathbf{J}_k^{-1} \mathbf{J}\|_2.$$

\Rightarrow (2.6.6) = $\|\cdot\|_2$ -minimale relative Korrektur von \mathbf{J}_{k-1} , die Sekantenbedingung (2.6.4) respektiert.

EXQUASINEWON

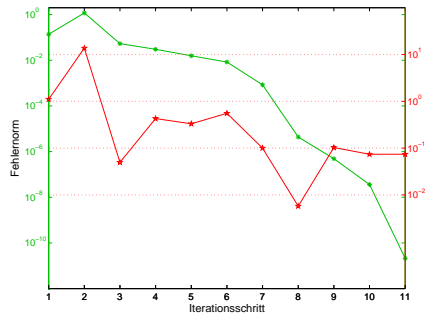
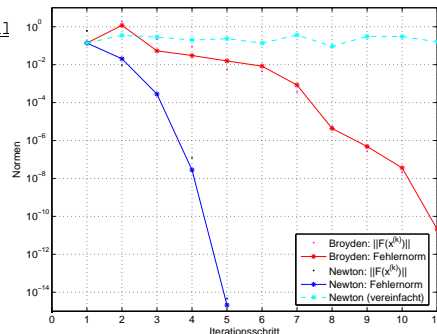
Beispiel 55 (Broydens Quasi-Newton-Verfahren: Konvergenz).

- Nullstellenbestimmung aus Bsp. 46, $n=2$
- $\mathbf{x}^{(0)} = (0.7, 0.7)^T$
- $\mathbf{J}_0 = DF(\mathbf{x}^{(0)})$

Langsamer als Newton-Verfahren
(2.5.1)

Besser als

vereinfachtes Newton-Verfahren (\rightarrow Bem. 48)



Konvergenzmonitor

Grösse, die Probleme mit der Konvergenz einer Iteration anzeigt

Hier:

$$\mu := \frac{\|\mathbf{J}_{k-1}^{-1} F(\mathbf{x}^{(k)})\|}{\|\Delta \mathbf{x}^{(k-1)}\|}$$

Heuristik: $\mu > 1 \Rightarrow$ keine Konvergenz

Bemerkung 56. Option: Gedämpftes Broyden-Verfahren (wie bei Newton-Verfahren, Abschn. 2.6.2)

Implementierung von (2.6.7): Mit Sherman-Morrison-Woodbury Update-Formel, Lemma 96

$$\mathbf{J}_{k+1}^{-1} = \left(\mathbf{I} - \frac{\mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)}) (\Delta \mathbf{x}^{(k)})^T}{\|\Delta \mathbf{x}^{(k)}\|_2^2 + \Delta \mathbf{x}^{(k)} \cdot \mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)})} \right) \mathbf{J}_k^{-1} = \left(\mathbf{I} + \frac{\Delta \mathbf{x}^{(k+1)} (\Delta \mathbf{x}^{(k)})^T}{\|\Delta \mathbf{x}^{(k)}\|_2^2} \right) \mathbf{J}_k^{-1}$$

(2.6.8)

Sinnvoll, falls $\|\mathbf{J}_k^{-1} F(\mathbf{x}^{(k+1)})\|_2 < \|\Delta \mathbf{x}^{(k)}\|_2$ „vereinfachte Quasi-Newton-Korrektur“

MATLAB-CODE: Broyden-Verfahren (2.6.7)

```
function x = broyden(F,x,J,tol)
k = 1;
[L,U] = lu(J); % Einmalige LU-Zerlegung !
s = U \ (L \ F(x)); sn = dot(s,s);
dx = [s]; dxn = [sn]; % Speicherung der Δx^(k), ||Δx^(k)||_2^2 (→ (2.6.8))
x = x - s; f = F(x);

while (sqrt(sn) > tol), k=k+1
    w = U \ (L \ f); % Löse gestaffelte LGS
    for l=2:k-1 % Abbruch, vgl. Abschn. 2.6.1
        w = w + dx(:,l) * (dx(:,l-1)' * w) ... % Aufbau von w := J_k^{-1} F(x^(k)) (→ Rekursion (2.6.8))
            / dxn(l-1);
    end
    if (norm(w) >= sn) % Konvergenzmonitor
        warning('Dubious step %d!',k);
    end
    z = s' * w; s = (1 + z / (sn - z)) * w; sn = s' * s; % Korrektur s = J_k^{-1} F(x^(k))
    dx = [dx,s]; dxn = [dxn,sn];
    x = x - s; f = F(x);
end
```

Rechenaufwand : $O(N^2 \cdot n)$ Operationen Vektorarithmetik, (Level-I)
 N Schritte
 1 LU-Zerlegung von J , $N \times$ Lösung gestaffeltes LGS \rightarrow Abschn. 3.2.3
 N F -Auswertungen !

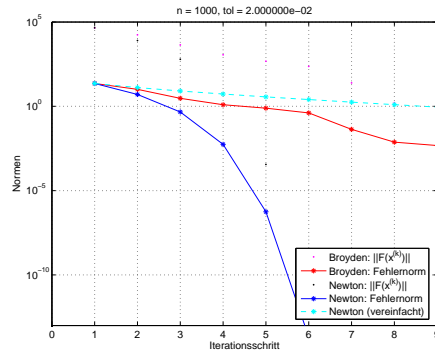
Speicheraufwand : LU-Faktoren von J + Hilfsvektoren $\in \mathbb{R}^n$
 N Schritte
 N Vektoren $\mathbf{x}^{(k)} \in \mathbb{R}^n$

Beispiel 57 (Broyden-Verfahren für grosses nichtlineares Gleichungssystem).

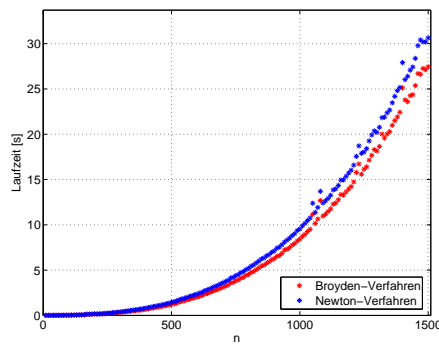
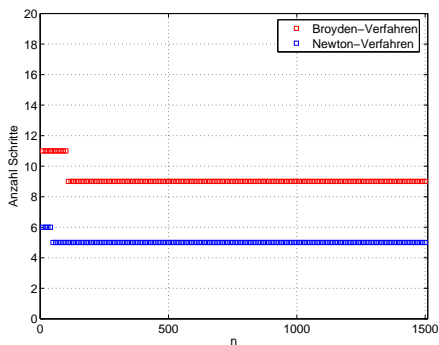
$$F(\mathbf{x}) = \begin{cases} \mathbb{R}^n \mapsto \mathbb{R}^n \\ \mathbf{x} \mapsto \text{diag}(\mathbf{x})\mathbf{A}\mathbf{x} - \mathbf{b}, \\ \mathbf{b} = (1, 2, \dots, n) \in \mathbb{R}^n, \\ \mathbf{A} = \mathbf{I} + \mathbf{a}\mathbf{a}^T \in \mathbb{R}^{n,n}, \\ \mathbf{a} = \frac{1}{\sqrt{1 \cdot \mathbf{b} - 1}}(\mathbf{b} - \mathbf{1}). \end{cases}$$

Auswertung wie in Bsp. 55

$h = 2/n$; $\mathbf{x}_0 = (2:h:4-h)'$;



Vergleich der Effizienzen: Broyden-Verfahren \longleftrightarrow Newton-Verfahren:
 (Problemgrösse $n \gg$ Toleranz $\text{tol} = 2n \cdot 10^{-5}$, $h = 2/n$; $\mathbf{x}_0 = (2:h:4-h)'$;)



Broyden-Verfahren konkurrenzfähig für $n \gg 1$ und geringe Genauigkeitsanforderungen

2.7 Nichtlineare Ausgleichsrechnung

[File: section-nichtlineare-ausgleichsrechnung.tex, SVN: section-nichtlineare-ausgleichsrechnung.tex 1240 2006-12-30 13:32:59Z hiptmair]

Gegeben: $F : D \subset \mathbb{R}^n \mapsto \mathbb{R}^m$, $m, n \in \mathbb{N}$, $m > n$

➤ Erweiterte Aufgabenstellung: **NICHTLINEARISCHES** **Nichtlineares Ausgleichsproblem**

$$\text{Suche } \mathbf{x}^* \in D: \mathbf{x}^* = \underset{\mathbf{x} \in D}{\text{argmin}} \Phi(\mathbf{x}), \quad \Phi(\mathbf{x}) := \frac{1}{2} \|F(\mathbf{x})\|_2^2. \quad (2.7.1) \quad \text{falsch}$$

(Existenz/Eindeutigkeit von \mathbf{x}^* jeweils für jeden Fall konkret nachzuprüfen !)

Terminologie: $D \hat{=}$ Parameterbereich, $x_1, \dots, x_n \hat{=}$ Parameter

Wir verlangen „Abhängigkeit von allen Parametern“

$$\exists \text{ Umgebung } \mathcal{U}(\mathbf{x}^*): DF(\mathbf{x}) \text{ hat vollen Rang } n \quad \forall \mathbf{x} \in \mathcal{U}(\mathbf{x}^*).$$

(bedeutet: Spalten der Jacobimatrix $DF(\mathbf{x})$ sind linear unabhängig.)

2.6
p. 117

1. Möglichkeit: (gedämpftes) **Newton-Verfahren**

$$\Phi(\mathbf{x}^*) = \min \Rightarrow \text{grad } \Phi(\mathbf{x}) = 0, \quad \text{grad } \Phi(\mathbf{x}) := \left(\frac{\partial \Phi}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial \Phi}{\partial x_n}(\mathbf{x}) \right)^T \in \mathbb{R}^n.$$

Idee: Newton-Verfahren zur Bestimmung einer Nullstelle von $\text{grad } \Phi : D \mapsto \mathbb{R}^n$

$$\text{Kettenregel} \Rightarrow \text{grad } \Phi(\mathbf{x}) = DF(\mathbf{x})^T F(\mathbf{x}),$$

$$\text{Produktregel} \Rightarrow H\Phi(\mathbf{x}) := D(\text{grad } \Phi)(\mathbf{x}) = DF(\mathbf{x})^T DF(\mathbf{x}) + \sum_{j=1}^m F_j(\mathbf{x}) D^2 F_j(\mathbf{x}).$$

($H\Phi(\mathbf{x}) =$ **HESSE** **Hesse-Matrix**)

➤ Lineares Gleichungssystem zur Bestimmung der Newton-Korrektur $\mathbf{s} \in \mathbb{R}^n$ zu $\mathbf{x}^{(k)}$:

$$\left(DF(\mathbf{x}^{(k)})^T DF(\mathbf{x}^{(k)}) + \sum_{j=1}^m F_j(\mathbf{x}^{(k)}) D^2 F_j(\mathbf{x}^{(k)}) \right) \mathbf{s} = -DF(\mathbf{x}^{(k)})^T F(\mathbf{x}^{(k)}). \quad (2.7.2) \quad \text{eq:NCNLSQ}$$

Newton-Verfahren = Minimierung einer lokalen **quadratischen Approximation** von Φ :

$$\Phi(\mathbf{x}) \approx Q(\mathbf{s}) := \Phi(\mathbf{x}^{(k)}) + \text{grad } \Phi(\mathbf{x}^{(k)})^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T H\Phi(\mathbf{x}^{(k)}) \mathbf{s}.$$

$$\text{grad } Q(\mathbf{s}) = 0 \Leftrightarrow H\Phi(\mathbf{x}^{(k)}) \mathbf{s} + \text{grad } \Phi(\mathbf{x}^{(k)}) = 0 \Leftrightarrow \text{eq:NCNLSQ} \quad (2.7.2).$$

2.7
p. 118

2.7
p. 119

2.7
p. 120

Bemerkung 58 (Differenzieren). → Wiederholung: Analysis Grundvorlesung

Seien V, W endlichdimensionale Vektorräume, $F : D \subset V \mapsto W$ hinreichend glatt. **Ableitung** $DF(\mathbf{x})$ von F in $\mathbf{x} \in V$

$$DF(\mathbf{x}) : V \mapsto W \text{ linear, } \|F(\mathbf{x} + \mathbf{h}) - F(\mathbf{x}) - DF(\mathbf{x})\mathbf{h}\| = o(\|\mathbf{h}\|) \quad \forall \mathbf{h}, \|\mathbf{h}\| < \delta.$$

• $F : V \mapsto W$ linear, d.h. $F(\mathbf{x}) = \mathbf{A}\mathbf{x}$, \mathbf{A} Matrix $\Rightarrow DF(\mathbf{x}) = \mathbf{A}$.

• **Kettenregel**: $F : V \mapsto W, G : W \mapsto U$ hinreichend glatt

$$D(G \circ F)(\mathbf{x})\mathbf{h} = DG(F(\mathbf{x}))(DF(\mathbf{x})\mathbf{h}), \quad \mathbf{h} \in V, \mathbf{x} \in D.$$

PRODUKTREGEL

• **Produktregel**: $F : D \subset V \mapsto W, G : D \subset V \mapsto U$ hinreichend glatt, $b : W \times U \mapsto Z$ bilinear

$$T(\mathbf{x}) = b(F(\mathbf{x}), G(\mathbf{x})) \Rightarrow DT(\mathbf{x})\mathbf{h} = b(DF(\mathbf{x})\mathbf{h}, G(\mathbf{x})) + b(F(\mathbf{x}), DG(\mathbf{x})\mathbf{h}), \quad \mathbf{h} \in V, \mathbf{x} \in D.$$

Für $F : D \subset V \mapsto \mathbb{R}$ sind der **Gradient** $\text{grad } F : D \mapsto \mathbb{R}^n, n := \dim V$, und die **Hessematrix** $HF(\mathbf{x}) : D \mapsto \mathbb{R}^{n,n}$ definiert durch

$$\text{grad } F(\mathbf{x})^T \mathbf{h} := DF(\mathbf{x})\mathbf{h}, \quad \mathbf{h}_1^T HF(\mathbf{x})\mathbf{h}_2 := D(DF(\mathbf{x})(\mathbf{h}_1))(\mathbf{h}_2), \quad \mathbf{h}, \mathbf{h}_1, \mathbf{h}_2 \in V.$$

△

2. Möglichkeit: **GAUSS-NEWTON**
Gauss-Newton-Verfahren

Idee: Lokale Linearisierung von F (anstatt lokal quadratischer Approximation von Φ)

Korrektur \mathbf{s} zur aktuellen Iterierten $\mathbf{x}^{(k)}$ so dass

$$\|F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})\mathbf{s}\|_2 \rightarrow \min \quad (\text{Lineares Ausgleichsproblem} \rightarrow \text{Abschnitt 3.5}) \quad \text{sec: numer-line-aus (2.7.3) GNS}$$

2.7
p. 122

Der Operator \backslash : die Wunderwaffe in MATLAB:

Für $\mathbf{A} \in \mathbb{R}^{m,n}$

$$\mathbf{x} = \mathbf{A} \backslash \mathbf{b}$$



\mathbf{x} Minimierer von $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$
mit minimaler 2-Norm,

Vorteil des Gauss-Newton-Verfahrens : Keine zweiten Ableitungen von F erforderlich

Nachteil des Gauss-Newton-Verfahrens : Keine lokal quadratische Konvergenz

MATLAB-Template für Gauss-Newton-Verfahren

```
function x = gn(x,f,J,tol)
s = J(x)\f(x);
x = x-s;
while (norm(s) > tol)
s = J(x)\f(x);
x = x-s;
end
```

3. Möglichkeit:

TRUSTREGION

Trust-Region-Verfahren (Levenberg-Marquardt-Verfahren)

Idee: Dämpfung der Gauss-Newton-Korrektur aus (2.7.3) durch **Strafterm** (engl. *penalty term*)

2.7
p. 121

$$\text{Statt } \|F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})\mathbf{s}\|_2^2 \text{ minimiere } \|F(\mathbf{x}^{(k)}) + DF(\mathbf{x}^{(k)})\mathbf{s}\|_2^2 + \lambda \|\mathbf{s}\|_2^2.$$

2.7
p. 123

$\lambda > 0 \hat{=}$ Strafparameter (wie zu wählen ? → Heuristik)

$$\lambda = \gamma \|F(\mathbf{x}^{(k)})\|_2, \quad \gamma := \begin{cases} 10 & , \text{ falls } \|F(\mathbf{x}^{(k)})\|_2 \geq 10, \\ 1 & , \text{ falls } 1 < \|F(\mathbf{x}^{(k)})\|_2 < 10, \\ 0.01 & , \text{ falls } \|F(\mathbf{x}^{(k)})\|_2 \leq 1. \end{cases}$$

► Modifizierte (regularisierte) Gleichung für Korrektur \mathbf{s} :

$$(DF(\mathbf{x}^{(k)})^T DF(\mathbf{x}^{(k)}) + \lambda \mathbf{I}) \mathbf{s} = -DF(\mathbf{x}^{(k)})F(\mathbf{x}^{(k)}). \quad (2.7.4) \quad \text{CMS}$$

NICHTLINREG

Beispiel 59 (Nichtlineare Regression, parametrische Statistik).

Parametrische Statistik: Fitten von Messwerten $(t_i, y_i), i = 1, \dots, m$ durch nichtlinear parameter-abhängige Funktion $f : \mathbb{R} \mapsto \mathbb{R}$, speziell $f(t) = x_1 + x_2 \exp(-x_3 t)$.

$$F(\mathbf{x}) = \begin{pmatrix} x_1 + x_2 \exp(-x_3 t_1) - y_1 \\ \vdots \\ x_1 + x_2 \exp(-x_3 t_m) - y_m \end{pmatrix} : \mathbb{R}^3 \mapsto \mathbb{R}^m, \quad DF(\mathbf{x}) = \begin{pmatrix} 1 & e^{-x_3 t_1} & -x_2 t_1 e^{-x_3 t_1} \\ \vdots & \vdots & \vdots \\ 1 & e^{-x_3 t_m} & -x_2 t_m e^{-x_3 t_m} \end{pmatrix}$$

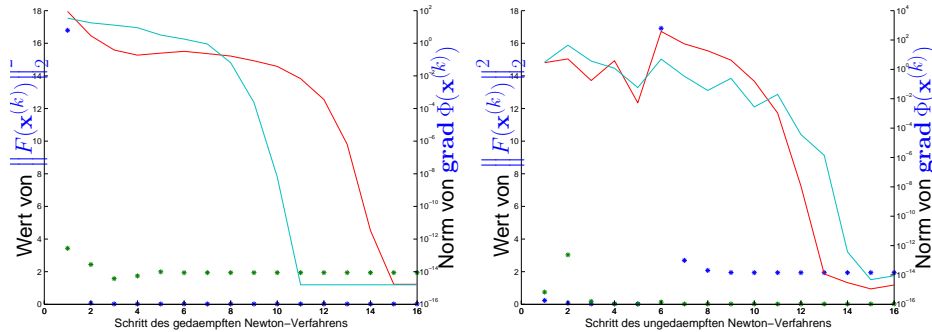
Numerisches Experiment:

Konvergenz des Newton-Verfahrens,
gedämpften Newton-Verfahrens
(→ Abschnitt 2.6.2) und Gauss-Newton-
Verfahrens für verschiedene Startwerte

2.7
p. 122

```
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);
```

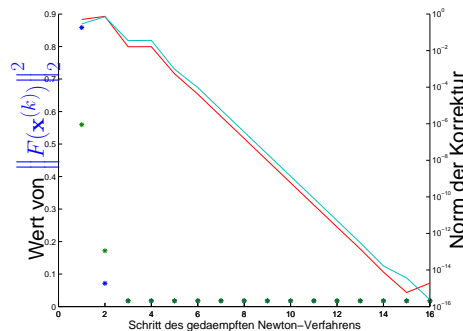
2.7
p. 124



Konvergenzverhalten des Newton-Verfahrens:

Anfangswert $(1.8, 1.8, 0.1)^T$ (rote Kurve) ➤ Newton-Verfahren gefangen in **lokalem Minimum**
 Anfangswert $(1.5, 1.5, 0.1)^T$ (cyan Kurve) ➤ Schnelle (lokal quadratische) Konvergenz

Gauss-Newton-Verfahren
 Anfangswert $(1.8, 1.8, 0.1)^T$ (rote Kurve)
 Anfangswert $(1.5, 1.5, 0.1)^T$ (cyan Kurve)
 Konvergenz in beiden Fällen
 Beobachtung: **Lineare Konvergenz**



3

[File: chapter-numerische-lineare-algebra.tex, SVN: chapter-numerische-lineare-algebra.tex 1055 2006-10-18 15:44:34Z kalai]

Numerische lineare Algebra

- ☛ Klasse von Problemen, beschreibbar durch elementare Operationen mit Matrizen und/oder Vektoren im \mathbb{R}^n oder \mathbb{C}^n , z.B.
 - Lineare Gleichungssysteme (engl. *linear systems of equations*)
 - Lineare Ausgleichsprobleme (engl. *linear least squares*)
 - Eigenwertprobleme (engl. *eigenproblems*)

Algorithmen/Konzepte der numerischen linearen Algebra sind Teil/Grundlage fast aller numerischer Methoden

2.7
p. 125

3.1 Grundbegriffe und -operationen

[File: section-grundbegriffe-und-operationen.tex, SVN: section-grundbegriffe-und-operationen.tex 1010 2006-09-11 15:44:43Z hiptmair]

Konventionen:

- Alle Vektoren sind **Spaltenvektoren** (engl. *column vectors*) [Zeilenvektor = \mathbf{x}^T],
 MATLAB Initialisierung: $\mathbf{x} = [1; 2];$ für Zeilenvektor $\mathbf{x}^t = [1, 2];$
 Notation: **fette** Kleinbuchstaben
 Komponenten $\mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{K}^n$, $\mathbb{K} = \mathbb{R}, \mathbb{C}$, bzw. $x_i = (\mathbf{x})_i$, $i = 1, \dots, n$.
- **Fette** Grossbuchstaben für Matrizen, z.B. **A**, Matriceinträge Kleinbuchstaben mit subscript Indizes, z.B. a_{ij} oder $(\mathbf{A})_{ij}$ ($i \rightarrow$ Zeilenindex, $j \rightarrow$ Spaltenindex)
- Für Vektor $\mathbf{x} \in \mathbb{R}^n / \mathbf{x} \in \mathbb{C}^n$, $n \in \mathbb{N}$:

$$\mathbf{x}^T = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}^T := (x_1 \cdots x_n), \quad \mathbf{x}^H = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}^H := (\bar{x}_1 \cdots \bar{x}_n), \quad \bar{} = \text{komplexe Konjugation}.$$

2.7
p. 126

3.1
p. 127

3.1
p. 128

- Für Matrix $A \in \mathbb{K}^{n,m} \rightarrow$ Transponierte Matrix (engl. *transposed matrix*)

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

- Spalten (engl. *columns*) einer Matrix $A \in \mathbb{K}^{m,n}$ | Zeilen (engl. *rows/lines*) einer Matrix $A \in \mathbb{K}^{m,n}$

$$A = [a_{\cdot,1} \dots a_{\cdot,n}]$$

MATLAB: `A(:,1) ... A(:,n)`

$$A = \begin{bmatrix} a_{1,\cdot}^T \\ \vdots \\ a_{n,\cdot}^T \end{bmatrix}$$

MATLAB: `A(1,:) ... A(m,:)`

3.1.1 Operationen

KOMPLEXITÄT

Komplexität \leftrightarrow Aufwand für Algorithmus in Abhängigkeit von Problemgrößenparameter n

Notation („Landau-O“): $f(n) = O(g(n)) \Leftrightarrow \exists C > 0, N > 0: |f(n)| \leq Cg(n)$ für alle $n > N$.

- LEVEL I
- (Euklidisches) Skalarprodukt: $x \cdot y = x^H y \in \mathbb{K}, x, y \in \mathbb{K}^n, \mathbb{K} = \mathbb{R}, \mathbb{C}$.
 - MATLAB operator `dot(x,y)`, Rechenaufwand = $n-1$ Additionen, n Multiplikationen.
 - Komplexität $O(n)$ (minimal)

- LEVEL I
- SAXPY: $z = \alpha x + y, x, y \in \mathbb{K}^n, \alpha \in \mathbb{K}$,
 - MATLAB `z=alpha*x+y`, Rechenaufwand = n Additionen, Multiplikationen.
 - Komplexität $O(n)$ (minimal)

- LEVEL II
- Vektor \times Matrix-Produkt $y = Ax, A \in \mathbb{K}^{n,m}$
 - MATLAB `y=A*x`, Rechenaufwand = m Skalarprodukte.
 - Komplexität $O(mn)$ (minimal)

- LEVEL III
- Matrix \times Matrix-Produkt $C = AB, A \in \mathbb{K}^{m,n}, B \in \mathbb{K}^{n,k}$
 - MATLAB `C=A*B`, Rechenaufwand = k -mal Matrix $A \times$ Vektor.
 - Komplexität $O(kmn)$ (nicht minimal !)

Bemerkung 60. Effiziente maschinennahe Implementierung der Grundoperationen:

BLAS
BLAS-Bibliothek (www.netlib.org/blas/) \rightarrow Vorlesung „Programmiertechniken“

△

Bemerkung 61. Blockweise Matrixmultiplikation: $A_{11} \in \mathbb{K}^{n,n}, A_{12} \in \mathbb{K}^{n,m}, A_{21} \in \mathbb{K}^{m,n}, A_{22} \in \mathbb{K}^{m,m}, B_{11} \in \mathbb{K}^{n,n}, B_{12} \in \mathbb{K}^{n,m}, B_{21} \in \mathbb{K}^{m,n}, B_{22} \in \mathbb{K}^{m,m}, m, n \in \mathbb{N}$:

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix} = \begin{pmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{pmatrix}. \quad (3.1.1) \quad \text{Eq:Block}$$

△

3.1.2 Matrix-Speicherformate

MATSPFORMAT

Zwei Hauptklassen von Matrizen:

VOLLBESMAT

Vollbesetzte (engl. *dense*) Matrizen



Dünnbesetzte (engl. *sparse*) Matrizen

Definition 3.1.1 (Dünnbesetzte Matrix). $A \in \mathbb{K}^{m,n}, m, n \in \mathbb{N}$, heisst dünnbesetzt, falls

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

3.1
p. 129

3.1
p. 131

- Speicherung vollbesetzter Matrizen („Standard“):

$A \in \mathbb{K}^{m,n} \rightarrow$ Lineares Array (Grösse mn) + Indexarithmetik
(Beachte führende Dimension (engl. *row major, column major*))

- Speicherung dünnbesetzter Matrizen, Beispiel CRS **Compressed-Row-Storage (CRS)**-Format:

Matrixinformation für $A \in \mathbb{K}^{n,n}$ wird in drei Arrays gespeichert:

```
real val      Grösse nnz(A) := # {(i,j) in {1,...,n}^2, a_ij != 0}
int col_ind   Grösse nnz(A)
int row_ptr   Grösse n+1 & row_ptr[N+1] = nnz(A) + 1
```

Zugriff auf Matrixeintrag $a_{ij} \neq 0, 1 \leq i, j \leq n$:

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

3.1
p. 130

3.1
p. 132

Beispiel 62 (CRS-Matrixspeicherformat).

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

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						

Variante: Diagonales CRS-Format (Matrixdiagonale wird separat gespeichert)

Dünnbesetzte Matrizen in MATLAB [18]:

Initialisierungsfunktionen: `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);`

Beispiel 63 (Effiziente Initialisierung von Sparse-Matrizen).

① : Initialisierung dünnbesetzter Matrizen I

```

A1 = sparse(n,n);
for i=1:n
    for j=1:n
        if (abs(i-j) == 1), A1(i,j) = A1(i,j) + 1; end;
        if (abs(i-j) == round(n/3)), A1(i,j) = A1(i,j) - 1; end;
    end; end

```

Direkte Initialisierung

② : Initialisierung dünnbesetzter Matrizen II

```

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

```

Akkumulation in dynamischem Vektor

③ : Initialisierung dünnbesetzter Matrizen III

```

dat = zeros(6*n,3); k = 0;
for i=1:n
    for j=1:n
        if (abs(i-j) == 1), k=k+1; dat(k,:) = [i,j,1.0]; end;
        if (abs(i-j) == round(n/3))
            k=k+1; dat(k,:) = [i,j,-1.0]; end;
    end; end;
A3 = sparse(dat(1:k,1),dat(1:k,2),dat(1:k,3),n,n);

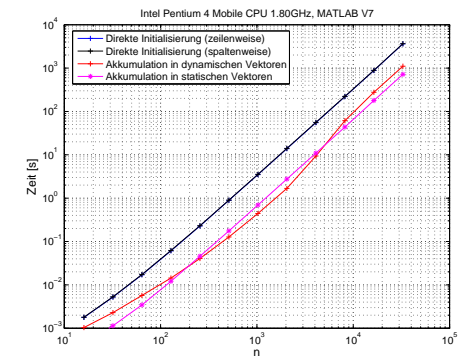
```

Akkumulation in statischem Vektor

Zeitmessung:

(MATLAB `tic,toc`-Funktionen)

(Mobile Intel Pentium 4 - M CPU 2.40GHz,
Linux MATLAB V6.5)



Beachte: Vollbesetzt * dünnbesetzt → vollbesetzt

dünnbesetzt +, -, * dünnbesetzt → dünnbesetzt

`S = [A,B;C,D]` dünnbesetzt, wenn ein Block dünnbesetzt

Beispiel 64 (Effiziente MATLAB-Implementierungen mit sparse-Matrizen).

Initialisierung:

`A = rand(n,n); v = rand(n,1);`

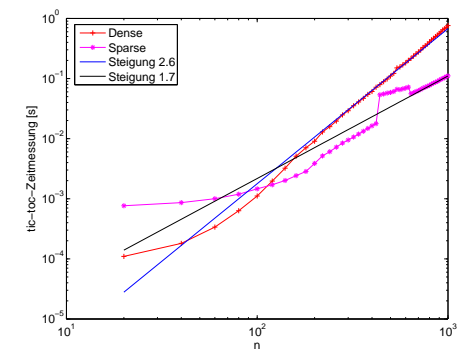
Aufgabe: Skalierung einer Matrix:

• DENSE-Matrixarithmetik:

`S = A*diag(v)+eye(n)`

• SPARSE-Matrixarithmetik:

`S = A*spdiags(v,0,n,n)+speye(n)`



3.1.3 Erinnerung: Normen

Begriffe **Kondition** (\rightarrow Abschn. 1.6), **Stabilität** (\rightarrow Abschn. 1.8), **Fehler** \rightarrow erfordern Metrik auf \mathbb{K}^n

Grundvorlesungen: Norm $\|\cdot\|$ auf $\mathbb{K}^n \Rightarrow$ Metrik auf \mathbb{K}^n

Definition. \rightarrow Def. 2.1.4
 $X = \mathbb{K}$ -Vektorraum, $\mathbb{K} = \mathbb{C}, \mathbb{R}$. Eine Abbildung $\|\cdot\| : X \rightarrow \mathbb{R}_0^+$ ist eine **Norm** auf X , falls

- (i) $\forall x \in X: x \neq 0 \Leftrightarrow \|x\| > 0$ (Definitheit),
- (ii) $\|\lambda x\| = |\lambda| \|x\| \quad \forall x \in X, \lambda \in \mathbb{K}$ (Homogenität),
- (iii) $\|x + y\| \leq \|x\| + \|y\| \quad \forall x, y \in X$ (Dreiecksungleichung).

Beispiele: (für Vektor $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{K}^n$)

Bezeichnung	:	Definition	MATLAB-Funktion
Euklidische Norm	:	$\ \mathbf{x}\ _2 := \sqrt{ x_1 ^2 + \dots + x_n ^2}$	<code>norm(x)</code>
1-Norm	:	$\ \mathbf{x}\ _1 := x_1 + \dots + x_n $	<code>norm(x, 1)</code>
Maximumnorm	:	$\ \mathbf{x}\ _\infty := \max\{ x_1 , \dots, x_n \}$	<code>norm(x, inf)</code>

Definition. \rightarrow Def. 2.2.6
Der Norm $\|\cdot\|$ auf \mathbb{R}^n zugeordnete **Matrixnorm**

$$\mathbf{M} \in \mathbb{R}^{n,n}: \|\mathbf{M}\| := \sup_{\mathbf{x} \in \mathbb{R}^n \setminus \{0\}} \frac{\|\mathbf{M}\mathbf{x}\|}{\|\mathbf{x}\|}.$$

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

Notation: $\|\mathbf{x}\|_2 \rightarrow \|\mathbf{M}\|_2, \|\mathbf{x}\|_1 \rightarrow \|\mathbf{M}\|_1, \|\mathbf{x}\|_\infty \rightarrow \|\mathbf{M}\|_\infty$

Beispiel 65 (Matrixnorm zur Maximumnorm und 1-Norm).

z.B. für $\mathbf{M} \in \mathbb{K}^{2,2}$:

$$\begin{aligned} \|\mathbf{M}\mathbf{x}\|_\infty &= \max\{|m_{11}x_1 + m_{12}x_2|, |m_{21}x_1 + m_{22}x_2|\} \\ &\leq \max\{|m_{11}| + |m_{12}|, |m_{21}| + |m_{22}|\} \|\mathbf{x}\|_\infty, \\ \|\mathbf{M}\mathbf{x}\|_1 &= |m_{11}x_1 + m_{12}x_2| + |m_{21}x_1 + m_{22}x_2| \\ &\leq \max\{|m_{11}| + |m_{21}|, |m_{12}| + |m_{22}|\} (|x_1| + |x_2|). \end{aligned}$$

\triangleright Matrixnorm zu $\|\cdot\|_\infty =$ **Zeilensummennorm** $\|\mathbf{M}\|_\infty := \max_{i=1,\dots,n} \sum_{j=1}^n |m_{ij}|, \quad (3.1.2)$

\triangleright Matrixnorm zu $\|\cdot\|_1 =$ **Spaltensummennorm** $\|\mathbf{M}\|_1 := \max_{j=1,\dots,m} \sum_{i=1}^n |m_{ij}|. \quad (3.1.3)$

3.2 Numerische Lösung linearer Gleichungssysteme

[File: section-numerische-loesung-linearer-gleichungssysteme.tex, SVN: section-numerische-loesung-linearer-gleichungssysteme.tex 1137 2006-11-13 09:27]

Gegeben: Matrix $\mathbf{A} \in \mathbb{K}^{n,n}$, Vektor $\mathbf{b} \in \mathbb{K}^n, n \in \mathbb{K}$

Gesucht: Lösungsvektor $\mathbf{x} \in \mathbb{K}^n$: $\mathbf{Ax} = \mathbf{b}$ \leftarrow **Lineares Gleichungssystem (LGS)**

Beispiel 66 (Knotenanalyse eines linearen elektrischen Netzwerks).

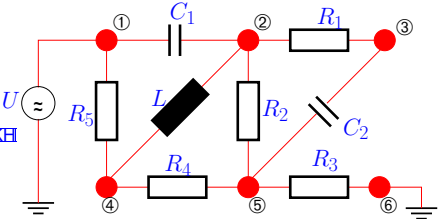
Numeriere Knoten $1, \dots, n$

I_{kj} : Strom Knoten $k \rightarrow$ Knoten $j, I_{kj} = -I_{jk}$

KIRCHHOFF
Kirchhoffsche Regel: Zweigstromsumme = 0

$$\forall k \in \{1, \dots, n\}: \sum_{j=1}^n I_{kj} = 0. \quad (3.2.1) \quad \text{KH}$$

Unbekannte: **Knotenpotentiale** $u_k, k = 1, \dots, n$.
(davon einige fixiert: Erdung, Quellen)



BAUELEMENTG
Bauelementgleichungen

(im Frequenzbereich mit Winkelfrequenz $\omega > 0$):

• **OHMSCHER** Widerstand: $I = R^{-1}U, [R] = 1\text{VA}^{-1}$

• **KONDENSATOR**: $I = i\omega CU$, Kapazität $[C] = 1\text{AsV}^{-1}$

• **SPULE**: $I = -i\omega^{-1}L^{-1}U$, Induktivität $[L] = 1\text{VsA}^{-1}$

Bauelementgleichungen + (3.2.1) \triangleright **Lineares Gleichungssystem**

$$\begin{aligned} \textcircled{2}: \quad & i\omega C_1(u_2 - u_1) + R_1^{-1}(u_2 - u_3) - i\omega^{-1}L^{-1}(u_2 - u_4) + R_2^{-1}(u_2 - u_5) = 0, \\ \textcircled{3}: \quad & R_1^{-1}(u_3 - u_2) + i\omega C_2(u_3 - u_5) = 0, \\ \textcircled{4}: \quad & R_5^{-1}(u_4 - u_1) - i\omega^{-1}L^{-1}(u_4 - u_2) + R_4^{-1}(u_4 - u_5) = 0, \\ \textcircled{5}: \quad & R_2^{-1}(u_5 - u_2) + i\omega C_2(u_5 - u_3) + R_4^{-1}(u_5 - u_4) + R_3(u_5 - u_6) = 0, \\ & u_1 = U, \quad u_6 = 0. \end{aligned}$$

$$\begin{pmatrix} i\omega C_1 + \frac{1}{R_1} - \frac{i}{\omega L} + \frac{1}{R_2} & -\frac{1}{R_1} & \frac{i}{\omega L} & -\frac{1}{R_2} \\ -\frac{1}{R_1} & \frac{1}{R_1} + i\omega C_2 & 0 & -i\omega C_2 \\ \frac{i}{\omega L} & 0 & \frac{1}{R_5} - \frac{i}{\omega L} + \frac{1}{R_4} & -\frac{1}{R_4} \\ -\frac{1}{R_2} & -i\omega C_2 & -\frac{1}{R_4} & \frac{1}{R_2} + i\omega C_2 + \frac{1}{R_4} \end{pmatrix} \begin{pmatrix} u_2 \\ u_3 \\ u_4 \\ u_5 \end{pmatrix} = \begin{pmatrix} i\omega C_1 U \\ 0 \\ \frac{1}{R_5} U \\ 0 \end{pmatrix}$$

3.2.1 Theorie und Kondition

Bekannt aus der linearen Algebra:

3.2

p. 139

3.2

p. 140

Definition 3.2.1. Der **RANG** einer Matrix $M \in \mathbb{K}^{m,n}$ (engl. rank) ist die maximale Anzahl linear unabhängiger Spalten/Zeilen von M .

Theorem 3.2.2 (Invertierbarkeit von Matrizen). Eine Matrix $A \in \mathbb{R}^{n,n}$ heisst **regulär** oder **invertierbar**, falls eine der folgenden äquivalenten Bedingungen erfüllt ist:

- $\exists B \in \mathbb{R}^{n,n}$: $BA = AB = I$, wobei I die $n \times n$ -Einheitsmatrix ist.
- $x \mapsto Ax$ definiert einen Isomorphismus des \mathbb{R}^n ($Ax = b$ hat genau eine Lösung $\forall b$)
- Die Spalten von A sind linear unabhängig (voller Spaltenrang).
- Die Zeilen von A sind linear unabhängig (voller Zeilenrang).
- $\det A \neq 0$.

Abschn. 3.1.3: Fixiere Norm $\|\cdot\|$ auf \mathbb{R}^n \Rightarrow induziert Matrixnorm (\rightarrow Def. 2.2.6) auf $\mathbb{R}^{n,n}$

Vorüberlegung: Relative Konditionszahl von $x \mapsto Ax$, A regulär (bzgl. Störungen in x):

$$A(x + \delta x) = y + \delta y \Rightarrow \frac{\|\delta y\|}{\|y\|} \leq \|A\| \|A^{-1}\| \frac{\|\delta x\|}{\|x\|}. \quad (3.2.2)$$

Herleitung der Abschätzung durch

$$\|\delta y\| \leq \|A\| \|\delta x\| \quad \text{and} \quad \|y\| = \|Ax\| \geq \|A^{-1}\|^{-1} \|x\|.$$

Bei LGS: Wie wirken sich (kleine) Störungen in Eingabedaten A, b auf x aus?

$$(A + \Delta A)(x + \Delta x) = b + \Delta b, \quad \Delta x = ?$$

Lemma 3.2.3 (Störungslemma).

$$B \in \mathbb{R}^{n,n}, \|B\| < 1 \Rightarrow I + B \text{ regulär} \wedge \|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|}.$$

Beweis. Dreiecksungleichung $\Rightarrow \|(I + B)x\| \geq (1 - \|B\|) \|x\|, \forall x \in \mathbb{R}^n \Rightarrow B$ regulär.

$$\Rightarrow \|(I + B)^{-1}\| = \sup_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|(I + B)^{-1}x\|}{\|x\|} = \sup_{y \in \mathbb{R}^n \setminus \{0\}} \frac{\|y\|}{\|(I + B)y\|} \leq \frac{1}{1 - \|B\|}$$

\Rightarrow Wegen Submultiplikativität der Matrixnorm: $\|AB\| \leq \|A\| \|B\|$:

$$A \text{ regulär: } (A + \Delta A) = A(I + A^{-1}\Delta A) \text{ regulär, falls } \|\Delta A\| < \|A^{-1}\|^{-1}.$$

$$Ax = b \Leftrightarrow (A + \Delta A)(x + \Delta x) = b + \Delta b \Rightarrow (A + \Delta A)\Delta x = \Delta b - \Delta Ax. \quad (3.2.3)$$

Theorem 3.2.4 (Kondition der Lösung linearer Gleichungssysteme).

Falls A regulär, $\|\Delta A\| < \|A^{-1}\|^{-1}$ und (3.2.3), dann

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

Relativer Fehler Relative Störungen

Lemma 3.2.3 $\Rightarrow \|(A + \Delta A)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\Delta A\|}$ & (3.2.3)

$$\Rightarrow \|\Delta x\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\Delta A\|} (\|\Delta b\| + \|\Delta Ax\|) \leq \frac{\|A^{-1}\| \|A\|}{1 - \|A^{-1}\| \|\Delta A\|} \left(\frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right) \|x\|$$

Definition 3.2.5 (Kondition). **Kondition** einer Matrix $A \in \mathbb{R}^{n,n}$: $\text{cond}(A) := \|A^{-1}\| \|A\|$

Beachte: $\text{cond}(A)$ abhängig von $\|\cdot\|$!

Merkregel:

Wenn $\text{cond}(A) \approx 10^s$ und relative Störungen $\|\Delta A\| : \|A\|, \|\Delta b\| : \|b\| \approx 10^{-k}, k > s$, dann können von der Lösung des LGS $Ax = b$ nur $k - s$ korrekte Dezimalstellen erwartet werden.

Beispiel 67 („Schlecht“ konditioniertes lineares Gleichungssystem).

$$Ax = b \text{ mit } A = \begin{pmatrix} 4.1 & 2.8 \\ 9.7 & 6.6 \end{pmatrix}, b = \begin{pmatrix} 4.1 \\ 9.7 \end{pmatrix}, x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Normen: 1-Norm $\|\cdot\|_1 \Rightarrow$ Matrixnorm = Spaltensummennorm (3.1.2)

$$\text{Störung: } \Delta b = \begin{pmatrix} 0.01 \\ 0 \end{pmatrix}, \|\Delta b\|_1 = 0.01 \Rightarrow \Delta x = \begin{pmatrix} -0.66 \\ -0.97 \end{pmatrix}, \|\Delta x\|_1 = 1.63 \text{ („gross“)}$$

$$\mathbf{A}^{-1} = \begin{pmatrix} -66 & 28 \\ 97 & 41 \end{pmatrix} \Rightarrow \text{cond}_1(\mathbf{A}) = \|\mathbf{A}\|_1 \cdot \|\mathbf{A}^{-1}\|_1 = 13.8 \cdot 163 = 2249.4.$$

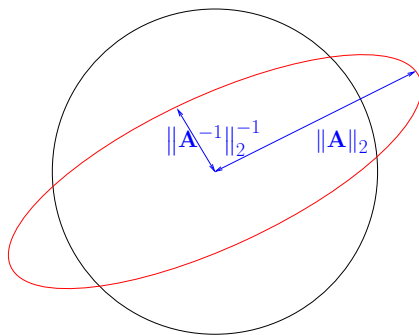
$$\updownarrow$$

$$\frac{\|\Delta \mathbf{x}\|}{\|\mathbf{x}\|} : \frac{\|\Delta \mathbf{b}\|}{\|\Delta \mathbf{x}\|} = 1.63 : 0.0007246 = 2249.4.$$

Kondition bzgl. Euklidischer Norm:

$$\text{cond}(\mathbf{A}) = \frac{\max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|}{\min_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|}$$

Für $\|\cdot\| = \|\cdot\|_2, n=2$: Verhältnis von maximalem zu minimalem Abstand des Bildes des Einheitskreises zum Ursprung.



MATLAB-Funktionen zur Konditionsberechnung/-abschätzung:

- $\text{cond}(\mathbf{A}) \rightarrow \text{cond}_2(\mathbf{A})$
- $\text{cond}(\mathbf{A}, 1)$ und $\text{cond}(\mathbf{A}, \text{inf}) \rightarrow \text{cond}_1(\mathbf{A})$ und $\text{cond}_\infty(\mathbf{A})$
- $\text{condest}(\mathbf{A})$ schätzt $\text{cond}_1(\mathbf{A})$, geeignet für grosse Matrizen

3.2.2 Die Gausselimination

GAUßELIM Ausnahmestellung linearer Gleichungssysteme (LGS):
☞ „exakt“ lösbar mit endlich vielen arithmetischen Grundoperationen

Standardtechnik: **Gausselimination** (→ Sekundarstufe, Kurs „Lineare Algebra“)



Idee: Transformation auf „einfachere“ äquivalente Gleichungssysteme durch sukzessive Zeilenumformungen (evtl. auch Spaltenumformungen)

Beispiel 68 (Gausselimination).

$$\begin{pmatrix} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 3 & -1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 4 \\ 1 \\ -3 \end{pmatrix} \longleftrightarrow \begin{cases} x_1 + x_2 = 4 \\ 2x_1 + x_2 - x_3 = 1 \\ 3x_1 - x_2 - x_3 = -3 \end{cases}$$

$$\begin{pmatrix} 1 & & \\ & 1 & \\ & & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 3 & -1 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \\ -3 \end{pmatrix} \xrightarrow{=L} \begin{pmatrix} 1 & & \\ & 1 & \\ & & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 0 & -1 & -1 \\ 0 & -4 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ -7 \\ -3 \end{pmatrix} \xrightarrow{=U}$$

$$\begin{pmatrix} 1 & & \\ & 1 & \\ & & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 0 & -1 & -1 \\ 0 & -4 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ -7 \\ -15 \end{pmatrix} \xrightarrow{=L} \begin{pmatrix} 1 & & \\ & 1 & \\ & & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 0 & -1 & -1 \\ 0 & 0 & 3 \end{pmatrix} \begin{pmatrix} 4 \\ -7 \\ 13 \end{pmatrix}$$

PIVOTELEM
= Pivotzeile, Pivotelement fett.

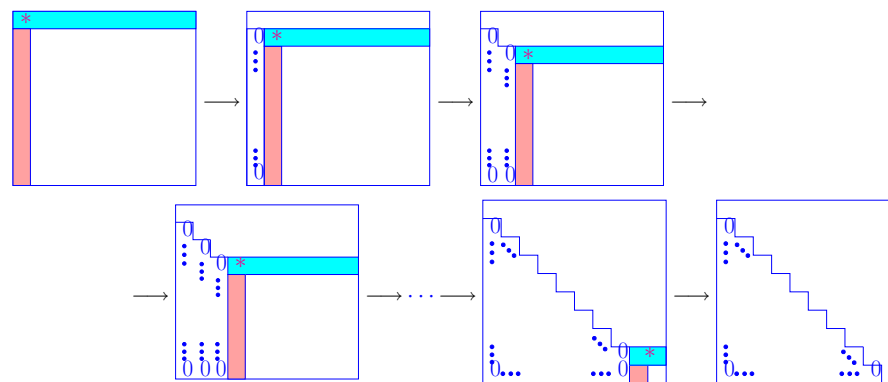
OBERE DREIECKSFORM
Transformation auf obere Dreiecksform

Anschließend: Lösung durch **RUECKSUBST**
Rücksubstitution:

$$\begin{cases} x_1 + x_2 = 4 \\ -x_2 - x_3 = -7 \\ 3x_3 = 13 \end{cases} \Rightarrow \begin{cases} x_3 = \frac{13}{3} \\ x_2 = 7 - \frac{13}{3} = \frac{8}{3} \\ x_1 = 4 - \frac{8}{3} = \frac{4}{3} \end{cases}$$

3.2
p. 145

3.2
p. 147



* $\hat{=}$ Pivotelement (notwendig $\neq 0$), = Pivotzeile

Im k . Schritt ($\mathbf{A} \in \mathbb{K}^{n,n}, 1 \leq k < n$, Pivotzeile \mathbf{a}_k^T): $\mathbf{A}\mathbf{x} = \mathbf{b} \xrightarrow{} \mathbf{A}'\mathbf{x} = \mathbf{b}'$ mit

$$a'_{ij} := \begin{cases} a_{ij} - \frac{a_{ik}a_{kj}}{a_{kk}} & \text{für } k < i, j \leq n, \\ 0 & \text{für } k < i \leq n, j = k, \\ a_{ij} & \text{sonst,} \end{cases} \quad b'_i := \begin{cases} b_i - \frac{a_{ik}b_k}{a_{kk}} & \text{für } k < i \leq n, \\ b_i & \text{sonst.} \end{cases} \quad (3.2.4)$$

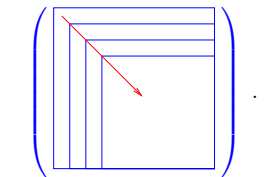
3.2
p. 146

3.2
p. 148

Block-Perspektive (erster Eliminationsschritt mit Pivotelement $\alpha \neq 0$):

$$A := \left(\begin{array}{c|c} \alpha & \mathbf{c}^T \\ \mathbf{d} & \mathbf{C} \end{array} \right) \rightarrow A' := \left(\begin{array}{c|c} \alpha & \mathbf{c}^T \\ \mathbf{0} & \mathbf{C}' := \mathbf{C} - \frac{\mathbf{d}\mathbf{c}^T}{\alpha} \end{array} \right). \quad (3.2.5) \quad \text{ur: block}$$

Sukzessive Durchführung:



Rechte Seite $\mathbf{b} \sim A(:, \text{end})$

Demonstrations- MATLAB-CODE: Rekursive Gausselimination

```
function A = blockgs(A)
% Recursive Gaussian elimination, no pivoting
n = size(A,1);
if (n ~= size(A,2))
    error('Size mismatch!'); end
if (n > 1)
    C = blockgs(A(2:end,2:end))-...
        A(2:end,1)*A(1,2:end)/A(1,1);
    A = [A(1,:) ; zeros(n-1,1), C];
end
```

Rechenaufwand (Anzahl elementarer Operationen) für die Gauss-Elimination:

$$\begin{aligned} \text{Elimination: } & \sum_{i=1}^{n-1} (n-i-1)(n-i+1) = \frac{1}{3}n^3 - \frac{5}{6}n + \frac{1}{2}n^2 + 1, \\ \text{Rücksubstitution: } & \sum_{i=1}^n (n-i+1) = \frac{1}{2}n(n+1). \end{aligned} \quad (3.2.6) \quad \text{eq: GEco}$$

Rechenaufwand für Gauss-Elimination (ohne Pivotsuche)
für generisches LGS $\mathbf{Ax} = \mathbf{b}$, $\mathbf{A} \in \mathbb{R}^{n,n}$ $= \frac{1}{3}n^3 + O(n^2)$

Programmiere niemals eine Gausselimination selbst!

Benutze Programmbibliotheken oder MATLAB (MATLAB operator: \)

3.2.3 Die LU-Zerlegung

LUZERL **MATEFAKTOR**
Interpretation der Gausselimination als Matrixfaktorisierung \rightarrow Bsp. 68: lex: GE

(Zeilenumformung = Multiplikation mit Eliminationsmatrix)

$$a_1 \neq 0 \Rightarrow \begin{pmatrix} 1 & 0 & \dots & 0 \\ -\frac{a_2}{a_1} & 1 & & 0 \\ -\frac{a_3}{a_1} & & \ddots & \\ \vdots & & & \ddots \\ -\frac{a_n}{a_1} & 0 & & 1 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} a_1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

► Nach $n-1$ Schritten der Gausselimination: ► Matrixfaktorisierung (\rightarrow Bsp. 68) lex: GE

$\mathbf{A} = \mathbf{L}_1 \cdots \mathbf{L}_{n-1} \mathbf{U}$ mit Eliminationsmatrizen \mathbf{L}_i , $i = 1, \dots, n-1$,
oberer Dreiecksmatrix $\mathbf{U} \in \mathbb{R}^{n,n}$.

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ l_2 & 1 & & 0 \\ l_3 & & \ddots & \\ \vdots & & & \ddots \\ l_n & 0 & & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & & 0 \\ 0 & h_3 & 1 & \\ \vdots & \vdots & & \ddots \\ 0 & h_n & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ l_2 & 1 & & 0 \\ l_3 & h_3 & 1 & \\ \vdots & \vdots & & \ddots \\ l_n & h_n & 0 & 1 \end{pmatrix}$$

$\mathbf{L}_1 \cdots \mathbf{L}_{n-1}$ sind **normalisierte untere Dreiecksmatrizen**
(Einträge = Faktoren $-\frac{a_{ik}}{a_{kk}}$ aus (3.2.4) \rightarrow Bsp. 68) lex: GE

3.2
p. 149

Definition 3.2.6 (Matrizentypen). Matrix $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{m,n}$ ist

- DIAGONAL** **Diagonalmatrix**, wenn $a_{ij} = 0$ für $i \neq j$,
- TRI U** **obere Dreiecksmatrix** (engl. upper triangular matrix), falls $a_{ij} = 0$ für $i > j$,
- TRI L** **untere Dreiecksmatrix** (engl. lower triangular matrix), falls $a_{ij} = 0$ für $i < j$.

Eine Dreiecksmatrix heisst **normalisiert**, wenn $a_{ii} = 1$, $i = 1, \dots, \min\{m, n\}$.

Die Vorwärtsrechnung in der Gausselimination für $\mathbf{Ax} = \mathbf{b}$, $\mathbf{A} \in \mathbb{R}^{n,n}$ ist algebraisch **und numerisch** (d.h. mit Berücksichtigung von Rundungsfehlern) äquivalent zur **LU-Zerlegung** (engl. LU decomposition) $\mathbf{A} = \mathbf{LU}$ in eine normalisierte untere Dreiecksmatrix \mathbf{L} und eine obere Dreiecksmatrix \mathbf{U} .

Eindeutigkeit der LU-Zerlegung:

Reguläre obere Dreiecksmatrizen/normalisierte untere Dreiecksmatrizen bilden eine Gruppe bzgl. Matrixmultiplikation. Der Schnitt beider Matrixmengen enthält nur die Einheitsmatrix.

$$\mathbf{L}_1 \mathbf{U}_1 = \mathbf{L}_2 \mathbf{U}_2 \Rightarrow \mathbf{L}_2^{-1} \mathbf{L}_1 = \mathbf{U}_2 \mathbf{U}_1^{-1} = \mathbf{I}.$$

3.2
p. 150

3.2
p. 151

3.2
p. 152

Direkte Herleitung der LU-Zerlegung

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

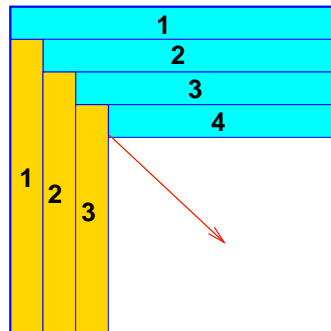
- Zeilenweise Berechnung von **U**
- Spaltenweise Berechnung von **L**

Einträge von **A** können überschrieben werden.

(Bezeichnung: Rechnung **in situ**)

CROUT

(Algorithmus von Crout)



Algorithmus 69.

in situ

LU-Zerlegung von $A \in \mathbb{R}^{n,n}$:

(ohne Pivotsuche)

A wird mit **L, U** überschrieben:

$B = \text{lurec}(A);$

$L = \text{tril}(B, -1) + \text{eye}(\text{size}(A));$

$U = \text{triu}(B);$

MATLAB-CODE: Rekursive LU-Zerlegung

```
function A = lurec(A)
% Recursive LU factorization
if (size(A,1) > 1)
    fac = A(2:end,1)/A(1,1);
    C = lurec(A(2:end,2:end)-...
              fac*A(1,2:end));
    A = [A(1,:) ; fac , C];
end
```

RECHNAUFWAND

Rechenaufwand für LU-Zerlegung von $A \in \mathbb{R}^{n,n} = \frac{1}{3}n^3 + O(n^2)$

(3.2.7) lab: lucos

Lösen eines linearen Gleichungssystems via LU-Zerlegung:

① LU-Zerlegung $A = LU$, Rechenaufwand $\frac{1}{3}n(n-1)(n+1)$

FORWARDSUBS

$Ax = b$: ② Vorwärtssubstitution, löse $Lz = b$, Rechenaufwand $\frac{1}{2}n(n-1)$

BACKWARDSUBS

③ Rücksubstitution, löse $Ux = z$, Rechenaufwand $\frac{1}{2}n(n+1)$

► Äquivalent zur Gaußelimination

► Aufwand zur Lösung von $Ax_k = b_k, b_k \in \mathbb{R}^n, k = 1, \dots, m, m \in \mathbb{N}$:

$1 \times$ LU-Zerlegung + $m \times$ Vorwärts- & Rückwärtssubstitution = $\frac{1}{3}n(n-1)(n+1) + mn^2$

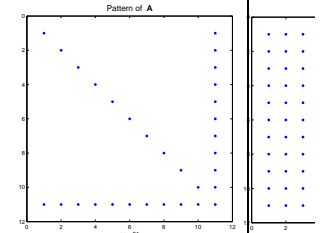
► Für $m \gg 1$ wesentlich „billiger“ als m -fache Gaußelimination!

BESETZMUSTER

Beispiel 70 (Besetzungsmuster der LU-Faktoren).

$A = [\text{diag}(1:10), \text{ones}(10,1); \text{ones}(1,10), 2];$

$[L,U] = \text{lu}(A); \text{spy}(A); \text{spy}(L); \text{spy}(U); \text{spy}(\text{inv}(A));$



Berechnung $A^{-1} = U^{-1}L^{-1}$ kann extrem unökonomisch sein!

3.2

p. 153

Bemerkung 71. „Teil-LU-Zerlegungen“ in linken oberen Blöcken:

$$\begin{pmatrix} \square & \square \\ \square & \square \end{pmatrix} = \begin{pmatrix} \square & \square \\ \square & \square \end{pmatrix} \begin{pmatrix} \square & \square \\ \square & \square \end{pmatrix}$$

Bemerkung 72. Blockweise Matrixmultiplikation (3.1.1) $\stackrel{\text{eq: Blockmul}}{=} \text{Block-LU-Zerlegung}$:

Mit $A_{11} \in \mathbb{K}^{n,n}$ regulär, $A_{12} \in \mathbb{K}^{n,m}$, $A_{21} \in \mathbb{K}^{m,n}$, $A_{22} \in \mathbb{K}^{m,m}$:

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} = \begin{pmatrix} I & 0 \\ A_{21}A_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} A_{11} & A_{12} \\ 0 & S \end{pmatrix}, \quad S := A_{22} - A_{21}A_{11}^{-1}A_{12}$$

SCHURKOMP

Schur-Komplement

(3.2.8) eq: Block

3.2

p. 154

3.2

p. 155

3.2

p. 156

3.2.4 Pivotsuche

PIVOTSUCHE

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

Zusammenbruch der Gauss-Elimination
Pivotelement = 0

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

Gauss-Elimination funktioniert

Idee: Vermeide sehr kleine Pivotelemente durch **Zeilenvertauschungen**

Terminologie: **Pivotsuche/Pivotstrategie** (engl. *pivoting*) = Wahl eines günstigen Pivotelements

Beispiel 73 (Pivotstrategie und Rundungsfehler).

$$A = \begin{pmatrix} \epsilon & 1 \\ 1 & 1 \end{pmatrix} \quad 0 < \epsilon < \text{eps}, \quad b = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow x = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

Einfache LU-Zerlegung:

$$\blacktriangleright L = \begin{pmatrix} 1 & 0 \\ \epsilon^{-1} & 1 \end{pmatrix}, \quad U = \begin{pmatrix} \epsilon & 1 \\ 0 & 1 - \epsilon^{-1} \end{pmatrix} = \tilde{U} := \begin{pmatrix} \epsilon & 1 \\ 0 & -\epsilon^{-1} \end{pmatrix} \quad \text{in } \mathbb{M}! \quad (3.2.9) \quad \text{u1}$$

$$\blacktriangleright \text{Lösung von } L\tilde{U}x = b: \quad x = \begin{pmatrix} 1 - \epsilon \\ 1 \end{pmatrix} \quad (\text{unbrauchbares Ergebnis!})$$

LU-Zerlegung nach Zeilenvertauschung:

$$A = \begin{pmatrix} 1 & 1 \\ \epsilon & 1 \end{pmatrix} \Rightarrow L = \begin{pmatrix} 1 & 0 \\ \epsilon & 1 \end{pmatrix}, \quad U = \begin{pmatrix} 1 & 1 \\ 0 & 1 - \epsilon \end{pmatrix} = \tilde{U} := \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \quad (3.2.10) \quad \text{u2}$$

$$\blacktriangleright \text{Lösung von } L\tilde{U}x = b: \quad x = \begin{pmatrix} \epsilon \\ 1 - \epsilon \end{pmatrix} \quad (\text{brauchbares Ergebnis!})$$

$$\text{Ohne Zeilenvertauschung, } \rightarrow \text{ (3.2.9): } L\tilde{U} = A + E \quad \text{mit} \quad E = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \quad \blacktriangleright \text{Instabil!}$$

$$\text{Nach Zeilenvertauschung, } \rightarrow \text{ (3.2.10): } L\tilde{U} = \tilde{A} + E \quad \text{mit} \quad E = \begin{pmatrix} 0 & 0 \\ 0 & \epsilon \end{pmatrix} \quad \blacktriangleright \text{Stabil!}$$

Geeignete Pivotstrategie entscheidend für numerische Stabilität der Gausselimination

Heuristik: „Explosion“ von Einträgen **L, U** \blacktriangleright Auslöschung bei Vorwärts- & Rückwärtssubstitution.

Algorithmus 74.

Rekursive in situ LU-Zerlegung von

$A \in \mathbb{R}^{n,n}$ mit **Spaltenpivotsuche**

Wahl des Pivotindex j :

$j \in \{i, \dots, n\}$ so, dass

$$\frac{|a_{ki}|}{\sum_{l=i}^n |a_{kl}|} \rightarrow \max \quad (3.2.11) \quad \text{u3}$$

für $k = j, k \in \{i, \dots, n\}$.

(relativ grösstes Pivotelement)

Demonstrations- MATLAB-CODE: Rekursive Gausselimination mit Spaltenpivotsuche

```
function A = gsrecpiv(A)
n = size(A,1);
if (n > 1)
    [p,j]=max(abs(A(:,1))./sum(abs(A)'))';
    if (p < eps*norm(A(:,1:n),1))
        error('Nearly Singular matrix'); end
    A([1,j],:) = A([j,1],:);
    A = gsrecpiv(A(2:end,2:end))-...
        A(2:end,1)*A(1,2:end)/A(1,1));
    A = [A(1,:) ; zeros(n-1,1), C];
end
```

Warum relativ grösstes Pivotelement in (3.2.11)? Skalieren des Gleichungssystems aus Bsp. 73: lex:pivdt

$$\begin{pmatrix} 2/\epsilon & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \epsilon & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 & 2/\epsilon \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2/\epsilon \\ 1 \end{pmatrix}$$

Keine Zeilenvertauschung, falls absolut grösstes Pivotelement genommen wird.

$$\begin{pmatrix} 2 & 2/\epsilon \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 2 & 2/\epsilon \\ 0 & 1 - 2/\epsilon \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 2 & 2/\epsilon \\ 0 & -2/\epsilon \end{pmatrix} \quad \text{in } \mathbb{M}.$$

3.2
p. 157

Konstruktiv bewiesen durch Algorithmus 74: lu:PS

Lemma 3.2.7 (LU-Zerlegung mit Pivotsuche). Zu jeder regulären Matrix $A \in \mathbb{R}^{n,n}$ existieren eine **Zeilenpermutation** P , eine **normalisierte untere Dreiecksmatrix** $L \in \mathbb{R}^{n,n}$ und eine **obere Dreiecksmatrix** $U \in \mathbb{R}^{n,n}$ (\rightarrow Def. 3.2.6) so dass def:LU

$$PA = LU.$$

MATLAB-Funktion: $[L,U,P] = \text{lu}(A)$ (P = Permutationsmatrix)

\blacktriangleright Gauss-Elimination mit Spaltenpivotsuche gemäss (3.2.11) ist „meistens“ numerisch stabil. lu:PS

Beispiel 75 (Wilkinsons Gegenbeispiel).

3.2
p. 158

3.2
p. 159

3.2
p. 160

$$n=10: \quad a_{ij} = \begin{cases} 1 & , \text{ falls } i=j \vee j=n, \\ -1 & , \text{ falls } i > j, \\ 0 & \text{ sonst} \end{cases}, \quad A = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & 1 & 0 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & 1 & 0 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & 1 & 0 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 & 0 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 \end{pmatrix}$$

Spaltenpivotsuche führt zu keinen Zeilenvertauschungen !

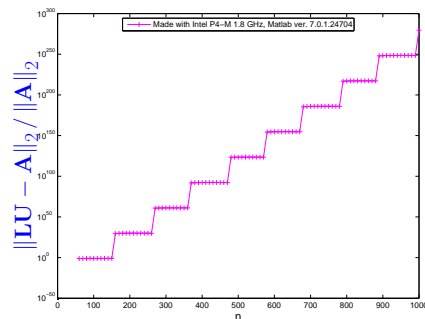
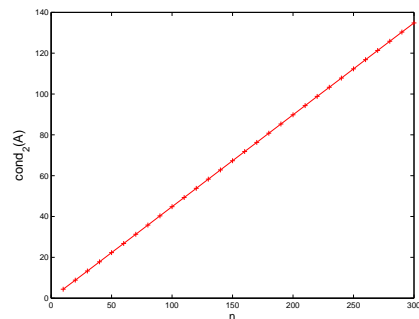
$$\blacktriangleright \quad A = LU, \quad l_{ij} = \begin{cases} 1 & , \text{ falls } i=j, \\ -1 & , \text{ falls } i > j, \\ 0 & \text{ sonst} \end{cases}, \quad u_{ij} = \begin{cases} 1 & , \text{ falls } i=j, \\ 2^{i-1} & , \text{ falls } j=n, \\ 0 & \text{ sonst} \end{cases}$$

Bsp. **Stabilität** der LU-Zerlegung ?

Ist $A - \tilde{L}\tilde{U}$ eine kleine relative Störung von A ?

```
res = [];
for n=10:10:300
    A = [tril(-ones(n,n-1))+2*eye(n-1);...
         zeros(1,n-1),ones(n,1)];
    [L,U,P] = lu(A);
    E = L*U-P*A;
    res = [res; n,cond(A),norm(E)/norm(A)];
end
```

MATLAB-CODE Stabilität der LU-Zerlegung



Bemerkung 76. **Totalpivotsuche** verwendet Zeilen- und Spaltenvertauschungen: verwende a_{ij} mit $|a_{ij}| \geq |a_{mk}|, i \leq m, k \leq n$ als Pivotelement im i . Schritt (bei in situ Rechnung !).

Gauss-Elimination mit Totalpivotsuche numerisch stabil (aber meist zu teuer!)

Beispiel 77 („Pfeilmatrix“).

$$A = \left(\begin{array}{c|c} \alpha & \mathbf{b}^T \\ \hline \mathbf{c} & \mathbf{D} \end{array} \right), \quad \begin{matrix} \alpha \in \mathbb{R}, \\ \mathbf{b}, \mathbf{c} \in \mathbb{R}^{n-1}, \\ \mathbf{D} \in \mathbb{R}^{n-1,n-1} \text{ reguläre Diagonalmatrix,} \end{matrix} \rightarrow \text{Def. 3.2.6}$$

(3.2.12) [Aarrow](#)

A ist Beispiel für **dünnbesetzte Matrix** (\rightarrow Def. 3.1.1), da

$$\#\{(i,j) \in \{1,\dots,n\}^2: a_{ij} \neq 0\} \ll n^2.$$

Anwendung von Algorithmus 74:

Faktormatrizen mit $O(n^2)$ Einträgen.

Rechenaufwand $O(n^3)$

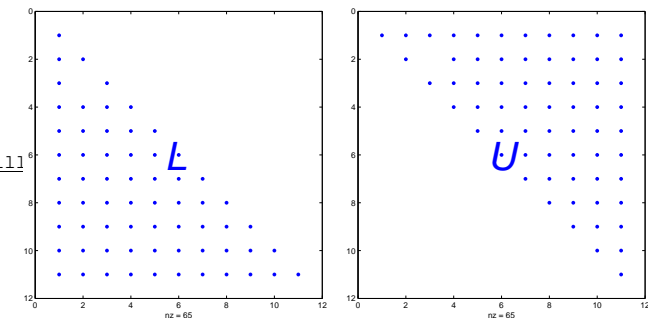
p. 161

MATLAB-CODE LU-Zerlegung von Pfeilmatrix

```
A = [ 1, ones(1,10); ...
      ones(10,1), eye(10,10)];
[L,U,P] = lu(A);
spy(L); spy(U);
```

3.2
p. 163

Auffüllen der LU-Faktoren
(engl. *Fill-in*) \rightarrow Def. 3.2.14



Rotiere Zeilen/Spalten:

- 1. Zeile/Spalte $\rightarrow n$. Zeile/Spalte,
- i . Zeile/Spalte $\rightarrow i-1$. Zeile/Spalte,
- $i = 2, \dots, n$

\blacktriangleright Alg. 69 kostet $O(n)$ Rechenoperationen:

$$A = \left(\begin{array}{c|c} & \mathbf{c} \\ \hline \mathbf{D} & \alpha \end{array} \right) \quad (3.2.13) \quad \text{u:Amo3}$$

3.2
p. 162

3.2
p. 164

Für permutiertes A aus (3.2.13), vgl. (3.2.8):

$$L = \left(\begin{array}{c|c} & \\ \hline & I \\ \hline b^T D^{-1} & 1 \end{array} \right), \quad U = \left(\begin{array}{c|c} & \\ \hline & D \\ \hline 0 & \sigma \end{array} \right), \quad \sigma := \alpha - b^T D^{-1} c.$$

➤ Kein Auffüllen bei Berechnung der LU-Zerlegung, Rechenaufwand $O(n)$

Lösen des LGS $Ax = y$ mit A aus 3.2.12: \overrightarrow{A} Zwei MATLAB-Implementierungen

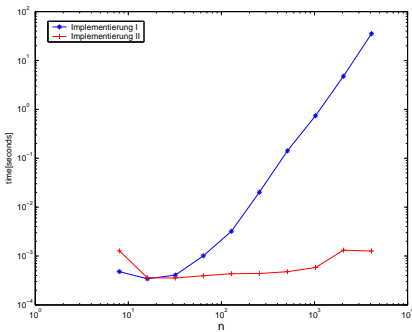
“Naive” Implementierung mittels “\”:

```
MATLAB-CODE: LGS mit Pfeilmatrix, Implementierung I
function x = sal(alpha,b,c,d,y)
A = [alpha, b'; c, diag(d)];
x = A\y;
```

„Strukturbewusste” Implementierung:

```
MATLAB-CODE: LGS mit Pfeilmatrix, Implementierung II
function x = sa2(alpha,b,c,d,y)
z = b./d;
xi = (y(1) - dot(z,y(2:end)))...
/(alpha-dot(z,c));
x = [xi; (y(2:end)-xi*c)./d];
```

```
Laufzeitmessung:
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 = sal(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-+');
```

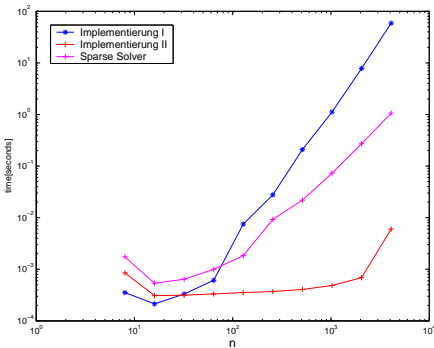


Resultat (Mobile Intel Pentium 4 - M CPU 2.40GHz, Linux MATLAB V6.5)

MATLAB kann es besser !

Mit dünnbesetzten Matrizen:
(Abschnitt 3.1.2)

```
MATLAB-CODE: Sparse Löser für Pfeilmatrix
function x = sa3(alpha,b,c,d,y)
n = length(d);
A = [alpha, b'; ...
c,spdiags(d,0,n,n)];
x = A\y;
```



Nütze Struktur eines linearen Gleichungssystems aus!



Vorsicht: Stabilität der Implementierung

3.2.5 Symmetrisch positiv definite Matrizen

SPDMATRIX

Definition 3.2.8 (Symmetrisch positiv definite Matrizen). $M \in \mathbb{K}^{n,n}$, $n \in \mathbb{N}$, ist **symmetrisch** (Hermiteisch) **positiv definit** (Abkürzung s.p.d.), falls

$$M = M^H \quad \wedge \quad x^H M x > 0 \quad \Leftrightarrow \quad x \neq 0.$$

Falls $x^H M x \geq 0$ für alle $x \in \mathbb{K}^n$, so nennt man M **positiv semidefinit**.

Bemerkung 78. (Lokale) Minimierung x^* von C^2 -Funktion $f : \mathbb{R}^n \mapsto \mathbb{R}$ ➤ Hessematrix $D^2 f(x^*)$ s.p.d.

Wie stelle ich fest, ob eine Matrix positiv definit ist ?

Notwendige Bedingungen:

Lemma 3.2.9 (Notwendige Bedingungen für s.p.d. Matrix). Für eine symmetrisch positiv definite Matrix $\mathbf{M} = \mathbf{M}^H \in \mathbb{K}^{n,n}$ gilt:

1. $m_{ii} > 0, i = 1, \dots, n$,
2. $m_{ii}m_{jj} - |m_{ij}|^2 > 0 \quad \forall i, j = 1, \dots, n$,
3. alle Eigenwerte von \mathbf{M} sind positiv. (\leftarrow auch hinreichend für symmetrisches \mathbf{M})

Beispiel 79 (Eigenschaften von Leitwertmatrizen). Siehe Bsp. ^{ex:elnet} 66

$\mathbf{A} \in \mathbb{R}^{n,n}$: Matrix, entstanden aus der Knotenanalyse eines linearen elektrischen Netzwerks mit ausschliesslich Ohmschen Widerständen

- $\mathbf{A} = \mathbf{A}^T$, $a_{kk} > 0$, $a_{kj} \leq 0$ für $k \neq j$,
- $\sum_{j=1}^n a_{kj} \geq 0$, $k = 1, \dots, n$.

◇

Definition 3.2.10 (Diagonaldominanz). $\mathbf{A} \in \mathbb{K}^{n,n}$ heisst **diagonaldominant**, falls

$$\forall k \in \{1, \dots, n\}: \sum_{j \neq k} |a_{kj}| \leq a_{kk}.$$

Lemma 3.2.11. Eine diagonaldominante Hermitesche Matrix ist positiv semidefinit.

Beweis. $\mathbf{A} = \mathbf{A}^H$ diagonaldominant, mit AGM-Ungleichung

$$\begin{aligned} \mathbf{x}^H \mathbf{A} \mathbf{x} &= \sum_{i=1}^n a_{ii} |x_i|^2 - \sum_{i \neq j} a_{ij} \bar{x}_i x_j \geq \sum_{i=1}^n a_{ii} |x_i|^2 - \sum_{i \neq j} |a_{ij}| |x_i| |x_j| \\ &\geq \sum_{i=1}^n a_{ii} |x_i|^2 - \frac{1}{2} \sum_{i \neq j} |a_{ij}| (|x_i|^2 + |x_j|^2) \\ &\geq \frac{1}{2} \left(\sum_{i=1}^n \{a_{ii} |x_i|^2 - \sum_{j \neq i} |a_{ij}| |x_i|^2\} \right) + \frac{1}{2} \left(\sum_{j=1}^n \{a_{ii} |x_j|^2 - \sum_{i \neq j} |a_{ij}| |x_j|^2\} \right) \\ &\geq \sum_{i=1}^n |x_i|^2 \left(a_{ii} - \sum_{j \neq i} |a_{ij}| \right) \geq 0. \end{aligned}$$

Theorem 3.2.12 (Gauss-Elimination für s.p.d. Matrizen). Zu jeder symmetrisch positiv definite Matrix (\rightarrow Def. 3.2.8) ^{def:spd} existiert eine LU-Zerlegung.

Beweisskizze. Gauss-Elimination ohne Pivotsuche durchführbar.

Induktionsbeweis: betrachte ersten Eliminationsschritt:

$$\mathbf{A} = \begin{pmatrix} a_{11} & \mathbf{b}^T \\ \mathbf{b} & \tilde{\mathbf{A}} \end{pmatrix} \xrightarrow[\text{Gauss-Elimination}]{1. \text{ Schritt}} \begin{pmatrix} a_{11} & \mathbf{b}^T \\ 0 & \tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^T}{a_{11}} \end{pmatrix}.$$

➤ Zu zeigen: $\tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^T}{a_{11}}$ s.p.d. (\rightarrow Induktionsschluss)

- Klar ist Symmetrie von $\tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^T}{a_{11}} \in \mathbb{R}^{n-1,n-1}$
- Da \mathbf{A} s.p.d. (\rightarrow Def. 3.2.8), ^{def:spd} für jedes $\mathbf{y} \in \mathbb{R}^{n-1} \setminus \{0\}$

$$0 < \begin{pmatrix} -\frac{\mathbf{b}^T \mathbf{y}}{a_{11}} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} a_{11} & \mathbf{b}^T \\ \mathbf{b} & \tilde{\mathbf{A}} \end{pmatrix} \begin{pmatrix} -\frac{\mathbf{b}^T \mathbf{y}}{a_{11}} \\ \mathbf{y} \end{pmatrix} = \mathbf{y}^T \left(\tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^T}{a_{11}} \right) \mathbf{y}.$$

► $\tilde{\mathbf{A}} - \frac{\mathbf{b}\mathbf{b}^T}{a_{11}}$ positiv definit. □

Beachte: keine Pivotsuche (\rightarrow Sect. 3.2.4) ^{sec:pivots-und-rund} erforderlich
(Spaltenpivotsuche wählt immer aktuelle Zeile)

3.2
p. 169

Lemma 3.2.13 (Cholesky-Zerlegung für s.p.d. Matrizen). Zu jedem s.p.d. $\mathbf{A} \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$, existiert eine **eindeutige obere Dreiecksmatrix** $\mathbf{R} \in \mathbb{R}^{n,n}$ mit $r_{ii} > 0, i = 1, \dots, n$, so dass $\mathbf{A} = \mathbf{R}^H \mathbf{R}$ (**Cholesky-Zerlegung**). ^{chol}

Thm. 3.2.12 \Rightarrow $\mathbf{A} = \mathbf{L}\mathbf{U}$ (LU-Zerlegung von \mathbf{A} , Sect. 3.2.3) ^{thm:GESpd} ^{sec:die-lr-zerlegung}

$$\mathbf{A} = \mathbf{L}\mathbf{D}\tilde{\mathbf{U}} \quad \begin{matrix} \mathbf{D} \hat{=} \text{Diagonale von } \mathbf{U}, \\ \tilde{\mathbf{U}} \hat{=} \text{normalisierte obere Dreiecksmatrix} \end{matrix} \rightarrow \text{Def. 3.2.6} \quad \text{def:LU}$$

Wegen Eindeutigkeit der LU-Zerlegung

$$\mathbf{A} = \mathbf{A}^T \Rightarrow \mathbf{U} = \mathbf{D}\mathbf{L}^T \Rightarrow \mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T,$$

mit eindeutigen \mathbf{L} , \mathbf{D} (Diagonalmatrix).

$$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 \quad \forall \mathbf{x} \neq 0 \Rightarrow \mathbf{y}^T \mathbf{D} \mathbf{y} > 0 \quad \forall \mathbf{y} \neq 0.$$

► \mathbf{D} hat positive Diagonaleinträge $\Rightarrow \mathbf{R} = \sqrt{\mathbf{D}}\mathbf{L}^T$. □

MATLAB-Kommando: `R = chol(A)`

3.2
p. 170

Rechnaufwand für Cholesky-Zerlegung: $\frac{1}{6}n^3 + O(n^2)$ \gg Halb so teuer wie Alg. ^{alg:LU} 69 ^{rechaufwchol}

3.2
p. 171

3.2
p. 172

Lösen eines LGS mit s.p.d. Koeffizientenmatrix via Cholesky-Zerlegung, Vorwärts- und Rückwärtssubstitution ist numerisch stabil.

Bemerkung 80 (LU-Zerlegung diagonaldominanter Matrizen).

A regulär, diagonaldominant $\Leftrightarrow \begin{cases} A \text{ LU-zerlegbar} \\ \Updownarrow \\ \text{Gauss-Elimination ohne Pivotsuche durchführbar} \end{cases}$

Beweis. $\xrightarrow{\text{lu:block (3.2.5)}}$ Induktion nach n : Nach 1. Eliminationsschritt:

$$a_{ij}^{(1)} = a_{ij} - \frac{a_{i1}}{a_{11}} a_{1j}, \quad i, j = 2, \dots, n \Rightarrow a_{ii}^{(1)} > 0.$$

$$\begin{aligned} \blacktriangleright |a_{ii}^{(1)}| - \sum_{\substack{j=2 \\ j \neq i}}^n |a_{ij}^{(1)}| &= \left| a_{ii} - \frac{a_{i1}}{a_{11}} a_{1i} \right| - \sum_{\substack{j=2 \\ j \neq i}}^n \left| a_{ij} - \frac{a_{i1}}{a_{11}} a_{1j} \right| \\ &\geq a_{ii} - \frac{|a_{i1}| |a_{1i}|}{a_{11}} - \sum_{\substack{j=2 \\ j \neq i}}^n |a_{ij}| - \frac{|a_{i1}|}{a_{11}} \sum_{\substack{j=2 \\ j \neq i}}^n |a_{1j}| \\ &\geq a_{ii} - \frac{|a_{i1}| |a_{1i}|}{a_{11}} - \sum_{\substack{j=2 \\ j \neq i}}^n |a_{ij}| - |a_{i1}| \frac{a_{11} - |a_{1i}|}{a_{11}} \geq a_{ii} - \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \geq 0. \end{aligned}$$

A

$\xrightarrow{\text{lu:PS 3.2.11}}$

i i

Bemerkung 81 (Übermitteln von Matriceigenschaften an MATLAB).

MATLAB-`\` kann einer allgemeinen (vollbesetzten) Matrix spezielle Eigenschaften (symmetrisch, s.p.d., Dreiecksmatrix) nicht „ansehen“.



Benutze `y = linsolve(A,b,opts)`

opts $\in \{$ LT \leftrightarrow A untere Dreiecksmatrix
UT \leftrightarrow A obere Dreiecksmatrix
UHES \leftrightarrow A obere Hessenbergmatrix
SYM \leftrightarrow A Hermitesche Matrix
POSDEF \leftrightarrow A positiv definite Matrix $\}$



3.2.6 Dünnbesetzte Gleichungssysteme

DÜNNBESGGLSYS

= LGS $Ax = b$ mit dünnbesetzter Koeffizientenmatrix A , \rightarrow Abschnitt 3.1.2, Def. 3.1.1 $\xrightarrow{\text{sec:matrix-def:parse}}$

Beispiel 82 (Dünnbesetzte Gleichungssysteme in der Schaltkreissimulation).

3.2
p. 173

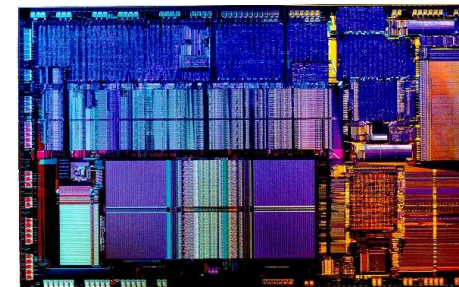
Moderne Schaltkreise (Chips):

$10^5 - 10^7$ Elemente

- Jedes Element mit nur wenigen Knoten verbunden
- Jeder Knoten mit nur wenigen Elementen verbunden

[Im Fall eines linearen Netzwerks]

Knotenanalyse \rightarrow dünnbesetzte
Schaltkreismatrix



LUFAKTDUENNMA

Beispiel 83 (LU-Faktorisierung dünnbesetzter Matrizen). \rightarrow Bsp. 77 $\xrightarrow{\text{ex:arrowmatrix}}$

$$A = \left(\begin{array}{ccc|ccc} 3 & -1 & & -1 & & \\ -1 & 3 & & & & \\ & & \ddots & & & \\ -1 & & & 3 & -1 & \\ & & & -1 & 3 & \\ & & & & & \ddots & -1 \\ & & & & & -1 & 3 \end{array} \right) \in \mathbb{R}^{n,n}, n \in 2\mathbb{N}$$

3.2
p. 174

MATLAB `lu`-Funktion ergibt (für $n = 10$):

3.2
p. 175

3.2
p. 176

$$A = \underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ -0.33 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -0.38 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -0.38 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -0.38 & 1 & 0 & 0 & 0 & 0 & 0 \\ -0.33 & -0.13 & -0.05 & -0.02 & -0.01 & 1 & 0 & 0 & 0 & 0 \\ 0 & -0.38 & -0.14 & -0.05 & -0.02 & -0.44 & 1 & 0 & 0 & 0 \\ 0 & 0 & -0.38 & -0.15 & -0.06 & -0.02 & -0.58 & 1 & 0 & 0 \\ 0 & 0 & 0 & -0.38 & -0.15 & -0.01 & -0.03 & -0.65 & 1 & 0 \\ 0 & 0 & 0 & 0 & -0.38 & -0.003 & -0.01 & -0.04 & -0.67 & 1 \end{pmatrix}}_L \underbrace{\begin{pmatrix} 3 & -1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & 2.67 & -1 & 0 & 0 & -0.33 & -1 & 0 & 0 & 0 \\ 0 & 0 & 2.63 & -1 & 0 & -0.13 & -0.38 & -1 & 0 & 0 \\ 0 & 0 & 0 & 2.62 & -1 & -0.05 & -0.14 & -0.38 & -1 & 0 \\ 0 & 0 & 0 & 0 & 2.62 & -0.02 & -0.05 & -0.15 & -0.38 & -1 \\ 0 & 0 & 0 & 0 & 0 & 2.62 & -1.15 & -0.06 & -0.02 & -0.01 \\ 0 & 0 & 0 & 0 & 0 & 0 & 2.06 & -1.19 & -0.07 & -0.02 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1.87 & -1.21 & -0.07 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1.78 & -1.19 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1.82 \end{pmatrix}}_U$$

A dünnbesetzt \nrightarrow LU -Faktoren dünnbesetzt (siehe Bsp. [lex:arrowmatrix](#) 77)

Definition 3.2.14 (Fill-in). Sei $A = LU$ eine LU -Faktorisierung von $A \in \mathbb{K}^{n,n}$. Wenn $l_{ij} \neq 0$ oder $u_{ij} \neq 0$ obwohl $a_{ij} = 0$, dann handelt es sich um **Fill-in** an Position (i, j) .

Definition 3.2.15 (Bandbreite). Für $A = (a_{ij})_{i,j} \in \mathbb{K}^{m,n}$ heisst

$\overline{m}(A) := \min\{k \in \mathbb{N} : j - i > k \Rightarrow a_{ij} = 0\}$ **Obere Bandbreite** (engl. upper bandwidth),

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

$m(A) := \overline{m}(A) + \underline{m}(A) + 1 =$ **Bandbreite** von A .

- $m(A) = 0 \triangleright A$ Diagonalmatrix, \rightarrow Def. [3.2.6](#)
- $\overline{m}(A) = \underline{m}(A) = 1 \triangleright A$ **Tridiagonalmatrix**

MATLAB-Funktion zum Erzeugen von Bandmatrizen:

Vollbesetzte Matrizen : `X=diag(v);`

Dünnbesetzte Matrizen : `X=spdiags(B,d,m,n);`

Tridiagonalmatrizen : `X=gallery('tridiag',c,d,e);` (\rightarrow sparse)

Definition 3.2.16 (Struktursymmetrie).

$A \in \mathbb{K}^{n,n}$ ist **struktursymmetrisch** (engl. structurally symmetric), falls

$$\forall i, j \in \{1, \dots, n\}: a_{ij} \neq 0 \Leftrightarrow a_{ji} \neq 0.$$

Definition 3.2.17 (Hülle einer Matrix). Für struktursymmetrisches, reguläres $A \in \mathbb{K}^{n,n}$ definie-

ZEILENBANDERE

ENV

Hülle (engl. envelope)

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

$$\text{env}(A) := \{(i, j), (j, i) : i - m_i(A) \leq j \leq i; 1 \leq i \leq n\}$$

Beispiel 84 (Hülle einer Matrix).

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

$\text{env}(A)$ = rote Einträge

$*$ \triangleq nichtverschwindender Eintrag $a_{ij} \neq 0$

Theorem 3.2.18 (Hülle & Fill-in). Falls $A \in \mathbb{K}^{n,n}$ regulär, struktursymmetrisch mit LU -Zerlegung $A = LU$, dann beschränkt sich der Fill-in auf $\text{env}(A)$.

Beweis. (induktiv) Im ersten Schritt der Gausselimination (ohne Pivotsuche):

$$A = \begin{pmatrix} a_{11} & \mathbf{b}^T \\ \mathbf{c} & \tilde{A} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 \\ -\frac{\mathbf{c}}{a_{11}} & \mathbf{I} \end{pmatrix}}_L \underbrace{\begin{pmatrix} a_{11} & \mathbf{b}^T \\ 0 & \tilde{A} - \frac{\mathbf{c}\mathbf{b}^T}{a_{11}} \end{pmatrix}}_U$$

$$\text{Wenn } (i, j) \notin \text{env}(A), i \geq j \Rightarrow \begin{cases} c_i = b_i = 0 \\ a_{ij} = a_{ji} = 0 \end{cases} \Rightarrow \tilde{l}_{ij} = 0 \ \& \ \tilde{u}_{ji} = 0.$$

Insbesondere $(i-1, j-1) \notin \text{env}(\tilde{U})$ für $i, j > 1$

► Speichere nur $a_{ij}, (i, j) \in \text{env}(A)$ bei in situ Berechnung der LU -Zerlegung von $A \in \mathbb{K}^{n,n}$

► Speicherbedarf: $n + 2 \sum_{i=1}^n m_i(A)$ Gleitpunktzahlen

► **hüllenorientierte Speichertechnik**

Beispiel 85 (Hüllenorientierte Speicherung).

Lineare hüllenorientierte Speicherung von s.p.d. $A = A^T \in \mathbb{R}^{n,n}$:

Zwei Arrays:

real val(1:P),
integer dptr(0:n)

$$P := n + \sum_{i=1}^n m_i(A).$$

(3.2.14) leg:Pdef

$$A = \begin{pmatrix} * & 0 & * & 0 & 0 & 0 & 0 \\ 0 & * & 0 & 0 & 0 & 0 & * \\ * & 0 & * & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & * & * & * & * \\ 0 & 0 & 0 & 0 & * & * & 0 \\ 0 & * & 0 & 0 & * & * & * \\ 0 & 0 & * & * & 0 & * & * \end{pmatrix}$$

Regel:

$$dptr[j] = k$$

$$\Updownarrow$$

$$val[k] = a_{jj}$$

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
val		a ₁₁	a ₂₂	a ₃₁	a ₃₂	a ₃₃	a ₄₄	a ₅₂	a ₅₃	a ₅₄	a ₅₅	a ₆₅	a ₆₆	a ₇₃	a ₇₄	a ₇₅	a ₇₆	a ₇₇
dptr	0	1	2	5	6	10	12	17										

◇

Bemerkung 86. Hüllenorientierte Cholesky-Zerlegung von s.p.d. $A \in \mathbb{R}^{n,n}$ überspringt $(i, j) \notin \text{env}(A)$

➤ Rechenaufwand: $O(P \cdot \max_{1 \leq i \leq N} m_i(A))$ elementare Operationen !

△

Algorithmus: Hüllenorientierte Cholesky-Zerlegung

```
for (j = 1; j ≤ N; j++) {  
  for (rjj = ajj, k = j - mj(A); k < j; k++)  
    rjj = rjj - rjk2;  
  rjj = √rjj;  
  for (i = j + 1; i ≤ N; i++) {  
    if ((i, j) ∈ env(A)) {  
      for (rij = aij, k = j - mi(A); k < j; k++)  
        rij = rij - rik * rjk;  
      rij = rij / rjj;  
    }  
  }  
}
```

Rechenaufwand:

$O(P \cdot \max_{1 \leq i \leq N} m_i(A))$ elementare Operationen

$$P := n + \sum_{i=1}^n m_i(A).$$

$A \in \mathbb{R}^{n,n}$ s.p.d. Bandmatrix

➤ Rechenaufwand für Cholesky-Zerlegung:

$$O(m(A) \cdot n)$$

BANDERRED

Bandbreitenreduktion:

Ziel: Reduziere $m_i(A)$, $A = (a_{ij}) \in \mathbb{R}^{N,N}$, durch Zeilen- und Spaltenvertauschungen.

Beispiel 87 (Bandbreitenreduktion durch Indexpermutation).

Spiegelung an Nebendiagonalen ➤ Reduktion von P aus leg:Pdef (3.2.14)

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

$i \leftarrow N + 1 - i$
 $P = 30$ $P = 22$

MATRIXGRAPH
Allgemeine Techniken benutzen **Matrixgraph**: (→ Graphentheorie)

3.2
p. 181

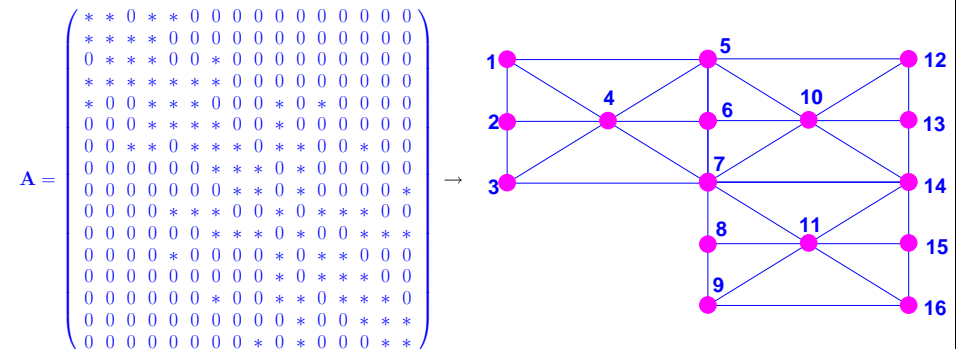
Definition 3.2.19 (Matrixgraph). Matrix-Graph $G(A) = (VG(A), EG(A))$ zu $A \in \mathbb{R}^{n,n}$:

- KNOTENMENGE Knotenmenge $VG(A)$ von $G(A) = \{1, \dots, n\}$
- KANTENMENGE Kantenmenge $EG(A)$ von $G(A)$: $EG(A) = \{(i, j) \in \{1, \dots, n\}^2; a_{ij} \neq 0\}$

A struktursymmetrisch $\Rightarrow G(A)$ ungerichtet.

$\iota \in VG(A)$: $\text{Grad } \deg \iota := \#\{\pi \in EG(A) : \exists \kappa \in VG(A), \pi = (\iota, \kappa)\}$

PLANARMATGRAPH
Beispiel 88 (Matrixgraph). Matrix mit planarem Matrixgraphen („Schaltkreis, Bsp. lex:elnet 66)



3.2
p. 182

3.2
p. 183

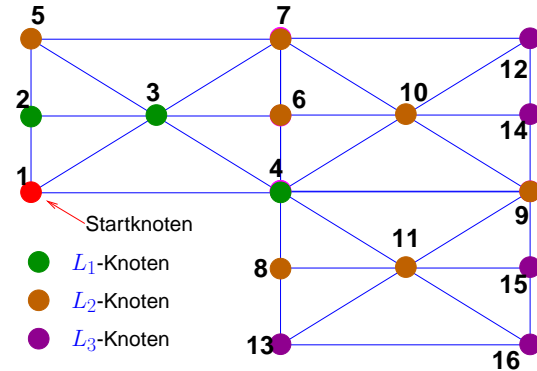
3.2
p. 184

Umordnungsalgorithmus von **Cuthill & McKee** [41] (für struktursymmetrische Matrix $A \in \mathbb{R}^{n,n}$):

- ① Startknoten minimalen Grades $u_0 \in VG(A) : \deg(u_0) = \min\{\deg(v) : v \in VG(A)\}$
- ② Partitionierung von $VG(A)$:

$$L_i := \{v \in VG(A) : \text{dist}(v, u_0) = i\}, \quad i \in \mathbb{N}.$$

- ③ Sortierung $VG(A) = \{u_0\} \cup L_1 \cup L_2 \cup \dots \cup L_N$,
interne Sortierung der L_i nach aufsteigendem Grad.



Beispiel 89 (Algorithmus von Cuthill & McKee).

Matrix aus Bsp. 88: Partitionierung und Numerierung erzeugt durch CMK-Algorithmus

Erweiterung: Reverse Cuthill-McKee-Algorithmus

(CMK-Alg. nach Spiegelung an Nebendiagonalen → Bsp. 87)

☞ MATLAB-Funktion: `symrcm(M)`

Bemerkung 90. Alternative Umsortierungsstrategie: Optimieren der Anordnung unter Berücksichtigung der „Fill-in“-Historie der Cholesky-Zerlegung, **minimum degree reordering** [2].

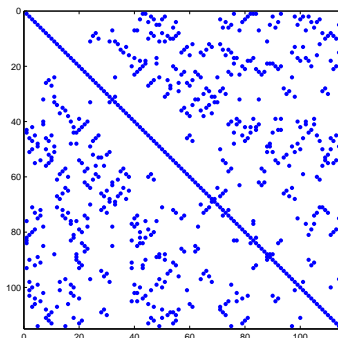
☞ MATLAB-Funktion: `symamd(M)`

Beispiel 91 (Fill-in-Minimierung durch Umordnung).

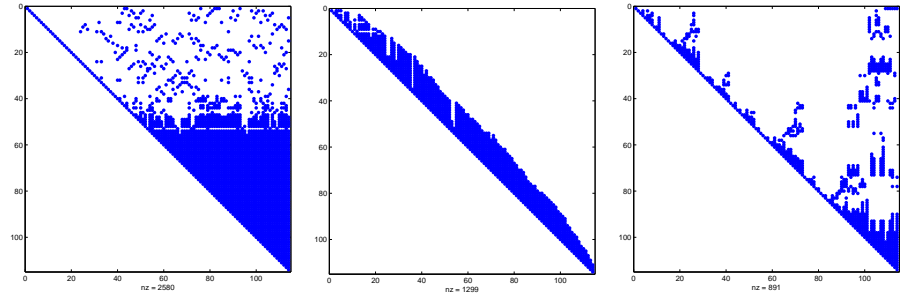
M: 114×114 s.p.d. Matrix aus FE-Diskretisierung einer PDE.

```
spy(M);
[R,P] = chol(M); spy(R);
r = symrcm(M);
[R,P] = chol(M(r,r)); spy(R);
m = symamd(M);
[R,p] = chol(M(m,m)); spy(R);
```

BESETZMUSTER
Besetzungsmuster von **M**



Untersuche Besetzungsmuster von Cholesky-Faktoren (→ Abschnitt 3.2.5) nach Umordnung:



keine Umordnung

Reverse Cuthill-McKee

Approximate Minimum Degree

Warnung: Immer Bibliotheksroutinen zur Lösung von allgemeinen LGS mit dünnbesetzten Koeffizientenmatrizen verwenden!

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

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

→ Pardiso (<http://www.computational.unibas.ch/cs/scicomp/software/pardiso>)

Matlab-\-Löser: verwendet Sparse-Techniken für **sparse**-Matrizen

3.2
p. 185

3.2.7 Die QR-Zerlegung

QRZERLEGUNG

LU-Probleme (Abschn. 3.2.4): **sec:pivot-und-rund** t.a. Pivotsuche notwendig (→ Strukturzerstörung)



Mögliche Instabilität (bei Spaltenpivotsuche)

Definition 3.2.20.

- $Q \in \mathbb{K}^{n,n}$, $n \in \mathbb{N}$, heisst **UNITÄR**, falls $Q^{-1} = Q^H$.
- $Q \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$, heisst **ORTHOGONAL**, falls $Q^{-1} = Q^T$.

Theorem 3.2.21 (Kriterium für Unitarität).

$$Q \in \mathbb{C}^{n,n} \text{ unitär} \Leftrightarrow \|Qx\|_2 = \|x\|_2 \quad \forall x \in \mathbb{K}^n.$$

► Q unitär $\Rightarrow \text{cond}(Q) = 1$ ► **leg:relcmv** (3.2.2) Unitäre Transformationen „stabilitätsfördernd“

3.2
p. 186

Falls $Q \in \mathbb{K}^{n,n}$ unitär, dann

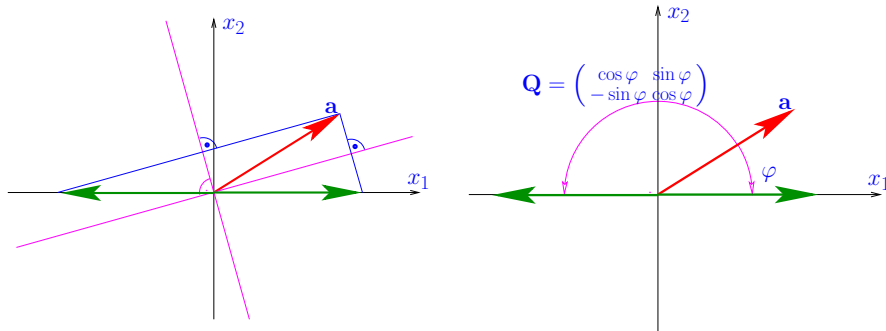
3.2
p. 187

3.2
p. 188

- haben alle Spalten- und Zeilenvektoren paarweise Euklidische Norm = 1,
- sind alle Spalten- und Zeilenvektoren paarweise orthogonal (bzgl. Euklidischem Skalarprodukt),
- ist $|\det Q| = 1$, alle Eigenwerte haben Betrag = 1.

Ziel: Finde orthogonale/unitäre Zeilenumformungen zum Annullieren von Matrixkomponenten.

In 2D: Zwei orthogonale Transformationen annullieren zweite Komponente von $\mathbf{a} \in \mathbb{R}^2$:



Spiegelung an Winkelhalbierenden

Rotation auf x_1 -Achse

Es gibt jeweils zwei mögliche Reflektionen/Rotationen

In n D: Zu $\mathbf{a} \in \mathbb{K}^n$ suche unitäre Matrix $Q \in \mathbb{K}^{n,n}$ mit $Q\mathbf{a} = \|\mathbf{a}\|_2 \mathbf{e}_1$, $\mathbf{e}_1 \hat{=} 1$. Einheitsvektor.

1. Möglichkeit: **HOUSEHOLDER**
Householder-Reflektionen

$$Q = H(\mathbf{v}) := I - 2 \frac{\mathbf{v}\mathbf{v}^H}{\mathbf{v}^H\mathbf{v}} \quad \text{mit} \quad \mathbf{v} = \frac{1}{2}(\mathbf{a} \pm \|\mathbf{a}\|_2 \mathbf{e}_1). \quad (3.2.15)$$

Zur Vermeidung von Auslöschung (\rightarrow Abschn. 11.7) wähle

$$\mathbf{v} = \begin{cases} \frac{1}{2}(\mathbf{a} + \|\mathbf{a}\|_2 \mathbf{e}_1) & , \text{ falls } \operatorname{Re} a_1 > 0, \\ \frac{1}{2}(\mathbf{a} - \|\mathbf{a}\|_2 \mathbf{e}_1) & , \text{ falls } \operatorname{Re} a_1 \leq 0. \end{cases}$$

2. Möglichkeit: **GIVENS**
Sukzessive Givens-Rotationen (\rightarrow 2D Fall)

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

Zur Vermeidung von Überlauf (in \mathbb{R}) \rightarrow Abschn. 11.5):

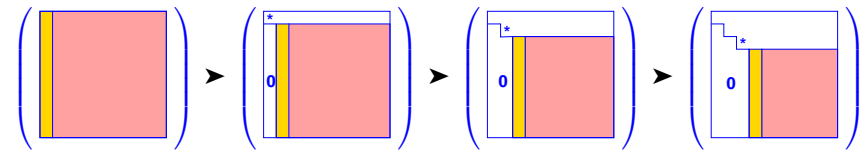
$$\begin{cases} \tau = \frac{a_1}{a_k}, & \sigma = \frac{1}{\sqrt{1+\tau^2}}, & \gamma = \sigma\tau, & \text{ falls } |a_k| > |a_1|, \\ \tau = \frac{a_k}{a_1}, & \gamma = \frac{1}{\sqrt{1+\tau^2}}, & \sigma = \gamma\tau, & \text{ falls } |a_k| \leq |a_1|. \end{cases}$$

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

MATLAB-CODE Givens-Rotation

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

Transformation auf obere Dreiecksgestalt durch sukzessive unitäre Transformationen:



3.2
p. 189

Yellow bar = „Zielspalte \mathbf{a} “ (bestimmt unitäre Transformation), Red bar = Zu modifizierender Restblock.

3.2
p. 191

$$Q_{n-1}Q_{n-2} \cdots Q_1 \mathbf{A} = \mathbf{R},$$

QR-Zerlegung von $\mathbf{A} \in \mathbb{C}^{n,n}$: $\mathbf{A} = \mathbf{Q}\mathbf{R}$, $\mathbf{Q} := Q_1^H \cdots Q_{n-1}^H$ unitäre Matrix, \mathbf{R} obere Dreiecksmatrix.

Verallgemeinerung auf $\mathbf{A} \in \mathbb{K}^{m,n}$

$$m > n: \quad \begin{pmatrix} \vdots \\ \mathbf{A} \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ \mathbf{Q} \\ \vdots \end{pmatrix} \begin{pmatrix} \mathbf{R} \\ \vdots \end{pmatrix}, \quad \mathbf{A} = \mathbf{Q}\mathbf{R}, \quad \begin{matrix} \mathbf{Q} \in \mathbb{K}^{m,n} \\ \mathbf{R} \in \mathbb{K}^{n,n} \end{matrix}$$

wobei $Q^H Q = I$ (orthonormierte Spalten), \mathbf{R} obere Dreiecksmatrix.

$$m < n: \quad \begin{pmatrix} \mathbf{A} \\ \vdots \end{pmatrix} = \begin{pmatrix} \mathbf{Q} \\ \vdots \end{pmatrix} \begin{pmatrix} \mathbf{R} \\ \vdots \end{pmatrix}, \quad \mathbf{A} = \mathbf{Q}\mathbf{R}, \quad \begin{matrix} \mathbf{Q} \in \mathbb{K}^{m,m} \\ \mathbf{R} \in \mathbb{K}^{m,n} \end{matrix}$$

3.2
p. 190

3.2
p. 192

wobei \mathbf{Q} unitär, \mathbf{R} obere Dreiecksmatrix.

Wann nimmt man welche Transformation ?

- ▶ Householder-Reflexionen vorzuziehen, wenn Zielspalten vollbesetzt
- ▶ Givens-Rotationen effizienter (weil selektiver), wenn Zielspalte nur wenige Nicht-Null-Komponenten besitzt.

MATLAB-Funktionen:

$$\begin{aligned} [\mathbf{Q}, \mathbf{R}] &= \text{qr}(\mathbf{A}) \quad \mathbf{Q} \in \mathbb{K}^{m,m}, \mathbf{R} \in \mathbb{K}^{m,n} \text{ für } \mathbf{A} \in \mathbb{K}^{m,n} \\ [\mathbf{Q}, \mathbf{R}] &= \text{qr}(\mathbf{A}, 0) \quad \mathbf{Q} \in \mathbb{K}^{m,n}, \mathbf{R} \in \mathbb{K}^{n,n} \text{ für } \mathbf{A} \in \mathbb{K}^{m,n}, m > n \\ &\text{(jede Funktion verfügbar für voll- und dünnbesetzte Matrizen)} \end{aligned}$$

Bemerkung 92 (QR-Orthogonalisierung).

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

Falls $m > n$, $\text{rank}(\mathbf{R}) = \text{rank}(\mathbf{A}) = n$ (voller Rang)

- ▶ $\{\mathbf{q}_{\cdot,1}, \dots, \mathbf{q}_{\cdot,n}\}$ ist **Orthonormalbasis** von $\text{Im}(\mathbf{A})$ mit $\text{Span}\{\mathbf{q}_{\cdot,1}, \dots, \mathbf{q}_{\cdot,k}\} = \text{Span}\{\mathbf{a}_{\cdot,1}, \dots, \mathbf{a}_{\cdot,k}\}, 1 \leq k \leq n$.

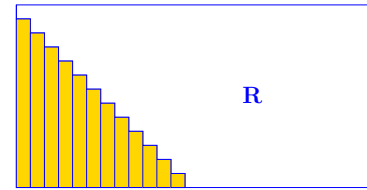
Bemerkung 93 (Akkumulation von unitären Transformationen).

Wie speichert man $\mathbf{G}_{i_1 j_1}(a_1, b_1) \cdots \mathbf{G}_{i_k j_k}(a_k, b_k), \mathbf{H}(\mathbf{v}_1) \cdots \mathbf{H}(\mathbf{v}_k)$?

☛ Für Householder-Reflexionen

$\mathbf{H}(\mathbf{v}_1) \cdots \mathbf{H}(\mathbf{v}_k)$: Speichere $\mathbf{v}_1, \dots, \mathbf{v}_k$

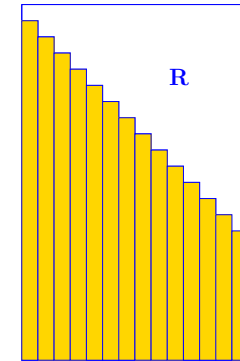
Für in situ QR-Zerlegung von $\mathbf{A} \in \mathbb{K}^{m,n}$: Speichere „Householder-Vektoren“ \mathbf{v}_j (abnehmende Längen !) im unteren Dreieck von \mathbf{A}



↑ Fall $m < n$

= Householder-Vektoren.

Fall $m > n \rightarrow$



3.2
p. 193

☛ Konvention für Givens-Rotationen ($\mathbb{K} = \mathbb{R}$)

$$\mathbf{G} = \begin{pmatrix} \bar{\gamma} & \bar{\sigma} \\ -\sigma & \gamma \end{pmatrix} \Rightarrow \text{Speichere } \rho := \begin{cases} 1 & , \text{ falls } \gamma = 0, \\ \frac{1}{2} \text{sign}(\gamma)\sigma & , \text{ falls } |\sigma| < |\gamma|, \\ 2 \text{sign}(\sigma)/\gamma & , \text{ falls } |\sigma| \geq |\gamma|. \end{cases}$$

$$\begin{cases} \rho = 1 & \Rightarrow \gamma = 0, \quad \sigma = 1 \\ |\rho| < 1 & \Rightarrow \sigma = 2\rho, \quad \gamma = \sqrt{1 - \sigma^2} \\ |\rho| > 1 & \Rightarrow \gamma = 2/\rho, \quad \sigma = \sqrt{1 - \gamma^2} \end{cases}$$

Speichere $\mathbf{G}_{ij}(a, b)$ als Tripel (i, j, ρ)

Speicherung von Orthogonaltransformationen als Matrizen meist ineffizient !

LOESEN LGS QR

Lösen eines linearen Gleichungssystems via QR-Zerlegung:

① QR-Zerlegung $\mathbf{A} = \mathbf{QR}$, Rechenaufwand $\frac{2}{3}n^3 + O(n^2)$ (ca. doppelt so aufwendig wie LU-Zerlegung ohne Pivotsuche)

$\mathbf{Ax} = \mathbf{b}$: ② Orthogonaltransformation $\mathbf{z} = \mathbf{Q}^H \mathbf{b}$, Rechenaufwand $4n^2 + O(n)$ bei kompakter Abspeicherung der Reflexionen/Rotationen

③ Rücksubstitution, löse $\mathbf{Rx} = \mathbf{z}$, Rechenaufwand $\frac{1}{2}n(n+1)$

3.2
p. 194

3.2
p. 195

3.2
p. 196

- Die Berechnung der (verallgemeinerten) QR-Zerlegung $A = QR$ mittels Householder-Reflektionen bzw. Givens-Rotationen ist für beliebiges $A \in \mathbb{C}^{m,n}$ stabil.
- Für beliebige reguläre Koeffizientenmatrizen lässt sich ein LGS mittels QR-Zerlegung + Orthogonaltransformation + Rücksubstitution stabil lösen.

Fill-in bei QR-Zerlegung ?

$$A \in \mathbb{C}^{n,n} \text{ mit QR-Zerlegung } A = QR \Rightarrow m(R) \leq m(A) \quad (\rightarrow \text{Def. 3.2.15})$$

Beispiel 94 (QR-basiertes Lösen eines tridiagonalen Gleichungssystems).

Elimination der Subdiagonalen durch $n-1$ sukzessive Givens-Rotationen:

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

MATLAB-Code (c, d, e, b = Spaltenvektoren der Länge n , $n \in \mathbb{N}$, $e(n)$, $c(n)$ nicht benutzt):

MATLAB-CODE : Lösen eines tridiagonalen LGS mit QR-Zerlegung

```
function y = tridiagqr(c,d,e,b)
n = length(d); t = norm(d)+norm(e)+norm(c);
for k=1:n-1
    [R,z] = planerot([d(k);e(k)]);
    if (abs(z(1))/t < eps), error('Matrix singular'); end;
    d(k) = z(1); b(k:k+1) = R*b(k:k+1);
    Z = R*[c(k), 0;d(k+1), c(k+1)];
    c(k) = Z(1,1); d(k+1) = Z(2,1);
    e(k) = Z(1,2); c(k+1) = Z(2,2);
end
A = spdiags([d,[0;c(1:end-1)],[0;0;e(1:end-2)]],[0 1 2],n,n);
y = A\b;
```

Rechenaufwand $O(n)$

Bemerkung 95 (Aufspüren fast singulärer Matrizen).

(Relativ) sehr kleine r_{ii} bei Berechnung der QR-Zerlegung $\leftrightarrow A$ „fast singulär“

3.2.8 Modifikationstechniken

MODIFTECH

Problem: Effizientes Aktualisieren von Faktorisierungen bei „geringfügiger“ Änderung der Matrix, [19, Abschnitt 12.6], [43, Abschnitt 4.9].

3.2.8.1 Rang-1-Modifikationen

RANGIMOD

$$A \in \mathbb{K}^{n,n} \mapsto \tilde{A} := A + uv^H, \quad u, v \in \mathbb{K}^n. \quad (3.2.16)$$

Allgemeine Rang-1-Matrix

Bemerkung 96 (Lösen eines Rang-1-modifizierten LGS).

SHERMORWOOD

Lemma 3.2.22 (Sherman-Morrison-Woodbury-Formel). Für reguläres $A \in \mathbb{K}^{n,n}$ und $U, V \in \mathbb{K}^{n,k}$, $n, k \in \mathbb{N}$, $k \leq n$, gilt

$$(A + UV^H)^{-1} = A^{-1} - A^{-1}U(I + V^HA^{-1}U)^{-1}V^HA^{-1},$$

wenn $I + V^HA^{-1}U$ invertierbar.

3.2
p. 197

3.2
p. 199

Aufgabe: Zu lösen: $\tilde{A}x = b$, wobei LU-Zerlegung $A = LU$ bekannt

Anwendung von Lemma 3.2.22 für $k = 1$:

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

Effiziente Implementierung !



MATLAB-CODE : Lösen eines Rang-1-modifizierten LGS

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

Aufgabe: Effiziente Berechnung der QR-Zerlegung (\rightarrow Abschn. 3.2.7) $A = QR$ von A aus (3.2.16), wenn eine QR-Zerlegung $A = QR$ bekannt.

① Mit $w := Q^Hu$: $A + uv^H = Q(R + ww^H)$

➔ Rechenaufwand $O(n^2)$ (abhängig von Abspeicherung von Q)

3.2
p. 198

3.2
p. 200

② Ziel: $\mathbf{w} \rightarrow \|\mathbf{w}\| \mathbf{e}_1$ durch $n-1$ Givens-Rotationen

$$\mathbf{w} = \begin{pmatrix} * \\ * \\ \vdots \\ * \\ * \\ * \\ * \end{pmatrix} \xrightarrow{\mathbf{G}_{n-1,n}} \begin{pmatrix} * \\ * \\ \vdots \\ * \\ * \\ * \\ 0 \end{pmatrix} \xrightarrow{\mathbf{G}_{n-2,n-1}} \begin{pmatrix} * \\ * \\ \vdots \\ * \\ * \\ 0 \\ 0 \end{pmatrix} \xrightarrow{\mathbf{G}_{n-3,n-2}} \dots \xrightarrow{\mathbf{G}_{12}} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

Beachte: Auswirkung auf \mathbf{R} !

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

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

OHESS

Obere Hessenbergmatrix: Eintrag $(i, j) = 0$, falls $i > j + 1$.

► $\mathbf{A} + \mathbf{u}\mathbf{v}^H = \mathbf{Q}\mathbf{Q}_1^H (\underbrace{\mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^H}_{\text{obere Hessenbergmatrix}})$ mit unitärem $\mathbf{Q}_1 := \mathbf{G}_{12} \dots \mathbf{G}_{n-1,n}$.

► Rechenaufwand $O(n^2)$

③ Sukzessive Givens-Rotationen: $\mathbf{R}_1 + \|\mathbf{w}\|_2 \mathbf{e}_1 \mathbf{v}^H \mapsto$ obere Dreiecksgestalt

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

$$\xrightarrow{\mathbf{G}_{n-2,n-1}} \begin{pmatrix} * & * & \dots & * & * & * & * \\ 0 & * & \dots & * & * & * & * \\ \vdots & & \ddots & & & & \vdots \\ 0 & \dots & 0 & * & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{n-1,n}} \begin{pmatrix} * & * & \dots & * & * & * & * \\ 0 & * & \dots & * & * & * & * \\ \vdots & & \ddots & & & & \vdots \\ 0 & \dots & 0 & * & * & * & * \\ 0 & \dots & 0 & 0 & * & * & * \\ 0 & \dots & 0 & 0 & 0 & * & * \\ 0 & \dots & 0 & 0 & 0 & 0 & * \end{pmatrix} =: \tilde{\mathbf{R}}. \quad (3.2.17) \quad \text{mod:OHI}$$

► Rechenaufwand $O(n^2)$

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

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

Symmetrische und positivitätserhaltende Rang-1-Modifikation:

$$\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n} \mapsto \tilde{\mathbf{A}} := \mathbf{A} + \alpha \mathbf{v}\mathbf{v}^H, \quad \mathbf{v} \in \mathbb{K}^n, \alpha > 0. \quad (3.2.18) \quad \text{mod:2} \quad 3.2$$

3.2
p. 201

p. 203

Aufgabe: Effiziente Berechnung der Cholesky-Zerlegung $\tilde{\mathbf{A}} = \tilde{\mathbf{R}}^H \tilde{\mathbf{R}}$ (\rightarrow Lemma 3.2.13) von $\tilde{\mathbf{A}}$ aus (3.2.18), wenn eine solche $\mathbf{A} = \mathbf{R}^H \mathbf{R}$ für \mathbf{A} bekannt.

① Mit $\mathbf{w} := \mathbf{R}^{-H} \mathbf{v}$: $\mathbf{A} + \alpha \mathbf{v}\mathbf{v}^H = \mathbf{R}^H (\mathbf{I} + \alpha \mathbf{w}\mathbf{w}^H) \mathbf{R}$.

► Rechenaufwand $O(n^2)$ (Rücksubstitution!)

② Idee: Formale Gauss-Elimination: mit $\tilde{\mathbf{w}} = (w_2, \dots, w_n)^T \rightarrow (3.2.5) \quad \text{lu:block}$

$$\mathbf{I} + \alpha \mathbf{w}\mathbf{w}^H = \left(\begin{array}{c|c} 1 + \alpha w_1^2 & \alpha w_1 \tilde{\mathbf{w}}^H \\ \hline \alpha w_1 \tilde{\mathbf{w}} & \mathbf{I} + \alpha \tilde{\mathbf{w}} \tilde{\mathbf{w}}^H \end{array} \right) \rightarrow \left(\begin{array}{c|c} 1 + \alpha w_1^2 & \alpha w_1 \tilde{\mathbf{w}}^H \\ \hline 0 & \mathbf{I} + \alpha^{(1)} \tilde{\mathbf{w}} \tilde{\mathbf{w}}^H \end{array} \right) \quad (3.2.19) \quad \text{mod:rec}$$

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

Analoge Struktur \rightarrow Rekursion

3.2
p. 202

3.2
p. 204

► Berechnung der Cholesky-Zerlegung

$$I + \alpha w w^H = R_1^H R_1.$$

Motiviert durch „Rekursion“
(3.2.19).

→ Rechenaufwand $O(n^2)$
($O(n)$ wenn nur Abspeicherung von $d, s \rightarrow$ (3.2.20))

MATLAB-CODE Cholesky-Zerlegung

```
function [d,s] = roid(alpha,w)
n = length(w);
d = []; s = [];
for i=1:n
    t = alpha*w(i);
    d = [d; sqrt(1+t*w(i))];
    s = [s; t/d(i)];
    % R=[R;zeros(1,i-1),d(i), ...
    %      s(i)*w(i+1:end)'];
    alpha = alpha - s(i)^2;
end
```

bei Rang-1-Modifikation

③ Spezielle Struktur von R_1 :

$$R_1 = \begin{pmatrix} d_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & d_n \end{pmatrix} + \begin{pmatrix} s_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & s_n \end{pmatrix} \begin{pmatrix} 0 & w_2 & w_3 & \cdots & \cdots & w_n \\ 0 & 0 & w_3 & \cdots & \cdots & w_n \\ \vdots & & \ddots & & & \vdots \\ \vdots & & & 0 & w_{n-1} & w_n \\ 0 & \cdots & \cdots & 0 & w_n \\ 0 & \cdots & \cdots & \cdots & 0 \end{pmatrix} \quad (3.2.20)$$

$$R_1 = \begin{pmatrix} d_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & d_n \end{pmatrix} + \begin{pmatrix} s_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & s_n \end{pmatrix} \begin{pmatrix} 0 & 1 & 1 & \cdots & \cdots & 1 \\ 0 & 0 & 1 & \cdots & \cdots & 1 \\ \vdots & & \ddots & & & \vdots \\ \vdots & & & 0 & 1 & 1 \\ 0 & \cdots & \cdots & 0 & 1 \\ 0 & \cdots & \cdots & \cdots & 0 \end{pmatrix} \begin{pmatrix} w_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & w_n \end{pmatrix}$$

► Geschickte Multiplikation

$$\tilde{R} := R_1 R$$

→ Rechenaufwand $O(n^2)$

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

MATLAB-CODE : roudchol

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

MATLAB-Funktion: `R = cholupdate(R,v);`

3.2.8.2 Hinzufügen einer Spalte

$$A \in \mathbb{K}^{m,n} \mapsto \tilde{A} = [a_1, \dots, a_{k-1}, v, a_k, \dots, a_n], \quad v \in \mathbb{K}^m. \quad (3.2.21)$$

Gegeben: QR-Zerlegung $A = QR$, $Q \in \mathbb{K}^{m,m}$ unitär, $R \in \mathbb{K}^{m,n}$ obere Dreiecksmatrix.

Aufgabe: Effiziente Berechnung der QR-Zerlegung $\tilde{A} = \tilde{Q} \tilde{R}$ von \tilde{A} aus (3.2.21), $\tilde{Q} \in \mathbb{K}^{m,m}$ unitär, $\tilde{R} \in \mathbb{K}^{m,n+1}$ obere Dreiecksmatrix.

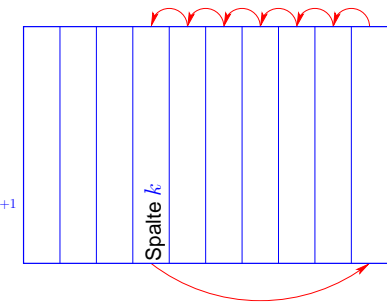
Idee: Aufgabe einfach, falls $k = n+1$ (Anfügen einer Spalte von rechts)

► \exists Teilzyklische Spaltenpermutation

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

~ Permutationsmatrix

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



$$\tilde{A} \rightarrow A_1 = \tilde{A} P = [a_1, \dots, a_n, v] = Q \begin{bmatrix} R & Q^H v \end{bmatrix} = Q \begin{pmatrix} R & Q^H v \end{pmatrix}$$

Spalte $Q^H v$
Fall $m > n+1$

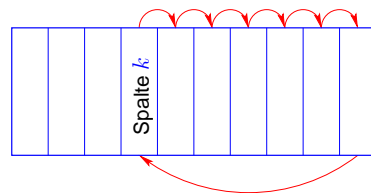
① Falls $m > n+1$: \exists orthogonale Transformation $Q_1 \in \mathbb{K}^{m,m}$ (Householder-Reflektion) mit

$$Q_1 Q^H v = \begin{pmatrix} * \\ \vdots \\ * \\ 0 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} n+1 \\ \\ m-n-1 \end{matrix} \quad \rightarrow \quad Q_1 Q^H A_1 = \begin{pmatrix} * & \cdots & \cdots & * \\ 0 & * & & * \\ \vdots & & \ddots & \vdots \\ & & & * & * \\ \vdots & & & & * \\ 0 & \cdots & \cdots & 0 & \\ \vdots & & & \vdots & \\ 0 & \cdots & \cdots & 0 & \end{pmatrix} \begin{matrix} n+1 \\ \\ m-n-1 \end{matrix}$$

→ Rechenaufwand $O(m-n)$ (Eine Reflektion)

Inverse Permutation:

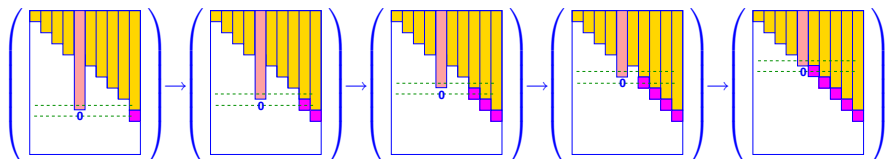
~ Multiplikation mit \mathbf{P}^H von rechts



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

② $n+1-k$ sukzessive Givens-Rotationen \Rightarrow obere Dreiecksmatrix $\tilde{\mathbf{R}}$

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



--- $\hat{=}$ Zielzeilen der Givens-Rotationen, \blacksquare $\hat{=}$ neue Einträge $\neq 0$

\Rightarrow Rechenaufwand $O((n-k)^2)$

3.2.8.3 Hinzufügen einer Zeile

$$\mathbf{A} \in \mathbb{K}^{m,n} \mapsto \tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{a}_{1,\cdot} \\ \vdots \\ \mathbf{a}_{k-1,\cdot} \\ \mathbf{v}^T \\ \mathbf{a}_{k,\cdot} \\ \vdots \\ \mathbf{a}_{m,\cdot} \end{bmatrix}, \quad \mathbf{v} \in \mathbb{K}^n. \quad (3.2.22) \text{ mod:row}$$

Gegeben: QR-Zerlegung $\mathbf{A} = \mathbf{Q}\mathbf{R}$, $\mathbf{Q} \in \mathbb{K}^{m+1,m+1}$ unitär, $\mathbf{R} \in \mathbb{K}^{m,n}$ obere Dreiecksmatrix.

Aufgabe: Effiziente Berechnung der QR-Zerlegung $\tilde{\mathbf{A}} = \tilde{\mathbf{Q}}\tilde{\mathbf{R}}$ von $\tilde{\mathbf{A}}$ aus (3.2.22), $\tilde{\mathbf{Q}} \in \mathbb{K}^{m+1,m+1}$ unitär, $\tilde{\mathbf{R}} \in \mathbb{K}^{m+1,n+1}$ obere Dreiecksmatrix.

- ① \exists Teilzyklische Zeilenpermutation $m+1 \leftarrow k, i \leftarrow i+1, i = k, \dots, m:$
 \rightarrow unitäre Permutationsmatrix $\mathbf{P} \in \{0,1\}^{m+1,m+1}$

$$\mathbf{P}\tilde{\mathbf{A}} = \begin{pmatrix} \mathbf{A} \\ \mathbf{v}^T \end{pmatrix} \blacktriangleright \begin{pmatrix} \mathbf{Q}^H & 0 \\ 0 & 1 \end{pmatrix} \mathbf{P}\tilde{\mathbf{A}} = \begin{pmatrix} \mathbf{R} \\ \mathbf{v}^T \end{pmatrix} = \begin{pmatrix} \text{yellow blocks} \\ \text{red block} \end{pmatrix}.$$

Fall $m = n$

3.2
p. 209

3.2
p. 211

② Obere Dreiecksgestalt durch sukzessive Givens-Rotationen:

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

$$\dots \xrightarrow{\mathbf{G}_{m-2,m}} \begin{pmatrix} * & \dots & * & \dots & * \\ 0 & * & * & & \vdots \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & * & * \\ 0 & \dots & 0 & * & * \end{pmatrix} \xrightarrow{\mathbf{G}_{m-1,m}} \begin{pmatrix} * & \dots & * & \dots & * \\ 0 & * & * & & \vdots \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & * & * \\ 0 & \dots & 0 & * & * \end{pmatrix} := \tilde{\mathbf{R}} \quad (3.2.23) \text{ mod:row}$$

③ Mit $\mathbf{Q}_1 = \mathbf{G}_{m-1,m} \dots \mathbf{G}_{1,m}$

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

3.2 \Rightarrow Ähnliche Update-Algorithmen für Streichen von Zeilen/Spalten
p. 210

3.3
p. 212

3.3 Numerische Berechnung von Eigenwerten und Eigenvektoren

[File: section-numerische-berechnung-von-eigenwerten-und-eigenvektoren.tex, SVN: section-numerische-berechnung-von-eigenwerten-und-eigenvektoren.te

NUMERISCH

Beispiel 97 (Resonanzen elektrischer Netzwerke).

Gegeben: Grosses elektrisches Netzwerk mit linearen Bauelementen
Gesucht: Resonanzfrequenzen

Bsp. 66: Kirchhoffsche Knotenregel + Bauelementgleichungen \Rightarrow Netzwerkmatrix

$A(\omega) := W + i\omega C - i\omega^{-1}S$, $W, C, S \in \mathbb{R}^{n,n}$ s.p.d. (\rightarrow Def. 3.2.8).

Resonanzfrequenzen = $\omega \in \{\omega \in \mathbb{R}: A(\omega) \text{ singular}\}$.

Wie können diese ω numerisch berechnet werden (\rightarrow quadratisches Eigenwertproblem).

3.3.1 Theorie und Kondition

Definition 3.3.1 (Eigenwerte und Eigenvektoren).

- $\lambda \in \mathbb{C}$ **Eigenwert** (engl. eigenvalue) von $A \in \mathbb{K}^{n,n}$ $\Leftrightarrow \det(\lambda I - A) = 0$
Charakteristisches Polynom $\chi(\lambda)$
- Spektrum** von $A \in \mathbb{K}^{n,n}$: $\sigma(A) := \{\lambda \in \mathbb{C}: \lambda \text{ Eigenwert von } A\}$
- Eigenraum** (engl. eigenspace) zum Eigenwert $\lambda \in \sigma(A)$: $\text{Eig}_A(\lambda) := \text{Ker}(\lambda I - A)$
- $x \in \text{Eig}_A(\lambda) \Rightarrow x$ ist **Eigenvektor**
- Geometrische **Vielfachheit** (engl. multiplicity) eines Eigenwerts $\lambda \in \sigma(A)$:
 $m(\lambda) := \dim \text{Eig}_A(\lambda)$

Lemma 3.3.2. Das Spektrum einer Matrix ist invariant unter Ähnlichkeitstransformationen:

$\forall A \in \mathbb{K}^{n,n}: \sigma(S^{-1}AS) = \sigma(A) \quad \forall \text{ regulären } S \in \mathbb{K}^{n,n}.$

Theorem 3.3.3 (Schursches Lemma).

$\forall A \in \mathbb{K}^{n,n}: \exists U \in \mathbb{C}^{n,n} \text{ unitär. } U^H A U = T \text{ mit } T \in \mathbb{C}^{n,n} \text{ obere Dreiecksmatrix.}$

Korollar 3.3.4 (Unitäre Diagonalisierbarkeit normaler Matrizen).

$A \in \mathbb{K}^{n,n}, AA^H = A^H A: \exists U \in \mathbb{C}^{n,n} \text{ unitär. } U^H A U = \text{diag}(\lambda_1, \dots, \lambda_n), \lambda_i \in \mathbb{C}.$

Beachte: $\lambda_1, \dots, \lambda_n$ = Eigenwerte von A
Spalten von U = Orthonormalbasis aus Eigenvektoren von A

Matrizen $A \in \mathbb{K}^{n,n}$ mit $AA^H = A^H A$ heissen normal.

Normal sind z.B.

- Hermiteische Matrizen: $A^H = A \Rightarrow \sigma(A) \subset \mathbb{R}$
- Unitäre Matrizen: $A^H = A^{-1} \Rightarrow |\sigma(A)| = 1$
- Schiefhermitesche Matrizen: $A = -A^H \Rightarrow \sigma(A) \subset i\mathbb{R}$

Eigenwertprobleme:

- Gegeben $A \in \mathbb{K}^{n,n}$ finde $\sigma(A)$.
- Gegeben $A \in \mathbb{K}^{n,n}$ finde $\sigma(A)$ und alle Eigenvektoren.
- Gegeben $A \in \mathbb{K}^{n,n}$ finde einige Eigenwerte und zugehörige Eigenvektoren.

Theorem 3.3.5 (Störungssatz für Eigenwerte). Zu $A \in \mathbb{K}^{n,n}$ gebe es $S \in \mathbb{C}^{n,n}$ regulär, so dass $S^{-1}AS = D, D$ diagonal. Dann

$\forall E \in \mathbb{K}^{n,n}, \lambda \in \sigma(A + E): \exists \lambda' \in \sigma(A): |\lambda - \lambda'| \leq \text{cond}_2(S) \|E\|_2$

Spalten von S = Eigenvektoren von A

Falls Eigenvektoren „fast“ linear abhängig ($\Rightarrow \text{cond}_2(S)$ gross), kann Eigenwertproblem sehr schlecht konditioniert sein

Falls A normal \Rightarrow EWP gut (absolut) konditioniert

Verfeinerte Analyse [19, Sect. 7.2.2], $\mathbf{A} \in \mathbb{K}^{n,n}$:

• $\lambda \in \mathbb{C}$ einfache Nullstelle von $\lambda \mapsto \det(\lambda \mathbf{I} - \mathbf{A})$ (einfacher Eigenwert)

• $\mathbf{x} \in \mathbb{C}^n$, $\|\mathbf{x}\|_2 = 1$ zugehöriger Eigenvektor (\rightarrow Def. 3.3.1),

• $\mathbf{y} \in \mathbb{C}^n$, $\|\mathbf{y}\|_2 = 1$, **LINKSEIGENVEKTOR** zu λ , d.h. $\mathbf{y}^H \mathbf{A} = \lambda \mathbf{y}^H \Leftrightarrow \mathbf{A}^H \mathbf{y} = \lambda \mathbf{y}$.

► absolute Konditionszahl von λ bzgl. Euklidischer Matrixnorm & Störungen von $\mathbf{A} =$

$$\kappa_{\text{abs}}(\lambda) := |\mathbf{y}^H \mathbf{x}|^{-1}.$$

Beispiel 98 (Störungsanfälligkeit von Eigenwerten).

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 4.001 \end{pmatrix} \rightarrow \begin{array}{l} \text{E.v. von } \mathbf{A}: \left\{ \begin{pmatrix} 1.0000 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0.5547 \\ 0.8321 \\ 0 \end{pmatrix}, \begin{pmatrix} 0.5547 \\ 0.8321 \\ 0.0002 \end{pmatrix} \right\}, \\ \text{E.v. von } \mathbf{A}^T: \left\{ \begin{pmatrix} 0.8285 \\ -0.5523 \\ 0.0920 \end{pmatrix}, \begin{pmatrix} 0 \\ 0.0002 \\ -1.0000 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1.0000 \end{pmatrix} \right\} \end{array}$$

► $\kappa_{\text{abs}}(\lambda_1) = 1.2069722$, $\kappa_{\text{abs}}(\lambda_2) = 6009.25224$, $\kappa_{\text{abs}}(\lambda_3) = 6009.19059$

$$\mathbf{E} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 10^{-3} & 0 & 0 \end{pmatrix} \rightarrow \sigma(\mathbf{A} + \mathbf{E}) = \{1.000111, 4.0581822, 3.9427066\}.$$

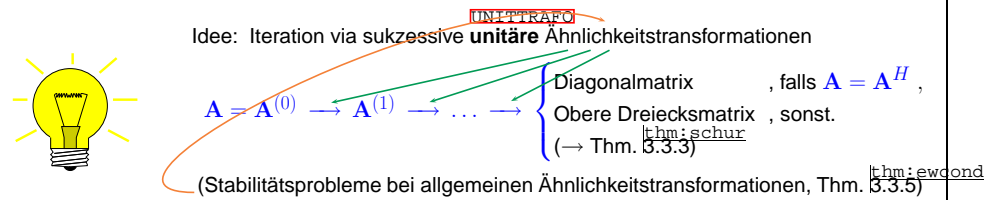
Bemerkung 99. Die (absolute) Kondition eines Eigenvektors einer diagonalisierbaren Matrix $\mathbf{A} \in \mathbb{K}^{n,n}$ mit n verschiedenen Eigenwerten ist schlecht, wenn

- die Kondition des zugehörigen Eigenwertes λ^* schlecht ist,
- oder wenn λ^* „schlecht separiert“ $\min\{|\lambda^* - \lambda|, \lambda \in \sigma(\mathbf{A}) \setminus \{\lambda^*\}\} \ll |\lambda^*|$.

3.3.2 Transformationsmethoden

TRANSFORMETE

Bemerkung 100. Alle Verfahren zur numerischen Berechnung von Eigenwerten/Eigenvektoren sind **Iterationsverfahren** \rightarrow Abschn. 2.1



(3.3.2) **QR-Algorithmus** (mit Shift): kubische Konvergenzordnung [19, Sect. 7.5,8.2]

Beispiel 101 (Unitäre Ähnlichkeitstransformation auf Tridiagonalgestalt).

$$\mathbf{A} = \begin{pmatrix} d_1 & b_1 & & & \\ b_1 & d_2 & & & \\ & & \ddots & & \\ & & & d_{n-1} & b_{n-1} \\ \alpha & & & b_{n-1} & d_n \end{pmatrix}, \quad \mathbf{d} \in \mathbb{R}^n, \mathbf{b} \in \mathbb{R}^{n-1}, \alpha \in \mathbb{R}.$$

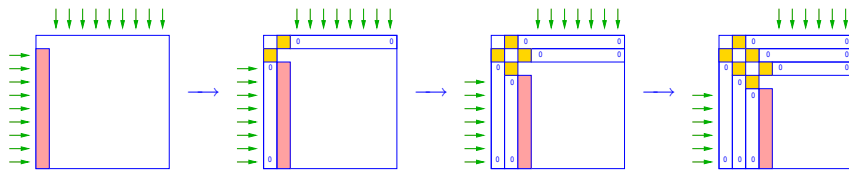
1. Schritt: Annulliere Matrixeinträge an Positionen $(n, 1)$ und $(1, n)$ durch Givens-Rotationen von rechts und links

► MATLAB-Animation

$$\mathbf{G}_{2,n} \mathbf{A} \mathbf{G}_{2,n}^T = \begin{pmatrix} * & * & & & 0 \\ * & * & & & * \\ * & * & & & * \\ & & \ddots & & \\ & & & * & * \\ 0 & * & * & & * \end{pmatrix} \rightarrow \begin{pmatrix} * & * & & & 0 \\ * & * & & & 0 \\ * & * & & & * \\ & & \ddots & & * \\ & & & * & * \\ 0 & 0 & * & * & * \end{pmatrix} \rightarrow \dots$$

► Rechenaufwand, Speicheraufwand = $O(n^3)$! (da Dünnbesetztheit zerstört)

Allgemeines Vorgehen: Sukzessive Householder-Ähnlichkeitstransformationen
 (→ $\hat{=}$ betroffene Zeilen/Spalten, $\hat{=}$ Zielvektor)



Endlich viele **unitäre** Householder/Givens-Ähnlichkeitstransformationen können jede Hermitesche Matrix auf Tridiagonalform bringen (Rechenaufwand $O(n^3)$)
 → Startwert $A^{(0)}$ für iteratives Verfahren

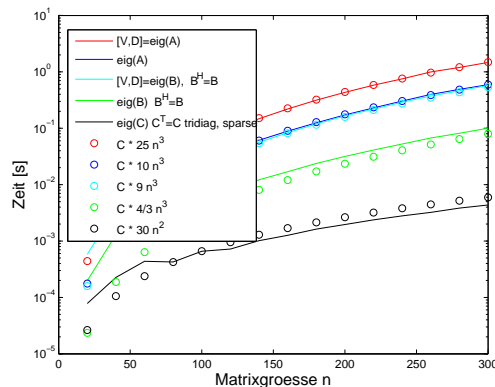
MATLAB-Funktionen: • Eigenwerte/normalisierte Eigenvektoren mit $[V,D] = \text{eig}(A)$;
 • Schur-Zerlegung mit $[U,T] = \text{schur}(A)$;

RECHNAUFWAND
 Rechenaufwand (elementare Operationen) für $\text{eig}()$ (Ergebnis auf Maschinengenauigkeit):

Eigenwerte und Eigenvektoren von $A \in \mathbb{K}^{n,n}$ $\sim 25n^3 + O(n^2)$
 Nur Eigenwerte von $A \in \mathbb{K}^{n,n}$ $\sim 10n^3 + O(n^2)$
 Eigenwerte und Eigenvektoren von $A = A^H \in \mathbb{K}^{n,n}$ $\sim 9n^3 + O(n^2)$
 Nur Eigenwerte von $A = A^H \in \mathbb{K}^{n,n}$ $\sim \frac{4}{3}n^3 + O(n^2)$
 Nur Eigenwerte von **tridiagonalem** $A = A^H \in \mathbb{K}^{n,n}$ $\sim 30n^2 + O(n)$
 Beachte: Keine MATLAB-Implementierung von eig für dünnbesetzte Matrizen!
 Ausnahme: $d = \text{eig}(A)$ für **Hermitesche** dünnbesetzte Matrizen

Beispiel 102. (Komplexität von MATLAB-Eigenlösern)

Experiment beschreiben und Verwe auf MATLAB-Code einfügen. Legen denbeschriftung kürzen und Erklärung im Text.



VERALLGEWPE

Verallgemeinertes Eigenwertproblem (VEWP) (engl. *generalized eigenvalue problem*):

Gegeben: $A, B \in \mathbb{K}^{n,n}$

Gesucht: Verallgemeinerte Eigenwerte $\lambda \in \mathbb{C}$: $Ax = \lambda Bx$
 Verallgemeinerte Eigenvektoren $x \in \mathbb{K}^n \setminus \{0\}$

MATLAB-Funktion: • $d = \text{eig}(A,B)$; (→ nur Eigenwerte)
 • $[V,D] = \text{eig}(A,B)$; (→ Eigenwerte + -vektoren)

Insiderinformation: **QZALGO** MATLAB benutzt QZ-Algorithmus **SOL89** [19, Sect. 7.7] (für vollbesetzte A, B)

Falls B regulär > VEWP formal äquivalent zu EWP für $B^{-1}A$
 (kein Rezept zur effizienten numerischen Implementierung!)

Bemerkung 103. Falls $B = B^H$ s.p.d. mit Cholesky-Zerlegung $B = R^H R$

$$Ax = \lambda Bx \Leftrightarrow \tilde{A}y = \lambda y \quad \text{mit } \tilde{A} := R^{-H} A R^{-1}, y := Rx.$$

3.3
 p. 221 → Algorithmisch relevant, da billiger als QZ.

Beispiel 104 (Quadratisches Eigenwertproblem). (Fortsetzung von Bsp. 97)

Finde $x \neq 0, \omega \neq 0$, so dass

$$A(\omega)x = (W + i\omega C - i\omega^{-1}S)x = 0. \quad (3.3.3) \quad \text{eq:netwo}$$

Substitution: $y = -i\omega^{-1}x$ [44, Sect. 3.4]:

$$(3.3.3) \Leftrightarrow \begin{pmatrix} W & S \\ I & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = -i\omega \begin{pmatrix} C & 0 \\ 0 & I \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

> Verallgemeinertes Eigenwertproblem

Übung 3.4. (i) Welche Bedingung muss der Vektor $a \in \mathbb{R}^{n-1}$ erfüllen, damit die Matrix

$$M := \begin{pmatrix} 1 & a^T \\ a & I \end{pmatrix} \in \mathbb{R}^{n,n}$$

positiv definit ist.

(ii) Schreibe eine effiziente MATLAB-Routine zur Berechnung der Eigenwerte für das verallgemeinerte Eigenwertproblem

$$\lambda \in \mathbb{R}, x = (x_1, \dots, x_n)^T \in \mathbb{R}^n \setminus \{0\}: x_1 e_1 = \lambda Mx,$$

wobei $e_1 := (1, 0, \dots, 0)^T \in \mathbb{R}^n$.

3.3
 p. 222

3.3
 p. 223

◇

3.4
 p. 224

3.4.1 Potenzmethoden

POTMETH

Spezielle Eigenwertprobleme für $A \in \mathbb{K}^{n,n}$:

- Finde den **betragsgrössten** Eigenwert (+ Eigenvektor) von A .
- Finde den **betragskleinsten** Eigenwert (+ Eigenvektor) von A .
- Zu $\alpha \in \mathbb{C}$ finde $\lambda \in \sigma(A)$ mit $|\alpha - \lambda| = \min\{|\alpha - \mu|, \mu \in \sigma(A)\}$.



Idee für diagonalisierbares $A \in \mathbb{K}^{n,n}$: $S^{-1}AS = \text{diag}(\lambda_1, \dots, \lambda_n)$

$$z = \sum_{j=1}^n \zeta_j s_{\cdot,j} \Rightarrow A^k z = \sum_{j=1}^n \zeta_j \lambda_j^k s_{\cdot,j}$$

Falls $|\lambda_1| < |\lambda_2| < \dots < |\lambda_n|$, $\|s_{\cdot,j}\|_2 = 1, j = 1, \dots, n$, $\zeta_n \neq 0$

$$\frac{A^k z}{\|A^k z\|} \rightarrow \pm s_{\cdot,n} = \text{Eigenvektor zu } \lambda_n \text{ für } k \rightarrow \infty. \quad (3.4.1) \quad \text{ewp:pm}$$

► Motiviert **direkte Potenzmethode** (engl. *power method*): Iterationsverfahren (\rightarrow Abschn. 2.1)

Startwert: $z^{(0)}$ „beliebig“,

Iterationsvorschrift: $w := Az^{(k-1)}$, $z^{(k)} := \frac{w}{\|w\|_2}$, $k = 1, 2, \dots$

Rechenaufwand: $1 \times \text{Matrix} \times \text{Vektor}$ pro Schritt \rightarrow **Billig für dünnbesetzte Matrizen**

Wie bekommt man den Eigenwert?

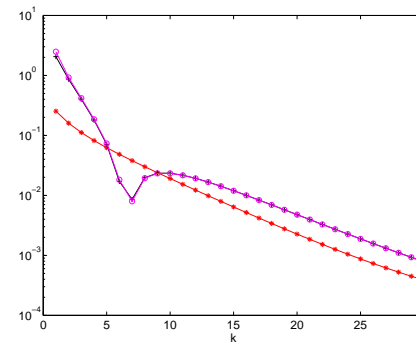
① Asymptotisch aus (3.4.1) $\xrightarrow{\text{ewp:pm}} Az^{(k)} \approx \lambda_n z^{(k)} \Rightarrow |\lambda_n| \approx \frac{\|Az^{(k)}\|}{\|z^{(k)}\|}$

② Idee (für $A = A^H$): $\lambda_n \approx \argmin_{\theta \in \mathbb{R}} \|Az^{(k)} - \theta z^{(k)}\|_2^2 \Rightarrow \lambda_n \approx \frac{(z^{(k)})^H Az^{(k)}}{\|z^{(k)}\|^2}$

Definition 3.4.1. Für $A \in \mathbb{K}^{n,n}$, $u \in \mathbb{K}^n$ ist der **Rayleigh-Quotient** gegeben durch

$$\rho_A(u) := \frac{u^H A u}{u^H u}$$

Beispiel 105 (Direkte Potenzmethode).



```
n = length(d);
S = triu(diag(n:-1:1,0))+...
    ones(n,n);
A = S*diag(d,0)*inv(S);
```

○ : Fehler $|\lambda_n - \rho_A(z^{(k)})|$
 * : Fehlernorm $\|z^{(k)} - s_{\cdot,n}\|$
 < : $\left| \lambda_n - \frac{\|Az^{(k-1)}\|_2}{\|z^{(k-1)}\|_2} \right|$

$z^{(0)}$ = Zufallsvektor

Experimente:

- 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 $d = 1 - 2.^{-(1:0.5:5)'}$; $\Rightarrow |\lambda_{n-1}| : |\lambda_n| = 0.9866$

3.4
p. 225

$$\rho_{EV}^{(k)} := \frac{\|z^{(k)} - s_{\cdot,n}\|}{\|z^{(k-1)} - s_{\cdot,n}\|},$$

$$\rho_{EW}^{(k)} := \frac{|\rho_A(z^{(k)}) - \lambda_n|}{|\rho_A(z^{(k-1)}) - \lambda_n|}.$$

	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

Vermutung:

Lineare Konvergenz (\rightarrow Def. 2.1.6)

Theorem 3.4.2 (Konvergenz der direkten Potenzmethode). $A \in \mathbb{K}^{n,n}$ habe einen betragsgrössten Eigenwert $\lambda_n > 0$ mit algebraischer Vielfachheit 1. Seien v, y Rechts- bzw. Linkseigenvektoren von A zu λ_n mit $\|y\|_2 = \|v\|_2 = 1$. Dann konvergieren

$$\|Az^{(k)}\|_2 \rightarrow \lambda_n, \quad z^{(k)} \rightarrow \pm v \quad \text{linear mit Rate } \frac{|\lambda_{n-1}|}{|\lambda_n|},$$

wobei $z^{(k)}$ die Iterierten der direkten Potenzmethode sind und $y^H z^{(0)} \neq 0$ vorausgesetzt wird.

3.4
p. 226

3.4
p. 227

3.4
p. 228

Bemerkung 106. Rundungsfehler $\triangleright \mathbf{y}^H \mathbf{z}^{(0)} \neq 0$ „numerisch“ immer erfüllt

Bemerkung 107. Abbruchkriterium für Potenziteration? (→ Abschn. 2.1) sec:iterationsverfahren

„Relative Änderung“ $\leq \text{tol}$:

$$\begin{cases} \|\mathbf{z}^{(k)} - \mathbf{z}^{(k-1)}\| \leq \text{tol}, \\ \left| \frac{\|\mathbf{A}\mathbf{z}^{(k)}\|}{\|\mathbf{z}^{(k)}\|} - \frac{\|\mathbf{A}\mathbf{z}^{(k-1)}\|}{\|\mathbf{z}^{(k-1)}\|} \right| \leq (1/L - 1)\text{tol} \end{cases} \quad \text{vgl. (2.1.2).}$$

Geschätzte Konvergenzrate \triangle

Wie findet man den **betragskleinsten** Eigenwert (+ Eigenvektor) von \mathbf{A} ?



Falls $\mathbf{A} \in \mathbb{K}^{n,n}$ regulär:

Betragskleinster EW von $\mathbf{A} = (\text{Betragsgrösster EW von } \mathbf{A}^{-1})^{-1}$

► Direkte Potenzmethode für $\mathbf{A}^{-1} = \text{Inverse Iteration}$
 Verallgemeinerung auf „Zu $\alpha \in \mathbb{C}$ finde $\lambda \in \sigma(\mathbf{A})$ mit $|\alpha - \lambda| = \min\{|\alpha - \mu|, \mu \in \sigma(\mathbf{A})\}$ “:
Inverse Iteration mit Shift:

$$\mathbf{z}^{(0)} \text{ beliebig, } \mathbf{w} = (\mathbf{A} - \alpha \mathbf{I})^{-1} \mathbf{z}^{(k-1)}, \quad \mathbf{z}^{(k)} := \frac{\mathbf{w}}{\|\mathbf{w}\|_2}, \quad k = 1, 2, \dots \quad (3.4.2)$$

Dabei: $(\mathbf{A} - \alpha \mathbf{I})^{-1} \mathbf{z}^{(k-1)} \hat{=}$ Löse $(\mathbf{A} - \alpha \mathbf{I}) \mathbf{w} = \mathbf{z}^{(k-1)}$ basierend auf Gaußelimination:
 einmalige LU-Zerlegung von $\mathbf{A} - \alpha \mathbf{I}$!

Und was passiert, wenn „versehentlich“ $\alpha \in \sigma(\mathbf{A})$ (→ $\mathbf{A} - \alpha \mathbf{I}$ singulär)?

„Magie der Gaußelimination“:

Die Gauss-Elimination mit Spaltenpivotsuche (bzw. LU-Zerlegung/Cholesky-Zerlegung + Vorwärtsrechnung + Rücksubstitution) angewandt auf das LGS $\mathbf{A}\mathbf{x} = \mathbf{b}$ berechnet (fast immer → Bsp. 175) ein $\tilde{\mathbf{x}}$ für das das **Residuum** $\mathbf{r} := \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$ „klein“ ist: $\|\mathbf{r}\| \approx \epsilon$!

Sei $\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n}$ mit $\mathbf{S}^H \mathbf{A} \mathbf{S} = \text{diag}(\lambda_1, \dots, \lambda_n)$, \mathbf{S} unitär, $0 < |\lambda_1| \ll |\lambda_j|, j = 2, \dots, n$

Eigenvektorentwicklung $\mathbf{b} = \sum_{j=1}^n \beta_j \mathbf{s}_{\cdot,j}, \quad \tilde{\mathbf{x}} = \sum_{j=1}^n \xi_j \mathbf{s}_{\cdot,j} \quad \triangleright \quad \mathbf{r} := \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}} = \sum_{j=1}^n (\beta_j - \lambda_j \xi_j) \mathbf{s}_{\cdot,j}$

$$\|\mathbf{r}\|_2 \text{ „klein“} \quad \triangleright \quad \begin{cases} \xi_j \approx \frac{\beta_j}{\lambda_j}, & j = 2, \dots, n, \\ \xi_1 \text{ „gross“}. \end{cases}$$

► Für $\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n}$: Rundungsfehler bei der Gauss-Elimination können nur grosse Fehler in Richtung des Eigenvektors zum betragskleinsten Eigenwert zur Folge haben.

Bemerkung 108. Falls tatsächlich einmal Pivotelement = 0 \triangleright ersetze es durch **eps**!

Thm. 3.4.2 \triangleright Für inverse Iteration (mit Shift):
 Asymptotisch lineare Konvergenz, Rayleigh-Quotient $\rightarrow \lambda_j$ mit Rate

$$\frac{|\lambda_j - \alpha|}{\min\{|\lambda_i - \alpha|, i \neq j\}} \quad \text{mit } \lambda_j \in \sigma(\mathbf{A}), \quad |\alpha - \lambda_j| \leq |\alpha - \lambda| \quad \forall \lambda \in \sigma(\mathbf{A}).$$

► Sehr schnell für $\alpha \approx \lambda_j$!



Idee:

A posteriori Anpassung des Shifts
 Benutze $\alpha := \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})$ im k . Schritt

3.4
p. 229

3.4
p. 231

RAYLEIGHQUOTIENT
 Rayleigh-
 Quotienten-
 Iteration

MATLAB-CODE: Rayleigh-Quotienten-Iteration

```
function [z,lmin] = rqui(A,tol,maxit)
alpha = 0; n = size(A,1);
z = rand(size(A,1),1); z = z/norm(z);
for i=1:maxit
    z = (A-alpha*speye(n))\z;
    z = z/norm(z); lmin=dot(A*z,z);
    if (abs(alpha-lmin) < tol), break; end;
    alpha = lmin;
end
```

(für normale $\mathbf{A} \in \mathbb{K}^{n,n}$)

erhält Dünnbesetztheit,
 siehe Abschnitt 3.1.2

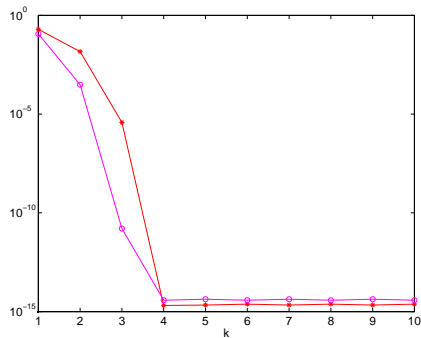
(3.4.3) **rqui**

Nachteil: Kein Recycling von LU-Zerlegungen möglich (im Vergleich zu (3.4.2)) sec:invitshift

Beispiel 109 (Rayleigh-Quotienten-Iteration).

3.4
p. 230

3.4
p. 232



```
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';
```

○ : $|\lambda_{\min} - \rho_A(\mathbf{z}^{(k)})|$
 * : $\|\mathbf{z}^{(k)} - \mathbf{x}_{\cdot,j}\|, \lambda_{\min} = \lambda_j$

k	$ \lambda_{\min} - \rho_A(\mathbf{z}^{(k)}) $	$\ \mathbf{z}^{(k)} - \mathbf{x}_{\cdot,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 3.4.3. Für $\mathbf{A} = \mathbf{A}^H$ konvergiert $\rho_A(\mathbf{z}^{(k)})$ für die iterierten $\mathbf{z}^{(k)}$ aus der Rayleigh-Quotienten-Iteration (3.4.3) lokal kubisch gegen den betragskleinsten Eigenwert.

Was tun, wenn man mehrere Eigenwerte $\lambda_{n-k}, \lambda_{n-k+1}, \dots, \lambda_n, k \ll n$, und zugehörige Eigenvektoren von $\mathbf{A} \in \mathbb{K}^{n,n}$ benötigt ?

DIRPOTITUR

Direkte Potenziteration im Unterraum (für Hermitesche Matrizen $\mathbf{A} \in \mathbb{K}^{n,n}$)

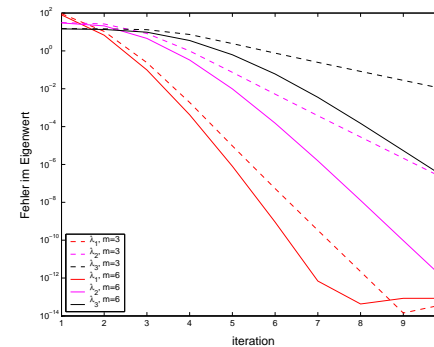
```
MATLAB-CODE : Potenzmethode im Unterraum
function [V,ev] = spowit(A,k,m)
n = size(A,1); V = rand(n,m); d = zeros(k,1);
for i=1:maxit
    V=A*V;
    [Q,R] = qr(V,0);
    T=Q'*A*Q;
    [S,D] = eig(T); [l,perm] = sort(-abs(diag(D)));
    V = Q*S(:,perm);
    if (norm(d+1(1:k)) < tol), break; end;
    d = -1(1:k);
end
V = V(:,1:k); ev = diag(D(perm(1:k),perm(1:k)));
```

RITZPROJ
Ritz-Projektion

Verallgemeinerung der Normierung auf $\|\mathbf{z}\|_2 = 1$

➔ Analog: Unterraumvarianten der inversen Iteration/Rayleigh-Quotienten-Iteration

Beispiel 110 (Direkte Potenziteration im Unterraum).



S.p.d. Testmatrix: $a_{ij} := \min\{\frac{i}{j}, \frac{j}{i}\}$
 $n=200$; $\mathbf{A} = \text{gallery('lehmer',n)}$;
 „Anfangsmatrix“:
 $\mathbf{V} = \text{eye}(n,m)$;

- Beobachtung:
Lineare Konvergenz der Eigenwerte
- Wahl $m > k$ beschleunigt Konvergenz der Eigenwerte

3.4.2 Vorkonditionierte inverse Iteration

VORKONDIINV11

3.4
p. 233

Aufgabe: Für symmetrisch positiv definites $\mathbf{A} \in \mathbb{K}^{n,n}$ (→ Def. 3.2.8) finde betragskleinsten Eigenwert (+ zugehörigen Eigenvektor).

➤ Optionen: Inverse Iteration, Rayleigh-Quotienten-Iteration (3.4.3)

Problem: Trotz sparse-Eliminationstechniken (→ 3.2.6): Direkte LU-Zerlegung (→ Abschn. 3.2.3) oft unmöglich für $n \gg 1$.

Beachte: Inverse Iteration ➤ Näherungslösung des LGS in (3.4.2) ausreichend !



Idee: (für) Inverse Iteration ohne Shift

Statt Lösung von $\mathbf{A}\mathbf{w} = \mathbf{z}^{(k-1)}$: $\mathbf{w} = \mathbf{B}^{-1}\mathbf{z}^{(k-1)}$ mit "billiger" s.p.d. approximativer Inversen $\mathbf{B}^{-1} \approx \mathbf{A}^{-1}$



Kann man in der inversen Iteration \mathbf{A}^{-1} einfach durch \mathbf{B}^{-1} ersetzen ?

NEIN, denn kleinster Eigenwert von \mathbf{B} interessiert nicht!



Ersetzung $\mathbf{A}^{-1} \rightarrow \mathbf{B}^{-1}$ nur möglich bei Anwendung auf ein Residuum
RESIDUUM
 Residuum = Grösse, die $\rightarrow 0$, wenn Konvergenz gegen gesuchte Lösung

3.4
p. 234

3.4
p. 235

3.4
p. 236

Kandidat: Residuum für Eigenwertproblem $\mathbf{Ax} = \lambda \mathbf{x}$:

$$\mathbf{r} := \mathbf{Az} - \rho_{\mathbf{A}}(\mathbf{z})\mathbf{z}, \quad \rho_{\mathbf{A}}(\mathbf{z}) = \text{Rayleigh-Quotient} \rightarrow \text{Def. 3.4.1.}^{\text{def:rq}}$$

Beachte: Nur Richtung von $\mathbf{A}^{-1}\mathbf{z}$ relevant in inverser Iteration $(3.4.2)^{\text{ewp:invitshift}}$

$$(\mathbf{A}^{-1}\mathbf{z}) \parallel (\mathbf{z} - \mathbf{A}^{-1}(\mathbf{Az} - \rho_{\mathbf{A}}(\mathbf{z})\mathbf{z})) \Rightarrow \text{definiert gleiche nächste Iterierte!}$$

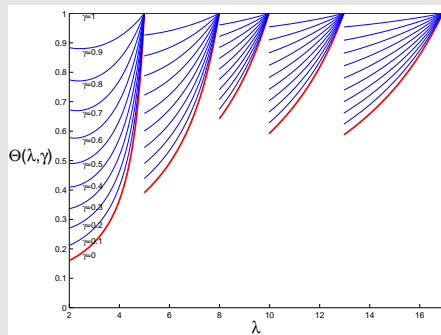
[Vorkonditionierte inverse Iteration für s.p.d. \mathbf{A}]

$$\mathbf{z}^{(0)} \text{ beliebig, } \mathbf{w} = \mathbf{z}^{(k-1)} - \mathbf{B}^{-1}(\mathbf{Az}^{(k-1)} - \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})\mathbf{z}^{(k-1)}), \quad k = 1, 2, \dots \quad (3.4.5)^{\text{eq:pcinv}}$$

$$\mathbf{z}^{(k)} = \frac{\mathbf{w}}{\|\mathbf{w}\|_2},$$

Theorie:

- Lineare Konvergenz von $(3.4.5)^{\text{eq:pcinv}}$
- Konvergenzgeschwindigkeit gross, wenn $\lambda_{\max}(\mathbf{B}^{-1}\mathbf{A}) : \lambda_{\min}(\mathbf{B}^{-1}\mathbf{A})$ klein



Beispiel: (\mathbf{A} aus Bsp. 112)
`gallery('poisson', 100)` (skaliert)
 $\lambda := \rho_{\mathbf{A}}(\mathbf{z}^{(k-1)})$
 $\lambda^* = \max\{\mu \in \sigma(\mathbf{A}), \mu \leq \lambda\}$
 Plot zeigt theoretische Schranke $[34, 33]$ für NE99c, NE99

$$\Theta(\lambda, \gamma) = \frac{|\rho_{\mathbf{A}}(\mathbf{z}^{(k)}) - \lambda^*|}{|\lambda - \lambda^*|}, \quad \rho := \|\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}\|_{\mathbf{A}}. \quad (3.4.6)^{\text{binv:rat}}$$

Notation: Die von einer s.p.d. Matrix $\mathbf{A} \in \mathbb{K}^{n,n}$ induzierte **Energienorm** ist

$$\|\mathbf{x}\|_{\mathbf{A}} := (\mathbf{x}^H \mathbf{A} \mathbf{x})^{1/2}, \quad \mathbf{x} \in \mathbb{K}^n.$$

Beachte: Matrix-Energienorm: $\|\mathbf{M}\|_{\mathbf{A}} := \sup_{\mathbf{x} \in \mathbb{K}^n \setminus \{0\}} \frac{\|\mathbf{M}\mathbf{x}\|_{\mathbf{A}}}{\|\mathbf{x}\|_{\mathbf{A}}} = \|\mathbf{A}^{1/2} \mathbf{M} \mathbf{A}^{-1/2}\|_2$

Aus Lemma 3.4.6 folgt:

$$0 < \gamma < 1: \quad \gamma(\mathbf{x}^H \mathbf{B} \mathbf{x}) \leq \mathbf{x}^H \mathbf{A} \mathbf{x} \leq (\mathbf{x}^H \mathbf{B} \mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{K}^n \Rightarrow \|\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}\|_{\mathbf{A}} \leq 1 - \gamma.$$

Konzept 3.4.4 (Vorkonditionierer). Eine s.p.d. Matrix $\mathbf{B} \in \mathbb{K}^{n,n}$ heisst **Vorkonditionierer** (engl. preconditioner) für die s.p.d. Matrix $\mathbf{A} \in \mathbb{K}^{n,n}$, falls

- $\exists 0 < \gamma < \Gamma$ „klein“: $\gamma(\mathbf{x}^H \mathbf{B} \mathbf{x}) \leq \mathbf{x}^H \mathbf{A} \mathbf{x} \leq \Gamma(\mathbf{x}^H \mathbf{B} \mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{K}^n$ und
- die Auswertung von $\mathbf{B}^{-1}\mathbf{x}$ erfordert ungefähr den gleichen Rechenaufwand wie die Matrix× Vektor-Multiplikation $\mathbf{A}\mathbf{x}$, $\mathbf{x} \in \mathbb{K}^n$.

Ein Vorkonditionierer liefert eine „billige“ approximative Inverse von \mathbf{A}

Definition 3.4.5 (Spektrale Kondition). **Spektrale Kondition** für reguläres $\mathbf{A} \in \mathbb{K}^{n,n}$

$$\kappa(\mathbf{A}) = \frac{\max\{|\sigma(\mathbf{A})|\}}{\min\{|\sigma(\mathbf{A})|\}}.$$

Lemma 3.4.6. Für s.p.d. $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{n,n}$

$$0 < \gamma < \Gamma: \quad \gamma(\mathbf{x}^H \mathbf{B} \mathbf{x}) \leq \mathbf{x}^H \mathbf{A} \mathbf{x} \leq \Gamma(\mathbf{x}^H \mathbf{B} \mathbf{x}) \quad \forall \mathbf{x} \in \mathbb{K}^n \Rightarrow \kappa(\mathbf{B}^{-1}\mathbf{A}) \leq \frac{\Gamma}{\gamma}.$$

3.4
p. 237

Beispiele für Vorkonditionierer (**Black-Box-Methoden**):

- Diagonalvorkonditionierer (Jacobi-Vorkonditionierer): $\mathbf{B} = \text{diag}(a_{11}, \dots, a_{nn})$

SYMGSVORK

- Symmetrischer Gauss-Seidel-Vorkonditionierer Idee: Approximatives Lösen von $\mathbf{Ax} = \mathbf{b}$ in zwei Stufen

① Approximiere $\mathbf{A}^{-1} \approx \text{tril}(\mathbf{A})$ (unterer Dreiecksanteil): $\tilde{\mathbf{x}} = \text{tril}(\mathbf{A})^{-1}\mathbf{b}$

② Approximiere $\mathbf{A}^{-1} \approx \text{triu}(\mathbf{A})$ (oberer Dreiecksanteil) und „löse“ damit die **Fehlgleichung** $\mathbf{A}(\mathbf{x} - \tilde{\mathbf{x}}) = \mathbf{r}$, mit **Residuum** $\mathbf{r} := \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$:

$$\mathbf{x} = \tilde{\mathbf{x}} + \text{triu}(\mathbf{A})^{-1}(\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}).$$

Mit $\mathbf{L}_{\mathbf{A}} := \text{tril}(\mathbf{A})$, $\mathbf{U}_{\mathbf{A}} := \text{triu}(\mathbf{A})$ findet man

$$\mathbf{x} = (\mathbf{L}_{\mathbf{A}}^{-1} + \mathbf{U}_{\mathbf{A}}^{-1} - \mathbf{U}_{\mathbf{A}}^{-1}\mathbf{A}\mathbf{L}_{\mathbf{A}}^{-1})\mathbf{b} \quad \triangleright \quad \mathbf{B} = \mathbf{L}_{\mathbf{A}}^{-1} + \mathbf{U}_{\mathbf{A}}^{-1} - \mathbf{U}_{\mathbf{A}}^{-1}\mathbf{A}\mathbf{L}_{\mathbf{A}}^{-1}.$$

Beispiel 111 (Jacobi-Vorkonditionierung).

3.4
p. 238

3.4
p. 239

3.4
p. 240

Matrix mit „schwerer“ Diagonalen

$$\mathbf{A} = \mathbf{A}^T \in \mathbb{R}^{n,n}: \exists \delta \in]0, 1[: \delta a_{ii} \geq \sum_{j \neq i}^n |a_{ij}|. \quad (3.4.7) \quad \text{Ladef}$$

$$\sigma(\mathbf{A}) \subset [(1 - \delta) \min_i \{a_{ii}\}, (1 + \delta) \max_i \{a_{ii}\}].$$

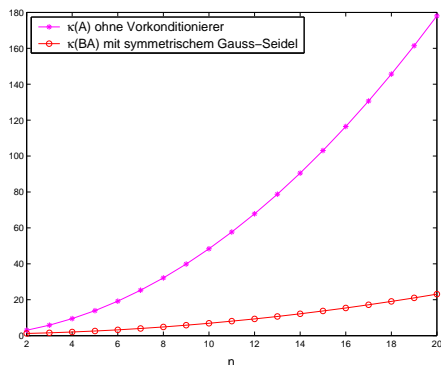
Lemma 3.4.7 (Gerschgorin-Kreise). Für $\mathbf{A} \in \mathbb{K}^{n,n}$ gilt

$$\sigma(\mathbf{A}) \subset \bigcup_{j=1}^n \{z \in \mathbb{C}: |z - a_{jj}| \leq \sum_{i \neq j} |a_{ji}|\}.$$

Für \mathbf{A} aus (3.4.7) (Jacobi-Vorkonditionierer = „Equilibrierung“)

$$\kappa(\text{diag}(a_{11}, \dots, a_{nn})^{-1} \mathbf{A}) \subset [1 - \delta, 1 + \delta].$$

Beispiel 112 (Gauss-Seidel-Vorkonditionierung).



$$\mathbf{A} = \begin{pmatrix} \mathbf{T} & -\mathbf{I} & 0 & \dots & \dots & 0 \\ -\mathbf{I} & \mathbf{T} & -\mathbf{I} & & & \vdots \\ 0 & \ddots & \ddots & \ddots & & \\ \vdots & & & -\mathbf{I} & \mathbf{T} & -\mathbf{I} \\ 0 & \dots & \dots & 0 & -\mathbf{I} & \mathbf{T} \end{pmatrix} \in \mathbb{R}^{n^2, n^2},$$

$$\mathbf{T} = \begin{pmatrix} 4 & -1 & 0 & \dots & \dots & 0 \\ -1 & 4 & -1 & & & \vdots \\ 0 & \ddots & \ddots & \ddots & & \\ \vdots & & & -1 & 4 & -1 \\ 0 & \dots & \dots & 0 & -1 & 4 \end{pmatrix} \in \mathbb{R}^{n,n}.$$

UNVOLLSTCHOL

- Unvollständige Cholesky-Zerlegung (engl. *incomplete Cholesky decomposition*)

„Cholesky-Zerlegung“ mit Unterdrückung von Fill-In (→ 3.2.14, Abschn. 3.2.6)

MATLAB-Funktion: `R=cholinc(A, '0')`; (nur sinnvoll für dünnbesetzte Matrizen !)

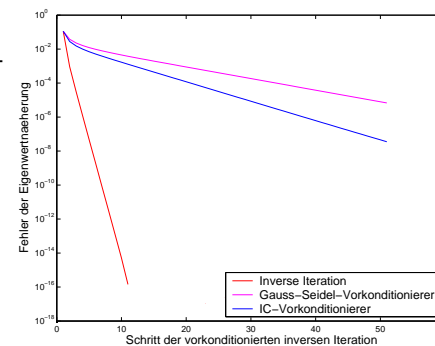
Beispiel 113 (Konvergenz der vorkonditionierten inversen Iteration).

Konvergenz der Eigenwertnäherung aus der vorkonditionierten inversen Iteration für

- `A = gallery('poisson', 30);` (s.p.d.)

und

- `B = A`
- `B ← Gauss-Seidel-Vorkonditionierer`
- `B ← cholinc(A, '0');`



3.4.3 Krylov-Unterraumverfahren

3.4 Gegeben: (Grosse, diagonalisierbare) Matrix $\mathbf{A} \in \mathbb{K}^{n,n}$: $\mathbf{S}^{-1} \mathbf{A} \mathbf{S} = \text{diag}(\lambda_1, \dots, \lambda_n)$
p. 241 (Eigenwerte sortiert: $|\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_n|$)

Eigenwertproblem: Finde $\lambda_1, \dots, \lambda_l$ (und Eigenvektoren $\mathbf{s}_{\cdot,1}, \dots, \mathbf{s}_{\cdot,l}$) bzw.
Finde $\lambda_{n-l+1}, \dots, \lambda_n$ (und $\mathbf{s}_{\cdot,n-l+1}, \dots, \mathbf{s}_{\cdot,n}$) , $l \ll n$

➤ Option: Potenzmethoden (direkt, invers), Abschnitte 3.4.1, 3.4.2 sec:directpotenzmethode-iter

Idee der Ritz-Projektion, vgl. (3.4.4). RITZ spowit

- $\mathbf{s}_{\cdot,m} \in V_l := \text{Im}(\mathbf{V}) \subset \mathbb{K}^n$, $1 \leq m \leq n$, $\mathbf{V} \in \mathbb{K}^{n,l}$, $\Rightarrow \lambda_m \in \sigma(\mathbf{V}^H \mathbf{A} \mathbf{V})$
 $1 \leq l \leq n$, $\{\mathbf{v}_{\cdot,1}, \dots, \mathbf{v}_{\cdot,l}\}$ Orthonormalbasis von V_l
- Falls V_l „signifikante Komponenten in Richtung von $\mathbf{s}_{\cdot,m}$ “ $\Rightarrow \exists \mu \in \sigma(\mathbf{V}^H \mathbf{A} \mathbf{V})$, $\mu \approx \lambda_m$
enthält

Aus Stabilitätsgründen !

Wie bekommt man Unterraum mit „signifikanten Komponenten in Richtung von $\mathbf{s}_{\cdot,n-l+1}, \dots, \mathbf{s}_{\cdot,n}$ “ ?

Tipp: Unterraumversion der direkten Potenzmethode (3.4.4) spowit

$$V_l^{(0)} \text{ beliebig}, \quad V_l^{(k)} := \mathbf{A} V_l^{(k-1)}, \quad k = 1, 2, \dots$$

3.4 Wie bekommt man Unterraum mit „signifikanten Komponenten in Richtung von $\mathbf{s}_{\cdot,1}, \dots, \mathbf{s}_{\cdot,l}$ “ ?
p. 242

3.4
p. 243

3.4
p. 244

Tipp:

$$\sigma(\|\mathbf{A}\|_2 \mathbf{I} - \mathbf{A}) = \|\mathbf{A}\|_2 - \sigma(\mathbf{A})$$

(\mathbf{A} s.p.d.: Betragsskleinste Eigenwerte von \mathbf{A} = betragsgrösste Eigenwerte von $\|\mathbf{A}\|_2 \mathbf{I} - \mathbf{A}$)

Definition 3.4.8. Für $\mathbf{A} \in \mathbb{K}^{n,n}$, $\mathbf{z} \in \mathbb{K}^n$, $\mathbf{z} \neq 0$, ist der l . *Krylov-Unterraum*

$$\mathcal{K}_l(\mathbf{A}, \mathbf{z}) := \text{Span} \{ \mathbf{z}, \mathbf{A}\mathbf{z}, \dots, \mathbf{A}^{l-1}\mathbf{z} \}.$$

Äquivalent: $\mathcal{K}_l(\mathbf{A}, \mathbf{z}) = \{ p(\mathbf{A})\mathbf{z} : p \text{ Polynom vom Grad } \leq l \}$

Sukzessive Konstruktion einer Orthonormalbasis (ONB) von $\mathcal{K}_l(\mathbf{A}, \mathbf{z})$: **ARNOLDI** *Arnoldi-Prozess*

Induktives Vorgehen: $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ ONB von $\mathcal{K}_m(\mathbf{A}, \mathbf{z}) \quad \forall 1 \leq m \leq l$

➤ $\mathbf{A}\mathbf{v}_l \in \mathcal{K}_{l+1}(\mathbf{A}, \mathbf{z})$ (warum ?)

$$\tilde{\mathbf{v}}_{l+1} := \mathbf{A}\mathbf{v}_l - \sum_{j=1}^l (\mathbf{v}_j^H \mathbf{A}\mathbf{v}_l) \mathbf{v}_j, \quad \mathbf{v}_{l+1} := \frac{\tilde{\mathbf{v}}_{l+1}}{\|\tilde{\mathbf{v}}_{l+1}\|_2} \Rightarrow \mathbf{v}_{l+1} \perp \mathcal{K}_l(\mathbf{A}, \mathbf{z}). \quad (3.4.8)$$

(vgl. **GRAMSCHMIDT** *Gram-Schmidt-Orthonormalisierung* → lineare Algebra)

orthogonal

$$\mathbf{V}_l = [\mathbf{v}_1, \dots, \mathbf{v}_l] : \mathbf{A}\mathbf{V}_l = \mathbf{V}_{l+1} \tilde{\mathbf{H}}_l, \quad \tilde{\mathbf{H}}_l \in \mathbb{K}^{l+1,l} \text{ mit } \tilde{h}_{ij} = \begin{cases} \mathbf{v}_i^H \mathbf{A}\mathbf{v}_j, & \text{falls } i \leq j, \\ \|\tilde{\mathbf{v}}_j\|_2, & \text{falls } i = j+1, \\ 0 & \text{sonst.} \end{cases}$$

➔ $\tilde{\mathbf{H}}_l$ = obere Hessenbergmatrizen

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} \begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_l \end{pmatrix} = \begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_{l+1} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{H}}_l \end{pmatrix}$$

Stabile Implementierung
des Arnoldi-Prozesses:

Modifizierter Gram-Schmidt

Im l . Schritt:

1× $\mathbf{A} \times$ Vektor

$l+1$ Skalarprodukte

l SAXPYs

n Divisionen

➤ Rechenaufwand für l Schritte,

wenn \mathbf{A} maximal k Einträge \neq

0 pro Zeile: $O(nkl^2)$

$H(l+1, l) = 0 \rightarrow$ Abbruch !

MATLAB-CODE: Arnoldi-Prozess

```
function [V,Ht] = arnoldi(A,k,v0)
V = [v0/norm(v0)];
Ht = zeros(k+1,k);
for l=1:k
    vt = A*V(:,l);
    for j=1:l
        Ht(j,l) = dot(V(:,j),vt);
        vt = vt - Ht(j,l)*V(:,j);
    end
    Ht(l+1,l) = norm(vt);
    V = [V, vt/Ht(l+1,l)];
end
```

(3.4.9) **arnoldi**

3.4
p. 245

Lemma 3.4.9. Für die Matrizen $\mathbf{V}_l \in \mathbb{K}^{n,l}$, $\tilde{\mathbf{H}}_l \in \mathbb{K}^{l+1,l}$ aus dem l . Schritt, $l \leq n$, des Arnoldi-Verfahrens gilt

(i) $\mathbf{V}_l^H \mathbf{V}_l = \mathbf{I}$,

(ii) $\mathbf{A}\mathbf{V}_l = \mathbf{V}_{l+1} \tilde{\mathbf{H}}_l$, $\tilde{\mathbf{H}}_l$ obere Hessenbergmatrix,

(iii) $\mathbf{V}_l^H \mathbf{A}\mathbf{V}_l = \mathbf{H}_l \in \mathbb{K}^{l,l}$, $h_{ij} = \tilde{h}_{ij}$ für $1 \leq i, j \leq l$,

(iv) Falls $\mathbf{A} = \mathbf{A}^H \rightarrow \mathbf{H}_l$ tridiagonal.

Abbruch des Arnoldi-Prozesses **arnoldi** (3.4.9)

im l . Schritt

➤ $\text{Span} \{ \mathbf{v}_1, \dots, \mathbf{v}_l \}$ invarianter Teilraum von \mathbf{A} .

(Praxis: rechne in diesem Fall mit „beliebigem“ \mathbf{v}_t weiter)

Beachte:

(iii) $\Leftrightarrow \mathbf{H}_l$ ist **Ritz-Projektion** von \mathbf{A} auf $\mathcal{K}_l(\mathbf{A}, \mathbf{z})$.

Arnoldi-Prozess zur Eigenwertapproximation:

Im l . Schritt: $\lambda_n \approx \mu_l^{(l)}$, $\lambda_{n-1} \approx \mu_{l-1}^{(l)}$, \dots , $\lambda_1 \approx \mu_1^{(l)}$,

$\sigma(\mathbf{H}_l) = \{ \mu_1^{(l)}, \dots, \mu_l^{(l)} \}$, $|\mu_1^{(l)}| \leq |\mu_2^{(l)}| \leq \dots \leq |\mu_l^{(l)}|$.

3.4
p. 246

3.4
p. 247

3.4
p. 248

MATLAB-CODE Arnoldi-Verfahren

```
function [dn,V,Ht] = arnolidev(A,v0,k,tol)
n = size(A,1); V = [v0/norm(v0)];
H = zeros(1,0); dn = zeros(k,1);
for l=1:n
    d = dn;
    Ht = [Ht, zeros(1,1); zeros(1,1)];
    vt = A*V(:,l);
    for j=1:l
        Ht(j,l) = dot(V(:,j),vt);
        vt = vt - Ht(j,l)*V(:,j);
    end
    ev = sort(eig(Ht(1:l,1:l)));
    dn(1:min(l,k)) = ev(end:-1:end-min(l,k)+1);
    if (norm(d-dn) < tol*norm(dn)), break; end;
    Ht(1+1,l) = norm(vt);
    V = [V, vt/Ht(1+1,l)];
end
```

Arnoldi-Verfahren
zur Berechnung der $k, k \leq n$
betragsgrössten Eigenwerte von
 $A \in \mathbb{K}^{n,n}$

Pro Schritt nur eine
Multiplikation mit Matrix A !
(Effizient für dünnbesetzte
Matrizen)

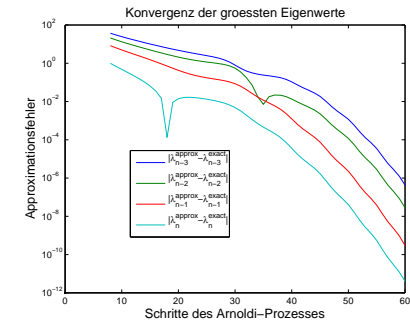
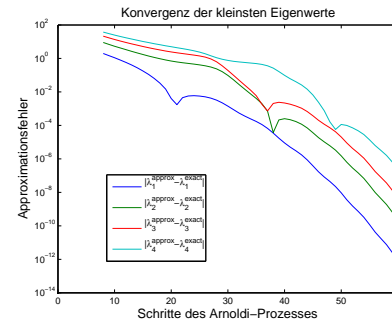
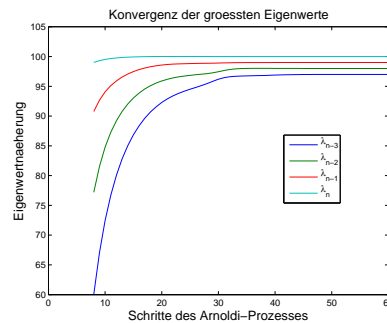
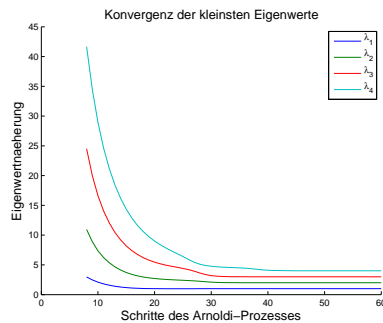
Monoton wachsender Speicherauf-
wand !

Beispiel 114 (Arnoldi-Verfahren zur Eigenwertberechnung).

```
n=100;
M=gallery('tridiag',-0.5*ones(n-1,1),2*ones(n,1),-1.5*ones(n-1,1));
[Q,R]=qr(M);
v=ones(n,1); v(2:2:end)=2*v(2:2:end);
S=diag(v)*Q; A=S*diag(1:n)*inv(S);
```

Startvektor: `ones(100,1)` :

$$\sigma(A) = \{1, 2, 3, 4, \dots, 97, 98, 99, 100\}$$



Im Fall $A = A^H \in \mathbb{K}^{n,n} \Rightarrow H_l = H_l^H \Rightarrow$ kurze Rekursionen: **LANCZOS**
Lanczos-Prozess

3.4
p. 249

Lemma 3.4.9 $\xrightarrow{\text{lem:arnoldi}}$ l . Schritt, $1 \leq l \leq n$ des Lanczos-Prozesses erzeugt $V_l \in \mathbb{K}^{n,l}$ mit orthonor-
mierten Spalten, so dass

$$V_l^H A V_l = \begin{pmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \beta_2 & \alpha_3 & \ddots & \\ & & \ddots & \ddots & \beta_{k-1} \\ & & & \beta_{k-1} & \alpha_k \end{pmatrix} =: T_l \in \mathbb{K}^{k,k} \quad [\text{Tridiagonalmatrix}]$$

3.4
p. 250

3.4
p. 251

3.4
p. 252

Lanczos-Prozess

Rechenaufwand pro Schritt:

- 1 \times $\mathbf{A} \times$ Vektor
- 2 Skalarprodukte
- 2 SAXPYs
- 1 Division

> Rechenaufwand für l Schritte, wenn

\mathbf{A} maximal k Einträge $\neq 0$ pro Zeile:

$O(nkl)$

$\beta(l+1) = 0 \rightarrow$ Abbruch!

Speicherbedarf: 2 Vektoren Länge n ,
wenn \mathbf{V} nicht berechnet wird.

MATLAB-CODE Lanczos-Prozess

```
function [V,alpha,beta] = lanczos(A,k,w)
V = [w/norm(w)];
alpha = zeros(k,1);
beta = zeros(k,1);
for l=1:k
    vt = A*V(:,l);
    if (l>1)
        vt = vt - beta(l)*V(:,l-1);
    end
    alpha(l) = dot(V(:,l),vt);
    vt = vt - alpha(l)*V(:,l);
    beta(l+1) = norm(vt);
    V = [V, vt/beta(l+1)];
end
beta = beta(2:end-1);
```

Lanczos-Prozess zur Eigenwertapproximation: (Theorie [12, Abschn. 8.5])

Im l . Schritt: $\lambda_n \approx \mu_l^{(l)}, \lambda_{n-1} \approx \mu_{l-1}^{(l)}, \dots, \lambda_1 \approx \mu_1^{(l)}$,
 $\sigma(\mathbf{T}_l) = \{\mu_1^{(l)}, \dots, \mu_l^{(l)}\}, |\mu_1^{(l)}| \leq |\mu_2^{(l)}| \leq \dots \leq |\mu_l^{(l)}|$.

Beispiel 115 (Rundungsfehlereinfluss beim Lanczos-Prozess).

$\mathbf{A} \in \mathbb{R}^{10,10}$, $a_{ij} = \min\{i, j\}$. $\mathbf{A} = \text{gallery('minij',10)};$

Berechnet durch $[V, \alpha, \beta] = \text{lanczos}(\mathbf{A}, n, \text{ones}(n, 1));$

$\mathbf{T} = \begin{pmatrix} 38.500000 & 14.813845 & & & & & & & & \\ 14.813845 & 9.642857 & 2.062955 & & & & & & & \\ & 2.062955 & 2.720779 & 0.776284 & & & & & & \\ & & 0.776284 & 1.336364 & 0.385013 & & & & & \\ & & & 0.385013 & 0.826316 & 0.215431 & & & & \\ & & & & 0.215431 & 0.582380 & 0.126781 & & & \\ & & & & & 0.126781 & 0.446860 & 0.074650 & & \\ & & & & & & 0.074650 & 0.363803 & 0.043121 & \\ & & & & & & & 0.043121 & 3.820888 & 11.991094 \\ & & & & & & & & 11.991094 & 41.254286 \end{pmatrix}$

$\sigma(\mathbf{A}) = \{0.255680, 0.273787, 0.307979, 0.366209, 0.465233, 0.643104, 1.000000, 1.873023, 5.048917, 44.766069\}$

$\sigma(\mathbf{T}) = \{0.263867, 0.303001, 0.365376, 0.465199, 0.643104, 1.000000, 1.873023, 5.048917, 44.765976, 44.766069\}$

► **Mysteriöser Cluster von Eigenwerten von \mathbf{T}** („ghost eigenvalues“ [19, Sect. 9.2.5])

$\mathbf{V}^H \mathbf{V} = \begin{pmatrix} 1.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000251 & 0.258801 & 0.883711 \\ 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000106 & 0.109470 & 0.373799 \\ 0.000000 & -0.000000 & 1.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000005 & 0.005373 & 0.018347 \\ 0.000000 & 0.000000 & 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000096 & 0.000328 \\ 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000001 & 0.000003 \\ 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 \\ 0.000000 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 & 0.000000 \\ 0.000251 & 0.000106 & 0.000005 & 0.000000 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 & 0.000000 \\ 0.258801 & 0.109470 & 0.005373 & 0.000096 & 0.000001 & 0.000000 & 0.000000 & -0.000000 & 1.000000 & 0.000000 & 0.000000 \\ 0.883711 & 0.373799 & 0.018347 & 0.000328 & 0.000003 & 0.000000 & 0.000000 & 0.000000 & 0.000000 & 1.000000 & 0.000000 \end{pmatrix}$

► **Orthogonalitätsverlust der Basisvektoren durch Rundungsfehlereinfluss**
(Lemma 3.4.9 nur von „theoretischem Wert“)

3.4
p. 253

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

► Gute Approximation der betragsgrössten Eigenwerte für $l < n$

Analytisch äquivalent: Arnoldi-Prozess angewandt auf \mathbf{A} , versuchen wir es mal:

3.4
p. 254

3.4
p. 255

3.4
p. 256

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	0.276159 0.349722 0.462449 0.643016 1.000000 1.873023 5.048917 44.766069									
9	0.263872 0.303009 0.365379 0.465199 0.643104 1.000000 1.873023 5.048917 44.766069									
10	0.255680 0.273787 0.307979 0.366209 0.465233 0.643104 1.000000 1.873023 5.048917 44.766069									

► Perfekt !!

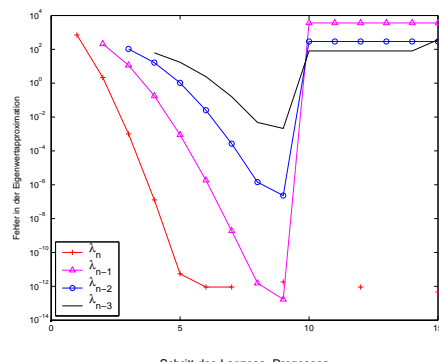
Beispiel 116 (Eigenwertberechnung mit Lanczos/Arnoldi-Prozess).

Lanczos-Verfahren:

```
A = gallery('minij',100);
Startvektor: ones(1,100)
```

$\lambda_l(\mathbf{T}_l), \lambda_{l-1}(\mathbf{T}_l), \dots$

Gute Konvergenz
bevor „ghost eigenvalues“ auftreten.

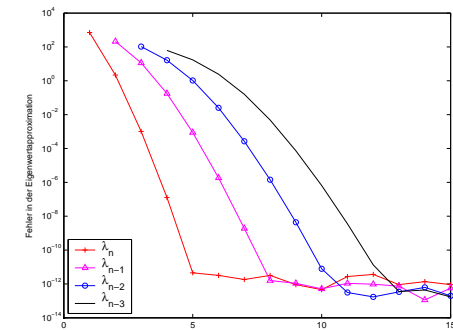


Arnoldi-Verfahren:

(Matrix/Startvektor wie oben)

$\lambda_l(\mathbf{T}_l), \lambda_{l-1}(\mathbf{T}_l), \dots$

Lanczos-Prozess rundungsfehleranfällig
Arnoldi-Prozess wesentlich stabiler
(Keine „ghost eigenvalues“)



Krylov-Unterraum-Verfahren basierend auf Arnoldi-/Lanczos-Prozess sind die Methoden der Wahl zur Berechnung von wenigen Eigenwerten/Eigenvektoren für *grosse dünnbesetzte Matrizen*

Bemerkung 117. Erweiterung der Krylov-Unterraum-Verfahren auf verallgemeinertes Eigenwertproblem $\mathbf{Ax} = \lambda \mathbf{Bx}$, \mathbf{B} s.p.d., durch Ersetzung des Euklidischen Skalarprodukts mit „**B-Skalarprodukt**“ $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x}^H \mathbf{B} \mathbf{y}$.

3.4

p. 257 MATLAB-Funktion:

`d = eigs(A,k,sigma)` : k grösste/kleinste Eigenwerte von \mathbf{A}
`d = eigs(A,B,k,sigma)` : k grösste/kleinste Eigenwerte für verallgemeinertes Eigenwertproblem $\mathbf{Ax} = \lambda \mathbf{Bx}$, \mathbf{B} s.p.d.
`d = eigs(Afun,n,k)` : Afun = Funktionshandle für Anwendung von $\mathbf{A}/\mathbf{A}^{-1}/\mathbf{A} - \alpha \mathbf{I}/(\mathbf{A} - \alpha \mathbf{B})^{-1}$ auf Vektor. Eigenschaften der durch Afun repräsentierten Matrix müssen durch Flaggen spezifiziert werden.

3.4

p. 259

3.4.4 Singulärwertzerlegungen

SINGWERTZERI

= Diagonalisierung durch separate unitäre Basistransformationen in Bild- und Urbildbereich

□ **Theorem 3.4.10.** Zu jedem $\mathbf{A} \in \mathbb{K}^{m,n}$ gibt es unitäre Matrizen $\mathbf{U} \in \mathbb{K}^{m,m}$, $\mathbf{V} \in \mathbb{K}^{n,n}$ und eine Diagonalmatrix $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p) \in \mathbb{R}^{m,n}$, $p := \min\{m, n\}$, $\sigma_1 \geq \sigma_2 \geq \sigma_p \geq 0$ mit

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^H.$$

3.4

p. 258

3.4

p. 260

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U} \end{pmatrix} \begin{pmatrix} \Sigma \end{pmatrix} \begin{pmatrix} \mathbf{V}^H \end{pmatrix}$$

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U} \end{pmatrix} \begin{pmatrix} \Sigma \end{pmatrix} \begin{pmatrix} \mathbf{V}^H \end{pmatrix}$$

Definition 3.4.11. Die Zerlegung $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$ aus Satz 3.4.10 heisst **Singulärwertzerlegung** (SVD) mit **Singulärwerten** σ_i .

Singulärwertzerlegung = Zerlegung in Rang-1-Komponenten:

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H = \sum_{j=1}^p \sigma_j \mathbf{u}_{:,j} \mathbf{v}_{:,j}^H. \quad (3.4.10)$$

Lemma 3.4.12. Die Quadrate σ_j^2 der nichtverschwindenden Singulärwerte von \mathbf{A} sind Eigenwerte von $\mathbf{A}^H\mathbf{A}$, $\mathbf{A}\mathbf{A}^H$ und $\begin{pmatrix} 0 & \mathbf{A}^H \\ \mathbf{A} & 0 \end{pmatrix}$ zu den Eigenvektoren $\mathbf{v}_{:,1}, \dots, \mathbf{v}_{:,p}$, $\mathbf{u}_{:,1}, \dots, \mathbf{u}_{:,p}$, bzw. $\begin{pmatrix} \mathbf{u}_{:,1} \\ \mathbf{v}_{:,1} \end{pmatrix}, \dots, \begin{pmatrix} \mathbf{u}_{:,p} \\ \mathbf{v}_{:,p} \end{pmatrix}$.

Thm. 3.3.5 $\Rightarrow |\sigma_k(\mathbf{A}) - \sigma_k(\mathbf{A} + \mathbf{E})| \leq \|\mathbf{E}\|_2$

Das Problem der Berechnung der Singulärwerte ist immer hervorragend konditioniert

Singulärwerte von \mathbf{A} sind eindeutig bestimmt

MATLAB-Funktionen:

`s = svd(A)` : Berechnet Singulärwerte der Matrix \mathbf{A}
`[U,S,V] = svd(A)` : Berechnet Singulärwertzerlegung gemäss Thm. 3.4.10
`[U,S,V] = svd(A,0)` : Sparsame Singulärwertzerlegung für $m > n$: $\mathbf{U} \in \mathbb{K}^{m,n}$, $\Sigma \in \mathbb{R}^{n,n}$, $\mathbf{V} \in \mathbb{K}^{n,n}$
`s = svds(A,k)` : k grösste Singulärwerte (sinnvoll für dünnbesetztes \mathbf{A})
`[U,S,V] = svds(A,k)` : Partielle Singulärwertzerlegung: $\mathbf{U} \in \mathbb{K}^{m,k}$, $\mathbf{V} \in \mathbb{K}^{n,k}$, $\Sigma \in \mathbb{R}^{k,k}$ diagonal, enthält k grösste Singulärwerte von \mathbf{A} .

Sparsame Singulärwertzerlegung:

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U} \end{pmatrix} \begin{pmatrix} \Sigma \end{pmatrix} \begin{pmatrix} \mathbf{V}^H \end{pmatrix}$$

3.4
p. 261

3.4
p. 263

SVD-Algorithmus ist in jedem Fall numerisch stabil
RECHAUFWSVD
 Rechenaufwand:
 $2mn^2 + 2n^3 + O(n^2) + O(mn)$ für `s = svd(A)`,
 $4m^2n + 22n^3 + O(mn) + O(n^2)$ für `[U,S,V] = svd(A)`,
 $O(mn^2) + O(n^3)$ für `[U,S,V] = svd(A,0)`, $m \gg n$.

NUMRANGBEST

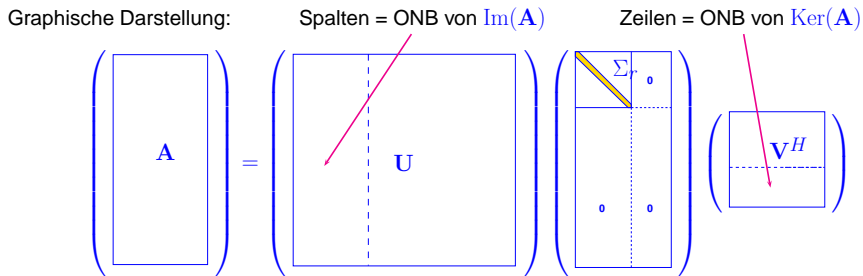
- Anwendung: Numerische Rangbestimmung

Lemma 3.4.13 (SVD und Rang einer Matrix). Erfüllen die Singulärwerte von \mathbf{A}

$$\begin{aligned} \sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0, \\ \text{rank}(\mathbf{A}) = r, \\ \text{Ker}(\mathbf{A}) = \text{Span} \{ \mathbf{v}_{:,r+1}, \dots, \mathbf{v}_{:,n} \}, \\ \text{Im}(\mathbf{A}) = \text{Span} \{ \mathbf{u}_{:,1}, \dots, \mathbf{u}_{:,r} \}. \end{aligned}$$

3.4
p. 262

3.4
p. 264



Anmerkung: MATLAB-Funktion `r=rank(A)` benutzt SVD zur Rangberechnung.

- Berechnung von Matrixnormen

Lemma 3.4.14 (SVD und Euklidische Matrixnorm).

$$\forall \mathbf{A} \in \mathbb{K}^{m,n}: \|\mathbf{A}\|_2 = \sigma_1(\mathbf{A}), \quad \forall \mathbf{A} \in \mathbb{K}^{n,n} \text{ regulär: } \text{cond}_2(\mathbf{A}) = \sigma_1/\sigma_n.$$

Anmerkung: MATLAB-Funktion `norm(A)` und `cond(A)` benutzen SVD.

- Anwendung: *Bestapproximation durch Niedrigrangmatrizen* (engl. *low rank approximation*)

Definition 3.4.15 (Frobeniusnorm). Die **Frobeniusnorm** von $\mathbf{A} \in \mathbb{K}^{m,n}$ ist definiert durch

$$\|\mathbf{A}\|_F^2 := \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2.$$

Frobeniusnorm und SVD:

$$\|\mathbf{A}\|_F^2 = \sum_{j=1}^p \sigma_j^2$$

Beachte: $\|\mathbf{A}\|_F$ invariant bei unitären Transformationen.

Notation: $\mathcal{R}_k(m, n) := \{\mathbf{A} \in \mathbb{K}^{m,n}: \text{rank}(\mathbf{A}) \leq k\}, m, n, k \in \mathbb{N}$

Theorem 3.4.16 (Niedrigrang-Bestapproximation). Sei $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^H$ die SVD von $\mathbf{A} \in \mathbb{K}^{m,n}$. Setze für $1 \leq k \leq \text{rank}(\mathbf{A})$ $\mathbf{U}_k := [\mathbf{u}_{\cdot,1}, \dots, \mathbf{u}_{\cdot,k}]$, $\mathbf{V}_k := [\mathbf{v}_{\cdot,1}, \dots, \mathbf{v}_{\cdot,k}]$, $\Sigma_k := \text{diag}(\sigma_1, \dots, \sigma_k)$. Dann gilt für $\|\cdot\| = \|\cdot\|_F$ und $\|\cdot\| = \|\cdot\|_2$

$$\|\mathbf{A} - \mathbf{U}_k \Sigma_k \mathbf{V}_k^H\| \leq \|\mathbf{A} - \mathbf{F}\| \quad \forall \mathbf{F} \in \mathcal{R}_k(m, n).$$

Beispiel 118 (Hauptkomponentenanalyse).

$\mathbf{A} \in \mathbb{R}^{m,n}, m \gg n$:

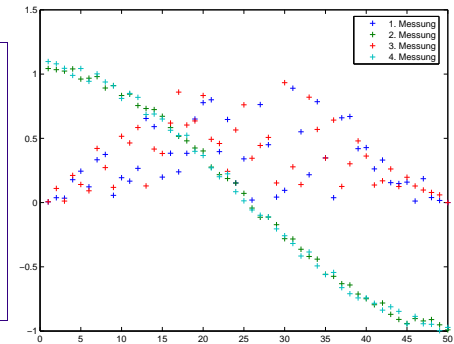
Zeilen von $\mathbf{A} \rightarrow$ Messung mehrerer Größen
Spalten von $\mathbf{A} \rightarrow$ Messgrößen

Ziel: Identifikation linearer Korrelationen

Im Jahresverlauf werden zwei verschiedene Größen an 10 verschiedenen Orten gemessen:

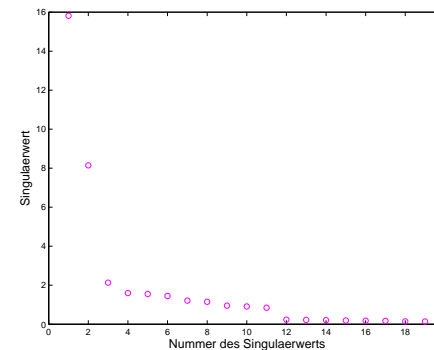
(Die Messungen sind fehlerbehaftet)

```
n = 10;
m = 50;
x = sin(pi*(1:m)'/m);
y = cos(pi*(1:m)'/m);
A = [];
for i = 1:n
    A = [A, x.*rand(m,1), ...
         y+0.1*rand(m,1)];
end
```



3.4
p. 265

3.4
p. 267

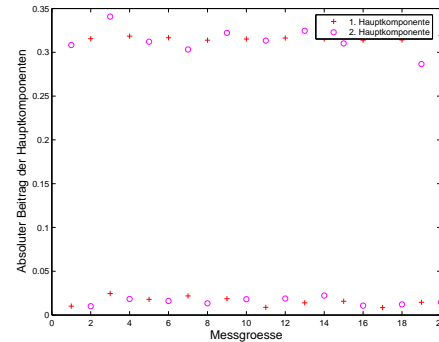
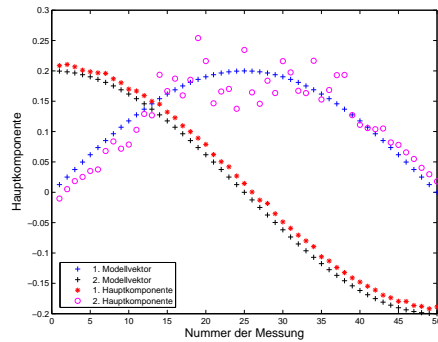


← Verteilung der Singulärwerte der Messmatrix
Zwei dominante Singulärwerte
Messungen linear korreliert mit **zwei** Hauptkomponenten.

Hauptkomponenten = $\mathbf{u}_{\cdot,1}, \mathbf{u}_{\cdot,2}$ (führende Spalten von \mathbf{U} in SVD)
Beiträge = $\mathbf{v}_{\cdot,1}, \mathbf{v}_{\cdot,2}$ (führende Spalten von \mathbf{V} in SVD)

3.4
p. 266

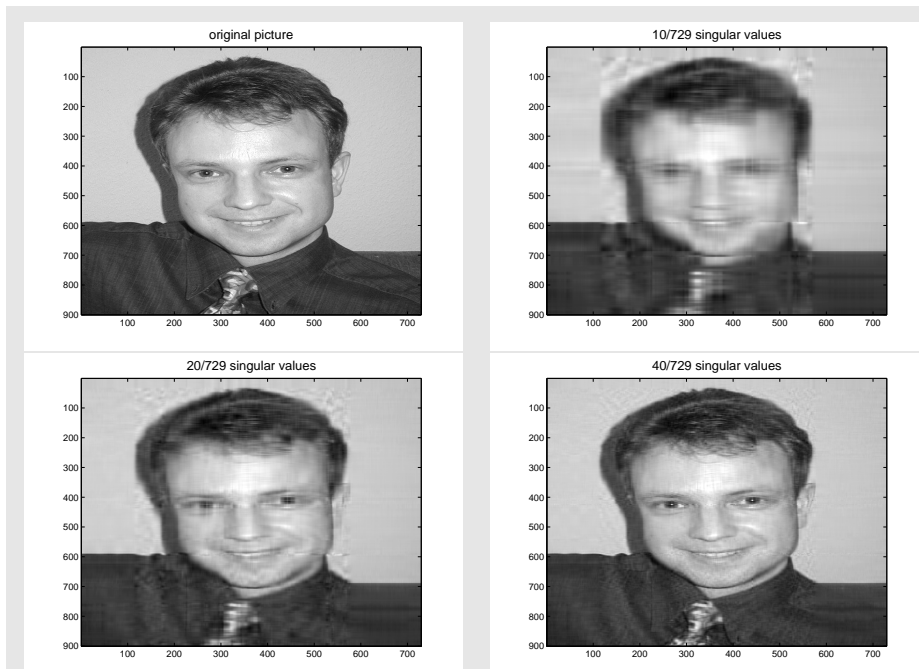
3.4
p. 268



BILDDATENKOMP

Beispiel 119 (Bildatenkompression). $m \times n$ -Pixelbild (Graustufen, BMP-Format) \leftrightarrow Matrix $A \in \mathbb{R}^{m,n}$, $a_{ij} \in \{0, \dots, 255\}$

► Niedrigrang-Bestapproximation des Bildes: $\tilde{A} = U_k \Sigma_k V^T$, Thm. 3.4.16 (Thm: rankapprox) (\Rightarrow Animation)



3.4
p. 269

3.5 Numerik linearer Ausgleichsprobleme

[File: section-numerik-linearer-ausgleichsprobleme.tex, SVN: section-numerik-linearer-ausgleichsprobleme.tex 1199 2006-12-05 06:48:09Z hiptmar]

LSQ

Lineares Ausgleichsproblem, „Methode der kleinsten Quadrate“ (engl. *linear least squares*):

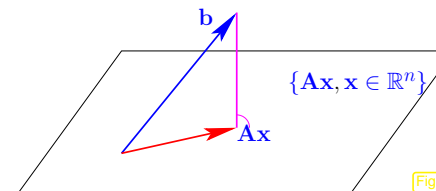
Gegeben: $A \in \mathbb{K}^{m,n}$, $m, n \in \mathbb{N}$, $b \in \mathbb{K}^m$

Gesucht: $x \in \mathbb{K}^n$ mit

$$(i) \|Ax - b\|_2 = \inf\{\|Ay - b\|_2 : y \in \mathbb{K}^n\},$$

(3.5.1) **eq:LSQ**

$$(ii) \|x\|_2 \text{ minimal unter Bedingung (i)}$$



[Fig. 3]

Geometrische Interpretation:

$x \triangleq$ Orthogonalprojektion von b auf Unterraum $\text{Span}\{a_1, \dots, a_n\}$

Lemma 3.5.1. Existenz & Eindeutigkeit von Lösungen des Linearen Ausgleichsproblems

3.4
p. 270

3.5
p. 271

3.5
p. 272

Definition 3.5.2 (Pseudoinverse). Die **Pseudoinverse** $\mathbf{A}^+ \in \mathbb{K}^{n,m}$ von $\mathbf{A} \in \mathbb{K}^{m,n}$ ist die Matrixdarstellung des (linearen) Lösungsoperators $\mathbb{R}^m \mapsto \mathbb{R}^n$, $\mathbf{b} \mapsto \mathbf{x}$ zum linearen Ausgleichsproblem (3.5.1) $\|\mathbf{Ax} - \mathbf{b}\| \rightarrow \min$, $\|\mathbf{x}\| \rightarrow \min$.

MATLAB: $\mathbf{x} = \mathbf{A} \backslash \mathbf{b}$ („backslash“) löst (3.5.1) für $\mathbf{A} \in \mathbb{K}^{m,n}$, $m \neq n$
 $\mathbf{p} = \text{pinv}(\mathbf{A})$ berechnet Pseudoinverse

Beispiel 120 (Lineares Fitten von Daten, lineare Regression). \rightarrow Bsp. 59

Gegeben: „Messpunkte“ $(t_i, y_i) \in \mathbb{K}^2$, $i = 1, \dots, m$, $t_i \in I \subset \mathbb{K}$
 Basisfunktionen $b_j : I \mapsto \mathbb{K}$, $j = 1, \dots, n$.

Gesucht: Koeffizienten $x_j \in \mathbb{K}$, $j = 1, \dots, n$, so dass

$$\sum_{i=1}^m |f(t_i) - y_i|^2 \rightarrow \min, \quad f(t) := \sum_{j=1}^n x_j b_j(t).$$

► Lineares Ausgleichsproblem mit $\mathbf{A} = (b_i(t_j))$, $1 \leq i \leq n$, $1 \leq j \leq m$.

Hier: $\|\mathbf{b} - \mathbf{Ax}\|_2$ misst Qualität des Fit-Modells.

Spezialfall: Polynomiales Fitten: $b_j(t) = t^{j-1}$

MATLAB-Funktion: $\mathbf{p} = \text{polyfit}(\mathbf{t}, \mathbf{y}, \mathbf{n})$; \mathbf{n} = Polynomgrad

Kondition des linearen Ausgleichsproblems:

Definition 3.5.3 (Verallgemeinerte Kondition einer Matrix, \rightarrow Def. 3.2.5). Seien $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0$, $p := \min\{m, n\}$, die singulären Werte (\rightarrow Def. 3.4.11) von $\mathbf{A} \in \mathbb{K}^{m,n}$. Dann ist

$$\text{cond}_2(\mathbf{A}) := \frac{\sigma_1}{\sigma_r}$$

die **verallgemeinerte Kondition** (bzgl. der 2-Norm) von \mathbf{A} .

Theorem 3.5.4. Für $m \geq n$, $\mathbf{A} \in \mathbb{K}^{m,n}$, $\text{rank}(\mathbf{A}) = n$, löse $\mathbf{x} \in \mathbb{K}^n$ das lineare Ausgleichsproblem $\|\mathbf{Ax} - \mathbf{b}\| \rightarrow \min$ und $\hat{\mathbf{x}}$ das lineare Ausgleichsproblem $\|(\mathbf{A} + \Delta\mathbf{A})\hat{\mathbf{x}} - \mathbf{b}\| \rightarrow \min$. Dann gilt

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_2}{\|\mathbf{x}\|_2} \leq \left(2 \text{cond}_2(\mathbf{A}) + \text{cond}_2^2(\mathbf{A}) \frac{\|\mathbf{r}\|_2}{\|\mathbf{A}\|_2 \|\mathbf{x}\|_2} \right) \frac{\|\Delta\mathbf{A}\|_2}{\|\mathbf{A}\|_2},$$

mit **Residuum** $\mathbf{r} = \mathbf{Ax} - \mathbf{b}$.

Falls $\|\mathbf{r}\|_2 \ll 1$ \triangleright Kondition(Lineares Ausgleichsproblem) $\approx \text{cond}_2(\mathbf{A})$

Falls $\|\mathbf{r}\|_2$ „gross“ \triangleright Kondition(Lineares Ausgleichsproblem) $\approx \text{cond}_2^2(\mathbf{A})$

3.5.1 Orthogonaltransformationsmethode

Annahme: $\mathbf{A} \in \mathbb{K}^{m,n}$, $m \geq n$, hat vollen Rang: $\text{rank}(\mathbf{A}) = n$

3.5
p. 273

QR-Zerlegung (\rightarrow Abschnitt 3.2.7): $\mathbf{A} = \mathbf{QR}$, $\mathbf{Q} \in \mathbb{K}^{m,m}$ unitär, $\mathbf{R} \in \mathbb{K}^{m,n}$ (reguläre) obere Dreiecksmatrix.

$$\|\mathbf{Ax} - \mathbf{b}\|_2 = \|\mathbf{Q}(\mathbf{Rx} - \mathbf{Q}^H \mathbf{b})\|_2 = \|\mathbf{Rx} - \tilde{\mathbf{b}}\|_2, \quad \tilde{\mathbf{b}} := \mathbf{Q}^H \mathbf{b}.$$

$$\|\mathbf{Ax} - \mathbf{b}\|_2 \rightarrow \min \Leftrightarrow \left\| \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} - \begin{pmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_m \end{pmatrix} \right\|_2 \rightarrow \min.$$

$$\mathbf{x} = \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_n \end{pmatrix}, \quad \text{Residuum } \mathbf{r} = \mathbf{Q} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \tilde{b}_{n+1} \\ \vdots \\ \tilde{b}_m \end{pmatrix}.$$

3.5
p. 274

3.5
p. 275

3.5
p. 276

Beachte:

$$\|\mathbf{r}\|_2 = \sqrt{\tilde{b}_{n+1}^2 + \dots + \tilde{b}_m^2}$$

Implementierung: Sukzessive Orthogonaltransformation (von links) von $\mathbf{A} \rightarrow \mathbf{R}$, $\mathbf{b} \rightarrow \tilde{\mathbf{b}}$

QR-basierter Algorithmus implementiert im Least-Squares-Löser des MATLAB-Operators \ (für vollbesetzte Matrizen)

Ausgleichsrechnung und SVD (für allgemeines $\mathbf{A} \in \mathbb{K}^{m,n}$, $\text{rank}(\mathbf{A}) = r \leq \min\{m, n\}$):

$$\mathbf{A} = [\mathbf{U}_1 \quad \mathbf{U}_2] \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{pmatrix}$$

$$\begin{pmatrix} \mathbf{A} \end{pmatrix} = \begin{pmatrix} \mathbf{U}_1 & \mathbf{U}_2 \end{pmatrix} \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{pmatrix},$$

$\mathbf{U}_1 \in \mathbb{K}^{m,r}$, $\mathbf{U}_2 \in \mathbb{K}^{m,m-r}$, $\Sigma_r = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{r,r}$, $\mathbf{V}_1 \in \mathbb{K}^{n,r}$, $\mathbf{V}_2 \in \mathbb{K}^{n,n-r}$, Spalten von $\mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1, \mathbf{V}_2$ orthonormiert.

$$\|\mathbf{Ax} - \mathbf{b}\|_2 = \left\| [\mathbf{U}_1 \quad \mathbf{U}_2] \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{pmatrix} \mathbf{x} - \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix} \right\|_2 = \left\| \begin{pmatrix} \Sigma_r \mathbf{V}_1^H \mathbf{x} \\ 0 \end{pmatrix} - \begin{pmatrix} \mathbf{U}_1^H \mathbf{b}_1 \\ \mathbf{U}_2^H \mathbf{b}_2 \end{pmatrix} \right\|_2$$

$$\mathbf{x} = (\mathbf{V}_1 \Sigma_r^{-1} \mathbf{U}_1^H) \mathbf{b}_1.$$

Praxis:

NUMRANK

Numerischer Rang:

$$r = \max\{i: \sigma_i / \sigma_1 > \text{tol}\}$$

MATLAB-CODE Lösung eines linearen

```
function y = lsqsvd(A,b)
[U,S,V] = svd(A,0);
sv = diag(S);
r = max(find(sv./sv(1) > eps));
y = V(:,1:r)*(diag(1./sv(1:r))*...
(U(:,1:r)'*b));
```

Ausgleichsproblems mit SVD

Theorem 3.5.5. Hat $\mathbf{A} \in \mathbb{K}^{m,n}$ die SVD (3.5.1), dann gilt $\mathbf{A}^+ = \mathbf{V}_1 \Sigma_r^{-1} \mathbf{U}_1^H$.

Anwendung der SVD: Minimierung auf der Einheitskugel:

Gegeben: $\mathbf{A} \in \mathbb{K}^{m,n}$, $m > n$: Suche $\mathbf{x} \in \mathbb{K}^n$, $\|\mathbf{x}\|_2 = 1$: $\|\mathbf{Ax}\|_2 \rightarrow \min$. (3.5.2)

Mit SVD $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^H$:

$$\min_{\|\mathbf{x}\|_2=1} \|\mathbf{Ax}\|_2^2 = \min_{\|\mathbf{x}\|_2=1} \|\mathbf{U} \Sigma \mathbf{V}^H \mathbf{x}\|_2^2 = \min_{\|\mathbf{V}^H \mathbf{x}\|_2=1} \|\mathbf{U} \Sigma (\mathbf{V}^H \mathbf{x})\|_2^2$$

$$= \min_{\|\mathbf{y}\|_2=1} \|\Sigma \mathbf{y}\|_2^2 = \min_{\|\mathbf{y}\|_2=1} (\sigma_1^2 y_1^2 + \dots + \sigma_n^2 y_n^2) \geq \sigma_n^2.$$

Das Minimum wird für $\mathbf{y} = \mathbf{e}_n$ angenommen $\Rightarrow \mathbf{x} = \mathbf{V} \mathbf{e}_n = \mathbf{v}_{\cdot,n}$.

Minimaler Wert = σ_n

Beispiel 121 (Fitten von Hyperebenen).

Hessesche Normalform einer Hyperebene im \mathbb{R}^d : $g = \{\mathbf{x} \in \mathbb{R}^d: c + \mathbf{n}^T \mathbf{x} = 0\}$, $\|\mathbf{n}\|_2 = 1$

Euklidischer Abstand $\text{dist}(g, \mathbf{y}) = |c + \mathbf{n}^T \mathbf{y}|$.

Ziel: Zu gegebenen Punkten $\mathbf{y}_1, \dots, \mathbf{y}_m$, $m > d$, finde $g \leftrightarrow c \in \mathbb{R}$, $\mathbf{n} \in \mathbb{R}^d$, $\|\mathbf{n}\|_2 = 1$, so dass

$$\sum_{j=1}^m \text{dist}(g, \mathbf{y}_j)^2 \rightarrow \min.$$

3.5
p. 277

3.5
p. 279

$$\|\mathbf{Ax}\|_2 \rightarrow \min \Leftrightarrow \left\| \begin{pmatrix} 1 & y_{1,1} & \dots & y_{1,n} \\ 1 & y_{2,1} & \dots & y_{2,n} \\ \vdots & \vdots & & \vdots \\ 1 & y_{m,1} & \dots & y_{m,n} \end{pmatrix} \begin{pmatrix} c \\ n_1 \\ \vdots \\ n_n \end{pmatrix} \right\|_2 \rightarrow \min \quad \text{wobei} \quad \|\mathbf{n}\|_2 = 1.$$

Mit QR-Zerlegung (\rightarrow Abschnitt 3.2.7)

$$\mathbf{A} := \begin{pmatrix} 1 & y_{1,1} & \dots & y_{1,n} \\ 1 & y_{2,1} & \dots & y_{2,n} \\ \vdots & \vdots & & \vdots \\ 1 & y_{m,1} & \dots & y_{m,n} \end{pmatrix} = \mathbf{QR}, \quad \mathbf{R} := \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1,d+1} \\ 0 & r_{22} & \dots & r_{2,d+1} \\ \vdots & \ddots & & \vdots \\ 0 & \dots & r_{d+1,d+1} & 0 \\ \vdots & & & \vdots \\ 0 & \dots & \dots & 0 \end{pmatrix} \in \mathbb{R}^{m,d+1}.$$

$$\|\mathbf{Ax}\|_2 \rightarrow \min \Leftrightarrow \|\mathbf{Rx}\|_2 = \left\| \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1,d+1} \\ 0 & r_{22} & \dots & r_{2,d+1} \\ \vdots & \ddots & & \vdots \\ 0 & \dots & r_{d+1,d+1} & 0 \\ \vdots & & & \vdots \\ 0 & \dots & \dots & 0 \end{pmatrix} \begin{pmatrix} c \\ n_1 \\ \vdots \\ n_d \end{pmatrix} \right\|_2 \rightarrow \min.$$

3.5
p. 280

3.5
p. 280

① Löse mittels SVD (Problem vom Typ $\begin{smallmatrix} \text{lsq:minconst} \\ (3.5.2) \end{smallmatrix}$):

$$\left\| \begin{pmatrix} r_{21} & r_{22} & \cdots & r_{2,d+1} \\ 0 & r_{33} & \cdots & r_{3,d+1} \\ \vdots & & \ddots & \vdots \\ 0 & & & r_{d+1,d+1} \end{pmatrix} \begin{pmatrix} n_1 \\ \vdots \\ n_d \end{pmatrix} \right\|_2 \rightarrow \min, \quad \|n\|_2 = 1.$$

② Da $r_{11} = \sqrt{d+1} \neq 0 \Rightarrow c = -r_{11}^{-1} \sum_{j=1}^d r_{1,j+1} n_j$.

MATLAB-Funktion:

Zu $A \in \mathbb{K}^{m,n}$ suche $\begin{smallmatrix} \text{MATLAB-CODE: Lineares Ausgleichsproblem mit Nebenbedingung} \\ \text{function [c,n] = clsq(A,dim);} \\ \text{[m,p] = size(A);} \\ \text{if p < dim+1, error('not enough unknowns'); end;} \\ \text{if m < dim, error('not enough equations'); end;} \\ \text{m = min(m, p);} \\ \text{R = triu(qr(A));} \\ \text{[U,S,V] = svd(R(p-dim+1:m,p-dim+1:p));} \\ \text{n = V(:,dim);} \\ \text{c = -R(1:p-dim,1:p-dim) \setminus R(1:p-dim,p-dim+1:p)*n;} \end{smallmatrix}$
 $n \in \mathbb{K}^{\dim}, c \in \mathbb{K}^{n-\dim}$
mit $\left\| A \begin{pmatrix} c \\ n \end{pmatrix} \right\|_2 \rightarrow \min$
mit Nebenbedingung: $\|n\|_2 = 1$.

3.5.2 Normalengleichungen

NORMGLEICH
 $A \in \mathbb{K}^{m,n}, m \geq n$, mit vollem Rang $\text{rank}(A) = n$. Zu lösen ist $\begin{smallmatrix} \text{lsq:LSO} \\ (3.5.1) \end{smallmatrix}$ für $b \in \mathbb{K}^m$

$$x \in \mathbb{K}^n: \|Ax - b\|_2 \rightarrow \min. \quad (3.5.3) \quad \text{sq:prb}$$

$$f(x) := \|Ax - b\|_2^2 = x^H (A^H A) x - 2b^H A x + b^H b.$$

$$\Rightarrow \text{grad } f(x) = 2(A^H A)x - 2A^H b.$$

$$\text{grad } f(x) \stackrel{!}{=} 0: \quad A^H A x = A^H b \quad = \text{Normalengleichung zu } \begin{smallmatrix} \text{lsq:prb} \\ (3.5.3) \end{smallmatrix}.$$

Beachte: $\text{rank}(A) = n \Rightarrow A^H A \in \mathbb{K}^{n,n}$ regulär



Vorsicht: Gefahr der Instabilität: Mit SVD $A = U \Sigma V^H$

$$\text{cond}_2(A^H A) = \text{cond}_2(V \Sigma^H U^H U \Sigma V^H) = \text{cond}_2(\Sigma^H \Sigma) = \frac{\sigma_1^2}{\sigma_n^2} = \text{cond}_2(A)^2.$$

Vorsicht: Informationsverlust bei Berechnung von $A^H A$, z.B.



$$A = \begin{pmatrix} 1 & 1 \\ \delta & 0 \\ 0 & \delta \end{pmatrix} \Rightarrow A^H A = \begin{pmatrix} 1+\delta^2 & 1 \\ 1 & 1+\delta^2 \end{pmatrix}$$

Wenn $\delta < \sqrt{\text{eps}}$ $\Rightarrow 1 + \delta^2 = 1$ in \mathbb{M} , d.h. $A^H A$ „numerisch singulär“, obwohl $\text{rank}(A) = 2$

```
>> A = [1 1;...
        sqrt(eps) 0;...
        0 sqrt(eps)];
>> rank(A)
ans = 2
>> rank(A'*A)
ans = 1
```

Vermeidung der Berechnung von $A^H A$: Residuum $r := Ax - b$ als neue Unbekannte:

$$A^H A x = A^H b \Leftrightarrow B \begin{pmatrix} r \\ x \end{pmatrix} := \begin{pmatrix} -I & A \\ A^H & 0 \end{pmatrix} \begin{pmatrix} r \\ x \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix}. \quad (3.5.4) \quad \text{sq:trafo}$$

Allgemeinere Substitution $z := \alpha x, r := \alpha^{-1} A x - b, \alpha > 0$ (zur Verbesserung der Kondition):

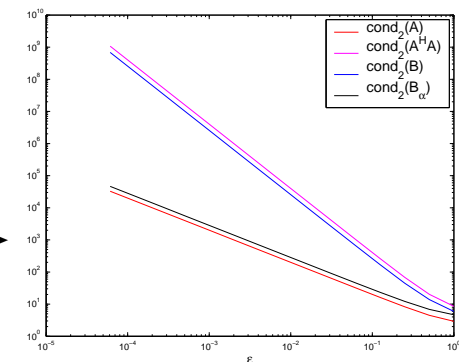
$$A^H A x = A^H b \Leftrightarrow B_\alpha \begin{pmatrix} r \\ z \end{pmatrix} := \begin{pmatrix} -\alpha I & A \\ A^H & 0 \end{pmatrix} \begin{pmatrix} r \\ z \end{pmatrix} = \begin{pmatrix} \alpha b \\ 0 \end{pmatrix} \quad (3.5.5) \quad \text{sq:trafo2}$$

Beispiel 122 (Kondition des erweiterten Systems).

3.5
p. 281

$$A = \begin{pmatrix} 1 + \epsilon & 1 \\ 1 - \epsilon & 1 \\ \epsilon & \epsilon \end{pmatrix}$$

Plot verschiedener Konditionszahlen in Abhängigkeit von ϵ ($\alpha = \|A\|_2 / \sqrt{2}$)



Bemerkung 123 (Normalengleichung vs. Orthogonaltransformationsmethode).

Vorteil der Normalengleichungen in der Form $\begin{smallmatrix} \text{lsq:trafo} \\ (3.5.4)/(3.5.5) \end{smallmatrix}$:

A sparse (\Rightarrow Def. 3.1.1) $\Rightarrow B/B_\alpha$ sparse $\begin{smallmatrix} \text{def:sparse} \end{smallmatrix}$

3.5
p. 282

3.5
p. 283

△ p. 284

3.5.3 Totales Ausgleichsproblem

TOTALAUSGLEICH

Gegeben: $\mathbf{A} \in \mathbb{K}^{m,n}$, $m \geq n$, $\text{rank}(\mathbf{A}) = n$, $\mathbf{b} \in \mathbb{K}^m$

Gesucht: $\hat{\mathbf{A}} \in \mathbb{K}^{m,n}$, $\hat{\mathbf{b}} \in \mathbb{K}^m$ mit

$$\left\| \underbrace{\begin{bmatrix} \mathbf{A} & \mathbf{b} \end{bmatrix}}_{=: \mathbf{C}} - \underbrace{\begin{bmatrix} \hat{\mathbf{A}} & \hat{\mathbf{b}} \end{bmatrix}}_{=: \hat{\mathbf{C}}} \right\|_F \rightarrow \min, \quad \hat{\mathbf{b}} \in \text{Im}(\hat{\mathbf{A}}).$$

$$\hat{\mathbf{b}} \in \text{Im}(\hat{\mathbf{A}}) \Rightarrow \text{rank}(\hat{\mathbf{C}}) = n \xrightarrow{(3.5.6) \Rightarrow \text{lsq:total}} \min_{\text{rank}(\hat{\mathbf{C}})=n} \left\| \mathbf{C} - \hat{\mathbf{C}} \right\|_F.$$

Thm. 3.4.16 $\xrightarrow{\text{thm:rankapprox}}$ Verwende SVD $\mathbf{C} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$

$$\hat{\mathbf{C}} = \sum_{j=1}^n \sigma_j \mathbf{u}_{:,j} \mathbf{v}_{:,j}^H \Rightarrow \hat{\mathbf{C}} \mathbf{v}_{:,n+1} = 0.$$

Wenn $\mathbf{v}_{n+1,n+1} \neq 0$, dann

$$\hat{\mathbf{A}} \mathbf{x} = \hat{\mathbf{b}} \quad \text{mit} \quad \mathbf{x} = \mathbf{v}_{n+1,n+1}^{-1} \mathbf{v}_{:,n+1}.$$

MATLAB-CODE Totales Ausgleichsproblem

```
function x = lsqtotal(A,b);
[m,n]=size(A);
[U, Sigma, V] = svd([A,b]);
s = V(n+1,n+1);
if s == 0,
    error('No solution')
end
x = -V(1:n,n+1)/s;
```

3.5.4 Ausgleichsrechnung mit linearen Nebenbedingungen

AUSGLEICHPROBLEME

Gegeben: $\mathbf{A} \in \mathbb{K}^{m,n}$, $m \geq n$, $\text{rank}(\mathbf{A}) = n$, $\mathbf{b} \in \mathbb{K}^m$,
 $\mathbf{C} \in \mathbb{K}^{p,n}$, $p < n$, $\text{rank}(\mathbf{C}) = p$, $\mathbf{d} \in \mathbb{K}^p$

Gesucht: $\mathbf{x} \in \mathbb{K}^n$ mit $\left\| \mathbf{A}\mathbf{x} - \mathbf{b} \right\|_2 \rightarrow \min$, $\mathbf{C}\mathbf{x} = \mathbf{d}$.

(3.5.7) $\xrightarrow{\text{lsq:lin}}$
p. 286

Lineare Nebenbedingung



Idee: Ankopplung der Nebenbedingung durch $\xrightarrow{\text{LagrMult}}$ Lagrange-Multiplikator $\mathbf{m} \in \mathbb{K}^p$

$$\mathbf{x} = \underset{\mathbf{x} \in \mathbb{K}^n}{\text{argmin}} \max_{\mathbf{m} \in \mathbb{K}^p} L(\mathbf{x}, \mathbf{m}), \quad L(\mathbf{x}, \mathbf{m}) := \frac{1}{2} \left\| \mathbf{A}\mathbf{x} - \mathbf{b} \right\|^2 + \mathbf{m}^H (\mathbf{C}\mathbf{x} - \mathbf{d}).$$

Notwendige (und hinreichende) Bedingungen für Lösung (\rightarrow Abschnitt 3.5.2) $\xrightarrow{\text{sec:normalengleichungen}}$

$$\frac{\partial L}{\partial \mathbf{x}}(\mathbf{x}, \mathbf{m}) = \mathbf{A}^H (\mathbf{A}\mathbf{x} - \mathbf{b}) + \mathbf{C}^H \mathbf{m} \stackrel{!}{=} 0, \quad \frac{\partial L}{\partial \mathbf{m}}(\mathbf{x}, \mathbf{m}) = \mathbf{C}\mathbf{x} - \mathbf{d} \stackrel{!}{=} 0;$$

$$\begin{pmatrix} \mathbf{A}^H \mathbf{A} & \mathbf{C}^H \\ \mathbf{C} & 0 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{m} \end{pmatrix} = \begin{pmatrix} \mathbf{A}^H \mathbf{b} \\ \mathbf{d} \end{pmatrix} \quad \xrightarrow{\text{ERNORMGL}} \text{Erweiterte Normalengleichungen (Sattelpunktproblem)}$$

(3.5.6) $\xrightarrow{\text{sq:total}}$ Algorithmus (basierend auf Block-LU-Zerlegung):

$$\begin{pmatrix} \mathbf{A}^H \mathbf{A} & \mathbf{C}^H \\ \mathbf{C} & 0 \end{pmatrix} = \begin{pmatrix} \mathbf{R}^H & 0 \\ \mathbf{G} & -\mathbf{S}^H \end{pmatrix} \begin{pmatrix} \mathbf{R} & \mathbf{G}^H \\ 0 & \mathbf{S} \end{pmatrix}, \quad \begin{matrix} \mathbf{R}, \mathbf{S} \in \mathbb{K}^{n,n} \text{ obere Dreiecksmatrizen} \\ \mathbf{G} \in \mathbb{K}^{p,n} \end{matrix}$$

3.5
p. 285

$\mathbf{R}^H \mathbf{R} = \mathbf{A}^H \mathbf{A} \rightarrow$ Cholesky-Zerlegung \rightarrow Sect. 3.2.5 $\xrightarrow{\text{sec:symm-posit-defin}}$
 $\mathbf{R}^H \mathbf{G}^H = \mathbf{C}^H \rightarrow n$ Vorwärtssubstitutionen \rightarrow Sect. 3.2.3 $\xrightarrow{\text{sec:die-lr-zerlegung}}$
 $\mathbf{G} \mathbf{G}^H = \mathbf{S}^H \mathbf{S} \rightarrow$ Cholesky-Zerlegung \rightarrow Sect. 3.2.5 $\xrightarrow{\text{sec:symm-posit-defin}}$

3.5
p. 287

Vorsicht:

Abschn. 3.5.2: Berechnung von $\mathbf{A}^H \mathbf{A}$ kann teuer und problematisch sein!
(Abhilfe durch Einführung einer neuen Unbekannten $\mathbf{r} = \mathbf{A}\mathbf{x} - \mathbf{b}$, vgl. (3.5.4)) $\xrightarrow{\text{lsq:trafo}}$

$$\xrightarrow{\text{lsq:trafo}} \begin{pmatrix} -\mathbf{I} & \mathbf{A} & 0 \\ \mathbf{A}^H & 0 & \mathbf{C}^H \\ 0 & \mathbf{C} & 0 \end{pmatrix} \begin{pmatrix} \mathbf{r} \\ \mathbf{x} \\ \mathbf{m} \end{pmatrix} = \begin{pmatrix} \mathbf{b} \\ 0 \\ \mathbf{d} \end{pmatrix}. \quad (3.5.8) \quad \xrightarrow{\text{sq:ctrl}}$$

SVD-basierter Algorithmus:

① Berechne Orthonormalbasis von $\text{Ker}(\mathbf{C})$ mit SVD (\rightarrow Abschnitt 3.5.1) $\xrightarrow{\text{sec:orth}}$

$$\mathbf{C} = \mathbf{U} \begin{bmatrix} \mathbf{\Sigma} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^H \\ \mathbf{V}_2^H \end{bmatrix}, \quad \mathbf{U} \in \mathbb{K}^{p,p}, \mathbf{\Sigma} \in \mathbb{R}^{p,p}, \mathbf{V}_1 \in \mathbb{K}^{n,p}, \mathbf{V}_2 \in \mathbb{K}^{n,n-p}$$

$$\xrightarrow{\text{lsq:trafo}} \text{Ker}(\mathbf{C}) = \text{Im}(\mathbf{V}_2).$$

und partikuläre Lösung

$$\mathbf{x}_0 := \mathbf{V}_1 \mathbf{\Sigma}^{-1} \mathbf{U}^H \mathbf{d}.$$

Darstellung für Lösung \mathbf{x} von (3.5.7) $\xrightarrow{\text{lsq:lin}}$ $\mathbf{x} = \mathbf{x}_0 + \mathbf{V}_2 \mathbf{y}, \quad \mathbf{y} \in \mathbb{K}^{n-p}$

② Einsetzen der Darstellung \rightarrow Standard lineares Ausgleichsproblem

$$\left\| \mathbf{A}(\mathbf{x}_0 + \mathbf{V}_2 \mathbf{y}) - \mathbf{b} \right\|_2 \rightarrow \min \Leftrightarrow \left\| \mathbf{A} \mathbf{V}_2 \mathbf{y} - (\mathbf{b} - \mathbf{A} \mathbf{x}_0) \right\|_2 \rightarrow \min.$$

3.5
p. 286

3.6
p. 288

3.6 Krylov-Verfahren für lineare Gleichungssysteme

[File: section-krylov-verfahren-fuer-lineare-gleichungssysteme.tex, SVN: section-krylov-verfahren-fuer-lineare-gleichungssysteme.tex 1062 2006-10-21 09:33:00]

Klasse von **Iterationsverfahren** (→ Abschnitt 2.1) zur approximativen Lösung linearer Gleichungssysteme $\mathbf{Ax} = \mathbf{b}$.

- Charakteristika:
- Nur Auswertung $\mathbf{A} \times \text{Vektor}$ erforderlich
 - (In vielen Fällen) Aufwand für einen Schritt \approx Kosten($\mathbf{A} \times \text{Vektor}$)

Iterationsverfahren zur Lösung linearer Gleichungssysteme !? Warum ?

- Eliminationsverfahren zu aufwändig (z.B. für \mathbf{A} sehr gross, dünnbesetzt), → (3.2.7)
- Direktes Lösen „overkill“, da nur geringe Genauigkeit erforderlich (z.B. im Newton-Verfahren → Abschnitt 2.5)
- Gute Näherung als Startwert verfügbar
- Koeffizientenmatrix nur als Prozedur $y = \text{evalA}(x) \leftrightarrow y = \mathbf{Ax}$ gegeben

Beobachtung: Lösungen vieler LGS aus Wissenschaft/Technik mit s.p.d. Koeffizientenmatrix weisen dominierende Anteile von Eigenvektoren zu betragsgrössten/betragskleinsten Eigenwerten auf.

[siehe Motivation für Krylovräume, Abschnitt 3.4.3]

Idee: Gegeben Näherungslösung $\mathbf{x}^{(0)}$, suche Lösung der **Fehlergleichung**

$$\mathbf{Ac} = \mathbf{r}, \quad \text{Residuum } \mathbf{r} := \mathbf{b} - \mathbf{Ax}^{(0)} \quad (3.6.1)$$

im Krylovraum $\mathcal{K}_l(\mathbf{A}, \mathbf{z})$ (→ Def. 3.4.8), „ \mathbf{z}, l „geeignet“

! Beachte: $\dim \mathcal{K}_l(\mathbf{A}, \mathbf{z}) \leq n$ → überbestimmtes Problem

3.6.1 Das Verfahren der konjugierten Gradienten (CG)

Annahme: $\mathbf{A} = \mathbf{A}^H \in \mathbb{K}^{n,n}$, $n \in \mathbb{N}$, Hermitesch, positiv definit (→ Def. 3.2.8)

A-Skalarprodukt $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x}^H \mathbf{Ay} \Rightarrow$ „**A**-Geometrie“

Definition 3.6.1 (Energienorm). Die von einer s.p.d. Matrix $\mathbf{A} \in \mathbb{K}^{n,n}$ induzierte **Energienorm** ist

$$\|\mathbf{x}\|_A := (\mathbf{x}^H \mathbf{Ax})^{1/2}, \quad \mathbf{x} \in \mathbb{K}^n.$$

3.6.1.1 Prinzip des CG-Verfahrens

Idee:

- Abbruch für exakte Lösung: $\mathbf{z} = \mathbf{r} := \mathbf{b} - \mathbf{Ax}^{(0)}$ (Residuum)
- Fixiere \mathbf{c} aus (3.6.1) durch **Galerkin-Bedingung**

$$\mathbf{r} - \mathbf{Ac} \perp \mathcal{K}_l(\mathbf{A}, \mathbf{r}) \Leftrightarrow \mathbf{r}^{(\text{neu})} := \mathbf{b} - \mathbf{Ax}^{(\text{neu})} \perp \mathcal{K}_l(\mathbf{A}, \mathbf{r}). \quad (3.6.2)$$

Euklidisch orthogonal

(Abstraktes) Verfahren der **konjugierten Gradienten (CG)** berechnet Näherungslösung $\mathbf{x}^{(l)}$ des LGS $\mathbf{Ax} = \mathbf{b}$ durch Korrektur der Startnäherung $\mathbf{x}^{(0)}$ in $\mathcal{K}_l(\mathbf{A}, \mathbf{r})$, $\mathbf{r} := \mathbf{b} - \mathbf{Ax}^{(0)}$, festgelegt durch **Galerkin-Bedingung**

falls $\{\mathbf{p}_1, \dots, \mathbf{p}_l\} = \text{Basis von } \mathcal{K}_l(\mathbf{A}, \mathbf{r})$

$$\mathbf{p}_j^H (\mathbf{r} - \mathbf{Ac}) = \mathbf{p}_j^H (\mathbf{b} - \mathbf{Ax}^{(\text{neu})}) = 0 \quad \forall j = 1, \dots, l. \quad (3.6.3)$$

Berechne Koeffizienten der Basisdarstellung $\mathbf{c} = \gamma_1 \mathbf{p}_1 + \dots + \gamma_l \mathbf{p}_l$ aus

$$\begin{pmatrix} \mathbf{p}_1^H \mathbf{Ap}_1 & \dots & \mathbf{p}_1^H \mathbf{Ap}_l \\ \vdots & & \vdots \\ \mathbf{p}_l^H \mathbf{Ap}_1 & \dots & \mathbf{p}_l^H \mathbf{Ap}_l \end{pmatrix} \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_l \end{pmatrix} = \begin{pmatrix} \mathbf{p}_1^H \mathbf{r} \\ \vdots \\ \mathbf{p}_l^H \mathbf{r} \end{pmatrix}. \quad (3.6.4)$$

Besonders einfach, wenn $\{\mathbf{p}_1, \dots, \mathbf{p}_l\}$ **A-orthogonale Basis** von $\mathcal{K}_l(\mathbf{A}, \mathbf{r})$: $\mathbf{p}_j^H \mathbf{Ap}_i = 0$ für $i \neq j$.

Bemerkung 124 (CG als direkter Löser).

Falls $\mathbf{A}^l \mathbf{r} \in \mathcal{K}_l(\mathbf{A}, \mathbf{r})$, dann $\mathbf{x}^{(l)} = \mathbf{A}^{-1} \mathbf{b}$ für die Näherungslösung aus dem abstrakten CG-Algorithmus.

CG terminiert nach spätestens n Schritten mit $\mathbf{x}^{(l)} = \mathbf{A}^{-1} \mathbf{b}$.

Sei l minimal, so dass $\mathbf{A}^l \mathbf{r} \in \mathcal{K}_l(\mathbf{A}, \mathbf{r})$. Nach dem Austauschaixiom $\mathbf{r} \in \text{Span}\{\mathbf{Ar}, \dots, \mathbf{A}^l \mathbf{r}\} = \mathcal{AK}_l(\mathbf{A}, \mathbf{r})$, also $\mathbf{A}(\mathbf{x}^* - \mathbf{x}^{(0)}) \in \mathcal{AK}_l(\mathbf{A}, \mathbf{r}) \Leftrightarrow \mathbf{x}^* - \mathbf{x}^{(0)} \in \mathcal{K}_l(\mathbf{A}, \mathbf{r})$, so dass die Korrektur die exakte Lösung liefert.

Bemerkung 125 (CG-Verfahren und quadratische Minimierung).

Lemma 3.6.2 (S.p.d. Gleichungssysteme und Minimierungsproblem). Für $\mathbf{A} \in \mathbb{K}^{n,n}$ s.p.d., $\mathbf{b} \in \mathbb{K}^n$ gilt

$$\mathbf{Ax} = \mathbf{b} \Leftrightarrow \mathbf{x} = \arg \min_{\mathbf{y} \in \mathbb{K}^n} J(\mathbf{y}), \quad J(\mathbf{y}) := \frac{1}{2} \mathbf{y}^H \mathbf{A} \mathbf{y} - \operatorname{Re}\{\mathbf{b}^H \mathbf{y}\}.$$

Unterraumkorrektur: Sei $V \subset \mathbb{K}^n$ Unterraum mit Basis $\{\mathbf{p}_1, \dots, \mathbf{p}_l\}$, $l \ll n$.

Approximatives Lösen der Fehlergleichung (3.6.1) durch „Minimierung in Richtung von V “:

(i) $\mathbf{c} = \operatorname{argmin}_{\mathbf{y} \in V} J(\mathbf{y})$, $J(\mathbf{y}) := \frac{1}{2} \mathbf{y}^H \mathbf{A} \mathbf{y} - \operatorname{Re}\{(\mathbf{b} - \mathbf{Ax}^{(0)})^H \mathbf{y}\}$, (3.6.5) GG:1

(ii) Neue Näherung $\mathbf{x}^{(\text{neu})} := \mathbf{x}^{(0)} + \mathbf{c}$. (3.6.6) GG:2

Lemma 3.6.3. \mathbf{c} , $\mathbf{x}^{(\text{neu})}$ aus (3.6.5), (3.6.6) eindeutig bestimmt durch die Galerkin-Bedingung (3.6.3). GG:1, GG:2, GG:3

Bemerkung 126. Für $\mathbf{b} \in \mathbb{R}^n$, $J(\mathbf{y}) := \frac{1}{2} \mathbf{y}^H \mathbf{A} \mathbf{y} - \operatorname{Re}\{\mathbf{b}^H \mathbf{y}\}$ ist $-\operatorname{grad} J(\mathbf{x}) = \mathbf{b} - \mathbf{Ax}$ (\Rightarrow Residuum = Abstiegsrichtung):

Wähle Residuum als Grundlage des Krylov-Raumes.

3.6.1.2 Implementierung des CG-Verfahrens

IMPLEMENTCG

Annahme: \mathbf{A} -orthogonale Basis $\{\mathbf{p}_1, \dots, \mathbf{p}_n\}$ von \mathbb{R}^n verfügbar, so dass

$$\operatorname{Span}\{\mathbf{p}_1, \dots, \mathbf{p}_l\} = \mathcal{K}_l(\mathbf{A}, \mathbf{r})$$

(Effiziente) sukzessive Berechnung der $\mathbf{x}^{(l)}$

Input : Anfangsnäherung $\mathbf{x}^{(0)} \in \mathbb{K}^n$
Output : Bessere Näherungslösung $\mathbf{x}^{(l)} \in \mathbb{K}^n$

$$\mathbf{r}^{(0)} := \mathbf{b} - \mathbf{Ax}^{(0)};$$

$$\text{for } j = 1 \text{ to } l \text{ do } \left\{ \mathbf{x}^{(j)} := \mathbf{x}^{(j-1)} + \frac{\mathbf{p}_j^H \mathbf{r}^{(0)}}{\mathbf{p}_j^H \mathbf{A} \mathbf{p}_j} \mathbf{p}_j; \right\} \quad (3.6.7) \quad \text{ALG: Kryl}$$

Für $1 \leq j \leq m \leq l$: Orthogonalität $\mathbf{p}_j \perp \mathbf{r}^{(m)} := \mathbf{b} - \mathbf{Ax}^{(m)}$

$$\mathbf{p}_j \cdot (\mathbf{b} - \mathbf{Ax}^{(m)}) = \mathbf{p}_j^H (\underbrace{\mathbf{b} - \mathbf{Ax}^{(0)}}_{=\mathbf{r}^{(0)}}) - \sum_{k=1}^m \frac{\mathbf{p}_k^H \mathbf{r}^{(0)}}{\mathbf{p}_k^H \mathbf{A} \mathbf{p}_k} \mathbf{A} \mathbf{p}_k = 0. \quad (3.6.8) \quad \text{KrylOrth}$$

Aufgabe: Konstruktion einer \mathbf{A} -orthogonalen Basis $\{\mathbf{p}_1, \dots, \mathbf{p}_l\}$ des Krylov-Raums $\mathcal{K}_l(\mathbf{A}, \mathbf{r}^{(0)})$ (\rightarrow Def. 3.4.8, Arnoldi-Prozess Abschnitt 3.4.3). def: Kryl, sec: kryl-unterr

(3.6.8) \Rightarrow Idee: Gram-Schmidt-Orthogonalisierung, \rightarrow (3.4.8), der intermediären Residuen $\mathbf{r}^{(j)} := \mathbf{b} - \mathbf{Ax}^{(j)}$: GRS, def: GRS, sec: GRS



$$\mathbf{p}_1 := \mathbf{r}^{(0)}, \mathbf{p}_{j+1} := (\underbrace{\mathbf{b} - \mathbf{Ax}^{(j)}}_{=\mathbf{r}^{(j)}}) - \sum_{k=1}^j \frac{\mathbf{p}_k^H \mathbf{A} \mathbf{r}^{(j)}}{\mathbf{p}_k^H \mathbf{A} \mathbf{p}_k} \mathbf{p}_k, \quad j = 1, \dots, l-1. \quad (3.6.9) \quad \text{GRS}$$

Lemma 3.6.4. Nichtverschwindende Vektoren \mathbf{p}_j , $1 \leq j \leq l$, und $\mathbf{r}^{(j)} := \mathbf{b} - \mathbf{Ax}^{(j)}$, $0 \leq j \leq l$, aus (3.6.7), (3.6.9) erfüllen ALG: KrylGRS

(i) $\{\mathbf{p}_1, \dots, \mathbf{p}_j\}$ ist \mathbf{A} -orthogonale Basis von $\mathcal{K}_j(\mathbf{A}, \mathbf{r}^{(0)})$,

(ii) $\{\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(j-1)}\}$ ist Orthogonalbasis von $\mathcal{K}_j(\mathbf{A}, \mathbf{r}^{(0)})$.

Beweis. \mathbf{A} -Orthogonalität der \mathbf{p}_j nach Konstruktion \rightarrow (3.6.9). GRS

$$(3.6.7) \text{ \& (3.6.9) } \Rightarrow \mathbf{p}_{j+1} = \mathbf{r}^{(0)} - \sum_{k=1}^j \frac{\mathbf{p}_k^H \mathbf{r}^{(0)}}{\mathbf{p}_k^H \mathbf{A} \mathbf{p}_k} \mathbf{A} \mathbf{p}_k - \sum_{k=1}^j \frac{\mathbf{p}_k^H \mathbf{A} \mathbf{r}^{(j)}}{\mathbf{p}_k^H \mathbf{A} \mathbf{p}_k} \mathbf{p}_k.$$

$$\Rightarrow \mathbf{p}_{j+1} \in \operatorname{Span}\{\mathbf{r}^{(0)}, \mathbf{p}_1, \dots, \mathbf{p}_j, \mathbf{A} \mathbf{p}_1, \dots, \mathbf{A} \mathbf{p}_j\}.$$

$$\mathbf{p}_1 := \mathbf{r}^{(0)}, \mathbf{p}_{j+1} := \underbrace{(\mathbf{b} - \mathbf{A}\mathbf{x}^{(j)})}_{=: \mathbf{r}^{(j)}} - \sum_{k=1}^j \frac{\mathbf{p}_k^H \mathbf{r}^{(j)}}{\mathbf{p}_k^H \mathbf{A} \mathbf{p}_k} \mathbf{p}_k, \quad j = 1, \dots, l-1, \quad \text{GRS (3.6.9)}$$

$$\mathbf{p}_1 := \mathbf{r}^{(0)}, \mathbf{p}_{j+1} := \underbrace{(\mathbf{b} - \mathbf{A}\mathbf{x}^{(j)})}_{=: \mathbf{r}^{(j)}} - \sum_{k=1}^j \frac{\mathbf{p}_k^H \mathbf{r}^{(j)}}{\mathbf{p}_k^H \mathbf{A} \mathbf{p}_k} \mathbf{p}_k, \quad j = 1, \dots, l-1. \quad \text{KrylOrth (3.6.8)}$$

$$\text{GRS (3.6.9) erglantz} \mathbf{r}^{(j)} \in \text{Span} \{ \mathbf{p}_1, \dots, \mathbf{p}_{j+1} \} \quad \& \quad \mathbf{p}_j \in \text{Span} \{ \mathbf{r}^{(0)}, \dots, \mathbf{r}^{(j-1)} \}.$$

$$\blacktriangleright \text{Span} \{ \mathbf{p}_1, \dots, \mathbf{p}_j \} = \text{Span} \{ \mathbf{r}^{(0)}, \dots, \mathbf{r}^{(j-1)} \} = \mathcal{K}_l(\mathbf{A}, \mathbf{r}^{(0)}).$$

$$\text{KrylOrth (3.6.8) resorth} \mathbf{r}^{(j)} \perp \text{Span} \{ \mathbf{p}_1, \dots, \mathbf{p}_j \} = \text{Span} \{ \mathbf{r}^{(0)}, \dots, \mathbf{r}^{(j-1)} \}. \quad \text{(3.6.12) resorth}$$

Bemerkung 127 (Lanczos-Prozess und CG).

$\{ \mathbf{r}^{(0)}, \dots, \mathbf{r}^{(j-1)} \} \leftrightarrow$ Orthonormalbasis von $\mathcal{K}_j(\mathbf{A}, \mathbf{r})$ erzeugt im Lanczos-Prozess \rightarrow Abschn. 3.4.3 sec:kryl,run p. 297

Analog zum Lanczos-Prozess: **Kurze Rekursion** für Basisvektoren \mathbf{p}_j :

$$\text{KrylOrth (3.6.8)} \Rightarrow \text{GRS (3.6.9) vereinfacht sich zu } \mathbf{p}_{j+1} := \mathbf{r}^{(j)} - \frac{\mathbf{p}_j^H \mathbf{A} \mathbf{r}^{(j)}}{\mathbf{p}_j^H \mathbf{A} \mathbf{r}^{(j)}} \mathbf{p}_j, \quad j = 1, \dots, l.$$

Rekursion für Residuen:

$$\text{ALG:Krylit (3.6.7)} \blacktriangleright \mathbf{r}^{(j)} = \mathbf{r}^{(j-1)} + \frac{\mathbf{p}_j^H \mathbf{r}^{(0)}}{\mathbf{p}_j^H \mathbf{A} \mathbf{p}_j} \mathbf{A} \mathbf{p}_j. \quad \text{(3.6.13)}$$

$$\text{Lemma 3.6.4, (i) Lem:CG} \blacktriangleright \mathbf{r}^{(j-1)} \cdot \mathbf{p}_j = \left(\mathbf{r}^{(0)} + \sum_{k=1}^{m-1} \frac{\mathbf{r}^{(0)} \cdot \mathbf{p}_k}{\mathbf{p}_k^H \mathbf{A} \mathbf{p}_k} \mathbf{A} \mathbf{p}_k \right) \cdot \mathbf{p}_j = \mathbf{r}^{(0)} \cdot \mathbf{p}_j. \quad \text{(3.6.14)}$$

Algorithmus: **CG-Verfahren** zur Lösung von $\mathbf{A}\mathbf{x} = \mathbf{b}$, \mathbf{A} s.p.d.

Input : Anfangsnäherung $\mathbf{x}^{(0)} \in \mathbb{R}^n$
Output : Bessere Näherungslösung $\mathbf{x}^{(l)} \in \mathbb{R}^n$

```

 $\mathbf{p}_1 = \mathbf{r}^{(0)} := \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)};$ 
for  $j = 1$  to  $l$  do {
   $\mathbf{x}^{(j)} := \mathbf{x}^{(j-1)} + \frac{\mathbf{p}_j \cdot \mathbf{r}^{(j-1)}}{\mathbf{p}_j \cdot \mathbf{A} \mathbf{p}_j} \mathbf{p}_j;$ 

   $\mathbf{r}^{(j)} = \mathbf{r}^{(j-1)} - \frac{\mathbf{p}_j \cdot \mathbf{r}^{(j-1)}}{\mathbf{p}_j \cdot \mathbf{A} \mathbf{p}_j} \mathbf{A} \mathbf{p}_{j-1};$ 

   $\mathbf{p}_{j+1} = \mathbf{r}^{(j)} - \frac{\mathbf{A} \mathbf{p}_j \cdot \mathbf{r}^{(j)}}{\mathbf{p}_j \cdot \mathbf{A} \mathbf{p}_j} \mathbf{p}_j;$ 
}

```

Input: Anfangsnäherung $\mathbf{x} \doteq \mathbf{x}^{(0)} \in \mathbb{R}^n$
Abbruchschranke $\tau > 0$
Output: Bessere Näherungslösung $\mathbf{x} \doteq \mathbf{x}^{(l)}$

```

 $\mathbf{p} := \mathbf{r}^{(0)} := \mathbf{r} := \mathbf{b} - \mathbf{A}\mathbf{x};$ 
for  $j = 1$  to  $l_{\max}$  do {
   $\beta := \mathbf{r} \cdot \mathbf{r};$ 
   $\mathbf{h} := \mathbf{A}\mathbf{p};$ 
   $\alpha := \frac{\beta}{\mathbf{p}^H \mathbf{h}};$ 
   $\mathbf{x} := \mathbf{x} + \alpha \mathbf{p};$ 
   $\mathbf{r} := \mathbf{r} - \alpha \mathbf{h};$ 
  if  $\|\mathbf{r}\| \leq \tau \|\mathbf{r}^{(0)}\|$  then stop;
   $\beta := \frac{\beta}{\mathbf{r} \cdot \mathbf{r}};$ 
   $\mathbf{p} := \mathbf{r} + \beta \mathbf{p};$ 
}

```

\blacktriangleright 1 Matrix×Vektor-Produkt, 3 Skalarprodukte, 3 **SAXPY-Operationen** pro Schritt:
Wenn \mathbf{A} dünnbesetzt, $\text{nnz}(\mathbf{A}) \sim n \blacktriangleright$ Rechenaufwand $O(n)$ für einen Schritt

MATLAB-Funktion:

$\mathbf{x} = \text{pcg}(\mathbf{A}, \mathbf{b}, \text{tol}, \text{maxit}, [], [], \mathbf{x}_0)$: Löse $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit maximal maxit CG Schritten:
terminiere wenn $\|\mathbf{r}^{(l)}\| : \|\mathbf{r}^{(0)}\| < \text{tol}$.
 $\mathbf{x} = \text{pcg}(\text{Afun}, \mathbf{b}, \text{tol}, \text{maxit}, [], [], \mathbf{x}_0)$: Afun = Handle auf Funktion, die Auswertung $\mathbf{A} \times \text{Vektor}$ realisiert.
 $[\mathbf{x}, \text{flag}, \text{relr}, \text{it}, \text{resv}] = \text{pcg}(\dots)$: Diagnostische Informationen über Verlauf der Iteration

Grundgerüst der Implementierung in MATLAB (type pcg):

MATLAB-CODE CG-Algorithmus

```

function x = cg(A,b,tol,maxit)
r = b - A * x; rho = 1;
for i = 1 : maxit
    rho1 = rho; rho = r' * r;
    if (i == 1), p = r;
    else beta = rho/rho1; p = r + beta * p; end
    q = A * p; alpha = rho/(p' * q);
    x = x + alpha * p;
    if (norm(b - A * x) <= tol*norm(b)) return; end
    r = r - alpha * q;
end

```

Bemerkung 128 (Abbruchkriterium für nicht-vorkonditioniertes CG-Verfahren).

Mit $\mathbf{r}^{(l)} := \mathbf{b} - \mathbf{A}\mathbf{x}^{(l)}$ gilt für beliebige Vektornormen und die zugehörige Matrixkondition (\rightarrow Def. 3.2.5)

$$\frac{1}{\text{cond}(\mathbf{A})} \frac{\|\mathbf{r}^{(l)}\|}{\|\mathbf{r}^{(0)}\|} \leq \frac{\|\mathbf{x}^{(l)} - \mathbf{x}\|}{\|\mathbf{x}^{(0)} - \mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{r}^{(l)}\|}{\|\mathbf{r}^{(0)}\|}. \quad (3.6.15) \quad \text{cg-term}$$

Relative Abnahme des Iterationsfehlers

Beispiel 129 (Orthogonalitätsverlust der Residuen).

$$\mathbf{R}^H \mathbf{R} = \begin{pmatrix} 1.000000 & -0.000000 & 0.000000 & -0.000000 & 0.000000 & -0.000000 & 0.016019 & -0.795816 & -0.430569 & 0.348133 \\ -0.000000 & 1.000000 & -0.000000 & 0.000000 & -0.000000 & 0.000000 & -0.012075 & 0.600068 & -0.520610 & 0.420903 \\ 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 & -0.000000 & 0.001582 & -0.078664 & 0.384453 & -0.310577 \\ -0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 & -0.000024 & 0.001218 & -0.024115 & 0.019394 \\ 0.000000 & -0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 & -0.000002 & 0.000151 & -0.000118 \\ -0.000000 & 0.000000 & -0.000000 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & 0.000000 & -0.000000 & 0.000000 \\ 0.016019 & -0.012075 & 0.001582 & -0.000024 & 0.000000 & -0.000000 & 1.000000 & -0.000000 & -0.000000 & 0.000000 \\ -0.795816 & 0.600068 & -0.078664 & 0.001218 & -0.000002 & 0.000000 & -0.000000 & 1.000000 & 0.000000 & -0.000000 \\ -0.430569 & -0.520610 & 0.384453 & -0.024115 & 0.000151 & -0.000000 & -0.000000 & 0.000000 & 1.000000 & 0.000000 \\ 0.348133 & 0.420903 & -0.310577 & 0.019394 & -0.000118 & 0.000000 & 0.000000 & -0.000000 & 0.000000 & 1.000000 \end{pmatrix}$$

> Orthogonalitätsverlust der Residuen als Folge von Rundungsfehlern, vgl. Lanczos-Prozess, Bsp. 115.

3.6.1.3 Konvergenzgeschwindigkeit

KONVERGENZGESCHW. Bsp. 129 > n CG-Schritte untauglich als „direkter Löser“.

3.6 p. 301 Hoffnung: Substantielle Fehlerreduktion bereits in wenigen $l \ll n$ Schritten

Frage (für grosse n): Wie gut ist die Näherungslösung $\mathbf{x}^{(l)}$, $l \ll n$?

Lemma 3.6.5 („Optimalität“ der CG-Iterierten). Mit der exakten Lösung $\mathbf{x} \in \mathbb{K}^n$ von $\mathbf{A}\mathbf{x} = \mathbf{b}$ gilt für die CG-Iterierten $\mathbf{x}^{(l)}$

$$\|\mathbf{x} - \mathbf{x}^{(l)}\|_A = \min\{\|\mathbf{y} - \mathbf{x}\|_A : \mathbf{y} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r})\} \quad , \quad \mathbf{r} := \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}.$$

► Monotone Abnahme der Energienorm des Fehlers während der CG-Iteration („Gedächtnis“ des CG-Verfahrens)

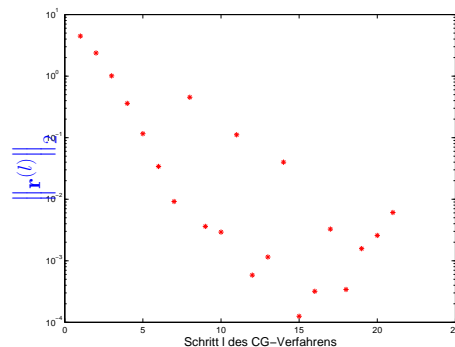
Beispiel 130 (Konvergenz des CG-Verfahrens).

Numerisches Experiment: $\mathbf{A} = \text{hilb}(20)$, $\mathbf{x}^{(0)} = \mathbf{0}$, $\mathbf{b} = (1, \dots, 1)^T$

Hilbert-Matrix: extrem schlecht konditioniert

Residuen während der CG-Iteration:

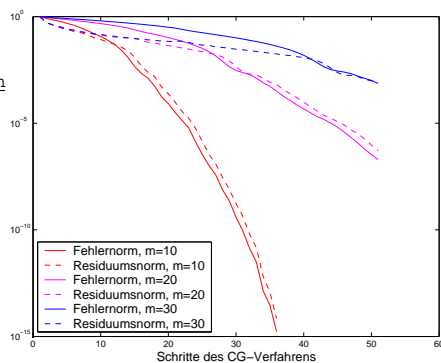
$$\mathbf{R} = [\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(10)}]$$



S.p.d. Matrix $A \in \mathbb{R}^{m^2, m^2}$, \rightarrow Bsp. 112 lex:gs:poisson

```
A = gallery('poisson',m);
x0 = (1:n)';
b = zeros(n,1);
```

Relative Abnahme von $\|r^{(l)}\|_2, \|x^{(l)}\|_2$



Quantitative Konvergenzabschätzung aus Lemma 3.6.5 lem:cgpt

$$y \in x^{(0)} + \mathcal{K}_l(A, r) \Leftrightarrow y = x^{(0)} + A p(A)(x - x^{(0)}), \quad p = \text{Polynom vom Grad } \leq l-1.$$

$$\blacktriangleright \quad x - y = q(A)(x - x^{(0)}) \quad \text{mit } q = \text{Polynom vom Grad } \leq l, \quad q(0) = 1.$$

$$\|x - x^{(l)}\|_A \leq \min \left\{ \max_{\lambda \in \sigma(A)} |q(\lambda)| : q \text{ Polynom vom Grad } \leq l, \quad q(0) = 1 \right\} \cdot \|x - x^{(0)}\|_A.$$

Abzuschätzen für $\lambda \in [\lambda_{\min}(A), \lambda_{\max}(A)]$ durch „Kandidatenpolynome“

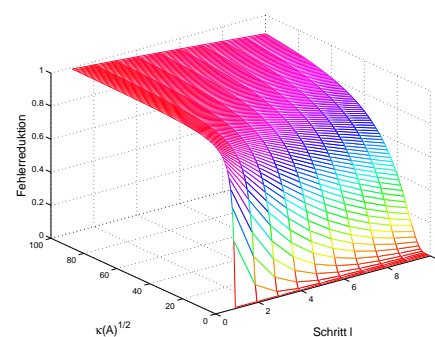
Mit Tschebyscheff-Polynomen (\rightarrow Abschnitt 4.1.3) zeigt man [21, Satz 9.4.2]: sec:tschbysch-Hat94p

Theorem 3.6.6 (Konvergenz des CG-Verfahrens). Für die Iterierten des CG-Verfahrens zur Lösung von $Ax = b$, $A = A^H$ positiv definit, gilt

$$\|x - x^{(l)}\|_A \leq \frac{2 \left(1 - \frac{1}{\sqrt{\kappa(A)}}\right)^l}{\left(1 + \frac{1}{\sqrt{\kappa(A)}}\right)^{2l} + \left(1 - \frac{1}{\sqrt{\kappa(A)}}\right)^{2l}} \|x - x^{(0)}\|_A.$$

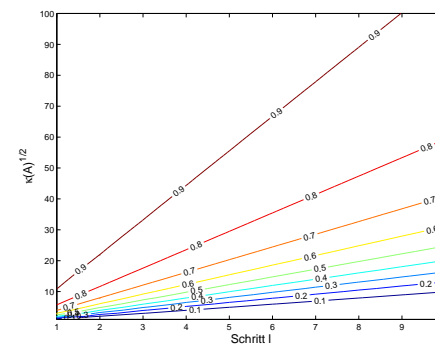
$(\kappa(A) = \text{spektrale Kondition} \rightarrow \text{Def. 3.4.5})$ von A , $\kappa(A) = \text{cond}_2(A)$ def:speccond

Theoretische Schranken aus Thm. 3.6.6 für Fehlerreduktion (in Energienorm) nach l CG-Schritten: thm:CGconv



Beispiel 131 (Konvergenzraten für CG-Verfahren).

```
A = gallery('poisson',m); n = size(A,1);
x0 = (1:n)'; b = ones(n,1); maxit = 30; tol = 0;
[x,flag,relres,iter,resvec] = pcg(A,b,tol,maxit,[],[],x0);
```

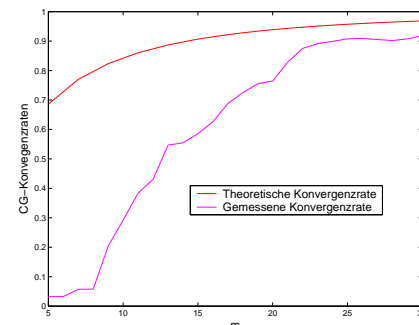
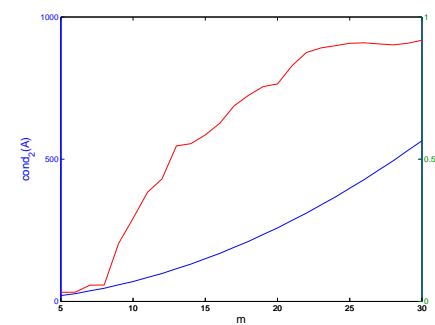


3.6
p. 305

Messen der Konvergenzrate:

$$\text{rate} = \sqrt[10]{\frac{\|r^{(30)}\|_2}{\|r^{(20)}\|_2}}.$$

3.6
p. 307



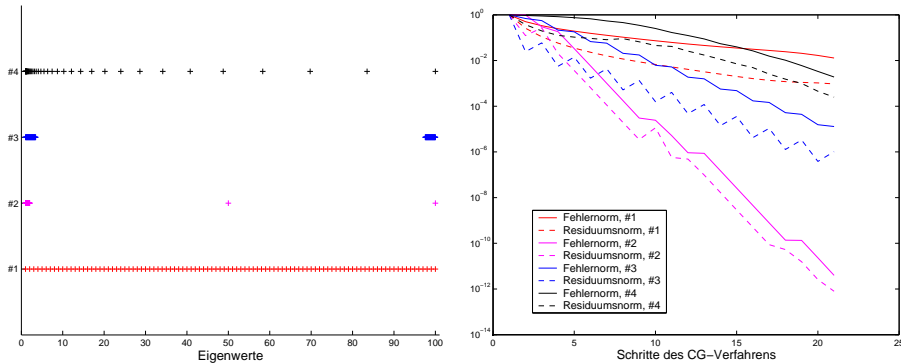
Beispiel 132 (CG-Konvergenz und Spektrum).

Testmatrix #1: $A = \text{diag}(d)$; $d = (1:100)$;
 Testmatrix #2: $A = \text{diag}(d)$; $d = [1 + (0:97)/97, 50, 100]$;
 Testmatrix #3: $A = \text{diag}(d)$; $d = [1 + (0:49)*0.05, 100 - (0:49)*0.05]$;
 Testmatrix #4: Spektrum exponentiell gehäuft um 1

```
x0 = cos((1:n)'); b = zeros(n,1);
```

3.6
p. 306

3.6
p. 308



- Verteilung der Eigenwerte hat entscheidenden Einfluss auf CG-Konvergenzgeschwindigkeit:
Wichtig: Schnellere Konvergenz bei gehäuften Eigenwerten

3.6.2 Vorkonditionierung

VORKOND Thm. 3.6.6 ^{thm:CGconv} CG konvergiert (unter Umständen) sehr langsam für $\kappa(\mathbf{A}) \gg 1$.
Idee: **Vorkonditionierung** (→ Abschnitt 3.4.2) ^{sec:vork-inverse-iter}



Wende CG an auf das **transformierte Gleichungssystem**

$$\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}, \quad \tilde{\mathbf{A}} := \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}, \quad \tilde{\mathbf{x}} := \mathbf{B}^{1/2}\mathbf{x}, \quad \tilde{\mathbf{b}} := \mathbf{B}^{-1/2}\mathbf{b}, \quad (3.6.16) \text{ cg:trf}$$

mit $\kappa(\tilde{\mathbf{A}}) = \text{„klein“}$, $\mathbf{B} = \mathbf{B}^H \in \mathbb{K}^{N,N}$ positiv definit **Vorkonditionierer** → Def. 3.4.4. ^{def:pc}

„Black-box“-Vorkonditionierer aus Abschnitt 3.4.2: Jacobi-, Gauss-Seidel-, IC-Vorkonditionierer. ^{sec:vork-inverse-iter}

Bemerkung 133 (Vorkonditionierung und Kontraktionszahl).

Für s.p.d. Vorkonditionierer \mathbf{B} von s.p.d. \mathbf{A} heisst $\rho := \|\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}\|_{\mathbf{A}}$ die **Kontraktionszahl** von \mathbf{B} bzgl. \mathbf{A} , siehe (3.4.6), Lemma 77. ^{binvit:rate lem:pcit}

$$\rho := \lambda_{\max}(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}) = \lambda_{\max}(\mathbf{I} - \mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}) < 1 \Rightarrow \begin{cases} \lambda_{\min}(\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}) \geq 1 - \rho, \\ \lambda_{\max}(\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}) \leq 1 + \rho, \end{cases}$$

$$\Rightarrow \kappa(\mathbf{B}^{1/2}\mathbf{A}\mathbf{B}^{1/2}) \leq \frac{1 + \rho}{1 - \rho}.$$

Implementierung: CG-Iteration in untransformierten Variablen \mathbf{x} (intrinsische Transformation)

Algorithmus: **VORKCG** **Vorkonditioniertes CG-Verfahren (PCG)**

Input: Anfangsnäherung $\mathbf{x} \in \mathbb{K}^n \hat{=} \mathbf{x}^{(0)} \in \mathbb{K}^n$, Abbruchschranke $\tau > 0$
Output: Bessere Näherung $\mathbf{x} \hat{=} \mathbf{x}^{(l)}$

```

p := r := b - Ax;  p := B^-1 r; q := p; tau_0 := p^H r;
for l = 1 to l_max do {
    beta := r^H q;  h := Ap;  alpha := beta / p^H h;
    x := x + alpha p;
    r := r - alpha h;
    q := B^-1 r;  beta := r^H q;
    if |q^H r| <= tau * tau_0 then stop;
    p := q + beta p;
}

```

(3.6.17) ^{alg:PCG}

➤ Rechenaufwand pro Schritt: 1 Auswertung $\mathbf{A} \times \text{Vektor}$, 1 Auswertung $\mathbf{B}^{-1} \times \text{Vektor}$, 3 Skalarprodukte, 3 **SAXPY-Operationen**

3.6 **Bemerkung 134** (Konvergenztheorie für PCG). Aussagen von Thm. 3.6.6 bleiben gültig mit $\tilde{\mathbf{A}}$ -Energienorm anstelle \mathbf{A} -Energienorm. ^{thm:CGconv} 3.6 p. 311

Bemerkung 135 (Abbruch des vorkonditionierten CG-Verfahrens).

Bem. 128, (3.6.15) ^{rem:det:term} ➤ Überwache transformiertes Residuum

$$\tilde{\mathbf{r}} = \tilde{\mathbf{b}} - \tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{B}^{-1/2}\mathbf{r} \Rightarrow \|\tilde{\mathbf{r}}\|_2^2 = \mathbf{r}^H \mathbf{B}^{-1} \mathbf{r}.$$

➤ Abschätzungen für Energienorm des Fehlers $\mathbf{e}^{(l)} := \mathbf{x} - \mathbf{x}^{(l)}$:

Benutze die Fehlergleichung (3.6.1) ^{eq:cfix} $\mathbf{A}\mathbf{e} = \mathbf{r}$:

$$\mathbf{r}^{(l)} \cdot \mathbf{B}^{-1} \mathbf{r}^{(l)} = \mathbf{B}^{-1} \mathbf{A} \mathbf{e}^{(l)} \cdot \mathbf{A} \mathbf{e}^{(l)} \leq \lambda_{\max}(\mathbf{B}^{-1} \mathbf{A}) \|\mathbf{e}^{(l)}\|_{\mathbf{A}}^2,$$

$$\|\mathbf{e}^{(l)}\|_{\mathbf{A}}^2 = \mathbf{A} \mathbf{e}^{(l)} \cdot \mathbf{e}^{(l)} = \mathbf{r}^{(l)} \cdot \mathbf{A}^{-1} \mathbf{r}^{(l)} = \mathbf{B}^{-1} \mathbf{r} \cdot \mathbf{B} \mathbf{A}^{-1} \mathbf{r}^{(l)} \leq \lambda_{\max}(\mathbf{B} \mathbf{A}^{-1}) \mathbf{B}^{-1} \mathbf{r}^{(l)} \cdot \mathbf{r}^{(l)}.$$

verfügbar in Algorithmus (3.6.17) ^{alg:PCG}

$$\frac{1}{\kappa(\mathbf{B}^{-1} \mathbf{A})} \frac{\|\mathbf{e}^{(l)}\|_{\mathbf{A}}^2}{\|\mathbf{e}^{(0)}\|_{\mathbf{A}}^2} \leq \frac{\mathbf{B}^{-1} \mathbf{r}^{(l)} \cdot \mathbf{r}^{(l)}}{\mathbf{B}^{-1} \mathbf{r}^{(0)} \cdot \mathbf{r}^{(0)}} \leq \kappa(\mathbf{B}^{-1} \mathbf{A}) \frac{\|\mathbf{e}^{(l)}\|_{\mathbf{A}}^2}{\|\mathbf{e}^{(0)}\|_{\mathbf{A}}^2} \quad (3.6.18) \text{ eq:pcgie}$$

3.6 $\kappa(\mathbf{B}^{-1} \mathbf{A})$ „klein“ ➤ \mathbf{B}^{-1} -Energienorm des Residuums $\approx \mathbf{A}$ -Energienorm des Fehlers!
 $(\mathbf{r}^{(l)} \cdot \mathbf{B}^{-1} \mathbf{r}^{(l)}) = \mathbf{q}^H \mathbf{r}$ in Algorithmus (3.6.17) ^{alg:PCG} 3.6 p. 312

MATLAB-Funktion: `[x,flag,relr,it,rv] = pcg(A,b,tol,maxit,B,[],x0);`
 (A, B können Handles auf Funktionen sein, die \mathbf{Ax} bzw. $\mathbf{B}^{-1}\mathbf{x}$ realisieren)

Interne MATLAB-Implementierung:

MATLAB-CODE: PCG-Algorithmus

```
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 - feval(Afun,x)) <= tol*b*norm(b)), return; end
    r = r - alpha * q;
end
```

Zweifelhaftes Abbruchkriterium !

3.6.3 Weitere Krylov-Unterraumverfahren

3.6.3.1 Residuenminimierende Verfahren

Idee: Ersetze Euklidisches Skalarprodukt \rightarrow \mathbf{A} -Skalarprodukt im CG-Algorithmus

$$\|\mathbf{x}^{(l)} - \mathbf{x}\|_{\mathbf{A}} \rightarrow \|\mathbf{A}(\mathbf{x}^{(l)} - \mathbf{x})\|_2 = \|\mathbf{r}^{(l)}\|_2$$



MINRES-Verfahren (für beliebige Hermitesche Matrizen !)

Theorem 3.6.7. Für $\mathbf{A} = \mathbf{A}^H$ erfüllen die Residuen $\mathbf{r}^{(l)}$ in der MINRES-Iteration

$$\|\mathbf{r}^{(l)}\|_2 = \min\{\|\mathbf{Ay} - \mathbf{b}\|_2 : \mathbf{y} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}^{(0)})\}$$

$$\|\mathbf{r}^{(l)}\|_2 \leq \frac{2 \left(1 - \frac{1}{\kappa(\mathbf{A})}\right)^l}{\left(1 + \frac{1}{\kappa(\mathbf{A})}\right)^{2l} + \left(1 - \frac{1}{\kappa(\mathbf{A})}\right)^{2l}} \|\mathbf{x} - \mathbf{x}^{(0)}\|_{\mathbf{A}}.$$

3.6 p. 313 Beachte: Konvergenzrate bestimmt durch $\kappa(\mathbf{A})$, nicht durch $\sqrt{\kappa(\mathbf{A})}$ wie bei CG !

MATLAB-Funktion: `• [x,flg,res,it,resv] = minres(A,b,tol,maxit,B,[],x0);`
`• [...] = minres(Afun,b,tol,maxit,Binvfun,[],x0);`

Rechenaufwand : 1 $\mathbf{A} \times$ Vektor, 1 $\mathbf{B}^{-1} \times$ Vektor pro Schritt, wenige Skalarprodukte & SAXPYs
 Speicheraufwand: Wenige Vektoren $\in \mathbb{K}^n$

Idee: Löse überbestimmtes Gleichungssystem



$$\mathbf{x}^{(l)} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}^{(0)}): \mathbf{Ax}^{(l)} = \mathbf{b}$$

im Sinn der Methode der kleinsten Quadrate, \rightarrow Abschnitt 3.5 sec: numer-line-ausgl

$$\mathbf{x}^{(l)} = \operatorname{argmin}\{\|\mathbf{Ay} - \mathbf{b}\|_2 : \mathbf{y} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}^{(0)})\}.$$



GMRES-Verfahren für allgemeine Matrizen $\mathbf{A} \in \mathbb{K}^{n,n}$

MATLAB-Funktion: `• [x,flag,relr,it,rv] = gmres(A,b,rs,tol,maxit,B,[],x0);`
`• [...] = gmres(Afun,b,rs,tol,maxit,Binvfun,[],x0);`

Rechenaufwand : 1 $\mathbf{A} \times$ Vektor, 1 $\mathbf{B}^{-1} \times$ Vektor pro Schritt,
 $O(l)$ Skalarprodukte & SAXPYs im l . Schritt
 Speicheraufwand: $O(l)$ Vektoren $\in \mathbb{K}^n$ im l . Schritt



3.6 p. 314

3.6 p. 315

3.6 p. 316

Bemerkung 136 (GMRES und Arnoldi-Prozess). GMRES \leftrightarrow Arnoldi-Prozess, \rightarrow Abschnitt 3.4.3 (vgl. Bem. 127) \rightarrow Lange Rekursionen \triangleright wachsender Speicher- und Rechenaufwand

Bemerkung 137 (Neustart von GMRES). Wenn Speicher- und Rechenaufwand zu gross wird: Starte GMRES-Iteration neu mit bisher erhaltener Näherung als Startwert \rightarrow rs-Parameter = Neustart nach jeweils rs Schritten (Vorsicht: Gefahr der Stagnation !)

3.6.3.2 Verfahren mit kurzen Rekursionen

Idee: Gegeben $\mathbf{x}^{(0)} \in \mathbb{K}^n$ bestimme verbesserte Näherung $\mathbf{x}^{(l)}$ aus **Petrov-Galerkin-Bedingung**

$$\mathbf{x}^{(l)} \in \mathbf{x}^{(0)} + \mathcal{K}_l(\mathbf{A}, \mathbf{r}^{(0)}): \mathbf{p}^H(\mathbf{b} - \mathbf{A}\mathbf{x}^{(l)}) = 0 \quad \forall \mathbf{p} \in W_l,$$

mit geeignetem **Testraum** W_l , $\dim W_l = l$, z.B. $W_l := \mathcal{K}_l(\mathbf{A}^H, \mathbf{r}^{(0)})$ (\rightarrow bi-conjugate gradients, BiCG)

► „Zoo“ von Verfahren mit kurzen Rekursionen (\rightarrow Lanczos-Prozess) für allgemeine Matrizen $\mathbf{A} \in \mathbb{K}^{n,n}$, siehe [3].

MATLAB-Funktion: `[x,flag,r,it,rv] = bicgstab(A,b,tol,maxit,B,[],x0)`
`[...] = bicgstab(Afun,b,tol,maxit,Binvfun,[],x0);`

Rechenaufwand : $2 \mathbf{A} \times \text{Vektor}$, $2 \mathbf{B}^{-1} \times \text{Vektor}$, 4 Skalarprodukte, 6 SAXPYs pro Schritt
 Speicheraufwand: 8 Vektoren $\in \mathbb{K}^n$

MATLAB-Funktion: `[x,flag,r,it,rv] = qmr(A,b,tol,maxit,B,[],x0)`
`[...] = qmr(Afun,b,tol,maxit,Binvfun,[],x0);`

Rechenaufwand : $2 \mathbf{A} \times \text{Vektor}$, $2 \mathbf{B}^{-1} \times \text{Vektor}$, 2 Skalarprodukte, 12 SAXPYs pro Schritt
 Speicheraufwand: 10 Vektoren $\in \mathbb{K}^n$

Keine Garantie für Konvergenz bei diesen Verfahren:
 Stagnation & Breakdowns „an der Tagesordnung“

Beispiel 138 (Versagen von Krylov-Raum basierten iterativen Lösern).

$$\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 & \cdots & \cdots & 0 \\ 0 & 0 & 1 & 0 & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 & 1 \\ 0 & & & \ddots & 0 & 1 \\ 1 & 0 & \cdots & \cdots & 0 & 0 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 0 \\ \vdots \\ \vdots \\ 0 \\ 1 \end{pmatrix} \quad \Rightarrow \quad \mathbf{x} = \mathbf{e}_1.$$

$$\mathbf{x}^{(0)} = 0 \quad \Rightarrow \quad \mathbf{r}^{(0)} = \mathbf{e}_n \quad \Rightarrow \quad \mathcal{K}_l(\mathbf{A}, \mathbf{r}^{(0)}) = \text{Span}\{\mathbf{e}_n, \mathbf{e}_{n-1}, \dots, \mathbf{e}_{n-l+1}\}$$

$$\Rightarrow \min\{\|\mathbf{y} - \mathbf{x}\|_2 : \mathbf{y} \in \mathcal{K}_l(\mathbf{A}, \mathbf{r}^{(0)})\} = \begin{cases} 1, & \text{falls } l \leq n, \\ 0, & \text{für } l = n. \end{cases}$$

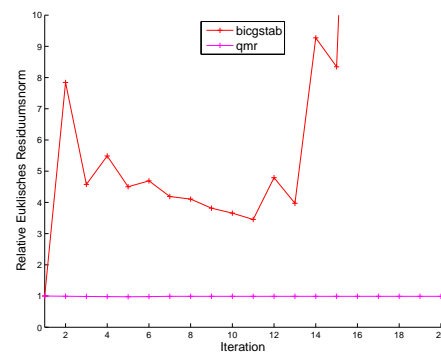
Was soll der Anwender tun? Ausprobieren („try and pray“)

Beispiel 139 (Konvergenz von Krylovverfahren für nichtsymmetrische Matrix).

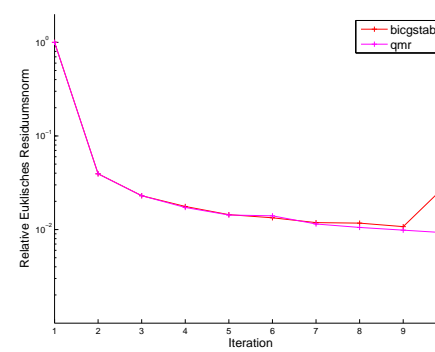
```
A = gallery('tridiag', -0.5*ones(n-1,1), 2*ones(n,1), -1.5*ones(n-1,1));
B = gallery('tridiag', 0.5*ones(n-1,1), 2*ones(n,1), 1.5*ones(n-1,1));
```

Plots zeigen $\|\mathbf{r}^{(l)}\|_2 : \|\mathbf{r}^{(0)}\|_2$:

3.6
p. 317



Tridiagonalmatrix **A**



Tridiagonalmatrix **B**

3.6
p. 319

3.7 Spezielle Matrizen

[File: section-specialmat.tex, SVN: section-specialmat.tex 1056 2006-10-18 15:52:41Z kalai]

3.6
p. 318

SPEZMAT Spezielle Algorithmen für $n \times n$ -Matrizen mit $O(n)$ Information (Bsp. \rightarrow Bandmatrizen)

3.7
p. 320

3.7.1 Diskrete Fouriertransformationen

DISKRETFOURIERTRAFO

$n \in \mathbb{N}$ fest \rightarrow n . Einheitswurzel $\omega_n := \exp(-2\pi i/n) = \cos(2\pi/n) - i \sin(2\pi/n)$

$\blacktriangleright \omega_n^k = \omega_n^{k+n} \forall k \in \mathbb{Z}, \omega_n^n = 1, \omega_n^{n/2} = -1,$ (3.7.1) [eq:comp1](#)

$\sum_{k=0}^{n-1} \omega_n^{kj} = \begin{cases} n, & \text{falls } j = 0, \\ 0, & \text{falls } j \neq 0. \end{cases}$ (3.7.2) [eq:comp1](#)

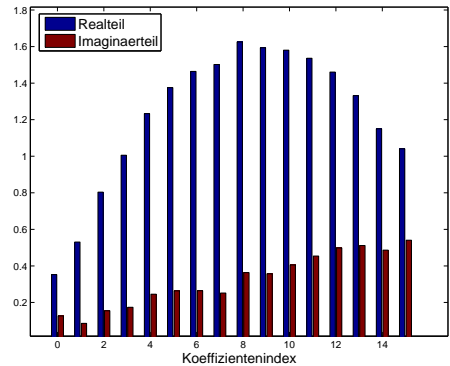
Eine Basistransformation in \mathbb{C}^n :

Standardbasis des \mathbb{K}^n TRIGBASIS
Trigonometrische Basis

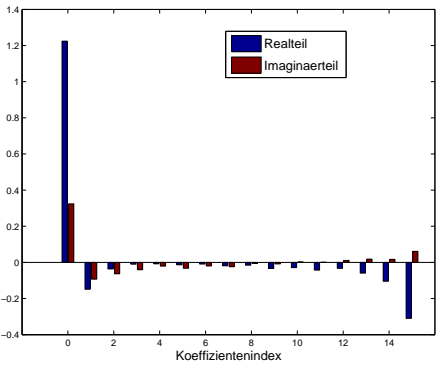
$$\left\{ \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \right\} \leftarrow \left\{ \begin{pmatrix} \omega_n^0 \\ \vdots \\ \omega_n^0 \end{pmatrix}, \begin{pmatrix} \omega_n^0 \\ \omega_n^1 \\ \vdots \\ \omega_n^{n-1} \end{pmatrix}, \dots, \begin{pmatrix} \omega_n^0 \\ \omega_n^{n-2} \\ \vdots \\ \omega_n^{(n-1)(n-2)} \end{pmatrix}, \begin{pmatrix} \omega_n^0 \\ \omega_n^{n-1} \\ \vdots \\ \omega_n^{(n-1)^2} \end{pmatrix} \right\}$$

Matrix der Basistransformation Trigonometrische Basis \rightarrow Standardbasis: Fourier-Matrix

$$F_n = \begin{pmatrix} \omega_n^0 & \omega_n^0 & \dots & \omega_n^0 \\ \omega_n^0 & \omega_n^1 & \dots & \omega_n^{n-1} \\ \omega_n^0 & \omega_n^2 & \dots & \omega_n^{2(n-2)} \\ \vdots & \vdots & \ddots & \vdots \\ \omega_n^0 & \omega_n^{n-1} & \dots & \omega_n^{(n-1)^2} \end{pmatrix} = (\omega_n^{ij})_{i,j=0}^{n-1} \in \mathbb{C}^{n,n}. \quad (3.7.3) \quad \text{[diff:FM](#)}$$



Vektor in Standardbasis



Vektor in trigonometrischer Basis

Definition 3.7.1 (Diskrete Fouriertransformation). Die lineare Abbildung $\mathcal{F}_n : \mathbb{C}^n \mapsto \mathbb{C}^n$, $\mathcal{F}_n(\mathbf{y}) := F_n \mathbf{y}$, $\mathbf{y} \in \mathbb{C}^n$, heisst diskrete Fouriertransformation (DFT), d.h. für $\mathbf{c} := \mathcal{F}_n(\mathbf{y})$

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{kj}, \quad k = 0, \dots, n-1. \quad (3.7.4) \quad \text{[eq:DFT](#)}$$

Konvention: Bei Diskussion der DFT: Vektorindizes laufen von 0 bis $n-1$.

Lemma 3.7.2 (Abbildungseigenschaft der diskreten Fourier-Transformation).

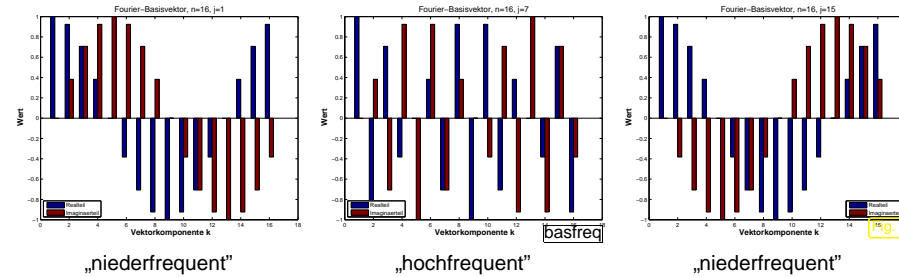
Die skalierte Fourier-Matrix $\frac{1}{\sqrt{n}} F_n$ ist unitär: $F_n^{-1} = \frac{1}{n} F_n^H = \frac{1}{n} \overline{F_n}^T$

MATLAB-Funktionen: $\mathbf{c} = \text{fft}(\mathbf{y}) \leftrightarrow \mathbf{y} = \text{ifft}(\mathbf{c})$;

Beispiel 140 (Frequenzanalyse mit DFT).

Einige Vektoren der Fourierbasis ($n=16$): \rightarrow Fig. 7 [speccircle](#)

3.7
p. 321



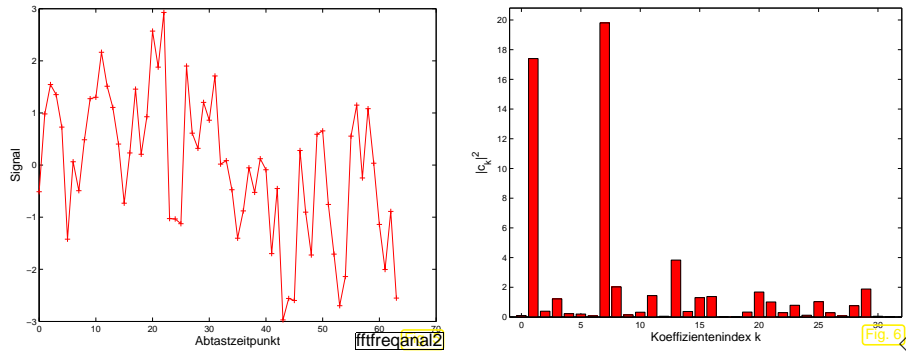
Extraktion charakteristischer Frequenzen aus gestörtem, zeitdiskret und periodischem Signal:

$$t = 0:63; \quad x = \sin(2\pi t/64) + \sin(7 \cdot 2\pi t/64);$$
$$y = x + \text{randn}(\text{size}(t));$$

3.7
p. 322

3.7
p. 323

3.7
p. 324



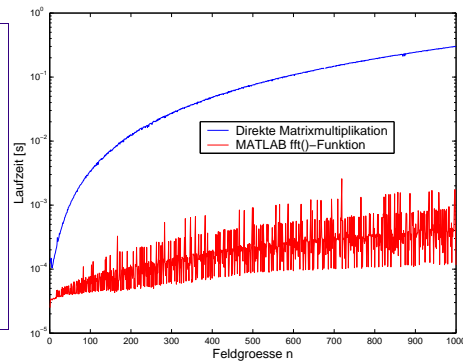
„Naive“ Auswertung von $\mathbf{F}\mathbf{y} \rightarrow O(n^2)$ Rechenoperationen

Beispiel 141 (Effizienz von fft).

tic-toc-Zeitmessung in MATLAB: Vergleich von **fft** und direkter Matrixmultiplikation (MATLAB V6.5, Linux, Mobile Intel Pentium 4 - M CPU 2.40GHz, Minimum über 100 Runs)

MATLAB-CODE Naive DFT-Implementierung

```
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
```



3.7.1.1 Schnelle Fouriertransformation

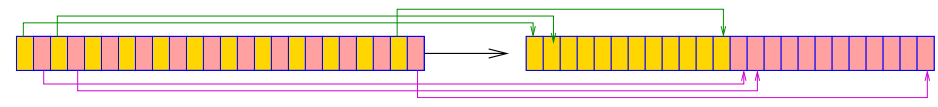
SCHELLFOURFA

„Milleniums-Algorithmus“ **Schnelle Fouriertransformation (FFT)** (C.F. Gauss, siehe [13])

http://orion.math.iastate.edu/burkardt/misc/algorithms_dongarra.html

Für $n = 2m$, $m \in \mathbb{N}$,

Permutation $P_m^{\text{OE}}(1, \dots, n) = (1, 3, \dots, n-1, 2, 4, \dots, n)$.



Wegen $\omega_n^{2j} = \omega_m^j$:

$$\text{Zeilenpermutation } P_m^{\text{OE}} \mathbf{F}_n = \begin{pmatrix} \mathbf{F}_m & \mathbf{F}_m \\ \mathbf{F}_m \begin{pmatrix} \omega_n^0 & \omega_n^1 & \dots & \omega_n^{n/2-1} \end{pmatrix} & \mathbf{F}_m \begin{pmatrix} \omega_n^{n/2} & \omega_n^{n/2+1} & \dots & \omega_n^{n-1} \end{pmatrix} \end{pmatrix} = \begin{pmatrix} \mathbf{F}_m & \\ & \mathbf{F}_m \end{pmatrix} \begin{pmatrix} \mathbf{I} & \mathbf{I} \\ \omega_n^0 & -\omega_n^0 \\ \omega_n^1 & -\omega_n^1 \\ \dots & \dots \\ \omega_n^{n/2-1} & -\omega_n^{n/2-1} \end{pmatrix}$$

Beispiel 142 (Zerlegung der Fourier-Matrix).

$$P_5^{\text{OE}} \mathbf{F}_{10} = \begin{pmatrix} \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 \\ \omega^0 & \omega^2 & \omega^4 & \omega^6 & \omega^8 & \omega^0 & \omega^2 & \omega^4 & \omega^6 & \omega^8 \\ \omega^0 & \omega^4 & \omega^8 & \omega^2 & \omega^6 & \omega^0 & \omega^4 & \omega^8 & \omega^2 & \omega^6 \\ \omega^0 & \omega^6 & \omega^2 & \omega^8 & \omega^4 & \omega^0 & \omega^6 & \omega^2 & \omega^8 & \omega^4 \\ \omega^0 & \omega^8 & \omega^6 & \omega^4 & \omega^2 & \omega^0 & \omega^8 & \omega^6 & \omega^4 & \omega^2 \\ \omega^0 & \omega^1 & \omega^2 & \omega^3 & \omega^4 & \omega^5 & \omega^6 & \omega^7 & \omega^8 & \omega^9 \\ \omega^0 & \omega^3 & \omega^6 & \omega^9 & \omega^2 & \omega^5 & \omega^8 & \omega^1 & \omega^4 & \omega^7 \\ \omega^0 & \omega^5 & \omega^0 & \omega^5 & \omega^0 & \omega^5 & \omega^0 & \omega^5 & \omega^0 & \omega^5 \\ \omega^0 & \omega^7 & \omega^4 & \omega^1 & \omega^8 & \omega^5 & \omega^2 & \omega^9 & \omega^6 & \omega^3 \\ \omega^0 & \omega^9 & \omega^8 & \omega^7 & \omega^6 & \omega^5 & \omega^4 & \omega^3 & \omega^2 & \omega^1 \end{pmatrix}, \quad \omega := \omega_{10}.$$

Rekursive Implementierung für $n = 2^p$, $p \in \mathbb{N}$

„Teile-und-herrsche“ (engl. *divide and conquer*)

Rechenaufwand $O(n \log_2 n)$
MATLAB-fft $\approx 5n \log_2 n$

(Auf Rekursionstiefe $0 \leq q < p$:
 2^q Aufrufe, jeder 2^{p-q} Additionen,
 2^{p-q-1} Multiplikationen)

MATLAB-CODE Rekursiv Implementierung von FFT

```
function c = myfft(y)
n = length(y);
if (n == 1), c = y; return; end
omega = exp(-2*pi*i/n);
z = (omega.^((0:n-1)')).*y;
c = [myfft(y(1:n/2)+y(n/2+1:n));...
     myfft(z(1:n/2)+z(n/2+1:n))];
c = reshape(reshape(c,n/2,2)',n,1);
```

FFT basierend auf allgemeiner Faktorisierung: $n = pq$, $p, q \in \mathbb{N}$ (Cooley-Tukey-Algorithmus)

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{jk} \stackrel{[j=lp+m]}{=} \sum_{m=0}^{p-1} \sum_{l=0}^{q-1} y_{lp+m} e^{-\frac{2\pi i}{pq}(lp+m)k} = \sum_{m=0}^{p-1} \omega_n^{mk} \sum_{l=0}^{q-1} y_{lp+m} \omega_q^{l(k \bmod q)}$$

(3.7.5) [fft:factor](#)

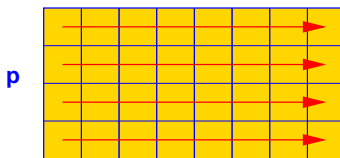
► p DFTs der Länge q berechnen $z_{m,k} := \sum_{l=0}^{q-1} y_{lp+m} \omega_q^{lk}$, $0 \leq m < p$, $0 \leq k < q$.

► für $k = rq + s$, $0 \leq r < p$, $0 \leq s < q$

$$c_{rq+s} = \sum_{m=0}^{p-1} e^{-\frac{2\pi i}{pq}(rq+s)m} z_{m,s} = \sum_{m=0}^{p-1} (\omega_n^{ms} z_{m,s}) \omega_p^{mr}$$

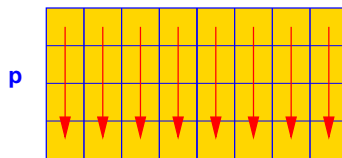
► q DFTs der Länge p berechnen alle c_k .

1. Durchlauf



q

2. Durchlauf



q

Aus dem 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 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 → Bsp. [ffttime](#) 141.

Bemerkung 143 (Zweidimensionale diskrete Fouriertransformation).

Zu $y_{j_1, j_2} \in \mathbb{C}$, $0 \leq j_1 < n_1$, $0 \leq j_2 < n_2$ berechne (geschachtelte DFTs !)

$$c_{k_1, k_2} = \sum_{j_1=0}^{n_1-1} \sum_{j_2=0}^{n_2-1} y_{j_1, j_2} \omega_{n_1}^{j_1 k_1} \omega_{n_2}^{j_2 k_2}, \quad 0 \leq k_1 < n_1, 0 \leq k_2 < n_2.$$

MATLAB-Kommando: `fft2(X)`

△

REPLACEMENT

Bemerkung 144 (Reelle DFT).

Aufgabe: Effiziente Berechnung der DFT (Def. [3.7.1](#)) (c_0, \dots, c_{n-1}) für reelle Koeffizienten $(y_0, \dots, y_{n-1})^T \in \mathbb{R}^n$, $n = 2m$, $m \in \mathbb{N}$

Redundanz:

$$\omega_n^{(n-k)j} = \overline{\omega_n^{kj}}, \quad k = 0, \dots, n-1 \quad \blacktriangleright \quad c_{n-k} = \overline{c_k}, \quad k = 1, \dots, n-1.$$

3.7
p. 329

3.7
p. 331

► Idee: Komplexifizierung & Rückführung auf DFT der Länge m .

$$h_k = \sum_{j=0}^{m-1} (y_{2j} + i y_{2j+1}) \omega_m^{jk} = \sum_{j=0}^{m-1} y_{2j} \omega_m^{jk} + i \cdot \sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}, \quad (3.7.6) \quad \text{fft1}$$

$$\overline{h_{m-k}} = \sum_{j=0}^{m-1} \overline{y_{2j} + i y_{2j+1}} \overline{\omega_m^{j(m-k)}} = \sum_{j=0}^{m-1} y_{2j} \omega_m^{jk} - i \cdot \sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}. \quad (3.7.7) \quad \text{fft2}$$

Spezialfall von (3.7.5) ($p=2$):

$$c_k = \sum_{j=0}^{n-1} y_j \omega_n^{jk} = \sum_{j=0}^{m-1} y_{2j} \omega_m^{jk} + \omega_n^k \cdot \sum_{j=0}^{m-1} y_{2j+1} \omega_m^{jk}. \quad (3.7.8) \quad \text{fft3}$$

$$\blacktriangleright \quad \begin{aligned} c_k &= \frac{1}{2}(h_k + \overline{h_{m-k}}) - \frac{1}{2}i \omega_n^k (h_k - \overline{h_{m-k}}), \quad k = 0, \dots, m-1, \\ c_m &= \text{Re}\{h_0\} - \text{Im}\{h_0\}, \\ c_k &= \overline{c_{n-k}}, \quad k = m+1, \dots, n-1. \end{aligned} \quad (3.7.9) \quad \text{fft4}$$

3.7
p. 330

3.7
p. 332

MATLAB-Implementierung
(auf der Grundlage
einer fft
der Länge $n/2$):

```
function c = fftreal(y)
n = length(y); m = n/2;
if (mod(n,2) ~= 0), error('n odd!'); end
y = y(1:2:n)+i*y(2:2:n);
h = fft(y); h = [h;h(1)];
c = 0.5*(h+conj(h(m+1:-1:1))) - ...
(0.5*i*exp(-2*pi*i/n).^((0:m)')).*...
(h-conj(h(m+1:-1:1)));
c = [c;conj(c(m:-1:2))];
```

Ineffizient!

FREQFILTDISKRSIG

Beispiel 145 (Frequenzfilterung diskreter Signale).

Frequenzfilterung diskreter periodischer Signale
durch Ausblenden von „Fourierkoeffizienten“

```
m = length(y)/2;
c = fft(y);
clow = c;
clow(m+1-k:m+1+k) = 0;
chigh = c-clow;
low = real(ifft(clow));
high = real(ifft(chigh));
```

(Optimierungspotential wegen $y_j \in \mathbb{R}$ und
 $c_{n/2-k} = \bar{c}_{n/2+k}$)

speccircle

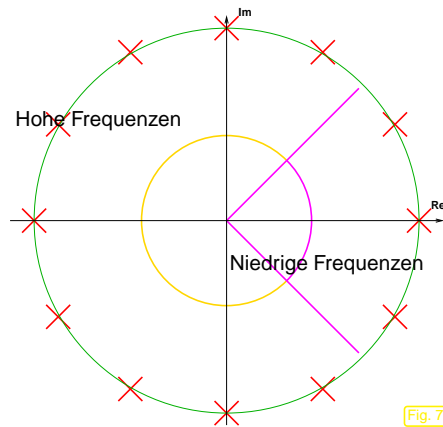
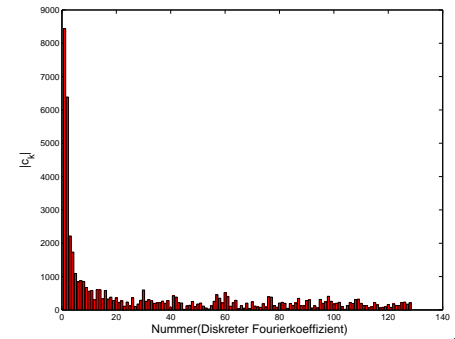
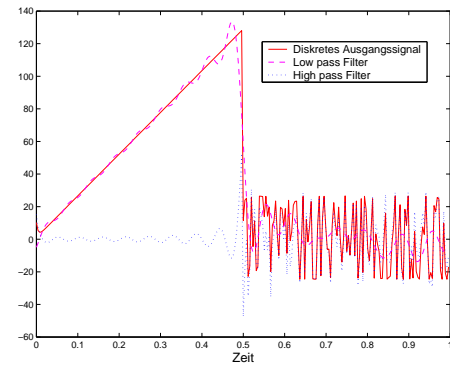


Fig. 7

Konkretes Signal:

```
k = 110; n = 256; m = n/2;
y = [(1:m)'; ones(m,1)] + n/10*sin(exp([- (1:m)'; (1:m)']));
```



Übung 3.8. Implement an efficient routine $x=\text{solvetrig}(m,y)$ for the solution of a linear system

3.7
p. 333

3.8
p. 335

of equations $\mathbf{M}\mathbf{x} = \mathbf{y}$ for $\mathbf{y} \in \mathbb{C}^p$, $p = 2m + 1$, $m \in \mathbb{N}$, and

$$\mathbf{M} := \begin{pmatrix} 1 & 1 & 1 & \cdots & \cdots & 1 \\ \cos(\varphi_1) & \cos(\varphi_2) & \cos(\varphi_3) & \cdots & \cdots & \cos(\varphi_p) \\ \sin(\varphi_1) & \sin(\varphi_2) & \sin(\varphi_3) & \cdots & \cdots & \sin(\varphi_p) \\ \cos(2\varphi_1) & \cos(2\varphi_2) & \cos(2\varphi_3) & \cdots & \cdots & \cos(2\varphi_p) \\ \sin(2\varphi_1) & \sin(2\varphi_2) & \sin(2\varphi_3) & \cdots & \cdots & \sin(2\varphi_p) \\ \vdots & \vdots & \vdots & & & \vdots \\ \cos(m\varphi_1) & \cos(m\varphi_2) & \cos(m\varphi_3) & \cdots & \cdots & \cos(m\varphi_p) \\ \sin(m\varphi_1) & \sin(m\varphi_2) & \sin(m\varphi_3) & \cdots & \cdots & \sin(m\varphi_p) \end{pmatrix},$$

where $\varphi_j = \frac{2\pi}{p}(j-1)$, $j = 1, \dots, p$.

Hint: First compute $\mathbf{M}\mathbf{M}^T$ and establish that this yields a diagonal matrix.

3.8.0.2 Sinustransformation

SINTRAF0

Eine weitere trigonometrische Basistransformation in \mathbb{R}^{n-1} , $n \in \mathbb{N}$:

3.7
p. 334

3.8
p. 336

Standardbasis des \mathbb{R}^{n-1}

$$\left\{ \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \right\} \leftarrow \left\{ \begin{pmatrix} \sin(\frac{\pi}{n}) \\ \sin(\frac{2\pi}{n}) \\ \vdots \\ \sin(\frac{(n-1)\pi}{n}) \end{pmatrix}, \begin{pmatrix} \sin(\frac{2\pi}{n}) \\ \sin(\frac{4\pi}{n}) \\ \vdots \\ \sin(\frac{2(n-1)\pi}{n}) \end{pmatrix}, \dots, \begin{pmatrix} \sin(\frac{(n-1)\pi}{n}) \\ \sin(\frac{2(n-1)\pi}{n}) \\ \vdots \\ \sin(\frac{(n-1)^2\pi}{n}) \end{pmatrix} \right\}$$

„Sinusbasis“

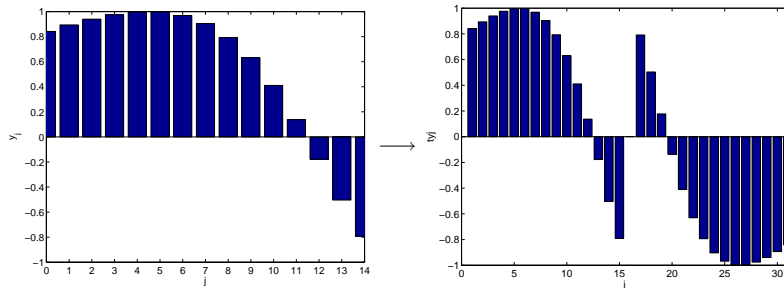
Matrix(Basistransformation Sinusbasis \rightarrow Standardbasis): $\mathbf{S}_n := (\sin(jk\pi/n))_{j,k=1}^{n-1} \in \mathbb{R}^{n-1, n-1}$

Lemma 3.8.1 (Eigenschaften der Sinusmatrix). $\sqrt{2/n} \mathbf{S}_n$ ist reell, symmetrisch und unitär.

Sinustransformation: $s_k = \sum_{j=1}^{n-1} y_j \sin(\pi j k / n)$, $k = 1, \dots, n-1$. (3.8.1)

DFT-basierter Algorithmus zur Auswertung der Sinustransformation ($\hat{=}$ $\mathbf{S}_n \times$ Vektor):

„Wrap around“: $\tilde{\mathbf{y}} \in \mathbb{R}^{2n}$: $\tilde{y}_j = \begin{cases} y_j & , \text{ falls } j = 1, \dots, n-1, \\ 0 & , \text{ falls } j = 0, n, \\ -y_{2n-j} & , \text{ falls } j = n+1, \dots, 2n-1. \end{cases}$ ($\tilde{\mathbf{y}}$ „punktsymmetrisch“)



$$\begin{aligned} (\mathbf{F}_{2n} \tilde{\mathbf{y}})_k &\stackrel{\text{Def. DFT (3.7.4)}}{=} \sum_{j=1}^{2n-1} \tilde{y}_j e^{-\frac{2\pi i}{2n} j k} \\ &= \sum_{j=1}^{n-1} y_j e^{-\frac{\pi i}{n} j k} - \sum_{j=n+1}^{2n-1} y_{2n-j} e^{-\frac{\pi i}{n} j k} \\ &= \sum_{j=1}^{n-1} y_j (e^{-\frac{\pi i}{n} j k} - e^{\frac{\pi i}{n} j k}) \\ &= -2i (\mathbf{S}_n \mathbf{y})_k, \quad k = 1, \dots, n-1. \end{aligned}$$

Wrap-around Implementierung

```
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 Sinustransformation

Bemerkung 146 (Sinustransformation per DFT halber Länge).

Schritt ①: Transformation der Koeffizienten

$$\tilde{y}_j = \sin(j\pi/n)(y_j + y_{n-j}) + \frac{1}{2}(y_j - y_{n-j}), \quad j = 1, \dots, n-1, \quad \tilde{y}_0 = 0.$$

3.8
p. 337

Schritt ②: Reelle DFT (\rightarrow Def. 3.7.1) von $(\tilde{y}_0, \dots, \tilde{y}_{n-1}) \in \mathbb{R}^n$: $c_k := \sum_{j=0}^{n-1} \tilde{y}_j e^{-\frac{2\pi i}{n} j k}$

$$\begin{aligned} \text{Re}\{c_k\} &= \sum_{j=0}^{n-1} \tilde{y}_j \cos(-\frac{2\pi i}{n} j k) = \sum_{j=1}^{n-1} (y_j + y_{n-j}) \sin(\frac{\pi j}{n}) \cos(\frac{2\pi j k}{n}) \\ &= \sum_{j=0}^{n-1} 2y_j \sin(\frac{\pi j}{n}) \cos(\frac{2\pi j k}{n}) = \sum_{j=0}^{n-1} y_j \left(\sin(\frac{2k+1}{n}\pi j) - \sin(\frac{2k-1}{n}\pi j) \right) \\ &= s_{2k+1} - s_{2k-1}. \\ \text{Im}\{c_k\} &= \sum_{j=0}^{n-1} \tilde{y}_j \sin(-\frac{2\pi i}{n} j k) = -\sum_{j=1}^{n-1} \frac{1}{2}(y_j - y_{n-j}) \sin(\frac{2\pi j k}{n}) = -\sum_{j=1}^{n-1} y_j \sin(\frac{2\pi j k}{n}) \\ &= -s_{2k}. \end{aligned}$$

Schritt ③: Extraktion der s_k

$$s_{2k+1}, \quad k = 0, \dots, \frac{n}{2} - 1 \quad \blacktriangleright \quad \text{Aus Rekursion } s_{2k+1} - s_{2k-1} = \text{Re}\{c_k\}, \quad s_1 = \sum_{j=1}^{n-1} y_j \sin(\pi j/n)$$

$$s_{2k}, \quad k = 1, \dots, \frac{n}{2} - 2 \quad \blacktriangleright \quad s_{2k} = -\text{Im}\{c_k\}.$$

MATLAB-Implementierung (basierend auf einer `fft` der Länge $n/2$):

3.8
p. 338

3.8
p. 339

3.8
p. 340

```
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
```

△

ANWSINUSTRAFO

Anwendung: Diagonalisierung lokaler translationsinvarianter linearer Operatoren

RUBENPDKNSTERN

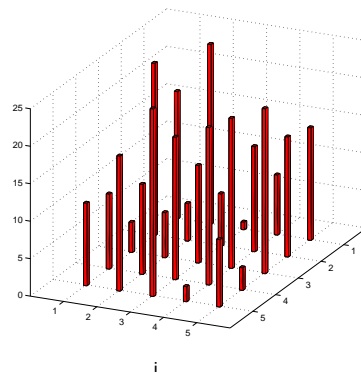
5-Punkt-Stern-Operator auf $\mathbb{R}^{n,n}$, $n \in \mathbb{N}$, in Gitterdarstellung:

$$T: \mathbb{R}^{n,n} \mapsto \mathbb{R}^{n,n}, \quad \mathbf{X} \mapsto T(\mathbf{X})$$

$$(T(\mathbf{X}))_{ij} := cx_{ij} + cyx_{i,j+1} + cyx_{i,j-1} + cx_{i+1,j} + cx_{i-1,j}$$

mit $c, cy, cx \in \mathbb{R}$, Konvention: $x_{ij} := 0$ für $(i,j) \notin \{1, \dots, n\}^2$.

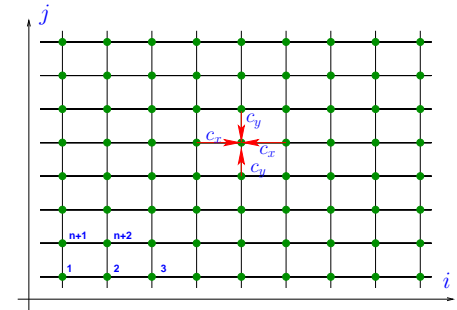
$\mathbf{X} \in \mathbb{R}^{n,n}$
 \Updownarrow
 Gitterfunktion $\in \{1, \dots, n\}^2 \mapsto \mathbb{R}$



Identifikation $\mathbb{R}^{n,n} \cong \mathbb{R}^{n^2}$, $x_{ij} \sim \tilde{x}_{(j-1)n+i} \rightarrow$ Matrixdarstellung $\mathbf{T} \in \mathbb{R}^{n^2, n^2}$ von T :

$$\mathbf{T} = \begin{pmatrix} \mathbf{C} & cy\mathbf{I} & 0 & \dots & \dots & 0 \\ cy\mathbf{I} & \mathbf{C} & cy\mathbf{I} & & & \vdots \\ 0 & \ddots & \ddots & \ddots & & \\ \vdots & & & cy\mathbf{I} & \mathbf{C} & cy\mathbf{I} \\ 0 & \dots & \dots & 0 & cy\mathbf{I} & \mathbf{C} \end{pmatrix} \in \mathbb{R}^{n^2, n^2},$$

$$\mathbf{C} = \begin{pmatrix} c & cx & 0 & \dots & \dots & 0 \\ cx & c & cx & & & \vdots \\ 0 & \ddots & \ddots & \ddots & & \\ \vdots & & & cx & c & cx \\ 0 & \dots & \dots & 0 & cx & c \end{pmatrix} \in \mathbb{R}^{n,n}.$$

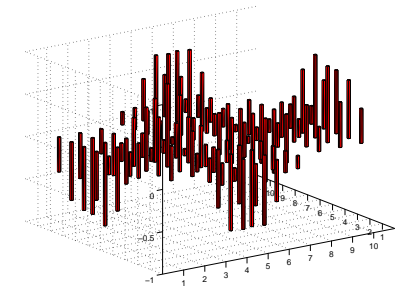


Sinusbasis von $\mathbb{R}^{n,n}$:

$$\mathbf{B}^{kl} = \left(\sin\left(\frac{\pi}{n+1}ki\right) \sin\left(\frac{\pi}{n+1}lj\right) \right)_{i,j=1}^n \quad (3.8.2)$$

$n = 10$: Gitterfunktion $\mathbf{B}^{2,3}$

➤



3.8
p. 341

3.8
p. 343

$$(T(\mathbf{B}^{kl}))_{ij} = c \sin\left(\frac{\pi}{n}ki\right) \sin\left(\frac{\pi}{n}lj\right) + cy \sin\left(\frac{\pi}{n}ki\right) \left(\sin\left(\frac{\pi}{n+1}l(j-1)\right) + \sin\left(\frac{\pi}{n+1}l(j+1)\right) \right) +$$

$$cx \sin\left(\frac{\pi}{n}lj\right) \left(\sin\left(\frac{\pi}{n+1}k(i-1)\right) + \sin\left(\frac{\pi}{n+1}k(i+1)\right) \right)$$

$$= \sin\left(\frac{\pi}{n}ki\right) \sin\left(\frac{\pi}{n}lj\right) (c + 2cy \cos\left(\frac{\pi}{n+1}l\right) + 2cx \cos\left(\frac{\pi}{n+1}k\right))$$

► \mathbf{B}^{kl} ist **Eigenvektor** von $T \leftrightarrow \mathbf{T}$ zum Eigenwert $c + 2cy \cos\left(\frac{\pi}{n+1}l\right) + 2cx \cos\left(\frac{\pi}{n+1}k\right)$

Algorithmus für Basistransformation:

$$\mathbf{X} = \sum_{k=1}^n \sum_{l=1}^n y_{kl} \mathbf{B}^{kl} \Rightarrow x_{ij} = \sum_{k=1}^n \sin\left(\frac{\pi}{n+1}ki\right) \sum_{l=1}^n y_{kl} \sin\left(\frac{\pi}{n+1}lj\right).$$

► Geschachtelte Sinustransformationen (\rightarrow Bem. 143) angewandt auf Zeilen/Spalten von $\mathbf{Y} = (y_{kl})_{k,l=1}^n$

Hier: Implementierung der Sinustransformation (3.8.1) mit „wrapping“-Technik

MATLAB-CODE Zweidimensionaler Sinustr.

```
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;
```

3.8
p. 342

3.8
p. 344

```

MATLAB-CODE FFT-basiertes Lösen aus lokalen
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));
translationsinvarianten linearen Operatoren

```

Diagonalisierung von T
durch 2D Sinustransformation

Effizienter Algorithmus
zum Lösen des LGS $T(X) = B$
Rechenaufwand $O(n^2 \log n)$

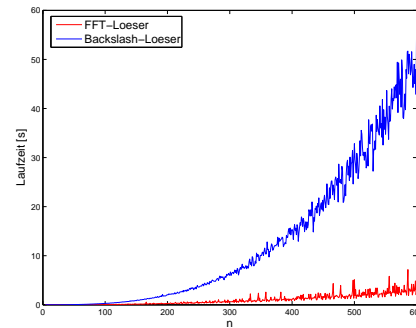
Beispiel 147 (Effizienz FFT-basierter Gleichungslöser).

tic-toc-Zeitmessung (MATLAB V7, Linux, Intel Pentium 4 Mobile CPU 1.80GHz)

```

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;

```



3.8.0.3 Kosinustransformation

COSTRANS

Noch eine trigonometrische Basistransformation im \mathbb{R}^n , $n \in \mathbb{N}$:

Standardbasis des \mathbb{R}^n „Kosinusbasis“

$$\left\{ \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \right\} \leftarrow \left\{ \begin{pmatrix} 2^{-1/2} \\ \cos(\frac{\pi}{2n}) \\ \cos(\frac{2\pi}{2n}) \\ \vdots \\ \cos(\frac{(n-1)\pi}{2n}) \end{pmatrix}, \begin{pmatrix} 2^{-1/2} \\ \cos(\frac{3\pi}{2n}) \\ \cos(\frac{6\pi}{2n}) \\ \vdots \\ \cos(\frac{3(n-1)\pi}{2n}) \end{pmatrix}, \dots, \begin{pmatrix} 2^{-1/2} \\ \cos(\frac{(2n-1)\pi}{2n}) \\ \cos(\frac{2(2n-1)\pi}{2n}) \\ \vdots \\ \cos(\frac{(n-1)(2n-1)\pi}{2n}) \end{pmatrix} \right\}$$

Matrix(Basistransformation Kosinusbasis \rightarrow Standardbasis):

$$C_n = (c_{ij}) \in \mathbb{R}^{n,n} \quad \text{mit} \quad c_{ij} = \begin{cases} 2^{-1/2} & , \text{ falls } i = 1, \\ \cos((i-1)\frac{2j-1}{2n}\pi) & , \text{ falls } i > 1. \end{cases}$$

Lemma 3.8.2 (Eigenschaften der Kosinusmatrix). $\sqrt{2/n} C_n$ ist reell und unitär (aber nicht symmetrisch).

Kosinustransformation: $c_k = \sum_{j=0}^{n-1} y_j \cos(k\frac{2j+1}{2n}\pi)$, $k = 1, \dots, n-1$, (3.8.3) **fft.cosfi**

$$c_0 = \frac{1}{\sqrt{2}} \sum_{j=0}^{n-1} y_j.$$

MATLAB-Implementierung von **Cy** („Wrapping“-Technik)

```

MATLAB-CODE Kosinustransformation
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-Implementierung von $C^{-1}y$ („Wrapping“-Technik):

```

MATLAB-CODE: Inverse Kosinustransformation
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);

```

3.8
p. 345

3.8.1 Zirkulante Matrizen

ZIRKMAT

Folge $(x_k)_{k \in \mathbb{Z}} = (\dots, x_{-2}, x_{-1}, x_0, x_1, x_2, \dots)$ ist n -periodisch, wenn $x_k = x_{n+k} \forall k \in \mathbb{Z}$.
(\rightarrow festgelegt durch x_0, \dots, x_{n-1})

DISCONV

Definition 3.8.3. Die **diskrete Faltung** (engl. discrete convolution) zweier n -periodischer Folgen $(x_k)_{k \in \mathbb{Z}}, (y_k)_{k \in \mathbb{Z}} \Rightarrow n$ -periodische Folge

$$z_k := (x_k) * (y_k) := \sum_{j=0}^{n-1} x_{k-j} y_j = (y_k) * (x_k) := \sum_{j=0}^{n-1} y_{k-j} x_j, \quad k \in \mathbb{Z}.$$

$$\begin{pmatrix} z_0 \\ \vdots \\ z_{n-1} \end{pmatrix} = \begin{pmatrix} x_0 & x_{n-1} & x_{n-2} & \cdots & \cdots & x_1 \\ x_1 & x_0 & x_{n-1} & & & \vdots \\ x_2 & x_1 & x_0 & \ddots & & \\ \vdots & \ddots & \ddots & \ddots & \ddots & \\ \vdots & & & \ddots & \ddots & x_{n-1} \\ x_{n-1} & \cdots & & & x_1 & x_0 \end{pmatrix} \begin{pmatrix} y_0 \\ \vdots \\ y_{n-1} \end{pmatrix}$$

3.8
p. 346

3.8
p. 347

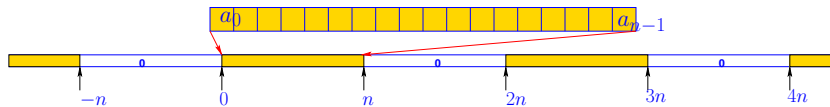
3.8
p. 348

Beispiel 148 (Multiplikation von Polynomen).

$$p(z) = \sum_{j=0}^{n-1} a_j z^j, \quad q(z) = \sum_{j=0}^{n-1} b_j z^j \quad \blacktriangleright \quad (pq)(z) = \sum_{j=0}^{2n-2} \left(\sum_{k=0}^j a_k b_{j-k} \right) z^j.$$

Konstruiere $2n - 1$ -periodische Koeffizientenfolgen

$$x_k := \begin{cases} a_k, & \text{falls } 0 \leq k < n, \\ 0, & \text{falls } n \leq k < 2n - 1, \end{cases} \quad y_k := \begin{cases} b_k, & \text{falls } 0 \leq k < n, \\ 0, & \text{falls } n \leq k < 2n - 1. \end{cases}$$



\blacktriangleright k . Koeffizient von pq : $c_j = ((x_k) * (y_k))_j, \quad j = 0, \dots, 2n - 2.$

◇

ZIRKULMAT

Definition 3.8.4 (Zirkulante Matrix).

Matrix $C = (c_{ij})_{i,j=1}^n \in \mathbb{K}^{n,n}$ = **zirkulant** $\Leftrightarrow \exists (u_k)_{k \in \mathbb{Z}}$ n -periodisch: $c_{ij} = u_{i-j}, 1 \leq i, j \leq n$.

- ☞ Zirkulante Matrix hat konstante (Neben-)Diagonalen $i - j = \text{const.}$
- ☞ Zeilen zyklisch verschoben.

$$\begin{pmatrix} u_0 & u_1 & u_2 & \cdots & \cdots & u_{n-1} \\ u_{n-1} & u_0 & & & & u_{n-2} \\ u_{n-2} & & & & & \vdots \\ \vdots & & & & & \vdots \\ \vdots & & & & & \vdots \\ u_1 & \cdots & \cdots & u_{n-1} & u_0 \end{pmatrix}$$

STRAHLTRANSP

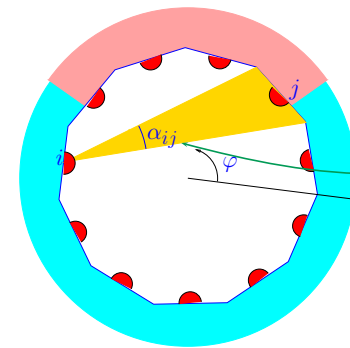
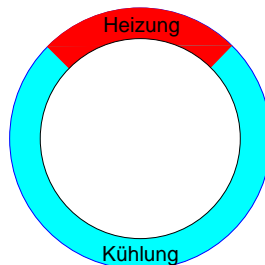
Beispiel 149 (Strahlungstransport).

Zylindrisches Rohr:

Heizung auf Teilen Γ_H des Umfangs (\rightarrow Wärmefluss vorgegeben)

Kühlung auf restlichem Umfang Γ_K (\rightarrow konstanter Wärmeabfluss)

Zu berechnen: Lokaler Wärmefluss



Modellierung:

- Approximation durch reguläres n -Eck, Seiten Γ_j
- Isotroper Halbkreisstrahler (Leistung I_j) auf Γ_j

$$\text{Strahlungsfluss } \Gamma_j \rightarrow \Gamma_i: P_{ji} := \frac{\alpha_{ij}}{\pi} I_j$$

Aperturwinkel: $\alpha_{ij} = \gamma_{|i-j|}, 1 \leq i, j \leq n$

$$\text{Gleichgewicht: } \sum_{i=1, i \neq j}^n P_{ji} - \sum_{i=1, i \neq j}^n P_{ij} = Q_j. \quad (3.8.4) \quad \text{eq: piper}$$

$Q_j \triangleq$ Wärmeeintritts-/abfluss durch Γ_j

$$Q_j := \int_{\frac{2\pi}{n}(j-1)}^{\frac{2\pi}{n}j} q(\varphi) d\varphi, \quad q(\varphi) := \begin{cases} \text{lokale Heizleistung} & \text{falls } \varphi \in \Gamma_H, \\ \frac{1}{|\Gamma_K|} \int_{\Gamma_H} q(\varphi) d\varphi & \text{falls } \varphi \in \Gamma_K. \end{cases}$$

$$\text{eq: piperadbal} \quad (3.8.4) \quad \Rightarrow \quad \text{LGS: } I_j - \sum_{i=1, i \neq j}^n \frac{\alpha_{ij}}{\pi} I_i = Q_j, \quad j = 1, \dots, n.$$

$$n = 8: \begin{pmatrix} 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 \\ -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 \\ -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 \\ -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 \\ -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 & -\gamma_3 \\ -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 & -\gamma_2 \\ -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 & -\gamma_1 \\ -\gamma_1 & -\gamma_2 & -\gamma_3 & -\gamma_4 & -\gamma_3 & -\gamma_2 & -\gamma_1 & 1 \end{pmatrix} \begin{pmatrix} I_1 \\ I_2 \\ I_3 \\ I_4 \\ I_5 \\ I_6 \\ I_7 \\ I_8 \end{pmatrix} = \begin{pmatrix} Q_1 \\ Q_2 \\ Q_3 \\ Q_4 \\ Q_5 \\ Q_6 \\ Q_7 \\ Q_8 \end{pmatrix}$$

\blacktriangleright Symmetrische, zirkulante, singuläre, positiv semidefinite (\rightarrow Lemma 3.4.7) Koeffizientenmatrix

◇

Effiziente Algorithmen für zirkulante Matrizen:

Sei $C \in \mathbb{C}^{n,n}$ zirkulante Matrix (\rightarrow Def. 3.8.4) $c_{ij} \stackrel{\text{def: zirmat}}{=} u_{i-j}, (u_k)_{k \in \mathbb{Z}} = n$ -periodische Folge
 $\mathbf{v}_k := (\omega_n^{-jk})_{j=0}^{n-1} \in \mathbb{C}^n, k = 0, \dots, n-1$ ($= k$. Spalte der Fourier-Matrix $\bar{\mathbf{F}}_n = n\mathbf{F}_n^{-1}$)

$$(\mathbf{C}\mathbf{v}_k)_j = \sum_{l=0}^{n-1} u_{j-l} \omega_n^{-lk} = \sum_{l=0}^{n-1} u_l \omega_n^{(l-j)k} = \omega_n^{-jk} \sum_{l=0}^{n-1} u_l \omega_n^{lk} = \lambda_k \omega_n^{-jk}.$$

unabhängig von j !

\mathbf{v}_k ist Eigenvektor von \mathbf{C} zum Eigenwert $\lambda_k = \sum_{l=0}^{n-1} u_l \omega_n^{lk} = (\mathbf{F}_n \mathbf{u})_k$.

Lemma 3.8.5 (Diagonalisierung zirkulanter Matrizen (\rightarrow Def. ^{Def:zirmat} 3.8.4)). Für jede zirkulante Matrix $\mathbf{C} \in \mathbb{K}^{n,n}$, $c_{ij} = u_{i-j}$, $(u_k)_{k \in \mathbb{Z}}$ n -periodische Folge, gilt

$$\mathbf{C} \bar{\mathbf{F}}_n = \bar{\mathbf{F}}_n \text{diag}(d_1, \dots, d_n) \quad , \quad \mathbf{d} = \mathbf{F}_n(u_0, \dots, u_{n-1})^T.$$

Beachte: $\bar{\mathbf{F}}_n = n \mathbf{F}_n^{-1} \rightarrow \mathbf{C} = \mathbf{F}_n^{-1} \text{diag}(d_1, \dots, d_n) \mathbf{F}_n$

Effizienter Algorithmus zur Berechnung der diskreten Faltung (\rightarrow Def. ^{Def:conv} 3.8.3) zweier n -periodischer Folgen $(x_k)_{k \in \mathbb{Z}}$, $(y_k)_{k \in \mathbb{Z}}$:

```
function z=conv(x,y)
n = length(y); z = zeros(n,1);
for i=1:n
    z(i)=dot(conj(x),...
        y([i:-1:1,n:-1:i+1]));
end
```

MATLAB-CODE Naive Diskrete Faltung

Gedankenlos: Aufwand $O(n^2)$

```
function z=conv(x,y)
z = ifft(fft(x).*fft(y));
```

MATLAB-CODE FFT-basierte diskrete Faltung

Professionell: Aufwand $O(n \log n)$

Effizienter FFT-basierter Algorithmus zur Invertierung einer zirkulanten Matrix

Bemerkung 150 (Primzahl-FFT).

Aus dem MATLAB-Manual:

When n is a prime number, the FFTW library first decomposes an n -point problem into three $(n-1)$ -point problems using Rader's algorithm ^{RAD68} [37]. It then uses the Cooley-Tukey decomposition described above to compute the $(n-1)$ -point DFTs.

Ein Satz aus der Zahlentheorie:

$$\forall p \in \mathbb{N} \text{ prim} \quad \exists g \in \{1, \dots, p-1\}: \{g^k \bmod p: k = 1, \dots, p-1\} = \{1, \dots, p-1\}.$$

Permutation $P_{p,g}: \{1, \dots, p-1\} \mapsto \{1, \dots, p-1\}$, $P_{p,g}(k) = g^k \bmod p$.
 Spiegelpermutation $P_k: \{1, \dots, k\} \mapsto \{1, \dots, k\}$, $P_k(i) = k - i + 1$.

Für Fouriermatrix $\mathbf{F} = (f_{ij})_{i,j=1}^p$: $P_{p-1} P_{p,g} (f_{ij})_{i,j=2}^p P_{p,g}^T$ ist zirkulant

3.8
p. 353

Beispiel für $p = 13$: $g = 2$, Permutation: $(2 \ 4 \ 8 \ 3 \ 6 \ 12 \ 11 \ 9 \ 5 \ 10 \ 7 \ 1)$.

$$\mathbf{F}_{13} \rightarrow \begin{matrix} \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 & \omega^0 \\ \omega^0 & \omega^2 & \omega^4 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 \\ \omega^0 & \omega^1 & \omega^2 & \omega^4 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} & \omega^7 \\ \omega^0 & \omega^7 & \omega^1 & \omega^2 & \omega^4 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} \\ \omega^0 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 & \omega^4 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 \\ \omega^0 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 & \omega^4 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 \\ \omega^0 & \omega^9 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 & \omega^4 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} \\ \omega^0 & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 & \omega^4 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} \\ \omega^0 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 & \omega^4 & \omega^8 & \omega^3 & \omega^6 \\ \omega^0 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 & \omega^4 & \omega^8 & \omega^3 \\ \omega^0 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 & \omega^4 & \omega^8 \\ \omega^0 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 & \omega^4 \\ \omega^0 & \omega^4 & \omega^8 & \omega^3 & \omega^6 & \omega^{12} & \omega^{11} & \omega^9 & \omega^5 & \omega^{10} & \omega^7 & \omega^1 & \omega^2 \end{matrix}$$

3.8
p. 355

Übung 3.9. Gegeben sei die Matrix $\mathbf{A} = (a_{lm})_{l,m=0}^{n-1} \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$. Zu berechnen sind die Werte

$$f_k := \sum_{l=0}^k a_{l,k-l} x_l y_{k-l}, \quad k = 0, \dots, n-1, \quad (3.9.1) \quad \text{Eq. 3k}$$

für beliebige Vektoren $\mathbf{x} := (x_0, \dots, x_{n-1})^T \in \mathbb{R}^n$, $\mathbf{y} := (y_0, \dots, y_{n-1})^T \in \mathbb{R}^n$.

3.8
p. 354

3.9
p. 356

- (i) Schreibe eine MATLAB-Funktion `function f = compf(A,x,y)`, die die f_k gemäss (3.9.1) berechnet.
- (ii) Was ist die asymptotische Komplexität von `compf` bzgl. n ?
- (iii) Entwickle eine *effiziente* Implementierung der Funktion `function f = compflr(p,q,x,y)` die `compf` für $A = pq^T$, $p, q \in \mathbb{R}^n$, realisiert (Hinweis: Beispiel 148)
- (iv) Was ist die asymptotische Komplexität von `compflr` bzgl. n ?

3.9.1 Toeplitz-Matrizen

TOEPLITZMAT Matrizen mit konstanten Diagonalen (Spezialfall: Zirkulante Matrizen \rightarrow Def. 3.8.4)

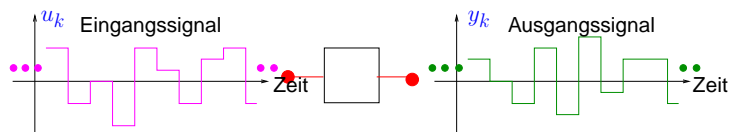
Definition 3.9.1 (Toeplitz-Matrix). $\mathbf{T} = (t_{ij})_{i,j=1}^n \in \mathbb{K}^{m,n}$ heisst *Toeplitz-Matrix*, wenn ein Vektor $\mathbf{u} = (u_{-m+1}, \dots, u_{n-1}) \in \mathbb{K}^{m+n-1}$ existiert, so dass $t_{ij} = u_{j-i}$, $1 \leq i \leq m$, $1 \leq j \leq n$.

$$\mathbf{T} = \begin{pmatrix} u_0 & u_1 & \cdots & \cdots & u_{n-1} \\ u_{-1} & u_0 & u_1 & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & u_1 \\ u_{1-m} & \cdots & \cdots & u_{-1} & u_0 \end{pmatrix}$$

Beispiel 151 (Lineare zeitinvariante Systeme).

- $(u_k)_{k \in \mathbb{Z}}$ m -periodisches zeitdiskretes Signal = input
- $(y_k)_{k \in \mathbb{Z}}$ m -periodisches Ausgangssignal eines *zeitinvarianten linearen Übertragungssystems*

$$\exists \mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{R}^n, \quad n \leq m : \quad y_k = \sum_{j=1}^n x_j u_{k-j+1}.$$



Identifikationsproblem = Gesucht: $\mathbf{x} \in \mathbb{R}^n$ so, dass

$$\|\mathbf{Ax} - \mathbf{y}\|_2 = \left\| \begin{pmatrix} u_0 & u_{-1} & \cdots & \cdots & u_{1-n} \\ u_1 & u_0 & u_{-1} & & \vdots \\ \vdots & u_1 & u_0 & \ddots & \\ \vdots & & \ddots & \ddots & u_{-1} \\ u_{n-1} & & & u_1 & u_0 \\ u_n & u_{n-1} & & u_1 & u_0 \\ \vdots & & & \vdots & \vdots \\ u_{m-1} & \cdots & \cdots & u_{m-n} & \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} - \begin{pmatrix} y_0 \\ \vdots \\ y_{m-1} \end{pmatrix} \right\|_2 \rightarrow \min.$$

> **Lineares Ausgleichsproblem**, \rightarrow Abschnitt 3.5, mit Toeplitz-Matrix \mathbf{A} : $a_{ij} = u_{i-j}$.

Matrix der Normalgleichungen (\rightarrow Abschnitt 3.5.2)

$$\mathbf{M} := \mathbf{A}^H \mathbf{A}, \quad m_{ij} = \sum_{k=1}^m u_{k-i} u_{k-j} = z_{i-j} \text{ wegen Periodizität der Folge } (u_k)_{k \in \mathbb{Z}}.$$

> $\mathbf{M} \in \mathbb{R}^{n,n}$ ist ebenfalls Toeplitz-Matrix & symmetrisch positiv semi-definit.

3.9
p. 357

Beispiel 152 (Lineare Regression bei stationären Markov-Ketten).

Folge von Zufallsvariablen: $(y_k)_{k \in \mathbb{Z}} =$ Markov-Kette mit zeitunabhängiger Korrelation

$$\text{Erwartungswert } E(y_{i-j} y_{i-k}) = u_{k-j} \quad \forall i, j, k \in \mathbb{Z}, \quad u_i = u_{-i}.$$

Modell: Linearer Zusammenhang

$$\exists \mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{R}^n : \quad y_k = \sum_{j=1}^n x_j y_{k-j} \quad \forall k \in \mathbb{Z}.$$

mit zu schätzenden Koeffizienten x_j , $j = 1, \dots, n$: Für festes $i \in \mathbb{Z}$

$$\mathbf{x} = \arg \max_{\mathbf{x} \in \mathbb{R}^n} E \left| y_i - \sum_{j=1}^n x_j y_{i-j} \right|^2 \quad \blacktriangleright \quad E|y_i|^2 - 2 \sum_{j=1}^n x_j u_k + \sum_{k,j=1}^n x_k x_j u_{k-j} \rightarrow \max.$$

$$\blacktriangleright \quad \mathbf{x}^T \mathbf{Ax} - \mathbf{b}^T \mathbf{x} \rightarrow \max \quad \text{mit } \mathbf{b} = (u_k)_{k=1}^n, \quad \mathbf{A} = (u_{i-j})_{i,j=1}^n.$$

Kovarianzmatrix \mathbf{A} ist s.p.d. Toeplitz-Matrix, \mathbf{x} Lösung der **Yule-Walker-Gleichung** $\mathbf{Ax} = \mathbf{b}$.

3.9.1.1 Toeplitz-Matrix-Arithmetik

TOEPMATHARITH
Gegeben: Toeplitz-Matrix $\mathbf{T} = (u_{j-i}) \in \mathbb{K}^{m,n}$ mit erzeugendem Vektor $\mathbf{u} = (u_{-m+1}, \dots, u_{n-1}) \in \mathbb{C}^{2n-1}$, Vektor $\mathbf{x} \in \mathbb{K}^n$

3.9
p. 359

3.9
p. 360

Aufgabe: Effiziente Realisierung des Matrix×Vektor-Produkts $\mathbf{T}\mathbf{x}$

Beobachtung: Folgende erweiterte Matrix ist zirkulant! (Beispiel für $m = n$)

$$\mathbf{C} = \begin{pmatrix} \mathbf{T} & \mathbf{S} \\ \mathbf{S} & \mathbf{T} \end{pmatrix} = \left(\begin{array}{cccc|cccc} u_0 & u_1 & \cdots & \cdots & u_{n-1} & 0 & u_{1-n} & \cdots & \cdots & u_{-1} \\ u_{-1} & u_0 & u_1 & \cdots & \vdots & u_{n-1} & 0 & \cdots & \cdots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots & \cdots & \vdots \\ \vdots & & \ddots & \ddots & \vdots & \vdots & & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & u_1 & \vdots & & \ddots & \ddots & u_{1-n} \\ u_{1-n} & \cdots & \cdots & \cdots & u_{-1} & u_0 & u_1 & \cdots & \cdots & u_{n-1} \\ 0 & u_{1-n} & \cdots & \cdots & u_{-1} & u_0 & u_1 & \cdots & \cdots & u_{n-1} \\ u_{n-1} & 0 & \cdots & \cdots & \vdots & u_{-1} & u_0 & u_1 & \cdots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots & \cdots & \vdots \\ \vdots & & \ddots & \ddots & \vdots & \vdots & & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & u_{1-n} & \vdots & & \ddots & \ddots & u_1 \\ u_1 & \cdots & \cdots & \cdots & u_{n-1} & 0 & u_{1-n} & \cdots & \cdots & u_{-1} & u_0 \end{array} \right)$$

Allgemein: $\mathbf{T} = \text{toeplitz}(u(0:-1:1-m), u(0:n-1));$
 $\mathbf{S} = \text{toeplitz}([0, u(n-1:-1:n-m+1)], [0, u(1-m:1:-1)]);$

$$\mathbf{C} \begin{pmatrix} \mathbf{x} \\ 0 \end{pmatrix} = \begin{pmatrix} \mathbf{T}\mathbf{x} \\ \mathbf{S}\mathbf{x} \end{pmatrix}$$

Rechenaufwand $O(n \log n)$ für Auswertung $\mathbf{T}\mathbf{x}$ (mit FFT, Abschn. 3.8.1)

3.9.1.2 Der Levinson-Algorithmus

Gegeben: S.p.d. Toeplitz-Matrix $\mathbf{T} = (u_{j-i})_{i,j=1}^n$
 mit erzeugendem Vektor $\mathbf{u} = (u_{-n+1}, \dots, u_{n-1}) \in \mathbb{C}^{2n-1}$, $u_{-k} = u_k$, o.B.d.A. $u_0 = 1$,
 Vektor $\mathbf{b} \in \mathbb{K}^n$

Gesucht: Effizientes Verfahren zur Lösung des LGS $\mathbf{T}\mathbf{x} = \mathbf{b}$

Induktive Konstruktion des Algorithmus:

Definiere: $\mathbf{T}_k := (u_{j-i})_{i,j=1}^k \in \mathbb{K}^{k,k}$ (linkes oberes Dreieck, ebenfalls s.p.d. & Toeplitz),
 $\mathbf{x}^k \in \mathbb{K}^k$: $\mathbf{T}_k \mathbf{x}^k = (b_1, \dots, b_k)^T$,
 $\mathbf{u}^k := (u_1, \dots, u_k)^T$

$$\mathbf{T}_{k+1} \mathbf{x}^{k+1} = \left(\begin{array}{c|c} \mathbf{T}_k & \begin{matrix} u_k \\ \vdots \\ u_1 \end{matrix} \\ \hline u_k & \cdots & u_1 & 1 \end{array} \right) \begin{pmatrix} \tilde{\mathbf{x}}^{k+1} \\ x_{k+1}^{k+1} \end{pmatrix} = \begin{pmatrix} b_1 \\ \vdots \\ b_k \\ b_{k+1} \end{pmatrix}$$

Spiegelungsp permutation: $P_k : \{1, \dots, k\} \mapsto \{1, \dots, k\}$, $P_k(i) := k - i + 1$

$$\begin{aligned} \tilde{\mathbf{x}}_{k+1} &= \mathbf{T}_k^{-1} (\tilde{\mathbf{b}}_{k+1} - x_{k+1}^{k+1} P_k \mathbf{u}^k) = \mathbf{x}^k - x_{k+1}^{k+1} \mathbf{T}_k^{-1} P_k \mathbf{u}^k, \\ x_{k+1}^{k+1} &= b_{k+1} - P_k \mathbf{u}^k \cdot \tilde{\mathbf{x}}^{k+1} = b_{k+1} - P_k \cdot \mathbf{x}^k + x_{k+1}^{k+1} P_k \cdot \mathbf{T}_k^{-1} P_k \mathbf{u}^k. \end{aligned}$$

Mit Hilfsvektoren $\mathbf{y}^k := \mathbf{T}_k^{-1} P_k \mathbf{u}^k$

$$\mathbf{x}^{k+1} = \begin{pmatrix} \tilde{\mathbf{x}}^{k+1} \\ x_{k+1}^{k+1} \end{pmatrix} \quad \text{mit} \quad \begin{aligned} x_{k+1}^{k+1} &= (b_{k+1} - P_k \mathbf{u}^k) / \sigma_k, \\ \tilde{\mathbf{x}}^{k+1} &= \mathbf{x}^k - x_{k+1}^{k+1} \mathbf{y}^k, \end{aligned} \quad \sigma_k := 1 - P_k \mathbf{u}^k \cdot \mathbf{y}^k.$$

MATLAB-Implementierung
 (rekursiv, u_{n+1} nicht verwendet!)

Lineare Rekursion:

Rechenaufwand $O(n-k)$ auf k . Rekursionsstufe, $k = 0, \dots, n-1$

Gesamtaufwand $O(n^2)$

```

MATLAB-CODE Levinson-Algorithmus
function [x,y] = levinson(u,b)
k = length(u)-1;
if (k == 0)
    x=b(1); y = u(1); return;
end
[xk,yk] = levinson(u(1:k),b(1:k));
sigma = 1-dot(u(1:k),yk);
t= (b(k+1)-dot(u(k:-1:1),xk))/sigma;
x= [ xk-t*yk(k:-1:1);t];
s= (u(k+1)-dot(u(k:-1:1),yk))/sigma;
y= [yk-s*yk(k:-1:1); s];
    
```

Bemerkung 153 („Superschnelle“ Toeplitz-Löser). FFT-basierte Algorithmen zur Lösung von $\mathbf{T}\mathbf{x} = \mathbf{b}$ benötigen Rechenaufwand $O(n \log^3 n)$

Übung 3.10. Gegeben eine grosse ganze Zahl $N \gg 1$ und zwei Polynome

$$p(z) = \sum_{j=0}^{2N} \alpha_j z^j, \quad q(z) = \sum_{j=0}^N \beta_j z^j, \quad \alpha_j, \beta_j \in \mathbb{C},$$

die Polynomdivision mit fragt nach einem Polynom s vom Grad N und einem Polynom r vom Grad $N-1$, so dass

$$p(z) = q(z)s(z) + r(z).$$

Implementiere eine effiziente MATLAB-Funktion `[s,r] = polydiv(a,b)`, der die Koeffizienten α_j und β_j übergeben werden und die die Koeffizienten von s und r berechnet. Dabei ist zu berücksichtigen, dass in MATLAB ein Polynom durch den Vektor seiner Koeffizienten zu *absteigenden* Potenzen kodiert wird, so etwa das Polynom q durch $(\alpha_{2N}, \alpha_{2N-1}, \dots, \alpha_0) \in \mathbb{C}^{2N+1}$.

Gib die Asymptotik des Rechenaufwandes in N an.

Übung 3.11. Für stetige Funktionen $f, u, h \in C^0([0, 1])$ betrachte die Integralgleichung

$$\int_0^t f(t-\tau)u(\tau) d\tau = h(t), \quad 0 \leq t \leq \tau. \quad (3.11.1)$$

Die Gleichung wird auf folgende Art *diskretisiert*: Man wählt ein äquidistantes Gitter $\{t_j := j/N\}_{j=0,\dots,N}$, $N \in \mathbb{N}$, betrachtet die Gleichung nur an den Gitterpunkten t_j , $j = 1, \dots, N$

(**Kollokation**), und approximiert die Integrale durch Mittelpunktsquadratur (die einfachste Gauss-Quadraturformel) auf den Intervallen $[t_{j-1}, t_j]$, $j = 1, \dots, N$.

So gelangt man zu der diskretisierten Integralgleichung

$$\frac{1}{N} \sum_{j=1}^n f(t_n - m_j) u(m_j) = h(t_n), \quad n = 1, \dots, N, \quad (3.11.2)$$

wobei $m_j := 1/2(t_j + t_{j-1})$. Implementiere eine effiziente MATLAB-Funktion `u = inteqsolve(N, h, f)`, die für Funktionshandles `h`, `f` die Werte $u(m_j)$, $j = 1, \dots, N$, aus (3.11.2) berechnet.

4 Interpolation und Approximation

Allgemeine „Funktion“ $f : D \subset \mathbb{R} \mapsto \mathbb{K}$ = Konzept der Analysis (unendlich viel Information)

Herausforderung: „Computergerechte“ Repräsentation von Funktionen



Idee: **Parametrisierung**

→ endlich viele Parameter $\alpha_1, \dots, \alpha_n$, $n \in \mathbb{N}$, charakterisieren Funktion.
(→ Bsp. [ex:NLLSQ](#))

Spezialfall: Darstellung durch (endliche) Linearkombination von **Basisfunktionen** $b_j : D \subset \mathbb{R} \mapsto \mathbb{K}$, $j = 1, \dots, n$:

$$f = \sum_{j=1}^n \alpha_j b_j, \quad \alpha_j \in \mathbb{K}.$$

► $f \in$ endlichdimensionaler **Funktionsraum** $V_n := \text{Span}\{b_1, \dots, b_n\}$

4.1 Polynomiale Techniken

[File: section-polynomiale-techniken.tex, SVN: section-polynomiale-techniken.tex 1299 2007-02-01 09:09:04Z kalai]

Notation: Vektorraum der Polynome vom Grad $\leq k$, $k \in \mathbb{N}$:

$$\mathcal{P}_k := \{t \mapsto \alpha_k t^k + \alpha_{k-1} t^{k-1} + \dots + \alpha_1 t + \alpha_0, \alpha_j \in \mathbb{K}\}. \quad (4.1.1)$$

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

Warum Polynome? → Einfach auszuwerten, zu integrieren, zu differenzieren

→ Vektorraumstruktur & Algebra

↔ Analysis: Taylorpolynome & Potenzreihen

Bemerkung 154. $t \mapsto \alpha_k t^k + \alpha_{k-1} t^{k-1} + \dots + \alpha_0$ = **monomiale Darstellung** eines Polynoms

MATLAB: $\alpha_k t^k + \alpha_{k-1} t^{k-1} + \dots + \alpha_0$ ► Vektor $(\alpha_k, \alpha_{k-1}, \dots, \alpha_0)$ (Anordnung !)

3.11
p. 365

- Auswertung eines Polynoms in monomialer Darstellung (→ Bem. [154](#)): **Hornerschema**

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

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

► Rechenaufwand $O(n)$

Benutze: MATLAB „built-in“-Funktion `polyval(x,p)`;

4.1.1 Polynominterpolation

Ziel: Rekonstruktion von Funktionen aus Wertepaaren, vgl. „Fitten“ Bsp. [ex:linfit](#)

4.1
p. 366

4.1
p. 367

4.1
p. 368

ALLGEMEINE LAGRANGE-POLYNOMINTERPOLATIONSAUFGABE

Zu (einfachen) Knoten (engl. nodes) $t_0, \dots, t_n, n \in \mathbb{N}, -\infty < t_0 < t_1 < \dots < t_n < \infty$ und Werten $y_0, \dots, y_n \in \mathbb{K}$ bestimme $p \in \mathcal{P}_n$ so, dass

$$p(t_j) = y_j \quad \text{für } j = 0, \dots, n. \quad (4.1.2) \quad \text{intp:agr}$$

ALLGEMEINE POLYNOMINTERPOLATIONSAUFGABE

Zu (evtl. mehrfachen) Knoten $t_0, \dots, t_n, n \in \mathbb{N}, -\infty < t_0 \leq t_1 \leq \dots \leq t_n < \infty$ und Werten $y_0, \dots, y_n \in \mathbb{K}$ bestimme $p \in \mathcal{P}_n$ so, dass

$$\frac{d}{dt^k} p(t_j) = y_j \quad \text{für } k = 0, \dots, l_j \quad \text{und } j = 0, \dots, n, \quad (4.1.3) \quad \text{intp:state}$$

mit $l_j := \max\{j - i : t_j = t_i, i = 0, \dots, n\}$ (Vielfachheit des Knotens t_j)

Doppelte Knoten $t_0 = t_1 < t_2 = t_3 < \dots < t_{2j} = t_{2j+1} < \dots < t_n$
~~HERMITE-INTERPOLATION~~
 Hermite-Interpolation

4.1.1.1 Theorie und Kondition

Für $t_0 < t_1 < \dots < t_n$ (\rightarrow Lagrange-Interpolation) betrachte

$$\text{LAGRANGE-POLYNOME} \quad L_i(t) := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j}, \quad i = 0, \dots, n. \quad (4.1.4) \quad \text{eq:lagrar}$$

$$L_i \in \mathcal{P}_n \quad \text{und} \quad L_i(t_j) = \delta_{ij}$$

Beispiel 155 (Lagrange-Polynome).

Lagrange-Polynome für äquidistante Knoten

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

Plot $n = 10$

lagrplot

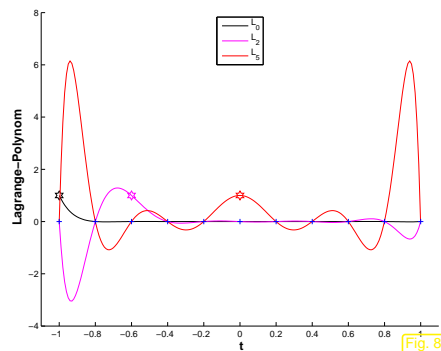


Fig. 8

$$\blacktriangleright \quad \text{Für Lagrange-Interpolation: Darstellung} \quad p(t) = \sum_{i=0}^n y_i L_i(t). \quad (4.1.5) \quad \text{eq:lagre}$$

Theorem 4.1.1 (Existenz & Eindeutigkeit des Interpolationspolynoms). Die Polynominterpolationsaufgabe (4.1.3) hat eine eindeutige Lösung $p \in \mathcal{P}_n$.

Beweis. für Lagrange-Interpolation:

$$p(t_j) = y_j \iff \sum_{i=0}^n a_i t_j^i = y_j, \quad j = 0, \dots, n \\ \iff (n+1) \times (n+1) \text{ lineares Gleichungssystem}$$

(4.1.5) \Rightarrow Existenz \Rightarrow Eindeutigkeit des Interpolationspolynoms \square

\blacktriangleright Polynominterpolation zu Knoten $\mathcal{T} := \{t_j\}_{j=0}^n =$ lineare Abbildung

$$I_{\mathcal{T}} : \mathbb{K}^{n+1} \mapsto \mathcal{P}_n, \quad (y_0, \dots, y_n)^T \mapsto \text{Interpolationspolynom } p. \quad (4.1.6) \quad \text{intp:point}$$

4.1
p. 369

Definition 4.1.2 (Verallgemeinerte Lagrange-Polynome).

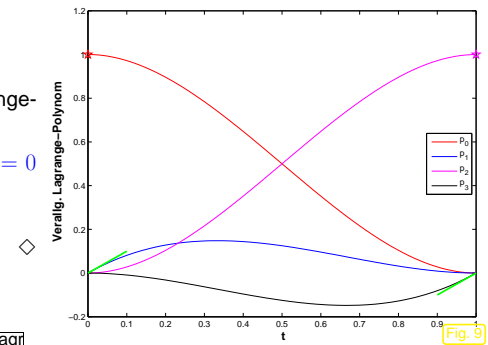
Die verallgemeinerten Lagrange-Polynome zu einer Knotenmenge $\mathcal{T} = \{t_j\}_{j=0}^n \subset \mathbb{R}$ sind $L_i := I_{\mathcal{T}}^{-1}(\mathbf{e}_{i+1}), i = 0, \dots, n$, wobei $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^T \in \mathbb{R}^{n+1} \hat{=} i$. Einheitsvektor.

Beispiel 156. (Verallgemeinerte Lagrange-Polynome für Hermite-Interpolation)

Doppelte Knoten $t_0 = 0, t_1 = 0, t_2 = 0, t_3 = 0$

$\blacktriangleright n = 3$

(Kubische Hermite-Interpolation)



hermlagr

Fig. 9

4.1
p. 370

4.1
p. 372

Erforderlich für Konditionsuntersuchung:

Normen auf $C^\infty(I)$, $I \subset \mathbb{R}$

NINP

Supremumsnorm $\|f\|_{L^\infty(I)} := \sup\{|f(t)| : t \in I\}$,

(4.1.7) **Ninf**

NLTWO

L^2 -Norm $\|f\|_{L^2(I)}^2 := \int_I |f(t)|^2 dt$,

(4.1.8) **NLTWO**

NLONE

L^1 -Norm $\|f\|_{L^1(I)} := \int_I |f(t)| dt$.

(4.1.9) **NLone**

Lemma 4.1.3 (Absolute Kondition der Polynominterpolation). Für eine Knotenmenge $\mathcal{T} \subset \mathbb{R}$ mit zugehörigen verallgemeinerten Lagrange-Polynomen L_i , $i = 0, \dots, n$, und festes $I \subset \mathbb{R}$ gilt

$$\|l_{\mathcal{T}}\|_{\infty \rightarrow \infty} := \sup_{\mathbf{y} \in \mathbb{K}^{n+1} \setminus \{0\}} \frac{\|l_{\mathcal{T}}(\mathbf{y})\|_{L^\infty(I)}}{\|\mathbf{y}\|_\infty} = \left\| \sum_{i=0}^n |L_i| \right\|_{L^\infty(I)}, \quad (4.1.10) \text{ eq:IPN1}$$

$$\|l_{\mathcal{T}}\|_{2 \rightarrow 2} := \sup_{\mathbf{y} \in \mathbb{K}^{n+1} \setminus \{0\}} \frac{\|l_{\mathcal{T}}(\mathbf{y})\|_{L^2(I)}}{\|\mathbf{y}\|_2} \leq \left(\sum_{i=0}^n \|L_i\|_{L^2(I)}^2 \right)^{\frac{1}{2}}. \quad (4.1.11) \text{ eq:IPN2}$$

Beweis. (für L^∞ -Norm) Mit \triangle -Ungleichung

$$\|l_{\mathcal{T}}(\mathbf{y})\|_{L^\infty(I)} = \left\| \sum_{j=0}^n y_j L_j \right\|_{L^\infty(I)} \leq \sup_{t \in I} \sum_{j=0}^n |y_j| |L_j(t)|,$$

Gleichheit in (4.1.10) für $\mathbf{y} := (\text{sgn}(L_j(t^*)))_{j=0}^n$, $t^* := \arg\max_{t \in I} \sum_{i=0}^n |L_i(t)|$. □

Beweis. (für L^2 -Norm) Mit \triangle -Ungleichung und Cauchy-Schwarz-Ungleichung

$$\|l_{\mathcal{T}}(\mathbf{y})\|_{L^2(I)} \leq \sum_{j=0}^n |y_j| \|L_j\|_{L^2(I)} \leq \left(\sum_{j=0}^n |y_j|^2 \right)^{\frac{1}{2}} \left(\sum_{j=0}^n \|L_j\|_{L^2(I)}^2 \right)^{\frac{1}{2}}. \quad \square$$

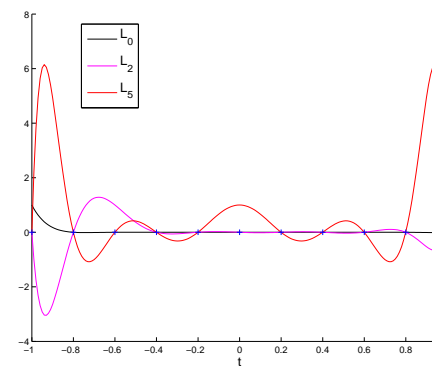
Terminologie:

LEBESGUE

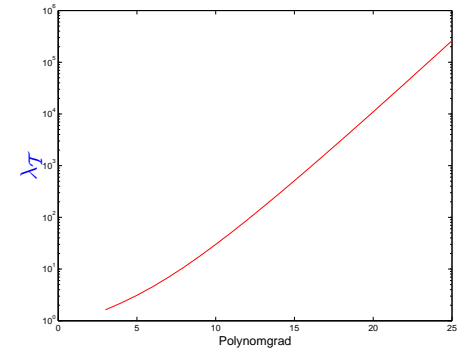
Lebesgue-Konstante zu \mathcal{T} : $\lambda_{\mathcal{T}} := \left\| \sum_{i=0}^n |L_i| \right\|_{L^\infty(I)}$

Beispiel 157 (Berechnung von Lebesgue-Konstanten).

$I = [-1, 1]$, $\mathcal{T} = \{-1 + \frac{k}{2n}\}_{k=0}^n$ (**äquidistante Knoten**)



Lagrange-Polynome für $n = 10$



Lebesgue-Konstante für äquidistante Knoten

Asymptotische Abschätzung (Stirlingsche Formel): für $n = 2m$, $I = [-1, 1]$, $\mathcal{T} = \{-1 + \frac{2k}{n}\}_{k=0}^n$

$$|L_m(1 - \frac{1}{n})| = \frac{\frac{1}{n} \cdot \frac{1}{n} \cdot \frac{3}{n} \dots \frac{n-3}{n} \cdot \frac{n+1}{n} \dots \frac{2n-1}{n}}{\left(\frac{2}{n} \cdot \frac{4}{n} \dots \frac{n-2}{n} \cdot 1\right)^2} = \frac{(2n)!}{(n-1)2^{2n}((n/2)!)^2 n!} \sim \frac{2^{n+1/2}}{\pi(n-1)n^{3/2}}$$

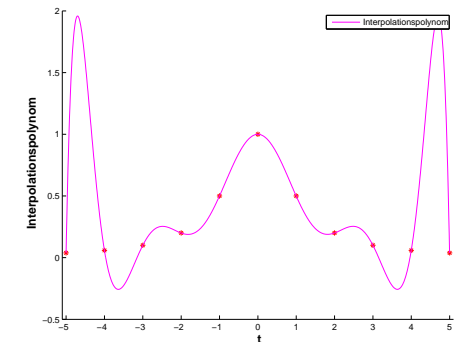
4.1 Theorie **COR92** [8]: für äquidistante Knoten $\lambda_{\mathcal{T}} \geq C e^{n/2}$ für $C > 0$ unabhängig von n . □
p. 373

Beispiel 158 (Oszillatorische Interpolationspolynome).

Polynominterpolation mit äquidistanten Knoten:

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

Plot $n = 10 \gg$



Gefahr: Starke Oszillationen von Interpolationspolynomen hohen Grades auf äquidistanten Knoten □

4.1.1.2 Algorithmen

- Vielfache Auswertung eines Interpolationspolynoms zu Knotenmenge $\mathcal{T} = \{t_j\}_{j=0}^n$, Stützwerten y_0, \dots, y_n :

$$p(t) = \sum_{i=0}^n \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j} y_i = \sum_{i=0}^n \lambda_i \prod_{\substack{j=0 \\ j \neq i}}^n (t - t_j) y_i = \prod_{j=0}^n (t - t_j) \cdot \sum_{i=0}^n \frac{\lambda_i}{t - t_i} y_i.$$

mit $\lambda_i = \frac{1}{(t_i - t_0) \cdots (t_i - t_{i-1})(t_i - t_{i+1}) \cdots (t_i - t_n)}$, $i = 0, \dots, n$.

$$\text{Aus } 1 = \prod_{j=0}^n (t - t_j) \sum_{i=0}^n \frac{\lambda_i}{t - t_i} \Rightarrow \prod_{j=0}^n (t - t_j) = \frac{1}{\sum_{i=0}^n \frac{\lambda_i}{t - t_i}}$$

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

Falls λ_i verfügbar:

Rechenaufwand(Auswertung von $p(t)$) = $O(n)$:

MATLAB-CODE: Auswertung eines Interpolationspolynoms

```
function p = intpolyval(t,y,x)
n = length(t); N = length(x);
for k = 1:n, lambda(k) = 1 / prod(t(k) - t([1:k-1,k+1:n])); end;
for i = 1:N
    z = (x(i)-t); j = find(z == 0);
    if (~isempty(j)), p(i) = y(j);
    else
        mu = lambda./z;
        p(i) = dot(mu,y)/sum(mu);
    end
end
```

Vermeide Division durch 0 !

- Auswertung eines Interpolationspolynoms an wenigen Stellen: **Aitken-Neville-Schema**

Gegeben: Knotenmenge $\mathcal{T} := \{t_j\}_{j=0}^n \subset \mathbb{R}$ paarweise verschieden, $t_i \neq t_j$ für $i \neq j$,
Stützwerte y_0, \dots, y_n , Auswertungsstelle $x \in \mathbb{R}$

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

p_{i_0, \dots, i_m} = Interpolationspolynom durch $(t_{i_0}, y_{i_0}), \dots, (t_{i_m}, y_{i_m})$:

$$\blacktriangleright p_i \equiv y_i, \quad i = 0, \dots, n, \quad p_{i_0, \dots, i_m}(t) = \frac{(t - t_{i_0})p_{i_1, \dots, i_m} - (t - t_{i_m})p_{i_0, \dots, i_{m-1}}}{t_{i_m} - t_{i_0}}. \quad (4.1.12)$$

Schema:

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

```
function v = ANipoleval(t,y,x)
for i=1:length(y)
    for k=i-1:-1:1
        y(k) = y(k+1) + (y(k+1)-y(k)) * ...
            (x-t(i)) / (t(i)-t(k));
    end
end
v = y(1);
```

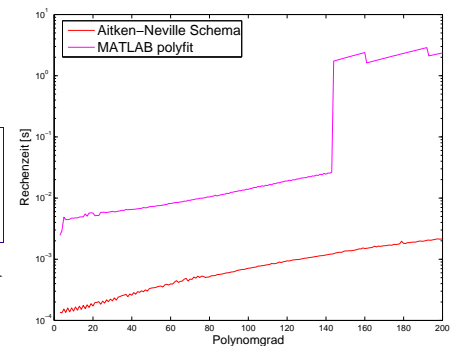
4.1
p. 377

Äquivalente MATLAB-Funktion:

(mit deutlich höherem Rechenaufwand !)

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

tic-toc-Rechenzeitmessung, 100 runs ➤



4.1
p. 379

Wichtige Anwendung:

Extrapolation nach Null

Problem: Gesucht ist $f(0)$ mit vorgeschriebener Genauigkeit, aber implementierte Funktion $y=f(h)$ ist für $|t| \ll 1$ nicht stabil auszuwerten.

Bekannt:

Existenz einer **asymptotischen Entwicklung** in h^2

$$f(h) = f(0) + A_1 h^2 + A_2 h^4 + \cdots + A_n h^{2n} + R(h), \quad A_k \in \mathbb{K},$$

mit **Restgliedabschätzung** $|R(h)| = O(h^{2n+2})$ für $h \rightarrow 0$.

4.1
p. 378

4.1
p. 380



Idee:

- ① Auswertung $f(t_i)$ für verschiedene $t_i, i = 0, \dots, n, |t_i| > 0$.
- ② $f(0) \approx p(0)$ mit Interpolationspolynom $p \in \mathcal{P}_n, p(t_i) = f(t_i)$.

Beispiel 159 (Numerische Differentiation durch Extrapolation, \rightarrow Bsp. ex:numdiff 27).

Für $2(n+1)$ -mal stetig differenzierbare Funktion $f : D \subset \mathbb{R} \mapsto \mathbb{R}, x \in D$ (Taylorsumme mit Lagrange-Restglied)

$$T(h) := \frac{f(x+h) - f(x-h)}{2h} \sim f'(x) + \sum_{k=1}^n \frac{1}{(2k)!} \frac{d^{2k}f}{dx^{2k}}(x) h^{2k} + \frac{1}{(2n+2)!} f^{(2n+2)}(\xi(x)) \cdot h^{2n+2}.$$

MATLAB-Code (f = handle auf Funktion):

MATLAB-CODE: Numerische Differentiation durch Interpolation

```
function d = diffex(f,x,h0,tol)
h = h0;
y(1) = (f(x+h0)-f(x-h0))/(2*h0);
for i=2:10
    h(i) = h(i-1)/2;
    y(i) = (f(x+h(i))-f(x-h(i)))/h(i-1);
    for k=i-1:-1:1
        y(k) = y(k+1)-(y(k+1)-y(k))*h(i)/(h(i)-h(k));
    end
    if (abs(y(2)-y(1)) < tol*abs(y(1))), break; end
end
d = y(1);
```

A posteriori Fehlerschätzung

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

Grad	Relativer Fehler	Grad	Relativer Fehler	h	Relativer Fehler
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	0.0000000589624
5	0.0000000746031	5	0.00000001253250	5	0.0000000009546
6	0.0000000001224	6	0.0000000004816	6	0.0000000000002
		7	0.0000000000021		

Vorteile:

Garantierte Genauigkeit > Effizienz

Zum Vergleich: Numerische Differentiation mit finiten (Vorwärts-) Differenzen \rightarrow Auslöschung
(\rightarrow Bsp. ex:numdiff 27)

```
x=1.1; h=2.^[-1:-5:-36];
atan_err = atan_err=abs(dirnumdiff(@atan,x,h)-1/(1+x^2))*(1+x^2);
sqrt_err=abs(dirnumdiff(@sqrt,x,h)-1/(2*sqrt(x)))*(2*sqrt(x));
exp_err=abs(dirnumdiff(@exp,x,h)-exp(x))/exp(x);
```

4.1
p. 381

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

4.1
p. 383

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

h	Relativer Fehler
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}	0.000000679016113

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

h	Relativer Fehler
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	Relativer Fehler
2^{-1}	0.29744254140026
2^{-6}	0.00785334954789
2^{-11}	0.00024418036620
2^{-16}	0.00000762943394
2^{-21}	0.00000023835113
2^{-26}	0.0000000429331
2^{-31}	0.00000012467100
2^{-36}	0.00000495453865

4.1.2 Interpolationsfehlerabschätzungen

INTERPASCHE

Wir betrachten Lagrangesche Polynominterpolation zu paarweise verschiedenen Knoten

$\mathcal{T} := \{t_0, \dots, t_n\} \subset I, I \subset \mathbb{R}$ Intervall mit Länge $|I|$.

4.1
p. 382

4.1
p. 384

Notation: Für eine stetige Funktion $f: I \mapsto \mathbb{K}$ definiere Polynominterpolationsoperator $l_{\mathcal{T}}(f) := l_{\mathcal{T}}(\mathbf{y}) \in \mathcal{P}_n$ mit $\mathbf{y} := (f(t_0), \dots, f(t_n))^T \in \mathbb{K}^{n+1}$.

Fragestellung: Genauigkeit der Approximation einer Funktion (= Interpoland) durch ihr Lagrangesches Interpolationspolynom zur Knotenmenge \mathcal{T} .

Spezielle Fragestellung: Asymptotisches Fehlerverhalten

$$\exists C \neq C(n): \|f - l_{\mathcal{T}}f\| \leq C T(n) \text{ für } n \rightarrow \infty.$$

Klassifikation (beste Schranken für $T(n)$):

$$\begin{aligned} \exists p > 0: \quad T(n) &\leq n^{-p} : \text{ALGCVG} \\ \exists 0 < q < 1: \quad T(n) &\leq q^n : \text{ALGCVG} \end{aligned}$$

Algebraische Konvergenz, Rate p ,
Exponentielle Konvergenz.

Beachte: Gleicher Begriff \leftrightarrow verschiedene Bedeutungen:

- Konvergenz einer Folge (von Iterierten $\mathbf{x}^{(k)} \rightarrow$ Abschn. 2.1) leg: iterationsverfahren
- Konvergenz einer Approximation (abhängig von einem Approximationsparameter)

Theorem 4.1.4 (Fehler der Polynominterpolation). Für $f \in C^{n+1}(I)$: $\forall t \in I$:

$$f(t) - l_{\mathcal{T}}(f)(t) = \int_0^1 \int_0^{\tau_1} \dots \int_0^{\tau_{n-1}} \int_0^{\tau_n} f^{(n+1)}(t_0 + \tau_1(t_1 - t_0) + \dots + \tau_n(t_n - t_{n-1}) + \tau(t - t_n)) d\tau d\tau_n \dots d\tau_1 \cdot \prod_{j=0}^n (t - t_j).$$

Beweis. Durch Induktion nach n , verwendet leg: ipolrec (4.1.12) und Hauptsatz der Integral- und Differentialrechnung RAN00 [38, Sect. 3.1]:

Bemerkung 160. thm: polintperr Thm. 4.1.4 gilt auch für Hermite-Interpolation.

Korollar 4.1.5. $f \in C^{n+1}(I)$: $\forall t \in I$: $\exists \tau_t \in [\min\{t, t_0, \dots, t_n\}, \max\{t, t_0, \dots, t_n\}]$:

$$f(t) - l_{\mathcal{T}}(f)(t) = \frac{f^{(n+1)}(\tau_t)}{(n+1)!} \cdot \prod_{j=0}^n (t - t_j).$$

Interpolationsfehlerabschätzungen erfordern Glattheit !

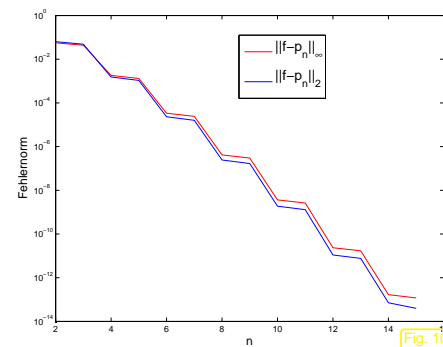
Beispiel 161 (Interpolationsfehler).

Interpolation von $f(t) = \sin t$, äquidistante Knoten in $I := [0, \pi]$: $\mathcal{T} := \{j\pi/n\}_{j=0}^n$: Interpolationspolynom $p := l_{\mathcal{T}}f \in \mathcal{P}_n$

$$\|f^{(k)}\|_{L^\infty(I)} \leq 1, \quad \forall k \in \mathbb{N}_0 \Rightarrow \|f - p\|_{L^\infty(I)} \leq \frac{1}{(1+n)!} \left\| \left(\cdot - 0 \right) \left(\cdot - \frac{\pi}{n} \right) \left(\cdot - \frac{2\pi}{n} \right) \dots \left(\cdot - \pi \right) \right\|_{L^\infty(I)} \leq \frac{1}{n+1} \left(\frac{\pi}{n} \right)^{n+1}.$$

➤ Gleichmäßige exponentielle Konvergenz der Interpolationspolynome
(Gilt für beliebige Knotenmengen \mathcal{T})

4.1
p. 385



MATLAB-Experiment:

Berechnung der Normen:

- L^∞ -Norm: Abtastung auf Gitter mit Maschenweite $\pi/1000$.
- L^2 -Norm: Numerische Quadratur (\rightarrow Abschnitt 4.4) mit Trapezregel auf Gitter mit Maschenweite $\pi/1000$. cha: numerische-quadratur

Beispiel 162 (Runge's Beispiel).

4.1
p. 386

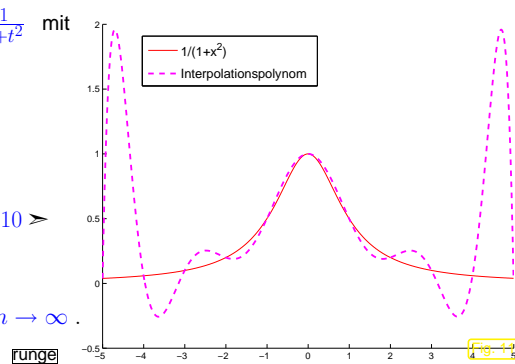
4.1
p. 387

4.1
p. 388

Polynominterpolation von $f(t) = \frac{1}{1+t^2}$ mit äquidistanten Knoten:

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

Plot $n = 10 \rightarrow$



Oszillationen von $\mathcal{I}_{\mathcal{T}}f$ am Intervallende

$$\|f - \mathcal{I}_{\mathcal{T}}f\|_{L^\infty([-5,5])} \rightarrow \infty \quad \text{für } n \rightarrow \infty.$$

Hier: $f(t) = \frac{1}{1+t^2} \rightarrow |f^{(n)}(t)| = 2^n n! \cdot O(|t|^{-2-n}).$

Schranke für Fehlerglied aus Cor. 4.1.5 $\rightarrow \infty$ für $n \rightarrow \infty$.

Bemerkung 163 (L^2 -Fehlerabschätzungen).

Cor. 4.1.5 \rightarrow Fehlerabschätzung für L^∞ -Norm

andere Normen ?

Aus Thm. 4.1.4 mit Cauchy-Schwarz-Ungleichung für Integrale:

$$\begin{aligned} \|f - \mathcal{I}_{\mathcal{T}}(f)\|_{L^2(I)}^2 &= \int_I \left| \int_0^1 \int_0^{\tau_1} \dots \int_0^{\tau_{n-1}} f^{(n+1)}(\dots) d\tau_n \dots d\tau_1 \cdot \prod_{j=0}^n (t - t_j) \right|^2 dt \\ &\leq \int_I |I|^{2n+2} \text{vol}(S_{n+1}) \int_{S_{n+1}} |f^{(n+1)}(\dots)|^2 d\tau \\ &\leq \int_I \frac{|I|^{2n+2}}{(n+1)!} \int_I \underbrace{\text{vol}(C_{t,\tau})}_{\leq 2^{(n-1)/2}/n!} |f^{(n+1)}(\tau)|^2 d\tau dt, \end{aligned}$$

wobei

$$S_{n+1} := \{\mathbf{x} \in \mathbb{R}^{n+1} : 0 \leq x_n \leq x_{n-1} \leq \dots \leq x_1 \leq 1\} \quad (\text{Einheitssimplex}), \\ C_{t,\tau} := \{\mathbf{x} \in S_{n+1} : t_0 + x_1(t_1 - t_0) + \dots + x_n(t_n - t_{n-1}) + x_{n+1}(t - t_n) = \tau\}.$$

$$\|f - \mathcal{I}_{\mathcal{T}}(f)\|_{L^2(I)} \leq \frac{2^{(n-1)/4} |I|^{n+1}}{\sqrt{(n+1)!} n!} \left(\int_I |f^{(n+1)}(\tau)|^2 d\tau \right)^{1/2}. \quad (4.1.13)$$

Beachte:

$$f \mapsto \|f^{(n)}\|_{L^2(I)} \quad \text{definiert eine Seminorm auf } C^{n+1}(I) \\ (\text{Sobolev-Seminorm, Mass für Glattheit einer Funktion})$$

4.1.3 Tschebyscheff-Interpolation

TSCHEBYSCHEFF

Knotenmenge: $\mathcal{T} := \{t_0 < t_1 < \dots < t_{n-1} < t_n\}, n \in \mathbb{N}, I := [t_0, t_n]$

Cor. 4.1.5 \rightarrow

$$\|f - p\|_{L^\infty(I)} \leq \frac{1}{(n+1)!} \|f^{(n+1)}\|_{L^\infty(I)} \|w\|_{L^\infty(I)}, \\ w(t) := (t - t_0) \dots (t - t_n).$$



Idee: Knoten t_0, \dots, t_n so, dass $\|w\|_{L^\infty(I)} \rightarrow \min!$

\leftrightarrow Finde $q \in \mathcal{P}_{n+1}$, führender Koeffizient = 1, $\|q\|_{L^\infty(I)} \rightarrow \min$.

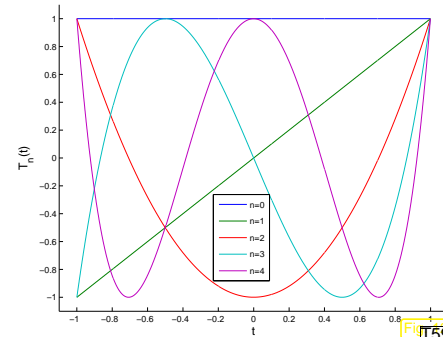
Wähle $t_0, \dots, t_n :=$ Nullstellen von q (Vorsicht: $t_j \in I$?)

Heuristik:

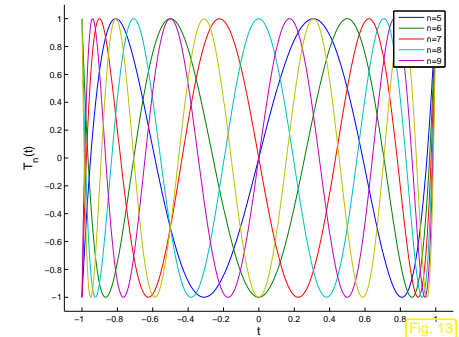
- q hat $n+1$ Nullstellen in I
- $|q(t_0)| = |q(t_n)| = \|q\|_{L^\infty(I)}$
- t^* Extremalstelle von $q \rightarrow |q(t^*)| = \|q\|_{L^\infty(I)}$

Definition 4.1.6 (Tschebyscheff-Polynome).

$$n. \text{ Tschebyscheff-Polynom} = T_n(t) := \cos(n \arccos t), \quad -1 \leq t \leq 1$$



Tschebyscheff-Polynome T_0, \dots, T_4



Tschebyscheff-Polynome T_5, \dots, T_9

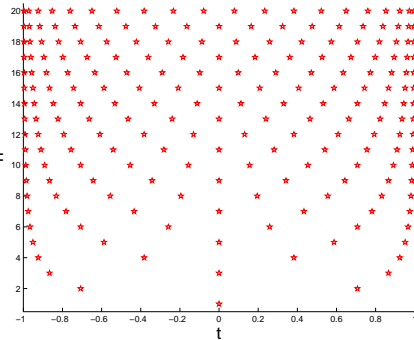
Nullstellen von T_n :

$$t_k = \cos\left(\frac{2k-1}{2n}\pi\right), \quad k = 1, \dots, n.$$

(4.1.14) \Leftarrow

Extrema (Alternanten) von T_n an

$$|T_n(\bar{t}_k)| = 1 \Leftrightarrow \exists k = 0, \dots, n: \bar{t}_k = \cos \frac{k\pi}{n}.$$



3-Term-Rekursion (\rightarrow Abschnitt 4.2) seg: dreitermrekursionen aus $\cos(n+1)x = 2\cos nx \cos x - \cos(n-1)x$ mit $\cos x = t$:

$$T_{n+1}(t) = 2tT_n(t) - T_{n-1}(t), \quad T_0 \equiv 1, \quad T_1(t) = t, \quad n \in \mathbb{N}. \quad (4.1.15)$$

eq:TPrek

$T_n \in \mathcal{P}_n$, führender Koeffizient 2^{n-1}

$\mathcal{P}_n = \text{Span}\{T_0, \dots, T_n\}$, $n \in \mathbb{N}_0$, T_n linear unabhängig (\rightarrow Basis)

Theorem 4.1.7 (Minimax Eigenschaft der Tschebyscheff-Polynome).

$$\|T_n\|_{L^\infty([-1,1])} = \inf\{\|p\|_{L^\infty([-1,1])} : p \in \mathcal{P}_n, p(t) = 2^{n-1}t^n + \dots\}, \quad \forall n \in \mathbb{N}.$$

DEH02
Zum Beweis siehe [12], Abschnitt 7.1.4.

Für $I = [-1, 1]$ \rightarrow „Optimale“ Interpolationsknoten $\mathcal{T} = \left\{\cos\left(\frac{2k+1}{2(n+1)}\pi\right), k = 0, \dots, n\right\}$,

$$w = 2^{-n}T_{n+1}, \quad \|w\|_{L^\infty(I)} = 2^{-n},$$

$$\|f - \mathcal{I}_{\mathcal{T}}(f)\|_{L^\infty([-1,1])} = \frac{2^{-n}}{(n+1)!} \|f^{(n+1)}\|_{L^\infty([-1,1])}.$$

SKALARG

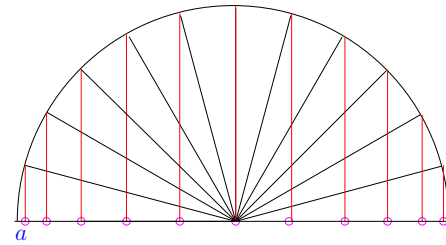
Ein Skalierungsargument: Intervalltransformation \leftrightarrow Transport von Funktionen

$$[-1, 1] \xrightarrow{\hat{t} \mapsto t := a + \frac{1}{2}(\hat{t}+1)(b-a)} [a, b] \leftrightarrow \hat{f}(\hat{t}) := f(t).$$

$$p \in \mathcal{P}_n \wedge p(t_j) = f(t_j) \Leftrightarrow \hat{p} \in \mathcal{P}_n \wedge \hat{p}(\hat{t}_j) = \hat{f}(\hat{t}_j).$$

Mit Transformationsformel für Integrale & $\frac{d^n \hat{f}}{d\hat{t}^n}(\hat{t}) = \left(\frac{1}{2}|I|\right)^n \frac{d^n f}{dt^n}(t)$:

$$\begin{aligned} \|f - \mathcal{I}_{\mathcal{T}}(f)\|_{L^\infty(I)} &= \|\hat{f} - \mathcal{I}_{\hat{\mathcal{T}}}(\hat{f})\|_{L^\infty([-1,1])} \leq \frac{2^{-n}}{(n+1)!} \left\| \frac{d^{n+1} \hat{f}}{d\hat{t}^{n+1}} \right\|_{L^\infty([-1,1])} \\ &\leq \frac{2^{-2n-1}}{(n+1)!} |I|^{n+1} \|f^{(n+1)}\|_{L^\infty(I)}. \end{aligned} \quad (4.1.16) \quad \text{eq:12intpr}$$



Für $I = [a, b]$: CHEBNODES Tschebyscheff-Knoten

$$t_k := a + \frac{1}{2}(b-a) \left(\cos\left(\frac{2k+1}{2(n+1)}\pi\right) + 1 \right), \quad k = 0, \dots, n. \quad (4.1.17) \quad \text{eq:CHEBNC}$$

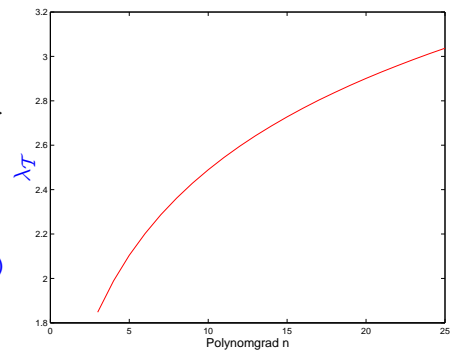
4.1
p. 393

Lebgue-Konstanten für Tschebyscheff-Knoten

Theorie CHB95 [6]:

$$\lambda_{\mathcal{T}} \sim \frac{2}{\pi} \log(1+n) + o(1),$$

$$\lambda_{\mathcal{T}} \leq \frac{2}{\pi} \log(1+n) + 1. \quad (4.1.18)$$



Beispiel 164 (Tschebyscheff-Interpolation).

Für $I = [a, b]$ seien $x_l := a + \frac{b-a}{N}l$, $l = 0, \dots, N$, $N = 1000$ und man approximiert

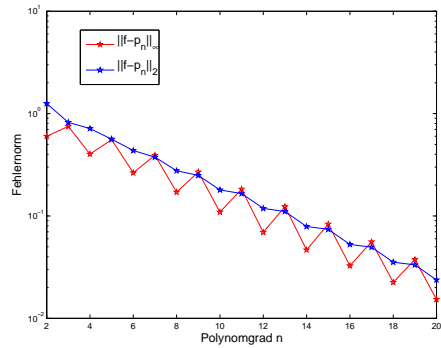
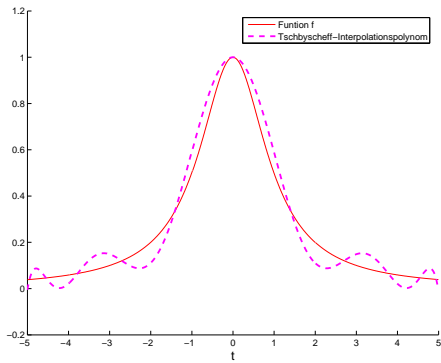
$$\|f - p\|_\infty \approx \max_{0 \leq l \leq N} |f(x_l) - p(x_l)|$$

$$\|f - p\|_2^2 \approx \frac{b-a}{2N} \sum_{0 \leq l < N} (|f(x_l) - p(x_l)|^2 + |f(x_{l+1}) - p(x_{l+1})|^2)$$

4.1
p. 394 $\textcircled{1}$ $f(t) = (1+t^2)^{-1} \rightarrow$ Bsp. 162, $I = [-5, 5]$, $n = 10$ im linken Plot. ex:runge

4.1
p. 395

4.1
p. 396



➤ Beobachtung: Exponentielle Konvergenz der Tschebyscheff-Interpolation

$$p_n \rightarrow f, \|f - I_n f\|_{L^2([-5,5])} \approx 0.8^n.$$

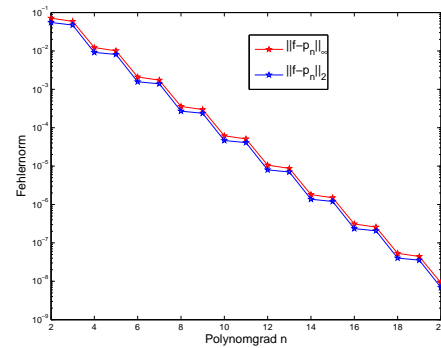
Nun: $I = [-1, 1]$

Weiterhin $f(t) = (1 + t^2)^{-1}$

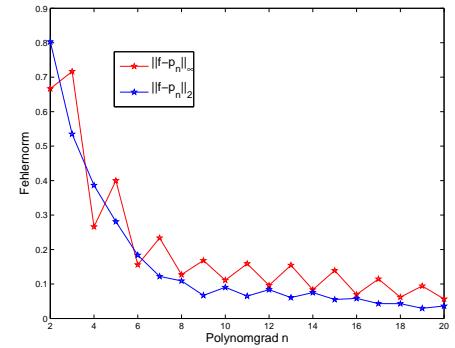
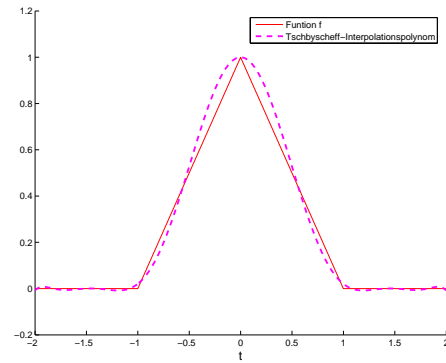
➤ Beobachtung:

Exponentielle Konvergenz:

$$\|f - I_n f\|_{L^2([-5,5])} \approx 0.42^n.$$



② $f(t) = \max\{1 - |t|, 0\}$, $I = [-2, 2]$, $n = 10$ im linken Plot.

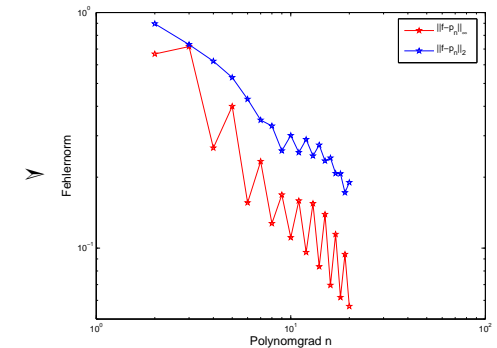


4.1
p. 397

Doppeltlogarithmischer Plot

➤ Beobachtung:

- keine exponentielle Konvergenz
- algebraische Konvergenz (?)

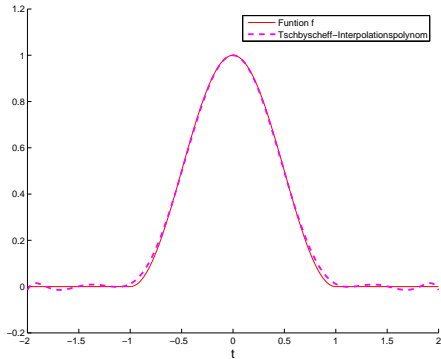


② $f(t) = 1 - \frac{1}{2}(1 + \cos \pi t)$ für $|t| < 1$, $f \equiv 0$ sonst, $I = [-2, 2]$, $n = 10$ im linken Plot.

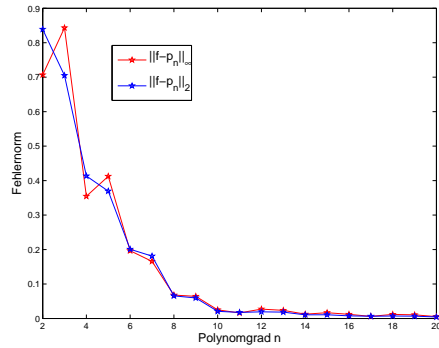
4.1
p. 398

4.1
p. 399

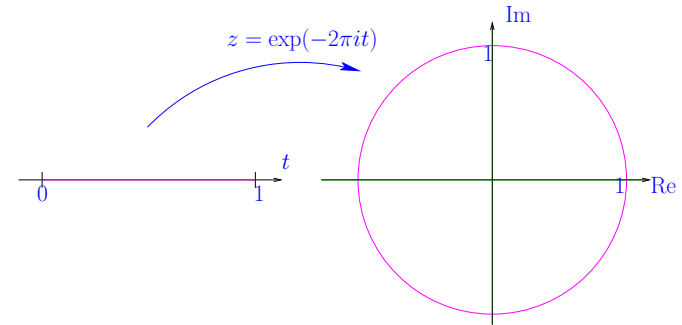
4.1
p. 400



➤ Beobachtung: Nur algebraische Konvergenz



◇



Definition 4.1.9 (Laurent-Polynome).

Der Raum \mathcal{P}_n^L der **Laurent-Polynome** vom Grad $\leq n$ über \mathbb{K} ist gegeben durch

$$\mathcal{P}_n^L := \begin{cases} \{z \mapsto \sum_{j=-m}^m \gamma_j z^j, \gamma_j \in \mathbb{K}\} & , \text{ falls } n = 2m, m \in \mathbb{N}, \\ \{z \mapsto \sum_{j=-m+1}^m \gamma_j z^j, \gamma_j \in \mathbb{K}\} & , \text{ falls } n = 2m - 1, m \in \mathbb{N}. \end{cases}$$

4.1.4 Trigonometrische Interpolation

4.1.4.1 Trigonometrische Polynome

Definition 4.1.8 (Trigonometrisches Polynom).

Der Raum \mathcal{P}_n^T der **trigonometrischen Polynome** vom Grad $\leq n$ ist gegeben durch

$$\mathcal{P}_n^T := \begin{cases} \{t \mapsto \sum_{j=-m}^m \gamma_j e^{-2\pi i j t}, \gamma_j \in \mathbb{C}\} & , \text{ falls } n = 2m, m \in \mathbb{N}, \\ \{t \mapsto \sum_{j=-m+1}^m \gamma_j e^{-2\pi i j t}, \gamma_j \in \mathbb{C}\} & , \text{ falls } n = 2m - 1, m \in \mathbb{N}. \end{cases}$$

Warum trigonometrische Polynome ?

$$p(t) := \sum_{j=-m+1}^m \gamma_j e^{-2\pi i j t} = \sum_{j=-m+1}^m \gamma_j z^j \hat{=} (\text{Laurent})\text{-Polynom auf } \mathbb{S}^1$$

Parametrisierung: $[0, 1[\xrightarrow[t \mapsto z]{z = \exp(-2\pi i t)} \mathbb{S}^1 := \{z \in \mathbb{C}: |z| = 1\}.$

4.1 Beachte: trigonometrische Polynome \subset 1-periodische Funktionen auf \mathbb{R} ($\leftrightarrow f(t) = f(t+1)$)
p. 401

4.1
p. 403

TIP

Trigonometrische Polynominterpolationsaufgabe

Zur äquidistanten Knoten $\mathcal{T} := \{j/n, j = 0, \dots, n-1\}$ und Stützwerten $y_0, \dots, y_{n-1} \in \mathbb{C}$ finde $p \in \mathcal{P}_{n-1}^T$, so dass

$$p(j/n) = y_j, \quad j = 0, \dots, n-1.$$

Theorem 4.1.10. Die trigonometrische Polynominterpolationsaufgabe ist eindeutig lösbar $\forall y_k$.

Beweis. siehe unten. □

▶ trigonometrische Polynominterpolationsaufgabe \Leftrightarrow lineare Abbildung $T_n : \mathbb{C}^n \mapsto \mathcal{P}_{n-1}^T$.

Beachte: 1-periodische Funktionen auf $\mathbb{R} \leftrightarrow$ Funktionen $\mathbb{S}^1 \mapsto \mathbb{K}$

$$f(t) = f(t+1) \quad \forall t \in \mathbb{R} \quad \leftrightarrow \quad g(z) = g(e^{2\pi i t}) := f(t) \quad \forall z \in \mathbb{S}^1.$$

4.1
p. 402

4.1
p. 404

► Trigonometrische Polynominterpolationsaufgabe \Leftrightarrow Äquidistante Polynominterpolation auf \mathbb{S}^1

\mathbb{S}^1

► Zu Knotenmenge $\mathcal{T}_n = \{\omega_n^j\}_{j=0}^{n-1}$ (komplexe Einheitswurzeln, Abschnitt 3.7.1) und Stützwerten y_0, \dots, y_{n-1} finde $p \in \mathcal{P}_{n-1}^L$ (komplexe Koeffizienten), so dass

$$p(\omega_n^k) = y_k, \quad k = 0, \dots, n-1.$$

► $\mathcal{T}_n \leftrightarrow$ Lineare Abbildung $\mathbb{C}^{n+1} \mapsto \mathcal{P}_n^L$.

Umformulierung der trigonometrischen Polynominterpolationsaufgabe:

❶ n gerade, $n = 2m, m \in \mathbb{N}$, $\Rightarrow p(t) = \sum_{j=-m+1}^m \gamma_j e^{-2\pi i j t} \in \mathcal{P}_{n-1}^L$

$$p(\frac{k}{n}) = y_k, \quad k = 0, \dots, n-1.$$

\Updownarrow

$$\sum_{j=-m+1}^m \gamma_j \omega_n^{jk} = \sum_{j=0}^m \gamma_j \omega_n^{jk} + \sum_{j=1}^{m-1} \gamma_{-j} \omega_n^{(n-j)k} = \sum_{j=0}^m \gamma_j \omega_n^{jk} + \sum_{j=m+1}^{n-1} \gamma_{j-n} \omega_n^{jk} = y_k,$$

$\omega_n \hat{=}$ komplexe Einheitswurzeln, Abschnitt 3.7.1.

❷ n ungerade, $n = 2m+1, m \in \mathbb{N}$, $\Rightarrow p(t) = \sum_{j=-m}^m \gamma_j e^{-2\pi i j t} \in \mathcal{P}_{n-1}^L$

$$p(\frac{k}{n}) = y_k, \quad k = 0, \dots, n-1.$$

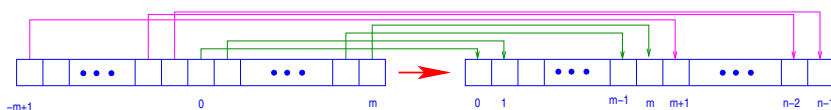
\Updownarrow

$$\sum_{j=-m}^m \gamma_j \omega_n^{jk} = \sum_{j=0}^m \gamma_j \omega_n^{jk} + \sum_{j=1}^m \gamma_{-j} \omega_n^{(n-j)k} = \sum_{j=0}^m \gamma_j \omega_n^{jk} + \sum_{j=m+1}^{n-1} \gamma_{j-n} \omega_n^{jk} = y_k.$$

$$\boxed{} \hat{=} \mathbf{F}_n \mathbf{g} = \mathbf{y}, \quad \mathbf{g} = (\gamma_0 \cdots \gamma_m \gamma_{m+1-n} \cdots \gamma_{n-1}), \quad \mathbf{y} = (y_0 \cdots y_{n-1}).$$

Fourier-Matrix (3.7.3)

► γ_j durch inverse diskrete Fouriertransformation (\rightarrow Def. 3.7.1) + Umsortieren der Koeffizienten:



Aufgabe: Auswertung eines trigonometrischen Polynoms $p \in \mathcal{P}_{n-1}^T$ (\rightarrow Def. 4.1.8) an äquidistanten Stellen $\frac{l}{N}, l = 0, \dots, N-1, N > n$, in $[0, 1]$:

Für $n = 2m$:

$$p(\frac{l}{N}) = \sum_{j=-m+1}^m \gamma_j \omega_N^{lj} = \sum_{j=0}^m \gamma_j \omega_N^{lj} + \sum_{j=1}^{m-1} \gamma_{-j} \omega_N^{(N-j)l} = \sum_{j=0}^{N-1} \tilde{\gamma}_j \omega_N^{lj},$$

$$\text{mit } \tilde{\gamma}_j := \begin{cases} \gamma_j & , \text{ falls } 0 \leq j \leq m-1, \\ \gamma_{j-N} & , \text{ falls } N-m+1 \leq j \leq N-1, \\ 0 & \text{sonst.} \end{cases}$$

MATLAB-CODE Auswertung eines

function v= evaliptrig(y,N)

N = length(y);

if (mod(n,2) ~= 0), error; end

c = ifft(y);

c = [c(1:n/2), zeros(1,N-n), ...

c(n/2+1:n)];

v = fft(c);

trigonometrischen Interpolationspolynom

► diskrete Fouriertransformation \rightarrow Def. 3.7.1

der Länge N

Koeffizienten in \mathbf{c}

umsortiert gespeichert !

4.1.4.2 Trigonometrische Interpolation 1-periodischer Funktionen:

Linearer trigonometrischer Interpolationsoperator: für 1-periodisches $f: \mathbb{R} \mapsto \mathbb{C}$

4.1
p. 405

$$\mathcal{T}_n(f) := p \in \mathcal{P}_{n-1}^T \quad \text{mit} \quad p(\frac{j}{n}) = y_j, \quad j = 0, \dots, n-1.$$

4.1
p. 407

Beispiel 165 (Trigonometrische Interpolation).

#1 Stufenfunktion: $f(t) = 0$ für $|t - \frac{1}{2}| > \frac{1}{4}$, $f(t) = 1$ für $|t - \frac{1}{2}| \leq \frac{1}{4}$

#2 Glatte periodische Funktion: $f(t) = \frac{1}{\sqrt{1 + \frac{1}{2} \sin(2\pi t)}}$.

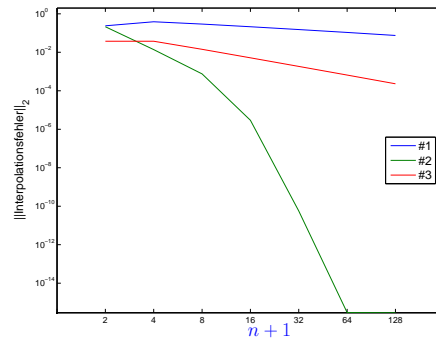
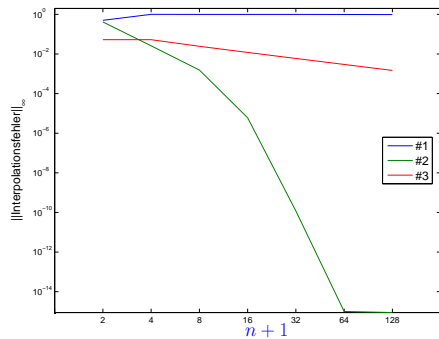
#3 Keilfunktion: $f(t) = |t - \frac{1}{2}|$

Berechnung der Normen der Interpolationsfehler:

```
t = 0:1/n:1; y = feval(f,t(1:end-1));
v = real(evaliptrig(y,4096));
v = [v,v(1)]; fv = feval(f,0:1/4096:1);
d = abs(fv-v); linf = max(d); l2 = 1/64*norm(d(1:end-1));
```

4.1
p. 406

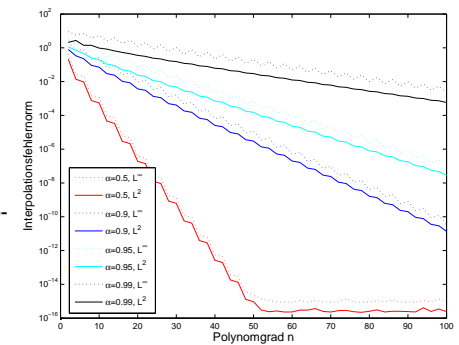
4.1
p. 408



Beobachtung: Fall #1, #3: Algebraische Konvergenz
Fall #2: Exponentielle Konvergenz

$$f(t) = \frac{1}{\sqrt{1 - \alpha \sin(2\pi t)}} \text{ auf } I = [0, 1].$$

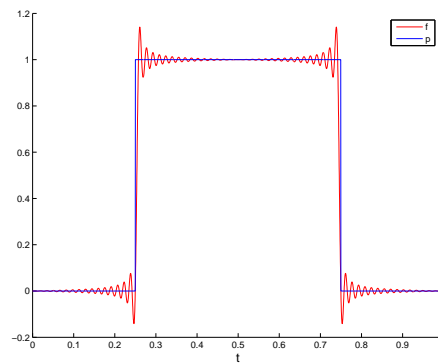
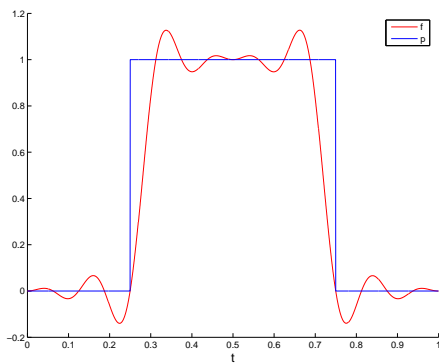
Approximative Messung der Normen: "oversampling" in 4096 Punkten



➤ Beobachtet: Exponentielle Konvergenz im Grad n , je schneller, je kleiner α

Analyse benutze: Trigonometrische Polynome = Partialsummen von **Fourierreihen**

Bekannt aus → Analysis, → Mathematische Methoden der Physik:



➤ $n = 16$ **Gibbsches Phänomen:** Überschwinger an Unstetigkeitsstellen $n = 128$

Beispiel 166 (Trigonometrische Interpolation analytischer Funktionen).

4.1
p. 409

Theorem 4.1.11 (L^2 -Konvergenz der Fourierreihe). Jede *quadratintegrierbare Funktion* $f \in L^2(]0, 1[) := \{f :]0, 1[\rightarrow \mathbb{C} : \|f\|_{L^2(]0, 1[)} < \infty\}$ ist der $L^2(]0, 1[)$ -Grenzwert ihrer **Fourierreihe**

$$f(t) = \sum_{k=-\infty}^{\infty} \hat{f}(k) e^{2\pi i k t} \text{ in } L^2(]0, 1[),$$

mit **Fourierkoeffizienten**

$$\hat{f}(k) = \int_0^1 f(t) e^{-2\pi i k t} dt, \quad k \in \mathbb{Z}.$$

Isometrieeigenschaft: $\sum_{k=-\infty}^{\infty} |\hat{f}(k)|^2 = \|f\|_{L^2(]0, 1[)}^2$ (4.1.19) [eq:Fouriso](#)

4.1
p. 410

4.1
p. 411

4.1
p. 412

□ **Lemma 4.1.12** (Ableitung und Fourierkoeffizienten).

$$f \in L^1([0, 1]) \text{ \& } f' \in L^1([0, 1]) \Rightarrow \widehat{f}'(k) = 2\pi i k \widehat{f}(k), \quad k \in \mathbb{Z}.$$

$$\blacktriangleright \left\| f^{(n)} \right\|_{L^2([0, 1])}^2 = (2\pi)^{2n} \sum_{k=-\infty}^{\infty} k^{2n} |\widehat{f}(k)|^2. \quad (4.1.20) \quad \text{eq:FN}$$

$$\text{Falls } |\widehat{f^{(n)}}(k)| \leq \int_0^1 |f^{(n)}(t)| dt < \infty \Rightarrow \widehat{f}(k) = O(k^{-n}) \text{ f\"ur } |k| \rightarrow \infty.$$

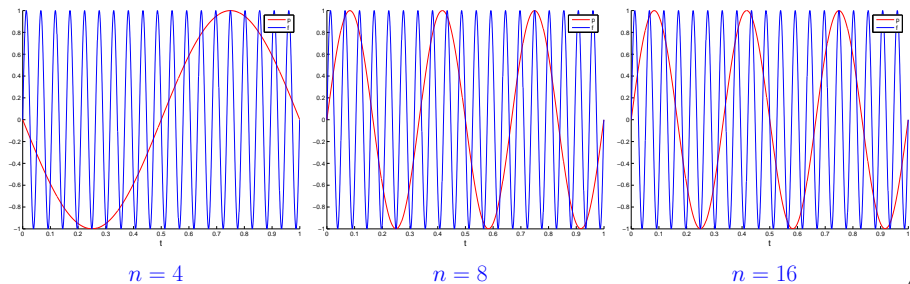
Je glatter eine 1-periodische Funktion, desto schneller der Abfall ihrer Fourierkoeffizienten

Bemerkung 167 (Aliasing).

$$\mathcal{T} = \left\{ \frac{j}{n} \right\}_{j=0}^{n-1} \blacktriangleright e^{\frac{2\pi i N t}{T}} = e^{\frac{2\pi i (N-n)t}{T}} \blacktriangleright \mathcal{T}_n(e^{2\pi i N \cdot}) = \mathcal{T}_n(e^{2\pi i (N \bmod n) \cdot}).$$

➤ Die trigonometrische Interpolation von $t \rightarrow e^{2\pi i N t}$ und $t \rightarrow e^{2\pi i (N \bmod n) t}$ mit Polynomgrad n liefert das gleiche trigonometrische Interpolationspolynom! \rightarrow **Aliasing**

Beispiel für $f(t) = \sin(2\pi \cdot 19t)$ ➤ $N = 38$:



[Linearität der trigonometrischen Interpolation]

Für $f \in C([0, 1])$ 1-periodisch, $n = 2m$:

$$f(t) = \sum_{j=-\infty}^{\infty} \widehat{f}(j) e^{2\pi i j t} \blacktriangleright \mathcal{T}_n(f)(t) = \sum_{j=-m+1}^m \gamma_j e^{2\pi i j t}, \quad \gamma_j = \sum_{l=-\infty}^{\infty} \widehat{f}(j + ln).$$

➤ Fourier-Koeffizienten der Interpolationsfehlerfunktion $e = f - \mathcal{T}_n(f)$:

$$\widehat{e}(j) = \begin{cases} -\sum_{|l|=1}^{\infty} \widehat{f}(j + ln) & , \text{ falls } -m+1 \leq j \leq m, \\ \widehat{f}(j) & , \text{ falls } j \leq -m \vee j > m. \end{cases}$$

$$\text{eq:Fourier} \quad (4.1.19) \quad \text{eq:fourierest} \quad \blacktriangleright \left\| f - \mathcal{T}_n(f) \right\|_{L^2([0, 1])}^2 \leq \left| \sum_{|l|=1}^{\infty} \widehat{f}(j + ln) \right|^2 + \sum_{|j| \geq m} |\widehat{f}(j)|^2. \quad (4.1.21) \quad \text{eq:fourie}$$

➤ Abschätzbar bei bekanntem Abfall der $\widehat{f}(j) \leftrightarrow$ Glattheit von f , vgl. (4.1.20) eq:FN

□ **Theorem 4.1.13** (Fehlerabschätzung für trigonometrische Interpolation). Für $k \in \mathbb{N}$:

$$f^{(k)} \in L^2([0, 1]) \Rightarrow \left\| f - \mathcal{T}_n f \right\|_{L^2([0, 1])} \leq \sqrt{1 + c_k} n^{-k} \left\| f^{(k)} \right\|_{L^2([0, 1])},$$

mit $c_k = 2 \sum_{l=1}^{\infty} (2l - 1)^{-2k}$.

4.1.4.3 Trigonometrische Interpolation analytischer Funktionen

Was passiert wenn $f: \mathbb{R} \mapsto \mathbb{C}$ 1-periodisch und $f \in C^\infty(\mathbb{R})$? (Aussage von Thm. 4.1.13) thm:triperrest

4.1
p. 413

4.1
p. 415

□ **Definition 4.1.14** (Analytische Funktion).

Eine Funktion $f: D \subset \mathbb{C} \mapsto \mathbb{C}$ heisst **analytisch** (\leftrightarrow holomorph) (in D), falls (im Sinne lokal gleichmässiger Konvergenz)

$$\forall z \in D: \exists r > 0, (\alpha_k)_{k=0}^{\infty} \in \mathbb{C}^{\mathbb{N}_0}: |w - z| < r \Rightarrow f(w) = \sum_{k=0}^{\infty} \alpha_k (w - z)^k.$$

Analytizität \leftrightarrow lokale Darstellbarkeit durch Potenzreihe

Noch glatter:

Ganze Funktionen \leftrightarrow global konvergente Potenzreihe

➤ analytisch in \mathbb{C} , z.B. Polynome, $\exp z, \sin z, \cos z, \dots$ + Kompositionen.

„Eine Funktion sucht sich ihr Analytizitätsgebiet selbst!“

Beispiel 168 (Analytizitätsgebiete von Funktionen).

• $f(t) = (1 + t^2)^{-1}$ (\rightarrow Bsp. 162): Interpretation als Funktion in $z \in \mathbb{C}$:

$$f(z) = \frac{1}{1 + z^2} = \frac{1}{(z - i)(z + i)} \Rightarrow f \text{ analytisch in } \mathbb{C} \setminus \{\pm i\}.$$

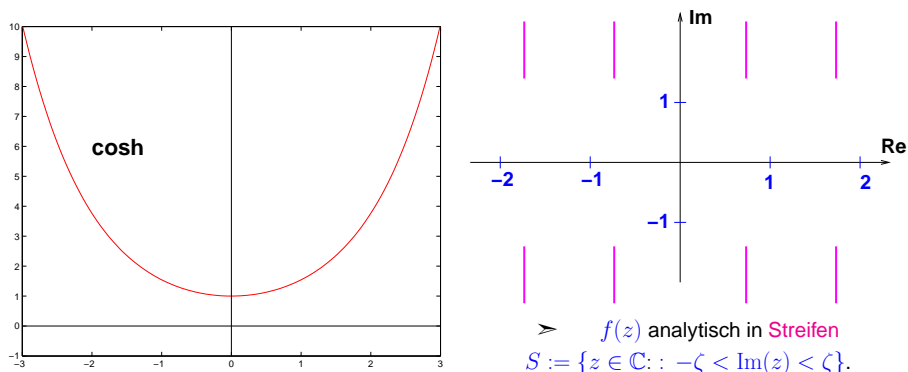
4.1
p. 414

4.1
p. 416

- $f(t) = \frac{1}{\sqrt{1 + \alpha \sin(2\pi t)}}$, $\alpha > 0$ (\rightarrow Bsp. 165): ^{ex:trigip}Fortsetzung in \mathbb{C} , $\sqrt{}$ = Hauptzweig, analytisch auf $\mathbb{C} \setminus \mathbb{R}_0^-$.

$$1 + \alpha \sin(2\pi z) \notin \mathbb{R}_0^- \Leftrightarrow \sin(2\pi z) = \sin(2\pi x) \cosh(2\pi y) + i \cos(2\pi x) \sinh(2\pi y) \notin [-\infty, -1 - \frac{1}{\alpha}]$$

► Analytizitätsgebiet: $\mathbb{C} \setminus \bigcup_{k \in \mathbb{Z}} (\frac{k}{2} + \frac{1}{4} + i(\mathbb{R} \setminus -\zeta, \zeta))$, $\zeta \in \mathbb{R}^+$, $\cosh(2\pi\zeta) = 1 + \frac{1}{\alpha}$.



► $f(z)$ analytisch in Streifen $S := \{z \in \mathbb{C} : -\zeta < \text{Im}(z) < \zeta\}$.

4.1
p. 417

Theorem 4.1.15 (Exponentielle Konvergenz der trigonometrischen Interpolation).

Falls $f : \mathbb{R} \mapsto \mathbb{C}$ 1-periodisch und ^{def:analytic}analytisch (\rightarrow Def. 4.1.14) im Streifen

$$\bar{S} := \{z \in \mathbb{C} : -\eta \leq \text{Im } z \leq \eta\}, \quad \eta > 0,$$

dann existiert ein nur von η abhängiges $C_\eta > 0$, so dass

$$\|f - T_n f\|_X \leq C_\eta e^{-\pi\eta n} \|f\|_{L^\infty(\bar{S})}, \quad n \in \mathbb{N}, \quad X = L^2([0, 1]), L^\infty([0, 1]).$$

Beweis auf der Grundlage des folgenden Satzes, der aus ^{eq:fourierest}(4.1.21) folgt:

Theorem 4.1.16. Falls $f \in L^2([0, 1])$ mit

$$\exists M > 0, \rho \in [0, 1]: |\hat{f}(k)| \leq M \rho^{|k|} \Rightarrow \begin{aligned} \|f - T_n(f)\|_{L^2([0, 1])} &\leq M \frac{\rho^{n/2}}{1 - \rho^n} \frac{2\sqrt{2}}{\sqrt{1 - \rho^2}}, \\ \|f - T_n(f)\|_{L^\infty([0, 1])} &\leq 4M \frac{\rho^{n/2}}{1 - \rho}. \end{aligned}$$

Proof. (Beweis von Theorem ^{thm:expcvg}4.1.15)

4.1
p. 418

Sei $f : \mathbb{R} \mapsto \mathbb{C}$ 1-periodisch, analytisch fortsetzbar in abgeschlossenen Streifen $S := \{z \in \mathbb{C} : -\eta^- \leq \text{Im}\{z\} \leq \eta^+\}$

► $g_r(t) := f(t + ir)$, 1-periodisch und glatt für $-\eta^- \leq r \leq \eta^+$

^{analstrip}

$$\hat{g}_r(k) = \int_0^1 f(t + ir) e^{-2\pi i k t} dt = \int_0^1 f(t) e^{2\pi i k (t - ir)} dt = e^{-2\pi r k} \hat{f}(k).$$

$$|\hat{f}(k)| \leq e^{-2\pi r k} \underbrace{\max\{|f(t + ir)|, 0 \leq t \leq 1\}}_{\|g_r\|_{L^\infty([0, 1])}} \quad \forall k \in \mathbb{Z}, \quad \forall -\eta^- \leq r \leq \eta^+.$$

Fig. 14

4.1
p. 419

$|\hat{f}(-k)| \leq \exp(-2\pi\eta^- k) \|f\|_{L^\infty(\bar{S})} \quad \forall k \in \mathbb{N},$
 $|\hat{f}(k)| \leq \exp(-2\pi\eta^+ k) \|f\|_{L^\infty(\bar{S})} \quad \forall k \in \mathbb{N}.$

Exponentieller Abfall $\hat{f}(k)$

^{thm:expcvg1}
4.1.16

(4.1.22) ^{eq:TIPstr}

Beispiel 169 (Exponentielle Konvergenz der trigonometrischen Interpolation).

$$f(t) = \frac{1}{\sqrt{1 + \alpha \sin(2\pi t)}}, \quad \alpha > 0 \quad (\rightarrow \text{Bsp. 165 \& Bsp. 168})$$

f analytisch in Streifen $S := \{x \in \mathbb{C} : |\text{Im } z| < \zeta, \cosh(2\pi\zeta) = 1 + \frac{1}{\alpha}\}$

Vorsicht: f unbeschränkt auf S ! ► betrachte schmäleren Streifen.

4.1
p. 420

4.1.4.4 Trigonometrische Interpolation und Tschebyscheff-Interpolation

$p \in \mathcal{P}_n \hat{=}$ Tschebyscheff-Interpolationspolynom auf $[-1, 1]$ zu $f : [-1, 1] \mapsto \mathbb{C}$ (\rightarrow Abschnitt 4.1.3)

$p(t_k) = f(t_k)$ für Tschebyscheff-Knoten, siehe (4.1.17), $t_k := \cos\left(\frac{2k+1}{2(n+1)}\pi\right)$, $k = 0, \dots, n$.

$$\begin{aligned} f : [-1, 1] &\mapsto \mathbb{C} &\leftrightarrow & g(s) := f(\cos 2\pi s), \\ p : [-1, 1] &\mapsto \mathbb{C} &\leftrightarrow & q(s) := p(\cos 2\pi s), \end{aligned} \quad \begin{matrix} 1\text{-periodisch,} \\ \text{achsensymmetrisch.} \end{matrix}$$

$$p(t_k) = f(t_k) \Leftrightarrow q\left(\frac{2k+1}{4(n+1)}\right) = g\left(\frac{2k+1}{4(n+1)}\right). \tag{4.1.23}$$

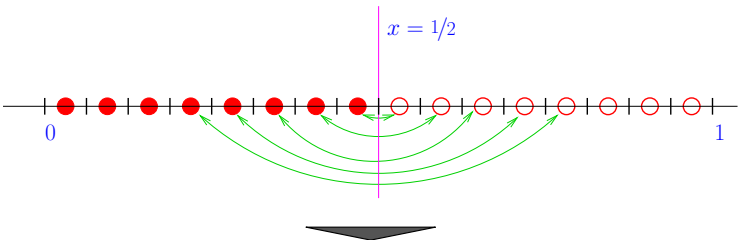
! Wir zeigen: $q(s)$ ist trigonometrisches Polynom: $q \in \mathcal{P}_{2n}^T$

$\mathcal{P}_n = \text{Span}\{T_0, \dots, T_n\}$, T_j = Tschebyscheff-Polynome \rightarrow Def. 4.1.6

$$p(t) = \gamma_0 T_0(t) + \gamma_1 T_1(t) + \dots + \gamma_n T_n(t), \quad \gamma_j \in \mathbb{C}.$$

$$q(s) = p(\cos 2\pi s) = \sum_{j=0}^n \gamma_j \cos(2\pi j s) = \sum_{j=-n}^n \frac{1}{2} \gamma_{|j|} e^{-2\pi i j s}, \quad s \in \mathbb{R}. \tag{4.1.24}$$

$$p(t_k) = f(t_k) \Leftrightarrow \sum_{j=0}^n \gamma_j \cos\left(2\pi j \cdot \frac{2k+1}{4(n+1)}\right) = g\left(\frac{2k+1}{4(n+1)}\right), \quad k = 0, \dots, n.$$



$$\begin{aligned} \tilde{g}(s) &:= g\left(s + \frac{1}{4(n+1)}\right), \\ \tilde{q}(s) &:= q\left(s + \frac{1}{4(n+1)}\right) \end{aligned} \Rightarrow \tilde{q}\left(\frac{k}{2(n+1)}\right) = \tilde{g}\left(\frac{k}{2n+2}\right), \quad k = 0, \dots, 2n+1.$$

$$\sum_{j=-n}^n \tilde{\gamma}_j \exp\left(2\pi i j \cdot \frac{k}{2(n+1)}\right) = \tilde{g}\left(\frac{k}{2n+2}\right), \quad k = 0, \dots, 2n+1,$$

$$\in \mathcal{P}_{2n}^T$$

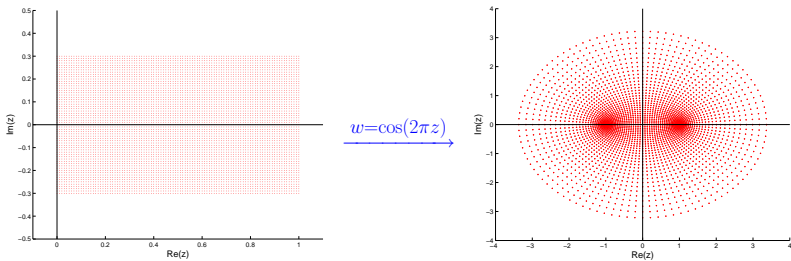
mit
$$\tilde{\gamma}_j := \begin{cases} \frac{1}{2} \exp\left(2\pi i \frac{j}{4(n+1)}\right) \gamma_{|j|}, & \text{falls } j \neq 0, \\ \gamma_0, & \text{für } j = 0. \end{cases}$$

► Transformiertes Tschebyscheff-Interpolationspolynom $p(\cos 2\pi t)$ = trigonometrisches Interpolationspolynom zu \tilde{g} ! $\Rightarrow \|f - p\|_{L^\infty([-1,1])} = \|\tilde{g} - \tilde{q}\|_{L^\infty([0,1])}$

► Übertragung der L^∞ Fehlerabschätzungen von Thm. 4.1.15

Wo muss f analytisch sein, damit $g(w) = f(\cos 2\pi z)$ analytisch in Streifen $S := \{w \in \mathbb{C} : -\zeta^- < \text{Im}\{z\} < \zeta^+, \zeta^-, \zeta^+ > 0\}$?

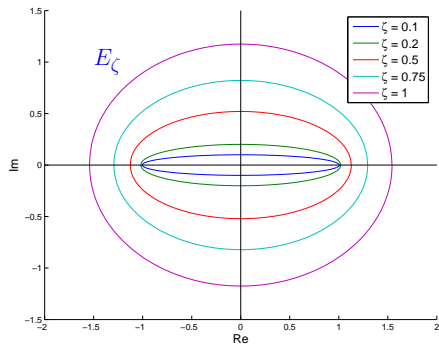
\Updownarrow
Gestalt von $\cos(S)$?



Mit $\cos(x + iy) = \cos x \cosh y - i \sin x \sinh y$: $\cos(S)$ = Inneres der Ellipse

$$E_\zeta := \{z = \cos \varphi \cosh(2\pi\zeta) - i \sin \varphi \sinh(2\pi\zeta) : 0 \leq \varphi \leq 2\pi\}, \quad \zeta := \max\{\zeta^-, \zeta^+\}.$$

E_ζ : horizontale Halbachse = $\cosh(2\pi\zeta)$, vertikale Halbachse = $\sinh(2\pi\zeta)$.



Analyse: Tschebyscheff-Interpolation von analytischen Funktionen

1. Transformation auf $[-1, 1]$
2. Analytizitätsgebiet $A \subset \mathbb{C}$ der transformierten Funktion
3. Maximale Ellipse $E_\zeta \subset A$ + Maximum auf Ellipse
4. Konvergenzabschätzung aus Thm. 4.1.15 thm:expcvq

4.1.5 Approximation durch Polynome

4.1.5.1 Bestapproximation

☞ Entscheidendes theoretisches Hilfsmittel:

\mathcal{T} = Knotenmenge in $I \subset \mathbb{R}$, $l_{\mathcal{T}} : C(I) \mapsto \mathcal{P}_n$ Polynominterpolationsoperator \rightarrow interp:pointpop (4.1.6)

$$l_{\mathcal{T}}(p) = p \quad \forall p \in \mathcal{P}_n \Rightarrow \|f - l_{\mathcal{T}}(f)\| = \|(f - p) - l_{\mathcal{T}}(f - p)\| \leq (1 + \|l_{\mathcal{T}}\|) \|f - p\|.$$

$$\blacktriangleright \|f - l_{\mathcal{T}}(f)\| \leq (1 + \|l_{\mathcal{T}}\|) \underbrace{\inf_{p \in \mathcal{P}_n} \|f - p\|}_{\text{Bestapproximationsfehler}}.$$

Speziell für $\|\cdot\| = \|\cdot\|_{L^\infty(I)}$:

$$\|f - l_{\mathcal{T}}(f)\|_{L^\infty(I)} \leq (1 + \lambda_{\mathcal{T}}) \inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty(I)} \quad \text{mit Lebesgue-Konstante } \lambda_{\mathcal{T}}. \quad (4.1.25) \quad \text{eq:lcrole}$$

Resultat von Jackson DAV75 [9, Thm. 13.3.7]:

⊗ **Theorem 4.1.17** (L^∞ -Bestapproximation durch Polynome).

$$f \in C^r([-1, 1]), r \geq 1: \inf_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty([-1, 1])} \leq (1 + \pi^2/2)^r \frac{(n-r)!}{n!} \|f^{(r)}\|_{L^\infty([-1, 1])}.$$

Thm. 4.1.17 + thm:polbestapprox (4.1.25) + Abschätzung $\lambda_{\mathcal{T}} \leq \frac{2}{\pi} \log(1+n) + 1$ RIV84b [39, Thm. 1.2]

➤ L^∞ -Fehlerabschätzung für Tschebyscheff-Interpolation (\mathcal{T} = Tschebyscheff-Knoten)

$$f \in C^r([-1, 1]), r \geq 1: \|f - l_{\mathcal{T}}(f)\|_{L^\infty([-1, 1])} \leq C(r) \left(1 + \frac{1}{\pi} \log n\right) \frac{(n-r)!}{n!} \|f^{(r)}\|_{L^\infty([-1, 1])},$$

mit $C(r) := 2(1 + \pi^2/2)^r$.

☞ Übertragung auf beliebige Intervalle mit Skalierungsargument.

4.1.5.2 Polynomiale Least-Squares Approximation

Gegeben: Knotenmenge $\mathcal{M} := \{\nu_0 < \nu_1 < \dots < \nu_m\}$ ($I := [\nu_0, \nu_m]$), $y_j \in \mathbb{C}$, $j = 0, \dots, m$

Gesucht: $p_n \in \mathcal{P}_n$ ($n < m$) mit
$$\sum_{j=0}^m |p(\nu_j) - y_j|^2 \rightarrow \min$$

➡ Lineares Ausgleichsproblem \rightarrow Abschnitt sec:numer-line-ausgl 3.5

4.1
p. 425

Abstrakt: Inneres Produkt auf \mathcal{P}_m :

$$\langle p, q \rangle := \sum_{j=0}^m p(\nu_j) \bar{q}(\nu_j)$$

➤ Norm auf \mathcal{P}_m : $\|p\|_{\mathcal{T}} := \sqrt{\langle p, p \rangle}$

Ausgleichsproblem: $p \in \mathcal{P}_n: \|p - y\|_{\mathcal{T}}^2 \rightarrow \min$,

wobei etwa $y := l_m(\mathbf{y})$.

4.1
p. 427

⊗ **Theorem 4.1.18.** Sei $(V, (\cdot, \cdot))$ ein Innenproduktraum (mit Norm $\|\cdot\|^2 := (\cdot, \cdot)$) mit Orthonormalbasis (ONB) $\{b_1, \dots, b_m\}$, $m := \dim V \in \mathbb{N}$, und $W_n := \text{Span}\{b_1, \dots, b_n\}$, $n \leq m$:

$$\forall y \in V: x^{(n)} := \arg \min_{v \in W_n} \|y - v\| \Leftrightarrow x^{(n)} = \sum_{j=1}^n (y, b_j) b_j.$$

Beweis: Benutze $v = \sum_{j=1}^m \xi_j b_j \Rightarrow \|v\|^2 = \sum_{j=1}^m |\xi_j|^2$
oder Normalengleichungen \rightarrow Abschnitt sec:normalengleichungen 3.5.2

4.1
p. 426

4.1
p. 428

Ziel: Finde $\langle \cdot, \cdot \rangle$ -ONB $\{r_0, r_1, \dots, r_m\}$ von \mathcal{P}_m mit

$$\text{Span}\{r_0, \dots, r_n\} = \mathcal{P}_n, \quad 0 \leq n \leq m. \quad (4.1.26) \quad \text{[p. 1]}$$



Idee: **Gram-Schmidt-Orthonormalisierung** wie im **Lanczos-Prozess**:
Gegeben ONB $\{r_0, \dots, r_n\}$, $r_j \in \mathcal{P}_j$, orthogonaliere $tr_n(t)$ gegen $\{r_0, \dots, r_n\}$.

Krylov-Connection: $C(I) \mapsto C(I)$, $f(t) \mapsto tf(t)$ ist $\langle \cdot, \cdot \rangle$ -selbstadjungierter Operator

Theorem 4.1.19 (Diskrete Orthogonalpolynome). Definiere $r_1 := 0$, $r_0 \equiv 1$ und

$$r_{k+1}(t) = (t - \alpha_{k+1})r_k(t) - \beta_k r_{k-1}(t), \quad k = 0, 1, \dots, m-1,$$

$$\text{mit } \alpha_{k+1} := \frac{\langle tr_k, r_k \rangle}{\|r_k\|_{\mathcal{T}}^2}, \quad \beta_k := \frac{\|r_k\|_{\mathcal{T}}^2}{\|r_{k-1}\|_{\mathcal{T}}^2}.$$

Dann ist $\{r_0, \dots, r_m\}$ eine $\langle \cdot, \cdot \rangle$ -orthogonale Basis von \mathcal{P}_m und erfüllt (4.1.26). Die r_j heissen (diskrete) **Orthogonalpolynome** ^{GGG05}

Beweis: (durch Induktion)

$$\begin{aligned} \langle r_{k+1}, r_j \rangle &= \langle (t - \alpha_{k+1})r_k - \beta_k r_{k-1}, r_j \rangle \\ &= \begin{cases} \langle r_k, (t - \alpha_{k+1})r_j \rangle - \beta_k \langle r_{k-1}, r_j \rangle = 0, & \text{falls } j < k-1, \\ \langle tr_k, r_{k-1} \rangle - \|r_k\|_{\mathcal{T}}^2 = 0, & \text{falls } j = k-1, \\ \langle tr_k, r_k \rangle - \alpha_{k+1} \langle r_k, r_k \rangle = 0, & \text{falls } j = k. \end{cases} \end{aligned}$$

MATLAB-Code aus ^{GGG05} [16, Sect. 4.2.5]:

Input:

ν : Vektor der m Knotenpunkte

y : Werte y_j , $j = 0, \dots, m$

n : Maximaler Wert für k

Output:

α : Rekursionskoeffizienten α_k

β : Rekursionskoeffizienten β_k

b : Koeffizienten $\langle y, r_j \rangle$, $j = 0, \dots, n$

P : $p_{kj} = r_k(\nu_j)$

MATLAB-CODE : Berechnung diskreter Orthogonalpolynome

```
function [alpha,beta,b,P]=dorthp(nu,y,n)
m = length(nu);
alpha(1) = sum(nu)/m;
p1 = ones(size(nu)); p2 = nu-alpha(1);
P = [p1;p2];
b(1) = dot(p1,y)/dot(p1,p1);
b(2) = dot(p2,y)/dot(p2,p2);
for k=1:min(n-1,m-2)
    p0 = p1; p1 = p2;
    alpha(k+1) = dot(p1,(nu.*p1))/norm(p1)^2;
    beta(k) = (norm(p1)/norm(p0))^2;
    p2 = (nu-alpha(k+1)).*p1-beta(k).*p0;
    P = [P;p2];
    b(k+2) = dot(p2,y)/norm(p2)^2;
end
```

➤ Aufwand zur Berechnung von $\alpha_k, \beta_k, \langle y, r_j \rangle$, $j = 0, \dots, n = O(nm)$

MATLAB-CODE Auswertung diskreter

```
function[Y,Z]= evlortho(alpha,beta,b,z)
m = max(size(b));
p1 = ones(size(z));
p2 = z-alpha(1);
Z = [ p1; p2];
Y = [p1*b(1); p1*b(1)+p2*b(2)];
for i = 2:m-1
    p0 = p1; p1 = p2;
    p2 = (z-alpha(i)).*p1...
        -beta(i-1).*p0;
    Z = [Z;p2];
    Y = [Y; Y(i,:)+p2*b(i+1)];
end
```

Orthogonalpolynome

Auswertung von diskreten Orthogonalpolynomen: (MATLAB-Code aus ^{GGG05} [16, Sect. 4.2.5])

Input:

☞ Ausgabe von **dorthp**

z : Auswertungsstellen

Output:

Y : Fitpolynomwerte an Stellen z

Z : Orthogonalpolynomwerte an Stellen z

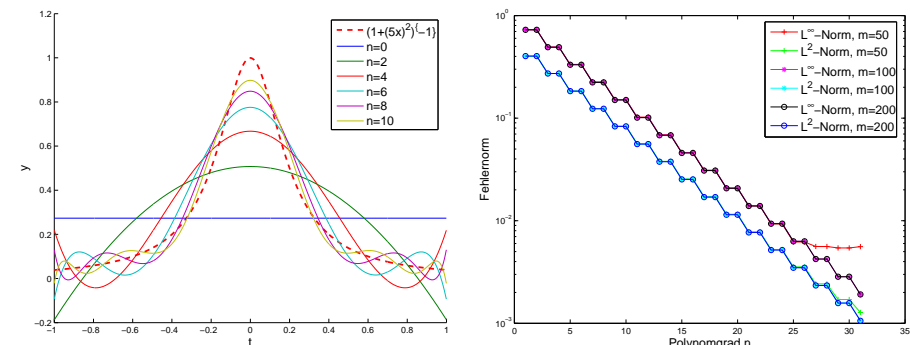
Beispiel 170 (Approximation durch Fitpolynome).

① $f(t) = (1 + (5t)^2)^{-1}$, $I = [-1, 1] \rightarrow$ Bsp. ^{GGG05} 162, analytische Funktion

Fitten in $t_k = -1 + k \frac{2}{m}$, $k = 0, \dots, m$, $m \in \mathbb{N}$

Approximative Berechnung der Fehlernormen (Sampling in $\xi_j = -1 + \frac{j}{500}$, $j = 0, \dots, 1000$)

4.1
p. 429



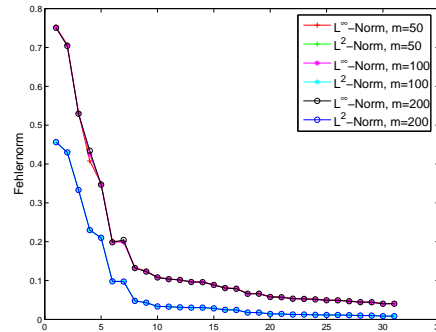
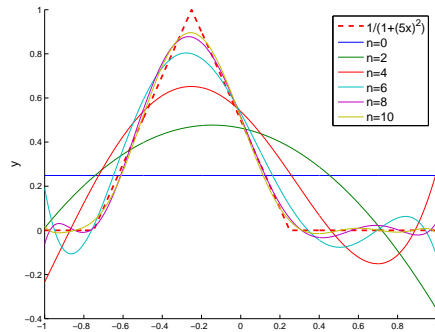
➤ Beobachtet: Exponentielle Konvergenz (im Polynomgrad n) für $n \ll m$

② $f(t) = \max\{0, 1 - 2 * |x + \frac{1}{4}|\}$, f nur $C^0([-1, 1])$, sonst wie oben

4.1
p. 430

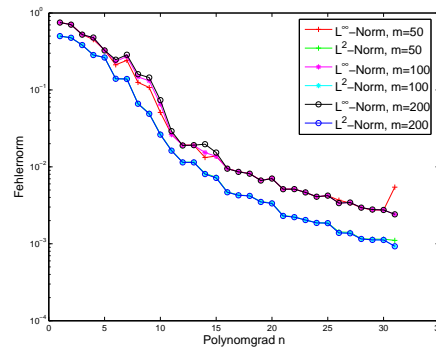
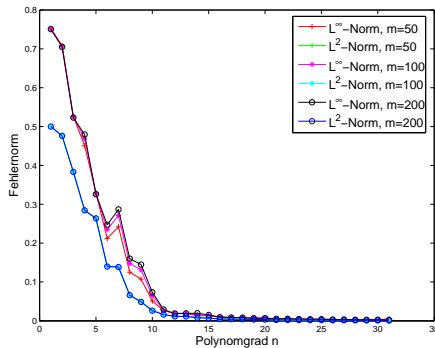
4.1
p. 431

4.1
p. 432



➤ Beobachtet: Nur algebraische Konvergenz (im Polynomgrad n) für $n \ll m$

③ $f(t) = \max\{\cos(4\pi|t + \frac{1}{4}|), 0\}$, nur $f \in C^1([-1, 1])$, sonst wie oben



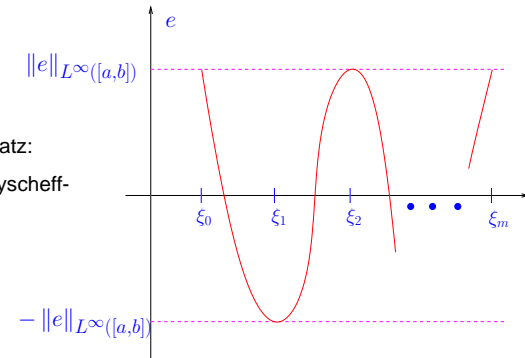
4.1.5.3 Tschebyscheff-Approximation

Gegeben: $f : I \subset \mathbb{R} \rightarrow \mathbb{C}$ beschränkt, I Intervall

Gesucht: L^∞ -bestapproximierendes Polynom: $q_n = \operatorname{argmin}_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty(I)}$

Theorem 4.1.20 (Tschebyscheffscher Alternantensatz). Für $f \in C[a, b]$, $a < b$, $n \in \mathbb{N}$:

$$q = \arg \min_{p \in \mathcal{P}_n} \|f - p\|_{L^\infty(I)} \Leftrightarrow \begin{aligned} &e(t) := f(t) - q(t) \text{ erfüllt:} \\ &\exists a \leq \xi_0 < \xi_1 < \dots < \xi_m \leq b, m \geq n+1, \\ &|e(\xi_j)| = \|e\|_{L^\infty([a,b])}, j = 0, \dots, m, \\ &e(\xi_{j-1}) = -e(\xi_{j+1}), j = 1, \dots, m. \end{aligned}$$



Veranschaulichung Alternantensatz:

Verhalten von $e \stackrel{\text{def: Tpol}}{\leftrightarrow}$ Tschebyscheff-Polynome \rightarrow Def. 4.1.6

4.1

p. 433 Terminologie: Extremalstellen ξ_j der Fehlerfunktion = **Alternanten**
Der Remes-Algorithmus [45, Abschn. II.5]:

Darstellung von Polynomen bzgl. Basis $\{u_0, \dots, u_n\}$ von \mathcal{P}_n : u_j = Monome, Tschebyscheff-Polynome, Def. 4.1.6

Idee: Iterative Aktualisierung einer Alternantenmenge: $\mathcal{A}^{(0)} \rightarrow \mathcal{A}^{(1)} \rightarrow \dots$, $\#\mathcal{A}^{(l)} = n+2$

① $\mathcal{A}^{(0)} := \{\xi_0^{(0)} < \xi_1^{(0)} < \dots < \xi_n^{(0)} < \xi_{n+1}^{(0)}\} \subset [a, b]$ „willkürlich“
(z.B. Extremalstellen des Tschebyscheff-Polynoms T_{n+1} , \rightarrow 4.1.6: **Tschebyscheff-Alternanten**)

② Falls $\mathcal{A}^{(l)} := \{\xi_0^{(l)} < \xi_1^{(l)} < \dots < \xi_n^{(l)} < \xi_{n+1}^{(l)}\} \subset [a, b]$ exakte Alternantenmenge (\rightarrow Thm. 4.1.20):

$$\sum_{j=0}^n \alpha_j u_j(\xi_k^{(l)}) + (-1)^k \delta = f(\xi_k^{(l)}), \quad k = 0, \dots, n+1, \quad (4.1.27)$$

liefert L^∞ -bestapproximierendes Polynom $q = \sum_{j=0}^n \alpha_j u_j$, $\delta = \pm \|f - q\|_{L^\infty([a,b])}$.

4.1
p. 434

➤ Löse (4.1.27) mit $\mathcal{A}^{(l)}$. Wenn $|\delta| - \|f - q\|_{L^\infty([a,b])} \leq \text{TOL}$ STOP

4.1
p. 435

4.1
p. 436

③ $\mathcal{A}^{(l+1)}$ = Extrempunkte von $f - q$ (zu finden durch Newton/Sekanten-Verfahren, Sampling).
 $l \leftarrow l + 1$ und GOTO ②.

Alternativ: $\mathcal{A}^{(l+1)} := \mathcal{A}^{(l)} \cup \{\text{Maximalstelle von } |f - q|\} \setminus \{\text{Minimalstelle von } |f - q|\}$
 (Austausch \rightarrow robusterer Algorithmus)

MATLAB-Code: Remes-Algorithmus zur
 Tschebyscheff-Approximation auf $[a, b]$:

f: Handle auf $f: [a, b] \mapsto \mathbb{R}$
 n: Polynomgrad
 N: Sampling-Auflösung

Tschebyscheff-Alternanten

Vandermonde-Matrix:

Die ersten $n + 1$ Spalten des LGS (4.1.27)

Suche: Maximum der Fehlerfunktion
 (Sampling-Technik)

```
nx = length(ind); xe = x(ind);
if (nx == d), xe=[a;xe;b];
elseif (nx == d+1)
    xmin = min(xe); xmax =
    max(xe);
    if ((xmin-a) > (b-xmax))
        xe = [a;xe];
    else xe = [xe;b];
end
```

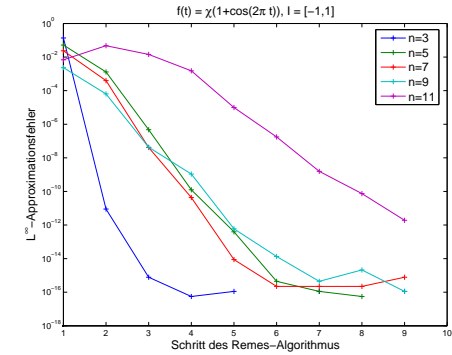
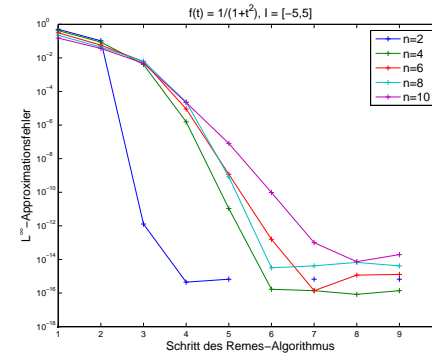
MATLAB-CODE: Remes-Algorithmus

```
function dev = remes(f,a,b,n,N,tol)
t = [a:(b-a)/N:b]'; y = feval(f,t);
xe=(a+b)/2+(a-b)/2*cos(pi/(n+1)*[0:d+1]');
fxe=feval(f,xe);
for k=1:maxit
    V=vander(xe);
    A=[V(:,2:d+2),(-1).^[0:d+1]'];
    c=A\fxe;
    delta = c(d+2);
    err = polyval(c(1:d+1),t) - y;
    delta = diff(err);
    s=[delta(1:n-1)].*[delta(2:n)];
    ind=find(s<0)+1;
    xe = ...;
    fxe=feval(f,xe);
    del = polyval(c(1:d+1),xe) - fxe;
    dev = max(abs(del));
    if (abs(dev-abs(e)) < tol), break; end
end;
```

Beispiel 171 (Konvergenz des Remes-Algorithmus).

• $f(t) = (1 + t^2)^{-1}$, $I = [-5, 5] \rightarrow$ Bsp. 162

• $f(t) = \begin{cases} \frac{1}{2}(1 + \cos(2\pi t)) & , \text{ falls } |t| < \frac{1}{2} \\ 0 & \text{sonst.} \end{cases}$, $I = [-1, 1]$

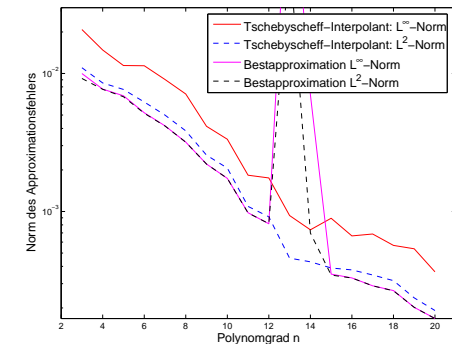
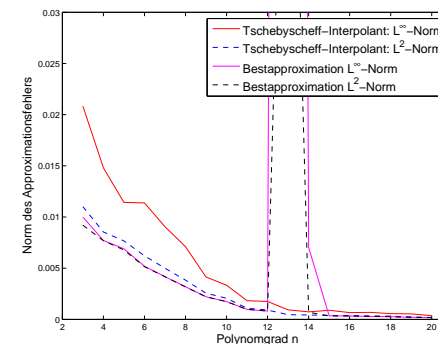


Beispiel 172 (Tschebyscheff-Interpolation vs. Bestapproximation).

$$f(t) = \begin{cases} \exp\left(\frac{1}{(t-0.1)^2 - \frac{1}{4}}\right) & , \text{ falls } -0.4 < t < 0.6 \\ 0 & \text{sonst.} \end{cases}, I = [-1, 1].$$

$f \in C^\infty([-1, 1])$, aber f nicht analytisch!

Approximative Berechnung der Fehlernormen für verschiedene Polynomgrade:



Beobachtungen: Indizien für Exponentielle Konvergenz

$n = 13$: Divergenz des Remes-Algorithmus !?

4.1.6 Clusteringapproximation

CLUSTAPPROX

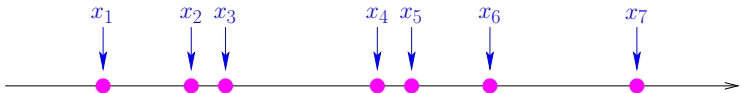
Aufgabenstellung:

KERNFN
Gegeben: **Kernfunktion** $G : I \times J \mapsto \mathbb{C}, I, J \subset \mathbb{R}$ Intervalle: $G(x, y)$ glatt für $x \neq y$
Kollokationspunkte $x_1 < x_2 < \dots < x_n, x_j \in I, y_1 < y_2 < \dots < y_m, y_j \in J$

KOLLMAT
Kollokationsmatrix: $\mathbf{M} \in \mathbb{C}^{n,m} \Leftrightarrow (\mathbf{M})_{ij} := G(x_i, y_j), \quad 1 \leq i \leq n, 1 \leq j \leq m.$
(4.1.28) **Eq:Model**

Gesucht: *Effizienter* Algorithmus für **approximative** Auswertung $\mathbf{M} \times \text{Vektor}$
(Naiv: Rechenaufwand $O(mn)$!)

Beispiel 173 (Interaktionsberechnung für Vielteilchensysteme).



n parallele stromdurchflossene Drähte, Draht $j \rightarrow$ Strom c_j , Position $x_j \in \mathbb{R}, c_j \in \mathbb{R}$

Zu berechnen: Magnetische Kraft auf *alle* Drähte

Kraft auf Draht j :
$$f_j = \sum_{\substack{k=1 \\ k \neq j}}^n \frac{1}{|x_j - x_k|} c_k c_j, \quad j = 1, \dots, n.$$

►
$$\mathbf{f} = \text{diag}(c_1, \dots, c_n) \mathbf{M} \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix}, \quad \mathbf{M} = (m_{ij})_{i,j=1}^n, \quad m_{ij} = \begin{cases} \frac{1}{|x_j - x_i|} & \text{falls } i \neq j, \\ 0 & \text{für } i = j. \end{cases}$$

➤ Kollokationsmatrix \mathbf{M} wird erzeugt durch **Kernfunktion** $G(x, y) = \frac{1}{|x - y|}$ ◇

Beispiel 174 (Gravitationskräfte in Galaxien).

Galaxie: $n (\approx 10^9)$ Sterne mit Positionen $\mathbf{x}_i \in \mathbb{R}^3$ und Massen $m_i, i = 1, \dots, n$.

Erforderlich für Simulation der Galaxiendynamik:

☞ Gravitationskraft auf *alle* Sterne !

$$f_j = \frac{G}{4\pi} \sum_{i \neq j} \frac{1}{\|\mathbf{x}_i - \mathbf{y}_j\|} m_i m_j, \quad j \in \{1, \dots, n\}.$$



$$\mathbf{f} = \text{diag}(m_1, \dots, m_n) \mathbf{M} \begin{pmatrix} m_1 \\ \vdots \\ m_n \end{pmatrix}, \quad m_{ij} := \begin{cases} \frac{G}{4\pi} \frac{1}{\|\mathbf{x}_i - \mathbf{y}_j\|} & \text{for } i \neq j, \\ 0 & \text{for } i = j, \end{cases} \quad 1 \leq i, j \leq n.$$

▷ dreidimensionale Verallgemeinerung der Aufgabenstellung ◇

4.1
p. 441

4.1.6.1 Separierte Kernapproximation

Beobachtung: Falls Kernfunktion **separiert**: $G(x, y) = g(x)h(y), \quad g : I \mapsto \mathbb{C}, h : J \mapsto \mathbb{C}$

$$\mathbf{M} = (g(x_j))_{j=1}^n \cdot \left((h(y_j))_{j=1}^m \right)^T \quad \blacktriangleright \quad \text{rank}(\mathbf{M}) = 1.$$



Rechenaufwand($\mathbf{M} \times \text{Vektor}$) = $m + n$

$$\left(\begin{array}{|c|} \hline \mathbf{M} \\ \hline \end{array} \right) \left(\begin{array}{|c|} \hline \vdots \\ \hline \end{array} \right) = \left(\begin{array}{|c|} \hline \vdots \\ \hline \end{array} \right) \underbrace{\left(\begin{array}{|c|} \hline \text{---} \\ \hline \end{array} \right)}_{\text{Skalarprodukt !}} \left(\begin{array}{|c|} \hline \vdots \\ \hline \end{array} \right)$$

Verallgemeinerung:
$$G(x, y) = \sum_{j=1}^q g_j(x) h_j(y), \quad g_j : I \mapsto \mathbb{C}, h_j : J \mapsto \mathbb{C}, q \in \mathbb{N}$$

$$\mathbf{M} = \mathbf{U} \mathbf{V}^T, \quad \begin{array}{l} \mathbf{U} \in \mathbb{R}^{n,q}, \quad u_{ij} = g_j(x_i), \quad j \in \{1, \dots, q\}, i \in \{1, \dots, n\}, \\ \mathbf{V} \in \mathbb{R}^{q,m}, \quad v_{ij} = h_i(y_j), \quad i \in \{1, \dots, q\}, j \in \{1, \dots, m\}. \end{array}$$

4.1
p. 442

4.1
p. 443

4.1
p. 444

$$\text{rank}(\mathbf{M}) = q$$

$$\text{Rechenaufwand}(\mathbf{M} \times \text{Vektor}) = q(m+n)$$

$$\begin{pmatrix} \text{M} \end{pmatrix} \begin{pmatrix} \text{Vektor} \end{pmatrix} = \begin{pmatrix} \text{U} \end{pmatrix} \left\{ \underbrace{\begin{pmatrix} \text{V}^T \end{pmatrix} \begin{pmatrix} \text{Vektor} \end{pmatrix}}_{q \text{ Skalarprodukte!}} \right\}$$



Idee: Globale Approximation von G durch (Summe) separierter Kernfunktionen:

$$G(x, y) \approx \tilde{G}(x, y) := \sum_{l=0}^d \sum_{k=0}^d \kappa_{l,k} g_l(x) h_k(y), \quad \kappa_{l,k} \in \mathbb{C}, \quad (x, y) \in I \times J.$$

(4.1.29) CA:1

$$\text{Ersetze } \mathbf{M} \rightarrow \tilde{\mathbf{M}} = (\tilde{m}_{ij}), \quad \tilde{m}_{ij} = \tilde{G}(x_i, y_j)$$

Bemerkung 175 (Qualitätsmass für Kernapproximation).

$$\|\mathbf{M} - \tilde{\mathbf{M}}\|_2^2 \leq \|\mathbf{M} - \tilde{\mathbf{M}}\|_F^2 = \sum_{i,j} (m_{ij} - \tilde{m}_{ij})^2 = \sum_{i,j} (G(x_i, y_j) - \tilde{G}(x_i, y_j))^2.$$

$$\text{Normalisiertes Qualitätsmass: } \frac{1}{mn} \sum_{i,j} (G(x_i, y_j) - \tilde{G}(x_i, y_j))^2. \quad (4.1.30) \text{ CA:quali}$$

△

Wie bekommt man eine separierte Kernapproximation?

Definition 4.1.21 (Tensorprodukt-Polynominterpolation). $L_j^x \in \mathcal{P}_n$ ($L_k^y \in \mathcal{P}_m$), $j = 0, \dots, n$ ($k = 0, \dots, m$), Lagrange-Polynome zu Knotenmenge $\mathcal{X} := \{x_j\}_{j=0}^n \subset I$ ($\mathcal{Y} := \{y_k\}_{k=0}^m \subset J$), $I, J \subset \mathbb{R}$ Intervalle. Zu stetigem $f: I \times J \rightarrow \mathbb{C}$ definiert

$$(I_{\mathcal{X} \times \mathcal{Y}} f)(x, y) := \sum_{j=0}^n \sum_{k=0}^m f(x_j, y_k) L_j^x(x) L_k^y(y), \quad x, y \in \mathbb{R},$$

das **Tensorprodukt-Interpolationspolynom**.

Abschnitt 4.1.3: Benütze **Tensorprodukt-Tschebyscheff-Polynominterpolation**

$$\tilde{G}(x, y) := \sum_{j=0}^d \sum_{k=0}^d G(t_j^x, t_k^y) L_j^x(x) L_k^y(y), \quad (4.1.31) \text{ CA:2}$$

$t_0^x, \dots, t_d^x, t_0^y, \dots, t_d^y$ Tschebyscheff-Knoten in I/J ,
 L_j^x, L_k^y zugeordnete Lagrange-Polynome.

Beispiel 176 (Globale separierbare Approximation bei glatter Kernfunktion).

- Glatte (analytische!) Kernfunktion $G(x, y) = \frac{1}{1 + |x - y|^2}$, Kollokationspunkte $x_i = y_i = \frac{i-1}{n}$, $i = 1, \dots, n$.
- \tilde{G} durch Tensorprodukt-Tschebyscheff-Polynominterpolation (178) vom Grad $d \in \mathbb{N}$.

$$\frac{1}{n^2} \sum_{i,j} (G(x_i, y_j) - \tilde{G}(x_i, y_j))^2$$

als Funktion von d für $n \in \{100, 200, 400\}$

smoothkernel

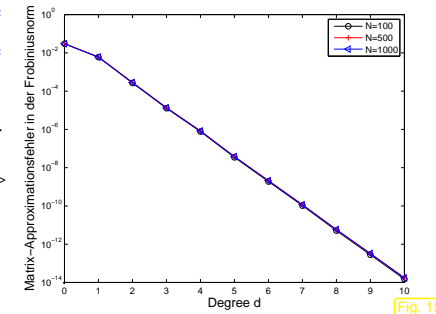


Fig. 15

4.1
p. 445

Beobachtet: **Exponentielle Konvergenz** $\|\mathbf{M} - \tilde{\mathbf{M}}\|_F \rightarrow 0$ in Abhängigkeit von d

4.1
p. 447

Beispiel 177 (Globale separierbare Approximation bei nichtglatter Kernfunktion).

Auswertungen wie in Bsp. 176, nun für Kernfunktion

$$G(x, y) = \begin{cases} \frac{1}{|x-y|} & \text{für } x \neq y, \\ 0 & \text{sonst,} \end{cases} \quad (4.1.32)$$

vgl. Bsp. 173.

nonsmoothkernel

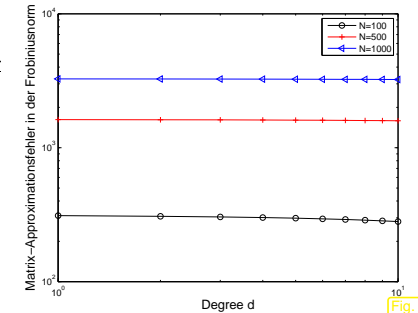


Fig. 16

Beobachtet: (Praktisch) keine Konvergenz $\|\mathbf{M} - \tilde{\mathbf{M}}\|_F \rightarrow 0$ für $d \rightarrow \infty$

Schuldiger:

Fehlende **globale** Glattheit

- Schlechte Approximierbarkeit von $G(x, y) := |x - y|^{-1}$ in Umgebung von $\{(x, y) \in I \times J: x = y\}$.

4.1
p. 446

4.1
p. 448

JEDOCH $G(x, y)$ aus (4.1.32) ist glatt (sogar analytisch \rightarrow Def. 4.1.14) in "grosserer Entfernung" von $\{(x, y) \in I \times J : x = y\}$.



Idee: Lokale Approximation von G durch (Summe) separierter Kernfunktionen:

$$G(x, y) \approx \sum_{l=0}^d \sum_{k=0}^d \kappa_{l,k} g_l(x) h_k(y), \quad \kappa_{l,j} \in \mathbb{C}, \quad (4.1.33)$$

$$(x, y) \in \tilde{I} \times \tilde{J}, \quad \tilde{I} \subset I, \quad \tilde{J} \subset J,$$

$$\text{mit } \tilde{I} \times \tilde{J} \cap \{(x, y) : x = y\} = \emptyset.$$

Konkret:

Lokale Approximation von G auf Rechtecken $\tilde{I} \times \tilde{J}$ einer Partitionierung von $I \times J$

- Möglichkeit der separierten Kernapproximation durch lokale Tensorprodukt-Tschebyscheff-Polynominterpolation (4.1.33)
- (zulässige Rechtecke)

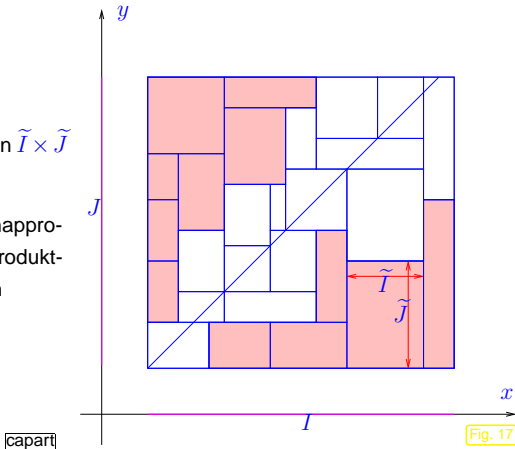


Fig. 17

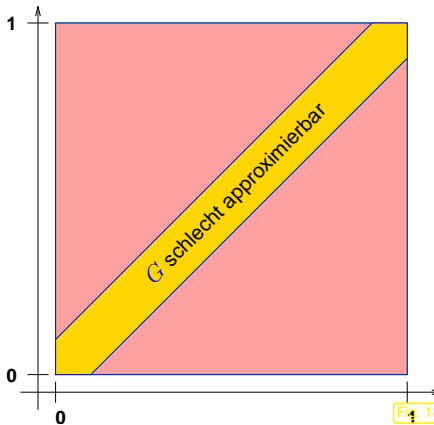


Fig. 18

Terminologie:

- Nahfeld:** G nicht approximierbar
- Fernfeld:** G gut approximierbar

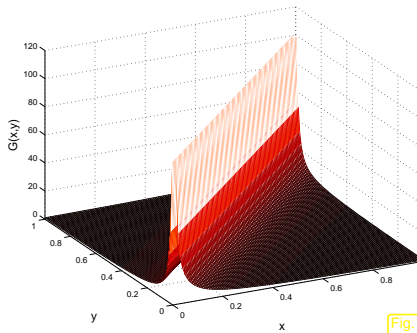


Fig. 19

Kernfunktion (4.1.32) auf $[0, 1]^2$

$$G(x, y) = \frac{1}{|x - y|}$$

- singulär für $x = y$
- Analytisch fern von $\{(x, y) : x = y\}$: Analytizitätsellipsen (\rightarrow Abschn. 4.1.4) wachsen mit $|x - y|$

Ähnliche Kernfunktionen: $G(x, y) = \log|x - y|$, $G(x, y) = \frac{\partial}{\partial x} \frac{1}{|x - y|}$, etc.

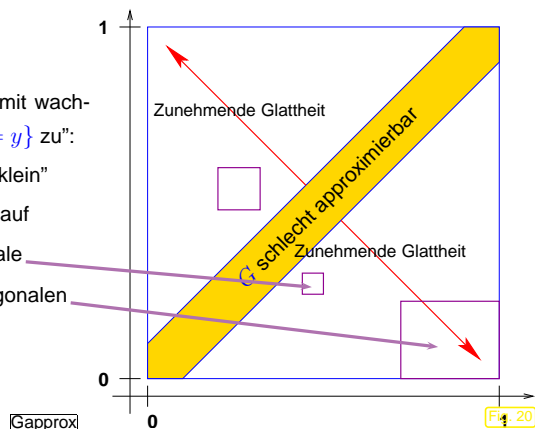
„Glattheit von $G(x, y) = \frac{1}{|x-y|}$ nimmt mit wachsender Entfernung von Diagonalen $\{x = y\}$ zu“:

➤ Keine Approximation, wenn $|x - y|$ „klein“

Tensorprodukt-Polynominterpolation auf

- kleinen Rechtecken nahe Diagonale
- grossen Rechtecken fern der Diagonalen

Wahl geeigneter Rechtecke ?

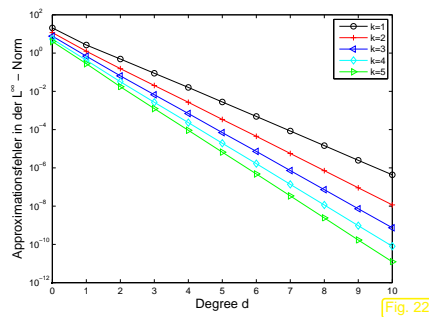
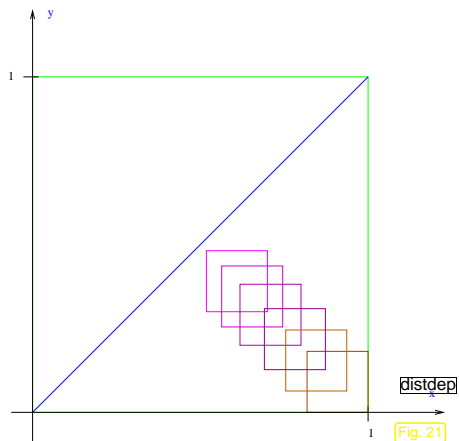


Beispiel 178 (Tensorprodukt-Tschebyscheff-Interpolation auf Rechtecken).

$I = J = [0, 1]^2$, $G(x, y) = |x - y|^{-1}$ from (4.1.32), (Genäherte) Supremumsnorm des Fehlers der lokalen Tensorprodukt-Tschebyscheff-Polynominterpolation von G auf

Rechtecken konstanter Grösse, aber mit wachsendem Abstand von $\{(x, y): x = y\}$.

$$\tilde{I}_k = [0.55 + k \cdot 0.05, 0.75 + k \cdot 0.05] \quad , \quad \tilde{J}_k = [0.25 - k \cdot 0.05, 0.45 - k \cdot 0.05] \quad , \quad k \in \{0, \dots, 5\}.$$



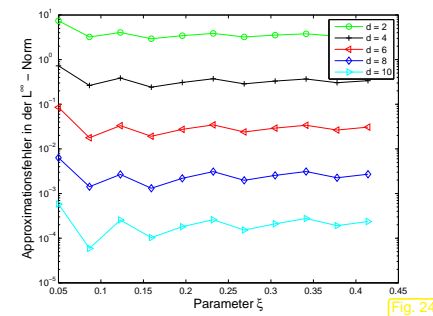
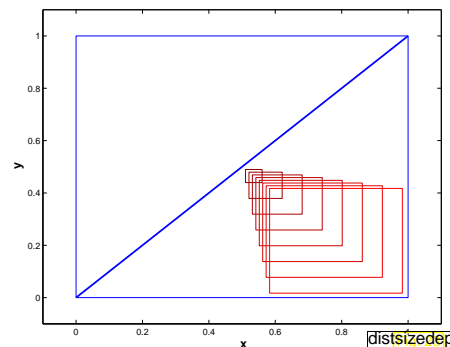
Abnahme des Interpolationsfehlers mit zunehmendem Abstand von $\{(x, y): x = y\}$

Beispiel 179 (Tensorprodukt-Tschebyscheff-Interpolation auf variablen Rechtecken).

siehe ➔ Bsp. 178: nun Supremumsnorm des Interpolationsfehlers auf Rechtecken

$$[\frac{1}{2}(\sqrt{2}-1)\xi + \frac{1}{2}, \frac{1}{2}(\sqrt{2}+1)\xi + \frac{1}{2}] \times [-\frac{1}{2}(\sqrt{2}-1)\xi + \frac{1}{2}, -\frac{1}{2}(\sqrt{2}+1)\xi + \frac{1}{2}] \quad , \quad 0.05 \leq \xi \leq \frac{1}{1+\sqrt{2}}$$

Grösse der Rechtecke \sim Abstand von Diagonalen



Abnahme des Interpolationsfehlers trotz Zunahme der Rechtecksgrösse

Zulässigkeit:

$[a, b] \times [c, d]$ heisst η -zulässig (engl. *admissible*), $\eta > 0$, falls

$$\eta \text{dist}([a, b], [c, d]) \geq \max\{b - a, d - c\}. \quad (4.1.34)$$

$$(\text{dist}([a, b], [c, d]) := \min\{|x - y| : x \in [a, b], y \in [c, d]\})$$

4.1.6.2 Clustertechnik

Aufgabe: Gegeben $\eta > 0$, finde (➔ **schneller** Algorithmus) Partitionierung von $I \times J$ „fern der Diagonalen“ in η -zulässige Rechtecke.

Bemerkung 180. Perspektiven:

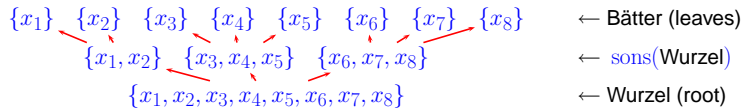
$$\text{Partitionierung von } I \times J \Rightarrow \text{Partitionierung von } \{x_i\}_{i=1}^n \times \{y_j\}_{j=1}^m \Leftrightarrow \text{Partitionierung der Indexmenge } \{1, \dots, n\} \times \{1, \dots, m\}$$

Clustertechnik:

Definition 4.1.22 (Clusterbaum). Ein Baum (➔ Informatik, Graphentheorie) T heisst **Clusterbaum** (engl. *cluster tree*) zu $\mathbb{P} := \{x_1, \dots, x_n\} \subset \mathbb{R}$: \Leftrightarrow

- Knoten von T (= **Cluster**) sind Teilmengen von \mathbb{P} .
- $\text{root}(T) = \mathbb{P}$.
- Für alle Knoten σ von T : $\{\sigma' : \sigma' \in \text{sons}(\sigma)\}$ ist Partitionierung von σ .

Bounding Box von Cluster $\sigma \in T$: $\Gamma(\sigma) := [\min\{x\}_{x \in \sigma}, \max\{x\}_{x \in \sigma}]$



Terminologie (→ Graphentheorie): Clusterbaum T : $\sigma \in T$: $\text{sons}(\sigma) = \emptyset \Leftrightarrow \sigma$ Blatt

Clusterbaum T :

Rekursive Konstruktion

MATLAB-Datenstruktur:

$N \times 6$ -Matrix, $N = \#T$,

Zeilen $\hat{=}$ Cluster

$T = [T; i, j, x_l, x_r, s_1, s_2]$

i, j : $\sigma = \{x_i, \dots, x_j\}$

x_l : $x_l = \min_{i \leq k \leq j} x_k$

x_r : $x_r = \max_{i \leq k \leq j} x_k$

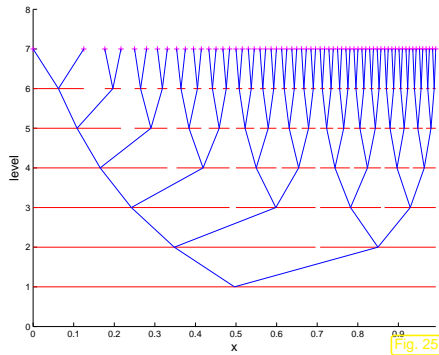
s_1, s_2 : Zeilenindizes der Söhne

$\#T \leq 2n$

MATLAB-CODE: Konstruktion eines Clusterbaumes

```
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);
```

Beispiel 181 (Clusterbaum).



← Binärer Clusterbaum:

$x = \text{sqrt}(0:1/64:(1-1/64));$

➤ 2^{l-1} Cluster auf Level l mit je 2^{7-l} Punkten.
 (ausgewogener Clusterbaum)



Idee: Wähle Approximationsrechtecke = Tensorprodukte von Bounding Boxen von Clustern $\in T$



Respektiere Zulässigkeitsbedingung (4.1.34) ^{CA:adm}

Cluster σ, μ zulässig \Leftrightarrow
 $\Gamma(\sigma), \Gamma(\mu)$ zulässig.

Rekursiver Aufbau einer
 zulässigen Partitionierung:

Pnear = Nahfeld:

Zone nahe der Diagonalen

(→ keine Approximation)

Pfar = Fernfeld:

Partitioniert mit η -zulässigen

Rechtecken

MATLAB-CODE: Rekursiver Aufbau einer zulässigen Partitionierung

```
function [Pnear,Pfar] = ...
    divide(Tx,Ty,sigma,mu,Pnear,Pfar,eta)
cls = Tx(sigma,:); clm = Ty(mu,:);
if (sigma == Blatt | mu == Blatt)
    Pnear = [Pnear; cls(2),cls(3),clm(2),clm(3)];
elseif ((sigma,mu) zulässig)
    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: Aufruf

```
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);
```

4.1
 p. 457

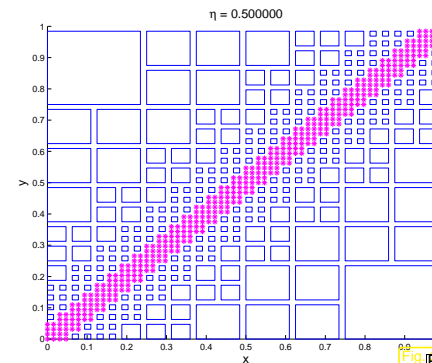
4.1
 p. 459

MATLAB Animation zum Aufbau einer Partitionierung in File

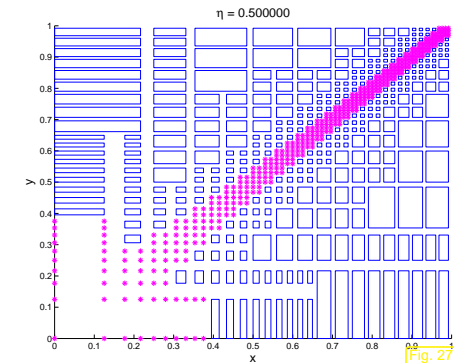
.MATLAB/Kapitel_4.ApproximationInterpolation/Clustering/clusteringanimation.m

Aufruf: clusteringapproximation(N,waittime)

N: Anzahl Kollokationspunkte in einer Dimension, waittime: Zeit (sec) zwischen 2 Animationsschritten



$x = (0:1/64:(1-1/64));$



$x = \text{sqrt}(0:1/64:(1-1/64));$

4.1
 p. 458

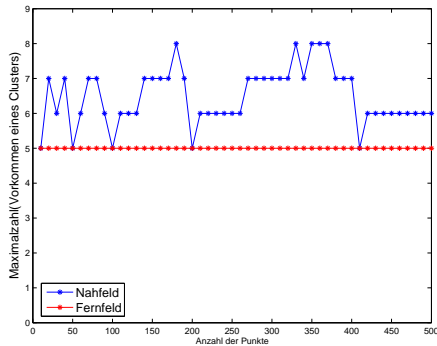
4.1
 p. 460

4.1
p. 464

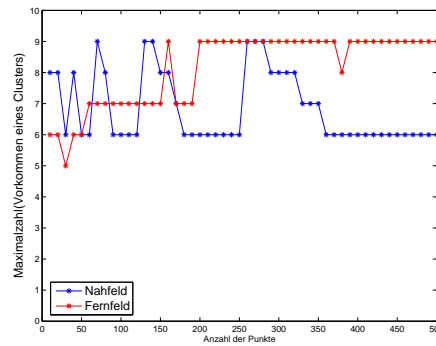
②: Berechne $\mathbf{X}_{\sigma,\mu}$ aus (4.1.35), $\sigma \in T_x, \mu \in T_y, (\sigma, \mu) \in P^{\text{far}}$.

Regelmässige Punktverteilung $\triangleright \max\{\#\{\mu \in T_y: (\sigma, \mu) \in P^{\text{far}}\}, \sigma \in T_x\} = O(1)$

Beispiel 182 (Vorkommen von Clustern in Partitionsrechtecken).



$$\mathbf{x} = (0:1/n:(1-1/n)), \eta = 1$$

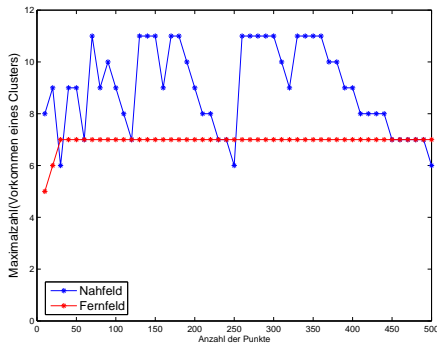


$$\mathbf{x} = \text{sqrt}(0:1/n:(1-1/n)), \eta = 1$$

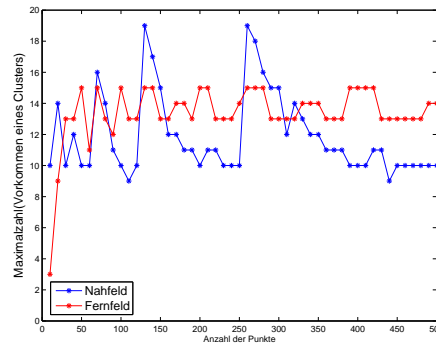
Bei „gleichmässiger“ Verteilung der x_i, y_j : Für alle $\sigma \in T_x$

$$\#\{\mu \in T_y: (\sigma, \mu) \in P^{\text{far}}\} = O(1),$$

\triangleright jeder Cluster trägt nur zu einer kleinen Anzahl von Partitionsrechtecken bei.



$$\mathbf{x} = (0:1/n:(1-1/n)), \eta = 0.5$$



$$\mathbf{x} = \text{sqrt}(0:1/n:(1-1/n)), \eta = 0.5$$

◇

\triangleright Rechenaufwand für Schritt ②: $O(n(d+1)^2)$, Speicheraufwand $O(n(d+1)^2)$

③: Fernfeldauswertung \triangleq Summe

$$\sum_{\substack{\sigma \in T_x \\ x_i \in \sigma}} \sum_{\substack{\mu \in T_y \\ (\sigma, \mu) \in P^{\text{far}}}} \mathbf{V}_\sigma \mathbf{X}_{\sigma,\mu} \mathbf{w}_\mu$$

foreach $\sigma \in T_x$ { Berechne \mathbf{V}_σ ; $\mathbf{s} := 0$
 foreach $\mu \in T_y, (\sigma, \mu) \in P^{\text{far}}$ { $\mathbf{s} \leftarrow \mathbf{s} + \mathbf{X}_{\sigma,\mu} \mathbf{w}_\mu$ }
 $\mathbf{f}_{|\sigma} \leftarrow \mathbf{f}_{|\sigma} + \mathbf{V}_\sigma \mathbf{s}$ }

\triangleright Rechenaufwand $O(n \log_2 n (d+1)^2)$, Speicheraufwand $O(n(d+1))$

④: Nahfeldberechnung

foreach $(\sigma, \mu) \in P^{\text{near}}$ { foreach $i: x_i \in \sigma$ { $\mathbf{f}_i \leftarrow \mathbf{f}_i + \sum_{j: y_j \in \mu} G(x_i, y_j) c_j$ } }

\triangleright Rechenaufwand $O(n)$

Gesamtrechenaufwand/Gesamtspeicherbedarf $O((d+1)^2 n \log_2 n)$

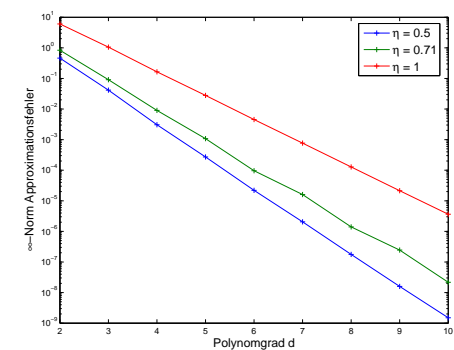
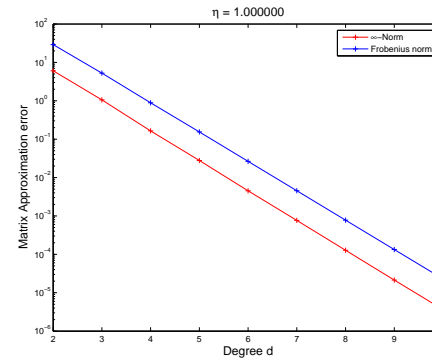
Bemerkung 183 (Konvergenz der Clusteringapproximation).

Falls $G(x, y)$ analytisch (\rightarrow Def. 4.1.14) in $\{(x, y): x \neq y\}$

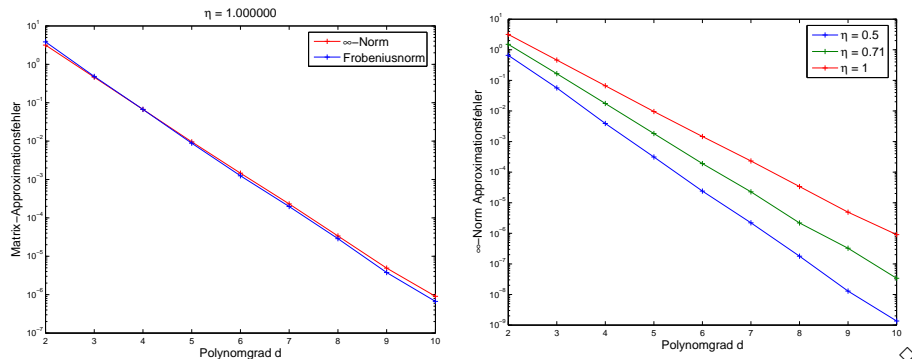
4.1 \triangleright Clusteringapproximation erbt **exponentielle Konvergenz** von Tschebyscheff-Interpolation, 4.1
 Bem. 4.1.4.4 \triangleright rem: TTPCS p. 465 \triangle p. 467

Beispiel 184 (Konvergenz der Clusteringapproximation einer Kollokationsmatrix).

• $\mathbf{x} = 0:1/128:1$; $G(x, y) = |x - y|^{-1}$ für $x \neq y$, $G(x, x) = 0$.



• $\mathbf{x} = \text{sqrt}(0:1/128:1)$; $G(x, y) = |x - y|^{-1}$ für $x \neq y$, $G(x, x) = 0$.



Beispiel 185 (Auswertung eines trigonometrischen Polynoms). ^{DUR95 [14]}

Gegeben: $\{t_0, \dots, t_{n-1}\} \subset [0, 1]$, $\alpha_{-m+1}, \dots, \alpha_m \in \mathbb{C}$, $n = 2m$, $m \in \mathbb{N}$, berechne

$$c_k := p(t_k), \quad j = 0, \dots, n-1 \quad \text{für} \quad p(t) := \underbrace{\sum_{j=-m+1}^m \alpha_j e^{-2\pi i j t}}_{\text{Trigonometrisches Polynom}}.$$

Diskrete Fouriertransformation (DFT) → Abschnitt ^{SEC: disk-four-transf} 5.7.1:

$$f_l := \sum_{j=-m+1}^m \alpha_j e^{-\frac{2\pi i}{n} j l} \quad \stackrel{\text{Lemma 5.7.2}}{\Leftrightarrow} \quad \alpha_j = \frac{1}{n} \sum_{l=-m+1}^m f_l e^{\frac{2\pi i}{n} l j}$$

FFT: f_l , $l = -m+1, \dots, m$, zu berechnen mit Rechenaufwand $O(n \log_2 n)$

$$\begin{aligned} c_k &= \sum_{j=-m+1}^m \alpha_j e^{2\pi i j t_k} = \sum_{j=-m+1}^m \frac{1}{n} \left(\sum_{l=-m+1}^m f_l e^{\frac{2\pi i}{n} l j} \right) e^{-2\pi i j t_k} \\ &= \frac{1}{n} \sum_{l=-m+1}^m f_l \sum_{j=-m+1}^m e^{-2\pi i j (t_k - l/n)} \\ &= \frac{1}{n} \sum_{l=-m+1}^m f_l e^{-2\pi i (t_k - l/n)(-m+1)} \frac{1 - e^{-2\pi i n (t_k - l/n)}}{1 - e^{-2\pi i (t_k - l/n)}} \\ &= \frac{1}{n} \sum_{l=-m+1}^m f_l e^{-\pi i t_k \sin(\pi n t_k)} \frac{1}{\sin(\pi (t_k - l/n))} (-1)^l e^{-\pi i l/n}. \end{aligned}$$

$$\blacktriangleright \mathbf{c} = \text{diag} \left(\frac{\sin(\pi n t_k)}{e^{\pi i t_k}} \right)_{k=-m+1}^m \mathbf{M} \text{diag} \left((-1)^l e^{-\pi i l/n} \right)_{l=-m+1, \dots, m} \mathbf{f}.$$

mit Kollokationsmatrix gemäss ^(4.1.28)

$$\mathbf{M} := \left(\frac{1}{\sin(\pi(t_k - l/n))} \right)_{\substack{k=-m+1, \dots, m \\ l=-m+1, \dots, m}} \in \mathbb{R}^{2n, 2n}.$$

➤ Approximative Auswertung mit Clustering-Algorithmus! → USFFT

Bemerkung 186. Clusteringapproximation Beispiel für schnelle approximative Realisierung eines Algorithmus der numerischen linearen Algebra (→ Kap. 3) → Trend in der numerischen linearen Algebra?

Problem	Exakte/direkte Methoden	Approximative/iterative Verfahren
Lineare Gleichungssysteme	Gauss-Elimination	CG-artige iterative Löser → Sect. 3.6
Eigenwertproblem	Transformationsmethoden	Krylov-Unterraumverfahren → Sect. 3.4.3
Kollokationsmatrix × Vektor	BLAS (SAXPY)	Clusteringtechniken

4.1
p. 469

Mit Clusteringapproximation verwandt, bzw. Verallgemeinerungen davon:

- \mathcal{H} -Matrix-Techniken ^{BGH03 [4]}
- \mathcal{H}^2 -Matrix-Techniken ^{HAB02 [22]}
- Multipol-Methoden ^{ROK85a, GRR97 [40, 20]} → weiterer „Milleniums-Algorithmus“
http://orion.math.iastate.edu/burkardt/misc/algorithms_dongarra.html

4.2 Dreitermrekursionen

[File: section-dreitermrekursion.tex, SVN: section-dreitermrekursion.tex 1011 2006-09-11 18:44:51Z hiptmair]

Technik zur Darstellung von Funktionen $I \subset \mathbb{R} \mapsto \mathbb{R}$ (→ Einleitung Kap. 4). ^{cha:interp-und-appr}

Gegeben: ^{FUNCSYS} Funktionensystem $b_k : I \mapsto \mathbb{R}$, $k \in \{1, \dots, n\}$, $n \in \mathbb{N}$
^{COEFFS} Entwicklungskoeffizienten α_j , $j \in \{1, \dots, n\}$

4.1
p. 470

4.2
p. 472

► Funktion $f(t) = \sum_{j=0}^{\infty} \alpha_j b_j(t)$, $t \in I$.

Aufgabe: Effiziente Auswertung von $f(x)$ für ein/mehrere $x \in I$

DREITERMREK

Dreitermrekursionen für viele relevante Funktionensysteme

$$b_0, b_1 \in \mathbb{R}, \quad b_{k+1}(x) = \kappa_k(x)b_k(x) + \sigma_k(x)b_{k-1}(x), \quad \kappa_k, \sigma_k \in \mathbb{R}, k \in \mathbb{N}. \quad (4.2.1) \quad \text{eq: dreiterm}$$

Beispiel 187 (Orthogonalpolynome). (→ Abschn. 4.1.5.2) |sec: least-squar-best

• **Legendre-Polynome** auf $I := [-1, 1] \rightarrow$ Def. 4.4.1, Fig. 38 |def: Legendre

$$\text{(4.4.6): } \text{leg: Legpol} \quad b_0(t) \equiv 1, b_1(t) = t, \quad b_{k+1}(t) = \frac{2k+1}{k+1}tb_k(t) - \frac{k}{k+1}b_{k-1}(t).$$

• **Tschebyscheff-Polynome** auf $I := [-1, 1] \rightarrow$ Def. 4.1.6, Figs. 12, 13 |def: Tpol T04T59

$$\text{(4.1.15): } \text{leg: TPrek} \quad b_0(t) \equiv 1, b_1(t) = t, \quad b_{k+1}(t) = 2tb_k(t) - b_{k-1}(t).$$

Immer gilt: Polynome vom Grad $\leq n$: $\mathcal{P}_n = \text{Span}\{b_0, \dots, b_n\}$

Beispiel 188 (Spezielle Funktionen).

• **Trigonometrische Funktionen** $b_k(t) = \cos(kt)$ bzw. $b_k(t) = \sin(kt)$

$$b_{k+1}(t) = 2 \cos t \cdot b_k(t) - b_{k-1}(t), \quad \begin{array}{l} b_0(t) \equiv 1, b_1(t) = \cos t \text{ für Cosinus,} \\ b_0(t) \equiv 0, b_1(t) = \sin t \text{ für Sinus.} \end{array}$$

Herleitung benutzt Additionstheorem

$$\cos x + \cos y = 2 \cos \frac{x+y}{2} \cos \frac{x-y}{2}, \quad x, y \in \mathbb{R},$$

für $x = (k+1)t, y = (k-1)t$.

• **Bessel-Funktionen** (erster Art) $b_k := J_k : \mathbb{R} \mapsto \mathbb{R}$

Reihenentwicklung:

$$J_k(t) := \sum_{l=0}^{\infty} \frac{(-1)^l}{2^{2l+k} l!(k+l)!} t^{2l+k}.$$

Dreitermrekursion:

$$b_{k+1}(t) = \frac{2k}{t}b_k(t) - b_{k-1}, \quad k \in \mathbb{N}, \quad (4.2.2) \\ b_0 := J_0, b_1 := J_1.$$



besselfunctions

Fig. 28

Wichtig z.B. in Fluidmechanik, Elektromagnetik: Lösungsdarstellung für Feldgleichungen bei Rotationssymmetrie

4.2
p. 473

Plotten von Funktionensystemen aus (F.2.1) |eq: plotbk

```
function plotbk(kappa,sigma,t,b0,b1,N)
% Plotten von Funktionensystemen
V = [b0;b1];
for k=1:N-1
    V = [V;kappa(k,t).*V(end,:) + ...
        sigma(k,t).*V(end-1,:)];
end
figure; hold on;
for k=0:N, plot(t,V(k+1,:),'-'); end
```

◁ MATLAB-Code zur graphischen Ausgabe der b_k
Funktionen σ, κ →
 $\sigma_k(x), \kappa_k(x)$
Für Besselfunktionen J_k :
• $\kappa = @(k,t) 2*k./t$;
• $\sigma = @(k,t) -1$;

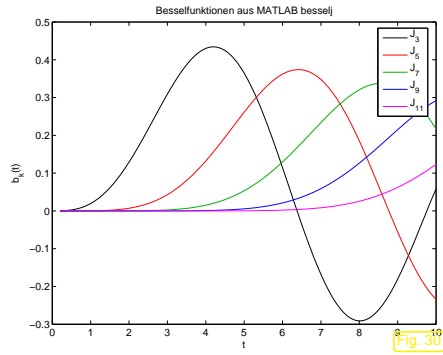
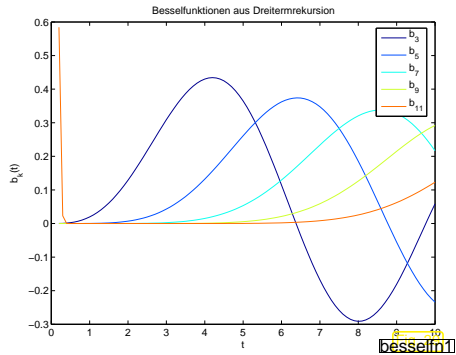
Beispiel 189 (Besselfunktionen aus Dreitermrekursion).

Beobachtung: "Instabilität" für kleine t , grosse k !?

4.2
p. 474

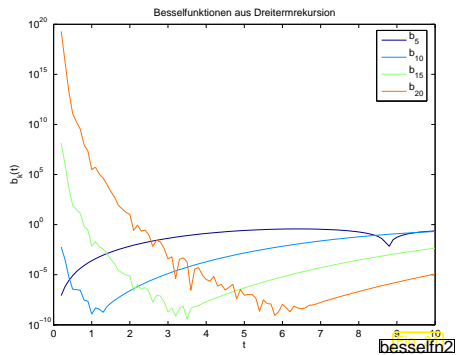
4.2
p. 475

4.2
p. 476



besseln1

Fig. 30



besseln2

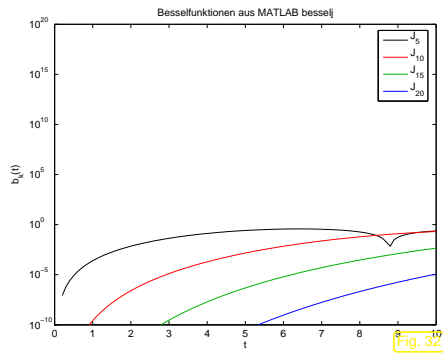


Fig. 30

Modellproblem:

Skalare Dreitermrekursion mit konstanten Koeffizienten:

$$x^{(k+1)} = \kappa \cdot x^{(k)} - x^{(k-1)}, \quad x^{(0)}, x^{(1)} \in \mathbb{R}.$$

(4.2.3) [eq:3tsim](#)

Ansatz:

$$x^{(k)} = \zeta^k, \quad \zeta \in \mathbb{R}$$

notwendig: $\zeta^2 - \kappa\zeta + 1 = 0.$

$$\zeta_{1,2} = \kappa/2 \pm \sqrt{(\kappa/2)^2 - 1} : \quad |\kappa| \leq 2 \Rightarrow |\zeta_{1,2}| = 1, \\ |\kappa| > 2 \Rightarrow |\zeta_1| > 1, |\zeta_2| < 1.$$

Allgemeine Lösung: $x^{(k)} = c_1 \zeta_1^k + c_2 \zeta_2^k, \quad c_1, c_2 \in \mathbb{R}.$

(4.2.4) [eq:3tgen](#)

Beispiel 190 (Instabile Dreitermrekursion).

[eq:3tsimple](#) (4.2.3) für $\kappa = 5, x^{(0)} = 1, x^{(1)} = 5/2 - \frac{1}{2}\sqrt{21}$

$$x^{(k)} = (5/2 - \frac{1}{2}\sqrt{21})^k.$$

Beobachtung:

„Umspringen“ zu $x^{(k)} = (5/2 + \frac{1}{2}\sqrt{21})^k$

? Warum denn das ?

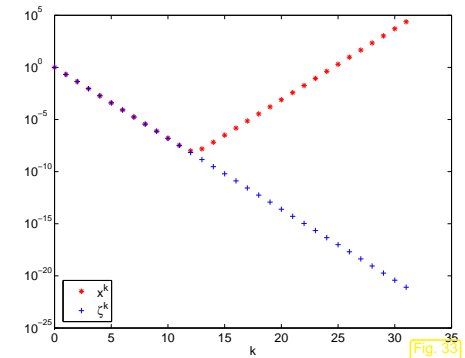


Fig. 33

4.2
p. 477

4.2
p. 479

Im Computerprogramm:

$$x^{(1)} = 0.20871215252208 \neq 5/2 - \frac{1}{2}\sqrt{21} \quad \text{Rundung !}$$

Gestörte $x^{(0)}, x^{(1)} \Rightarrow$ Dominanz exponentiell wachsender Lösung für $k \gg 1$

Beispiel 191 (Instabile Dreitermrekursion für Besselfunktionen).

[eq:3tbess](#) Dreitermrekursion (4.2.2):

$$t = 1/2, b_0 = J_0(1/2), b_1 = J_1(1/2)$$

Beobachtung:

„Wuchern“ einer exponentiell wachsenden Lösung

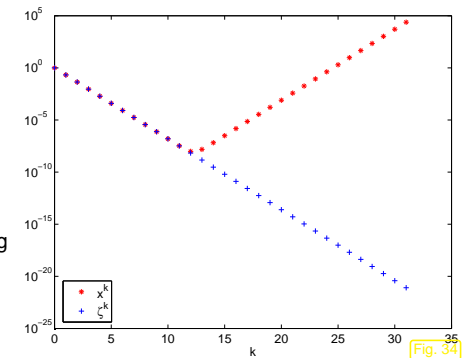


Fig. 34

4.2
p. 478

4.2
p. 480



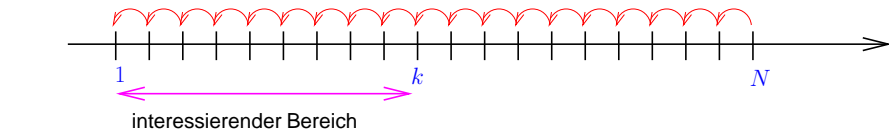
Idee:

Rückwärtsrekursion

Fixiere $N \gg n$ und berechne $b_k(t)$, $k = 0, \dots, n$, aus (für $\sigma_k(t) \neq 0$)

$$b_{k-1}(t) = -\frac{\kappa_k(t)}{\sigma_k(t)} b_k(t) + \frac{1}{\sigma_k(t)} b_{k+1}(t) \quad , \quad b_N(t), b_{N-1}(t) \text{ "beliebig"}$$

(4.2.5) [eq:inv3t](#)



Beispiel 192 (Approximation von Besselfunktionen mittels Rückwärtsrekursion).

Für $t = 10, 20, 30, 40$ Näherungen $\tilde{J}_k(t)$ durch Rückwärtsrekursion für $k = 0, \dots, 20$

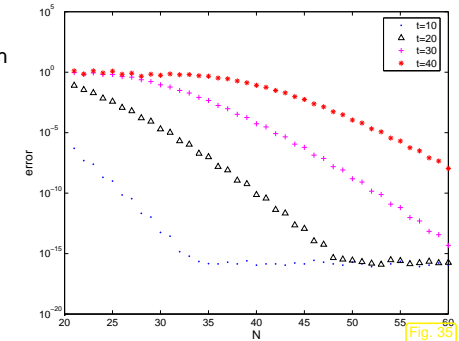
Fehlermass:

$$\max_{0 \leq k \leq 20} |J_k(t) - \tilde{J}_k(t)|$$

Beobachtung

Asymptotisch **exponentielle Konvergenz**

[revrecbess](#)



[Fig. 39](#)

4.2 Bemerkung 193.

p. 481

4.2
p. 483

Adaptive Rückwärtsrekursion für Besselfunktionen

```
function y = besscomp(t,n,tol)
d = realmax; N = n + 10;
x = revrecbess(t,n,N);
while(d > tol)
    N = N + 10;
    y = revrecbess(t,n,N);
    d = max(abs(x-y));
end
```

Adaptive Wahl von N :

Erhöhe N solange bis sich "Approximation stabilisiert"

(d.h. Fehler < Toleranz tol)

4.3 Stückweise Polynome

[File: section-stueckweise-polynome.tex, SVN: section-stueckweise-polynome.tex 1270 2007-01-15 10:44:02Z hiptmair]

STKWPOLY

Aufgabenstellung: Modellieren eines funktionalen Zusammenhangs $f: I \subset \mathbb{R} \mapsto \mathbb{R}$ basierend auf (genauen) Messungen (t_i, y_i) , $i = 0, \dots, n$:

➤ **Interpolationsbedingung** $f(t_i) = y_i$, $i = 0, \dots, n$

Beispiel 194 (Versagen der Polynominterpolation).

4.2
p. 482

4.3
p. 484



Idee:

Nutze zusätzliche **Normierungsbedingung**

$$\sum_{k=0}^{\infty} \omega_k b_k(t) = \text{bekannt}, \quad \omega_k \in \mathbb{R}.$$

(4.2.6) [eq:normt](#)

Normierungsbedingungen:

• Für Modellproblem [eq:3tsimple](#) (4.2.3): $\frac{x(t)}{x(0)} = 1$

• Für Besselfunktionen [ABS70](#) [1]: $J_0(t) + 2 \sum_{k=1}^{\infty} J_{2k}(t) = 1$

Rückwärtsrekursion für Besselfunktionen

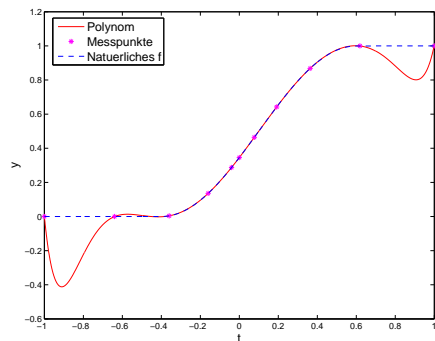
```
function x = revrecbess(t,n,N)
x = [1,0];
for k=N:-1:1
    x = [2*k/t*x(1)-x(2),x];
end
G = x(1) + 2*sum(x(3:2:length(x)));
x = x(1:n+1)/G;
```

◁ MATLAB-Funktion: Auswertung von Besselfunktionen J_k , $k = 0, \dots, n$ für Argument $t > 0$

Technik:

Rückwärtsrekursion
+
Normierungsbedingung

t_i	-1.0000	-0.6400	-0.3600	-0.1600	-0.0400	0.0000	0.0770	0.1918	0.3631	0.6187	1.0000
y_i	0.0000	0.0000	0.0039	0.1355	0.2871	0.3455	0.4639	0.6422	0.8678	1.0000	1.0000



← Interpolationspolynom, Grad = 10

Oszillationen an Intervallenden (→ Bsp. 194)

- Keine Lokalität
- Keine Monotonie
- Keine Krümmungserhaltung

$$f(t) = \begin{cases} 0 & , \text{ falls } t < -\frac{2}{5}, \\ \frac{1}{2}(1 + \cos(\pi(t - \frac{3}{5}))) & , \text{ falls } -\frac{2}{5} < t < \frac{3}{5}, \\ 1 & , \text{ sonst.} \end{cases}$$



Idee: Benutze **stückweise Polynome** bzgl. einer **Partition** (Gitter, engl. *mesh*) $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$ des Intervalls $I := [a, b]$, $a < b$.

Beachte: $x_j \neq t_j$ erlaubt!

4.3.1 Stückweise Polynominterpolation

Für Polynomgrad $d \in \mathbb{N}$ fasse jeweils $d+1$ Stützstellen zusammen ➤ Gitter \mathcal{M}

	t_0	t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	t_9	t_{10}	t_{11}	t_{12}	t_{13}	t_{14}	t_{15}	t_{16}	t_{17}	t_{18}	t_{19}	t_{20}
$d = 2$:	x_0	x_1	x_2	x_3	x_4	x_5					x_6	x_7	x_8	x_9	x_{10}						
$d = 4$:	x_0		x_1		x_2						x_3					x_4					x_5

➤ Stückweises Interpolationspolynom $s : [x_0, x_n] \rightarrow \mathbb{K}$:

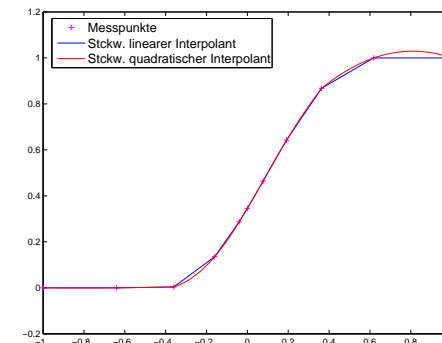
$$s_j := s|_{[x_{j-1}, x_j]} \in \mathcal{P}_d \quad \text{und} \quad s_j(t_i) = y_i \quad i = 0, \dots, n, \quad j = 1, \dots, m.$$

Beispiel 195 (Stückweise Polynominterpolation von Messpunkten).

Messpunkte aus Bsp. 194

Stückweise lineare/quadratische Interpolation

Lokalität



4.3.1.1 Stückweise lineare Interpolation

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

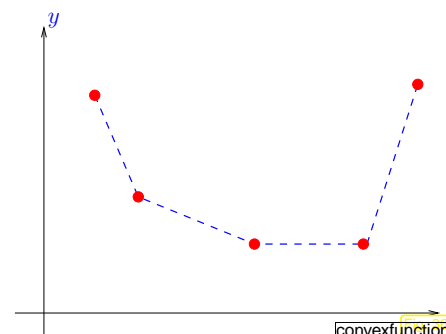
4.3
p. 485

Definition 4.3.1 (Monotone Daten). Daten *monoton*, wenn $y_i \geq y_{i-1}$ oder $y_i \leq y_{i-1}$, $i = 1, \dots, n$.

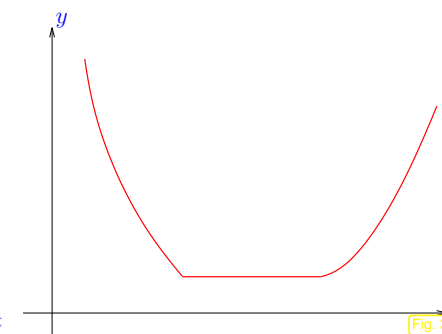
4.3
p. 487

Definition 4.3.2 (Konvexe/konkave Daten). Daten $\{(t_i, y_i)\}_{i=0}^n$ *konvex (konkav)*, falls

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



Konvexe Daten



Konvexe Funktion

4.3
p. 486

4.3
p. 488

Definition 4.3.3 (Konvexe/konkave Funktion).

$$f: I \subset \mathbb{R} \mapsto \mathbb{R} \quad \begin{array}{l} \text{konvex} \\ \text{konkav} \end{array} \quad :\Leftrightarrow \quad \begin{array}{l} f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y) \quad \forall 0 \leq \lambda \leq 1, \\ f(\lambda x + (1-\lambda)y) \geq \lambda f(x) + (1-\lambda)f(y) \quad \forall x, y \in I. \end{array}$$

SHAPEPRESERV

Theorem 4.3.4 (Formerhaltung bei stückweise linearer Interpolation).

Für stückweise linearen Interpolanten $s \in C([t_0, t_n])$ von $(t_i, y_i) \in \mathbb{R}^2, i = 0, \dots, n$:

$$\begin{array}{ll} (t_i, y_i) \text{ monoton (steigend/fallend)} & \Rightarrow s \text{ monoton (steigend/fallend)} \\ (t_i, y_i) \text{ konvex/konkav} & \Rightarrow s \text{ konvex/konkav} \end{array}$$

Lineare Interpolation ist linear (als Abbildung $y \mapsto s$) und lokal:

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

4.3.1.2 Stückweise polynomiale Interpolation von Funktionen

Lokale Lagrange-Interpolation von $f \in C(I)$:

Zu Gitter \mathcal{M} wähle Knotenmenge $\mathcal{T}^j := \{t_0^j, \dots, t_{n_j}^j\} \subset I_j$ für jede Gitterzelle $I_j := [x_{j-1}, x_j]$.

► Stückweises Interpolationspolynom $s: [x_0, x_m] \rightarrow \mathbb{K}$:

$$s_j := s|_{I_j} \in \mathcal{P}_{n_j} \quad \text{und} \quad s_j(t_i^j) = f(t_i^j) \quad i = 0, \dots, n_j, \quad j = 1, \dots, m.$$

Beachte: $t_{n_j-1}^{j-1} = x_{j-1} = t_0^j, j = 2, \dots, m \Leftrightarrow s \in C^0(I)$

Fragestellung: Asymptotik des Interpolationsfehlers für n fest, $m \rightarrow \infty$
(Interpolationsfehler $\leq T(h)$, Maschenweite $h := \max\{|x_j - x_{j-1}|: j = 1, \dots, m\}$)

Beispiel 196 (Stückweise Polynominterpolation).

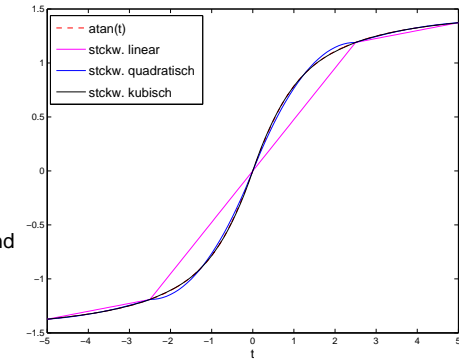
vgl. Bsp. 194: `ex:logpolyintpol`

$$f(t) = \arctan t, I = [-5, 5]$$

$$\text{Gitter } \mathcal{M} := \{-5, -\frac{5}{2}, 0, \frac{5}{2}, 5\}$$

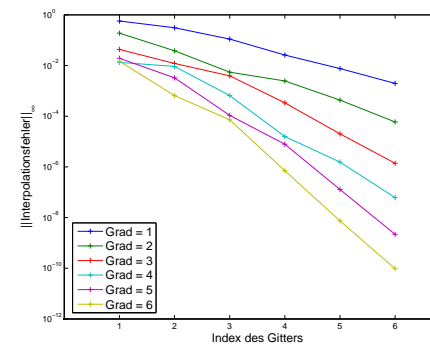
\mathcal{T}^j äquidistant in I_j

Plots der stückweise linearen, quadratischen und kubischen Polynominterpolanten



Interpolationsfehlernormen auf Gittern $\mathcal{M}_i := \{-5 + j 2^{-i} 10\}_{j=0}^{2^i}, i = 1, \dots, 6$, äquidistante Knotenmengen:

4.3
p. 489

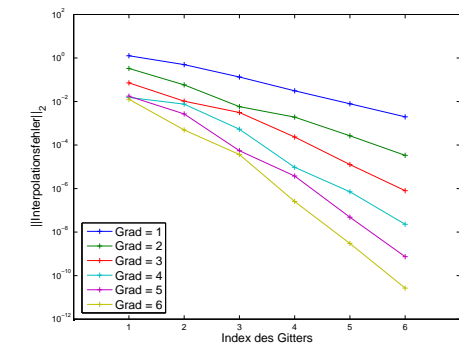


► Algebraische Konvergenz in Maschenweite

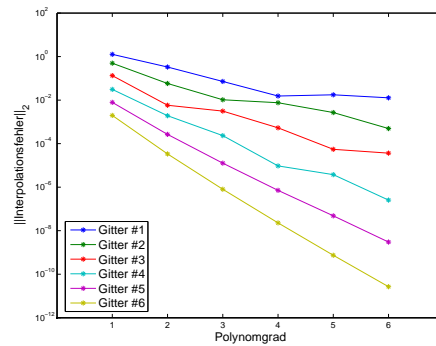
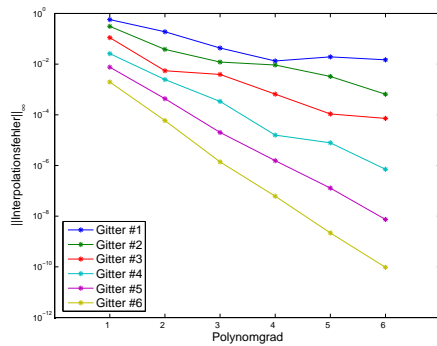
► Exponentielle Konvergenz im Polynomgrad

4.3
p. 490

4.3
p. 491



4.3
p. 492



◇

Interpolationsfehlerabschätzungen:

- Für festen Polynomgrad $n = n_j$, $j = 1, \dots, m$:

Anwendung von Cor 4.1.5 auf jedes Teilintervall: falls $f \in C^{n+1}([x_0, x_m])$

$$\|f - s\|_{L^\infty([x_0, x_m])} \leq \frac{h^{n+1}}{(n+1)!} \|f^{(n+1)}\|_{L^\infty([x_0, x_m])},$$

mit Maschenweite $h := \max\{|x_j - x_{j-1}| : j = 1, \dots, m\}$.

- Für festes Gitter: Situation wie bei Standard-Polynominterpolation → Abschnitt 4.1.2.



4.3.1.3 Kubische Hermite-Interpolation

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

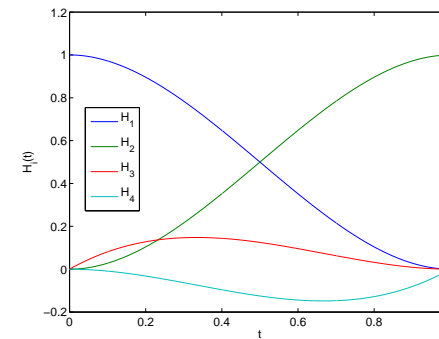
Gesucht: Modellfunktion $f \in C^1([x_0, x_n])$, $f(t_i) = y_i$, $i = 0, \dots, n$

- Stückweise **kubisches Hermite-Interpolationspolynom** $s \in C^1([t_0, t_n])$:
zu gegebenen **Steigungen** $c_i \in \mathbb{R}$, $i = 0, \dots, n$:

$$s|_{[t_{i-1}, t_i]} \in \mathcal{P}_3, \quad i = 1, \dots, n, \quad s(t_i) = y_i, \quad i = 0, \dots, n, \quad s'(t_i) = c_i, \quad i = 0, \dots, n.$$

$$\blacktriangleright s(t) = y_{i-1}H_1(t) + y_iH_2(t) + c_{i-1}H_3(t) + c_iH_4(t), \quad t \in [t_{i-1}, t_i], \quad (4.3.1)$$

[eq:kub](#)



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

Stückweise kubisches Polynom s auf

t_1, t_2 mit $s(t_1) = y_1$, $s(t_2) = y_2$,

$s'(t_1) = c_1$, $s'(t_2) = c_2$:

Effiziente lokale Auswertung

```
function s = kubeval(t,t1,t2,y1,y2,c1,c2)
h = t2-t1; t = (t-t1)/h;
a1 = y2-y1; a2 = a1-h*c1;
a3 = h*c2-a1-a2;
s = y1+(a1+(a2+a3*t).*(t-1)).*t;
```

Wahl der Steigungen c_i ?

☞ Mittelwert von Abschnittssteigungen:

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

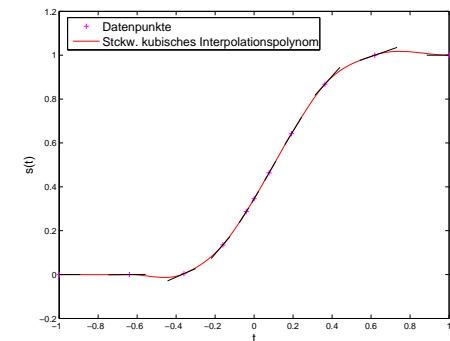
► Linearer lokaler Interpolationsoperator

Beispiel 197 (Stückweise kubische Hermite-Interpolation).

Daten aus Bsp. 194

Verwendung gewichtet gemittelter Steigungen

Keine Erhaltung von Monotonie und Krümmung



◇

Wahl der c_i mit „Limiter“ \rightarrow Monotonieerhaltung ^{FRC80} [15]

$$c_i = \begin{cases} 0 & , \text{ falls } \operatorname{sgn}(\Delta_i) \neq \operatorname{sgn}(\Delta_{i+1}) , \\ \text{gewichtetes Mittel von } \Delta_i, \Delta_{i+1} & \text{sonst} \end{cases} , \quad i = 1, \dots, n-1 .$$

Beispiel 198 (Monotonieerhaltende stückweise kubische Polynominterpolation).

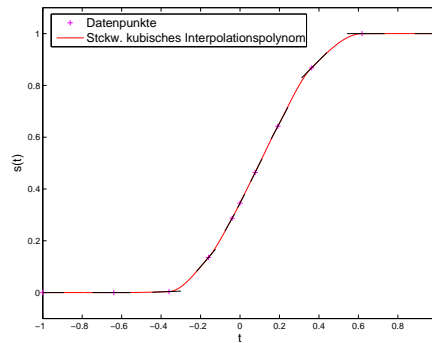
Daten aus Bsp. 194 ^{ex:logpolyintpol}

MATLAB-Funktion:

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

t: Stützstellen
y: Stützwerte
x: Auswertungsstellen
v: Vektor $s(x_i)$

Lokaler Interpolationsoperator
! Nichtlinearer Interpolationsoperator



Berechnung der c_i in pchip (Details in ^{FRC80} [15]):

```
function c = crfslopes(x,y)
n = length(x); h = diff(x); del = diff(y)./h;
if n==2, c = [del(1),del(1)];
return;
end
k = find(sign(del(1:n-2)).*sign(del(2:n-1)) > 0);
c = zeros(size(y)); hs = h(k)+h(k+1);
w1 = (h(k)+hs)/(3*hs); w2 = (hs+h(k+1))/(3*hs);
dmax = max(abs(del(k)), abs(del(k+1))); dmin = min(abs(del(k)), abs(del(k+1)));
c(k+1) = dmin./conj(w1.*(del(k)./dmax) + w2.*(del(k+1)./dmax));
c(1) = ((2*h(1)+h(2))*del(1) - h(1)*del(2))/(h(1)+h(2));
if (sign(d(1)) ~= sign(del(1)))
d(1) = 0;
elseif (sign(del(1)) ~= sign(del(2))) && (abs(d(1)) > abs(3*del(1)))
c(1) = 3*del(1);
end
c(n) = ((2*h(n-1)+h(n-2))*del(n-1) - h(n-1)*del(n-2))/(h(n-1)+h(n-2));
if (sign(d(n)) ~= sign(del(n-1)))
d(n) = 0;
elseif (sign(del(n-1)) ~= sign(del(n-2))) && (abs(d(n)) > abs(3*del(n-1)))
c(n) = 3*del(n-1);
end
```

4.3.2 Splines

SPLINES

Definition 4.3.5 (Splineräum). für ein Intervall $I := [a, b] \subset \mathbb{R}$, eine **Knotenmenge** (Gitter) $\mathcal{M} := \{a = t_0 < t_1 < \dots < t_{n-1} < t_n = b\}$ ist der Vektorraum $\mathcal{S}_{d,\mathcal{M}}$ der **Splinefunktionen** vom Grad d und der Ordnung $d+1$ definiert durch

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

$$\blacktriangleright s \in \mathcal{S}_{d,\mathcal{M}} \Rightarrow s' \in \mathcal{S}_{d-1,\mathcal{M}} \wedge \int_a^b s(\tau) d\tau \in \mathcal{S}_{d+1,\mathcal{M}} .$$

- $d = 0$: \mathcal{M} -stückweise konstante *unstetige* Funktionen
- $d = 1$: \mathcal{M} -stückweise lineare *stetige* Funktionen
- $d = 2$: *Stetig differenzierbare* \mathcal{M} -stückweise quadratische Funktionen

Abzählargument (Heuristik):

$$\dim \mathcal{S}_{d,\mathcal{M}} = n \cdot \dim \mathcal{P}_d - (n-1) \cdot d = n + d .$$

4.3
p. 497

4.3
p. 499

4.3.2.1 Splineinterpolation

Spezialfall: Interpolation in $\mathcal{S}_{1,\mathcal{M}}$ = Stückweise lineare Interpolation

Weiterer Spezialfall: **Kubische Splineinterpolation** $d = 3$ (\rightarrow Abschnitt ^{sec:kubische-herm-interp} 4.3.1.3)

^{leg:kub}
(4.3.1)

$$\blacktriangleright s|_{[t_{j-1}, t_j]}(t) = s(t_{j-1})(1 - 3\tau^2 + 2\tau^3) + s(t_j)(3\tau^2 - 2\tau^3) + h_j s'(t_{j-1})(\tau - 2\tau^2 + \tau^3) + h_h s'(t_j)(-\tau^2 + \tau^3) ,$$

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

\blacktriangleright Falls $s(t_j), s'(t_j), j = 0, \dots, l$ bekannt \rightarrow einfache Auswertung von s (\rightarrow Abschnitt ^{sec:kubische-herm} 4.3.1.3)

Interpolationsbedingung $s(t_j) = y_j, j = 0, \dots, n$: Steigungen $c_j := s'(t_j) = ?$

! $s \in C^2([t_0, t_n]) \blacktriangleright n-1$ Stetigkeitsbedingung für $s''(t) \blacktriangleright n-1$ Gleichungen

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

\blacktriangleright Zwei zusätzliche Bedingungen erforderlich:

4.3
p. 498

4.3
p. 500

① **Vollständige kubische Splineinterpolation:** $s'(t_0) = c_0, s'(t_n) = c_n$ vorgeschrieben

② **Natürliche kubische Splineinterpolation:** $s''(t_0) = s''(t_n) = 0$

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

► Jeweils LGS mit tridiagonaler s.p.d. (→ Def. 3.2.8, Lemma 3.2.11) Koeffizientenmatrix $\rightarrow c_0, \dots, c_n$

Thm. 3.2.18 \Rightarrow Rechenaufwand zur Lösung $= O(n)$

③ **Periodische kubische Splineinterpolation:** $s'(t_0) = s'(t_n), s''(t_0) = s''(t_n)$

► $n \times n$ -LGS mit s.p.d. Koeffizientenmatrix

$$A := \begin{pmatrix} a_1 & b_2 & 0 & \cdots & 0 & b_n \\ b_2 & a_2 & b_3 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & 0 \\ 0 & & & \ddots & \ddots & b_{n-1} \\ b_n & 0 & \cdots & 0 & b_{n-1} & a_{n-1} \end{pmatrix}, \quad \begin{aligned} b_i &:= \frac{1}{h_{i-1}}, \quad i = 2, 3, \dots, n, \\ a_i &:= \frac{2}{h_i} + \frac{2}{h_{i+1}}, \quad i = 1, 2, \dots, n-1. \end{aligned}$$

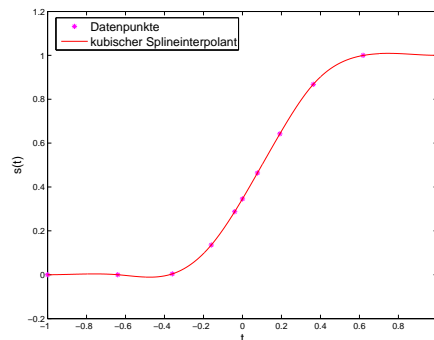
► Lösung durch Rang-1-Modifikationstechniken (→ Abschnitt 3.2.8.1, Lemma 3.2.22) + tridiagonale Elimination, Rechenaufwand $O(n)$

Beispiel 199 (Vollständige kubische Splineinterpolation).

Daten aus Bsp. 194

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

Kubischer Splineinterpolant nicht monotonie- oder krümmungserhaltend
(kubische Splineinterpolation linear !)



MATLAB-Funktion: `v = spline(t,y,x)`: Natürliche/vollständige Splineinterpolation
(→ Spline-Toolbox in MATLAB)

Beispiel 200 (Lokalität der natürlichen kubischen Splineinterpolation).

Gegeben: Gitter $\mathcal{M} := \{t_0 < t_1 < \dots < t_n\}$

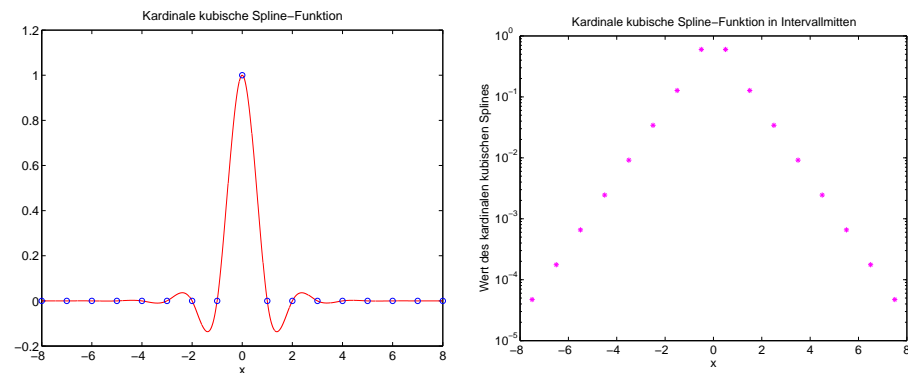
$$L_i \in \mathcal{S}_{3,\mathcal{M}} = i. \text{ natürlicher Kardinalspline} \Leftrightarrow L_i(t_j) = \delta_{ij}, \quad L_i''(t_0) = L_i''(t_n) = 0.$$

► Natürlicher Spline-Interpolant: $s(t) = \sum_{j=0}^n y_j L_j(t).$

Abfall der $L_i \leftrightarrow$ Lokalität der kubischen Splineinterpolation

Hier: $\mathcal{M} = \{-30, -29, \dots, 29, 30\}$

4.3
p. 501



► Exponentieller Abfall der Kardinalsplines \rightarrow kubische Splineinterpolation „fast lokal“

Beispiel 201 (Approximation durch vollständigen kubischen Splineinterpolanten).

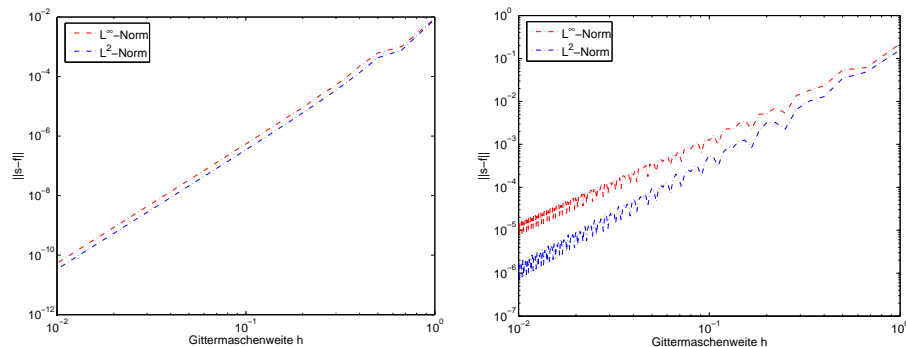
Gitter $\mathcal{M} := \{-1 + \frac{2}{n}j\}_{j=0}^n, n \in \mathbb{N} \rightarrow$ Maschenweite $h = 2/n, I = [-1, 1]$

$$f_1(t) = \frac{1}{1 + e^{-2t}} \in C^\infty(I), \quad f_2(t) = \begin{cases} 0 & \text{falls } t < -\frac{2}{5}, \\ \frac{1}{2}(1 + \cos(\pi(t - \frac{3}{5}))) & \text{falls } -\frac{2}{5} < t < \frac{3}{5}, \\ 1 & \text{sonst} \end{cases} \in C^1(I).$$

4.3
p. 502

4.3
p. 503

4.3
p. 504



$\|f_1 - s\|_{L^\infty([-1,1])} = O(h^4)$ $\|f_2 - s\|_{L^\infty([-1,1])} = O(h^2)$
 Theorie [27]: $f \in C^4([t_0, t_n]) \Rightarrow \|f - s\|_{L^\infty([t_0, t_n])} \leq \frac{5}{384} h^4 \|f^{(4)}\|_{L^\infty([t_0, t_n])}$

Glattheit des kubischen Splineinterpolanten

Für $f : [a, b] \mapsto \mathbb{R}$, $f \in C^2([a, b])$: $\frac{1}{2} \int_a^b |f''(t)|^2 dt$ = elastische Krümmungsenergie

Zu Gitter $\mathcal{M} := \{t_0 < t_1 < \dots < t_n\}$: $s \in \mathcal{S}_{3,\mathcal{M}}$ = natürlicher kubischer Splineinterpolant von $(t_i, y_i) \in \mathbb{R}^2$, $i = 0, \dots, n$.

$k \in C^2([t_0, t_n])$ erfülle $k(t_i) = 0$, $i = 0, \dots, n$:

$$h(\lambda) := \frac{1}{2} \int_a^b |s'' + \lambda k|^2 dt \Rightarrow \frac{dh}{d\lambda}|_{\lambda=0} = \int_a^b s''(t) k''(t) dt = \sum_{j=1}^n \int_{t_{j-1}}^{t_j} s''(t) k''(t) dt$$

Zweimalige partielle Integration:

$$\Rightarrow \int_a^b s''(t) k''(t) dt = s''(t_n) k'(t_n) - s''(t_0) k'(t_0) = 0.$$

Theorem 4.3.6 (Optimalität des natürlichen kubischen Splineinterpolanten). *Unter allen interpolierenden zweimal stetig differenzierbaren Funktionen minimiert der natürliche kubische Splineinterpolant die elastische Krümmungsenergie.*

4.3.2.2 Formerhaltende Splineinterpolation

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

Gesucht: Gitter $\mathcal{M} \subset [t_0, t_n]$ & formerhaltende interpolierende quadratische Splinefunktion $s \in \mathcal{S}_{2,\mathcal{M}}$, $s(t_i) = y_i$, $i = 0, \dots, n$ ($\mathcal{M} \neq \{t_j\}_{j=0}^n$!)

① „Formtreue“ Wahl von Steigungen c_i , $i = 0, \dots, n$ [MAR81, DQ001] \rightarrow Abschnitt 4.3.1.3

$$\text{Limiter } c_i := \begin{cases} \frac{2}{\Delta_i^{-1} + \Delta_{i+1}^{-1}}, & \text{falls } \text{sign}(\Delta_i) = \text{sign}(\Delta_{i+1}) \\ 0 & \text{sonst,} \end{cases} \quad i = 1, \dots, n-1.$$

$$c_0 := 2\Delta_1 - c_1, \quad c_n := 2\Delta_n - c_{n-1}.$$

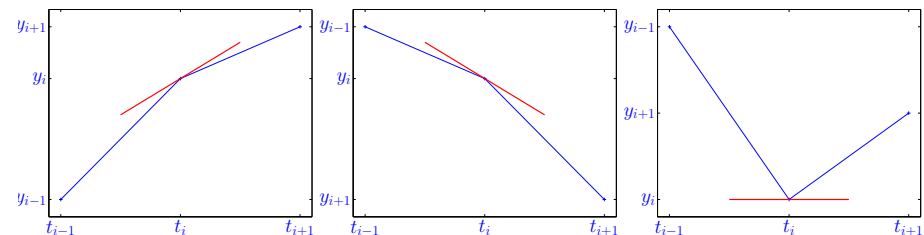
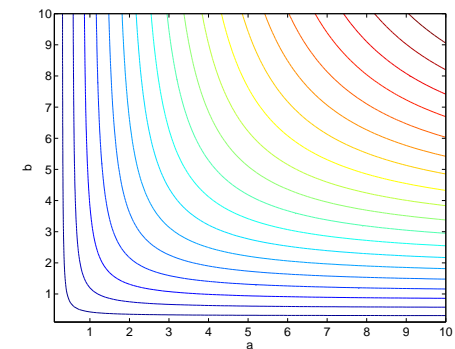
4.3
p. 505

$$c_i := \frac{2}{\Delta_i^{-1} + \Delta_{i+1}^{-1}}$$

= Harmonisches Mittel

Harmonisches Mittel von a und b

„Geglättete $\min(\cdot, \cdot)$ -Funktion“



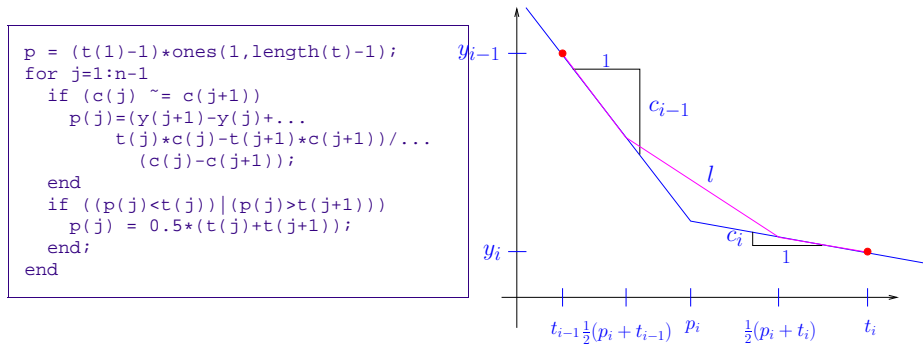
4.3
p. 506

4.3
p. 507

4.3
p. 508

② Wahl von „Zwischenpunkten“ $p_i \in]t_{i-1}, t_i], i = 1, \dots, n$:

$$p_i = \begin{cases} \text{Schnittpunkt der Geraden durch } (t_{i-1}, y_{i-1}), (t_i, y_i) \\ \text{mit Steigungen } c_{i-1}, c_i & , \text{ falls Schnittpunkt } \in]t_{i-1}, t_i] , \\ \frac{1}{2}(t_{i-1} + t_i) & \text{sonst.} \end{cases}$$



➤ Gitter für quadratischen Spline: $\mathcal{M} = \{t_0 < p_1 < t_1 < p_2 < \dots < p_n < t_n\}$

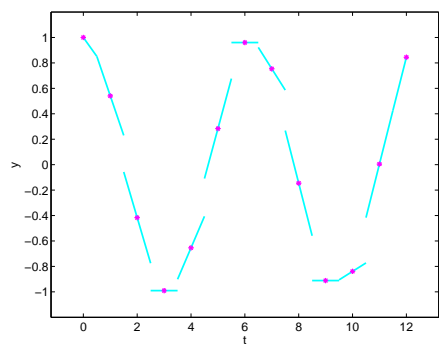
③ Setze l = Linearer Spline auf Gitter

$$\{t_0 < \frac{1}{2}(t_0 + p_1) < \frac{1}{2}(p_1 + t_1) < \frac{1}{2}(t_1 + p_2) < \dots < \frac{1}{2}(t_{n-1} + p_n) < \frac{1}{2}(p_n + t_n) < t_n\},$$

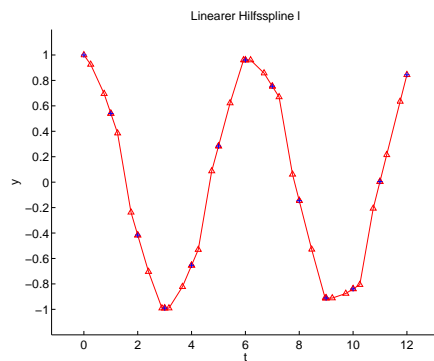
$$\text{mit } l(t_i) = y_i, \quad l'(t_i) = c_i.$$

Beispiel 202 (Hilfskonstruktion für formerhaltende quadratische Splineinterpolation).

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



Lokale Steigungen $c_i, i = 0, \dots, n$



Linearer Hilfsspline l

◇ p. 510

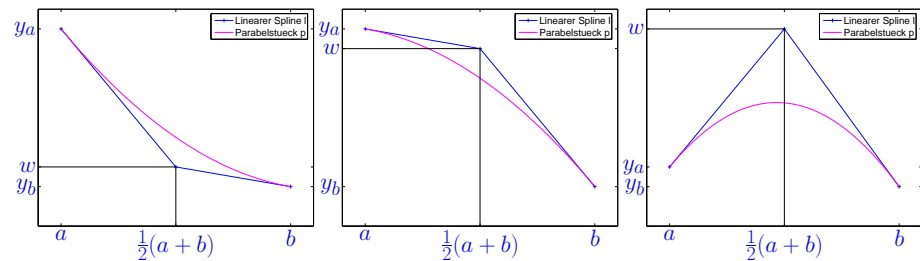
l „erbt“ lokale Monotonie/Krümmung der Daten

③ Lokale quadratische Approximation/Interpolation von l :

g : Streckenzug (linearer Spline) durch $(a, y_a), (\frac{1}{2}(a+b), w), (b, y_b), a < b, y_a, y_b, w \in \mathbb{R}$.

$$\text{Parabel: } p(t) := (y_a(b-t)^2 + 2w(t-a)(b-t) + y_b(t-a)^2)/(b-a)^2, \quad a \leq t \leq b.$$

$$p(a) = y_a, \quad p(b) = y_b, \quad p'(a) = g'(a), \quad p'(b) = g'(b).$$



4.3
p. 509

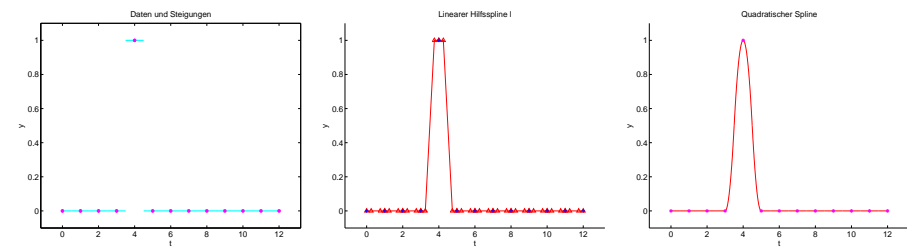
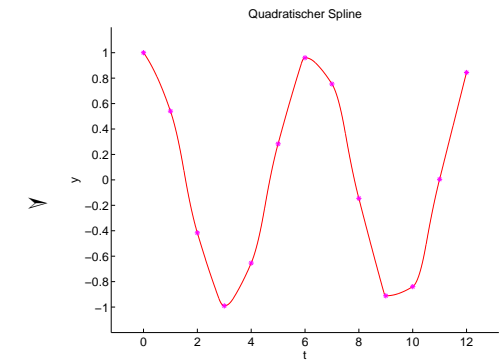
g monoton wachsend/fallend $\Rightarrow p$ monoton wachsend/fallend

g konvex/konkav $\Rightarrow p$ konvex/konkav

4.3
p. 511

Fortsetzung von Bsp. 202: ex:formspline

Interpolierender quadratischer Spline



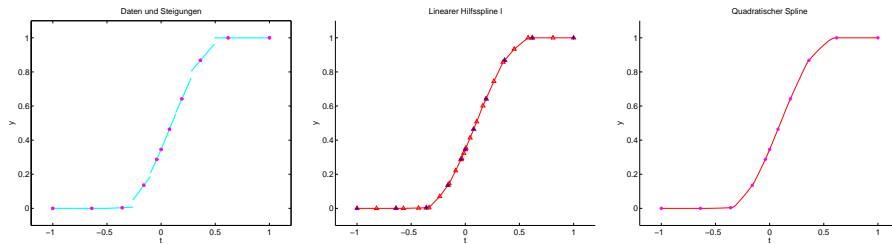
4.3
p. 510

4.3
p. 512

Formerhaltende quadratische Splineinterpolation = lokal + nichtlinear

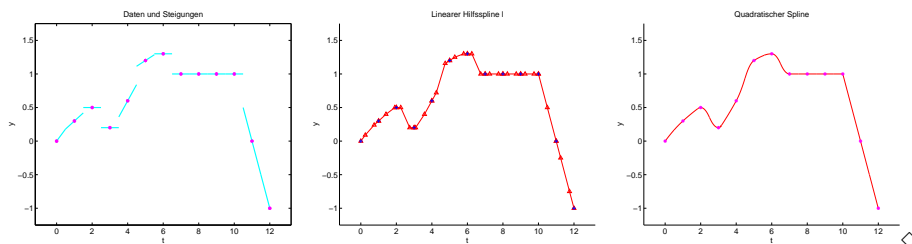
Beispiel 203 (Formerhaltende quadratische Splineinterpolation).

Daten aus Bsp. 194: `ex:logpolyintpol`



Daten aus [31]:

t_i	0	1	2	3	4	5	6	7	8	9	10	11	12
y_i	0	0.3	0.5	0.2	0.6	1.2	1.3	1	1	1	0	-1	



4.3.3 Bezier-Techniken

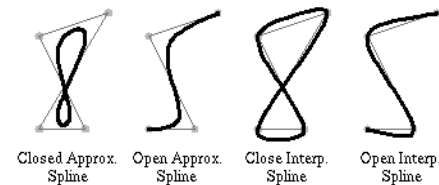
Ziel: Kurvenapproximation (nicht Interpolation) durch stückweise Polynome

Aus dem XFIG-Manual (<http://www.xfig.org/>):

About Spline Curves

A Spline curve is a smooth curve controlled by specified points.

- **CLOSED APPROXIMATING SPLINE**: Smooth closed curve which approximates specified points.
- **OPEN APPROXIMATING SPLINE**: Smooth curve which approximates specified points.
- **CLOSED INTERPOLATING SPLINE**: Smooth closed curve which passes through specified points.
- **OPEN INTERPOLATING SPLINE**: Smooth curve which passes through specified points.



Using splines, curves such as the following may be easily drawn.



4.3
p. 513

Ein Beweis des Weierstrassschen Approximationssatzes (\rightarrow Analysis, [9, Sekt. 6.2]):

Theorem 4.3.7 (Approximation durch Bernstein-Polynome). Falls $f \in C([0, 1])$ und

$$p_n(t) := \sum_{j=0}^n f(j/n) \binom{n}{j} t^j (1-t)^{n-j}, \quad 0 \leq t \leq 1,$$

dann $p_n \rightarrow f$ gleichmäßig für $n \rightarrow \infty$. Wenn $f \in C^m([0, 1])$, dann $p_n^{(k)} \rightarrow f^{(k)}$ gleichmäßig für $n \rightarrow \infty$.

Beispiel 204 (Bernstein-Approximation).

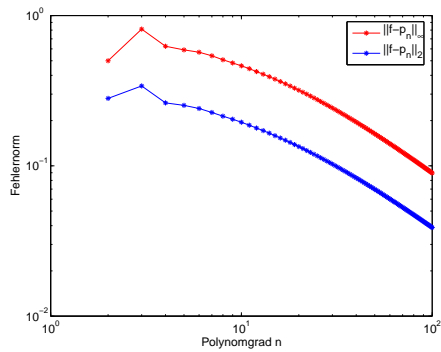
$$f_1(t) := \begin{cases} 0 & , \text{ falls } |2t - 1| > \frac{1}{2} \\ \frac{1}{2}(1 + \cos(2\pi(2t - 1))) & \text{sonst.} \end{cases}, \quad f_2(t) := \frac{1}{1 + e^{-12(x-1/2)}}.$$

Normen des Approximationsfehlers $f - p_n, p_n$ aus Thm. 4.3.7

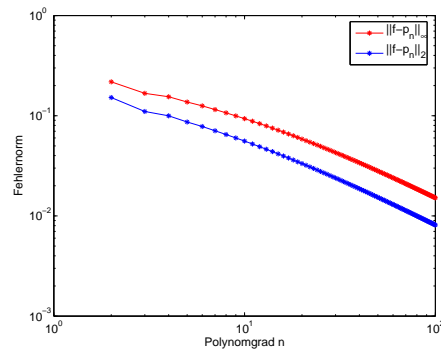
4.3
p. 514

4.3
p. 515

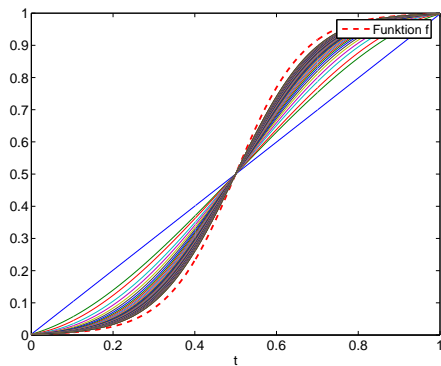
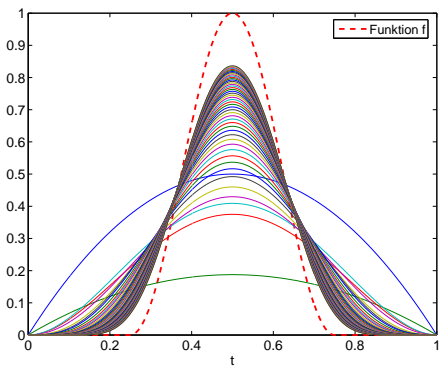
4.3
p. 516



Bernstein-Approximation von f_1



Bernstein-Approximation von f_2



➤ Schlechte Approximation, aber gute „Formwiedergabe“ durch p_n

Aus [9, Sect. 6.3]: Monotonic and convex functions yield monotonic and convex approximants, respectively. In a work, the Bernstein approximants mimic the the behavior of the function to a remarkable degree. There is a price that must be paid for these beautiful approximation properties: the convergence of Bernstein polynomials is very slow.

It is far slower than what can be achieved by other means. If f is bounded, then at a point t where $f'(t)$ exists and does not vanish, $p_n(t)$ converges to $f(t)$ precisely like C/n . This fact seems to

have precluded any numerical application of Bernstein polynomials from having been made (1965!). Perhaps they will find application when the properties of the approximant in the large are of more importance than closeness of the approximation.

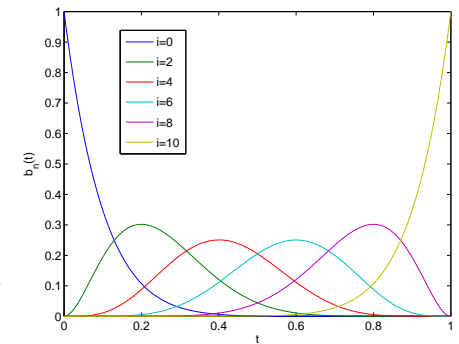
Definition 4.3.8 (Bernstein-Polynome).

Bernstein-Polynome vom Grad n :

$$b_{i,n}(t) := \binom{n}{i} t^i (1-t)^{n-i}, \quad i = 0, \dots, n.$$

$$(b_{i,n} \equiv 0 \text{ für } i < 0, i > n)$$

Bernstein-Polynome für $n = 10$



4.3
p. 517

4.3
p. 519

Lemma 4.3.9 (Eigenschaften von Bernstein-Polynomen).

- (i) $t = 0 = i$ -fache Nullstelle von $b_{i,n}$, $t = 1 = n - i$ -fache Nullstelle von $b_{i,n}$,
- (ii) $\sum_{i=0}^n b_{i,n} \equiv 1$ (Zerlegung der Eins),
- (iii) $b_{i,n}(t) = (1-t)b_{i,n-1} + tb_{i-1,n-1}(t)$,
- (iv) $b_{i,n}(t) \geq 0 \quad \forall 0 \leq t \leq 1$.

Lemma 4.3.9 (iii)

➤ Effiziente Auswertung von $b_{i,n}(t)$ (Rekursion)

```
function V = bernstein(n,x)
V = [ones(1,length(x));...
     zeros(n,length(x))];
for j=1:n
    for k=j+1:-1:2
        V(k,:) = x.*V(k-1,:) + (1-x).*V(k,:);
    end
    V(1,:) = (1-x).*V(1,:);
end
```

4.3
p. 518

4.3
p. 520

4.4 Numerische Quadratur

[File: section-numerische-quadratur.tex, SVN: section-numerische-quadratur.tex 1331 2007-03-07 09:16:14Z hiptmai]

NUMQUAD

= Näherungsweise Auswertung von $\int_{\Omega} f(x) dx$ (\neq Quadratur des Kreises)

$f : \Omega \subset \mathbb{R}^n \mapsto \mathbb{R} \hat{=} \text{function } y = f(x)$ (Nur Punktauswertungen verfügbar)

Spezialfall $n = 1$: $\Omega = [a, b]$ (Intervall)

QUADFORM

Quadraturformel: $\int_a^b f(t) dt \approx Q_n(f) := \sum_{j=0}^n \omega_j f(\xi_j)$.

ω_j : **Quadraturgewichte** $\in \mathbb{R}$ (engl. *weights*)
 ξ_j : **Quadraturknoten** $\in [a, b]$ (engl. *nodes*)

Fragestellung: (analog zur Interpolationsfehlerabschätzung, Abschnitt ^{sec:interp}4.1.2)

4.3
p. 521

Asymptotisches Verhalten $\int_a^b f(t) dt - Q_n(f)$ für $n \rightarrow \infty$ \triangleright Algebraische Konvergenz
 \triangleright Exponentielle Konvergenz

4.4
p. 523

4.4.1 Polynomiale Quadraturformeln



Idee: Ersetze $f \rightarrow p_n, p_n \in \mathcal{P}_n =$ Interpolationspolynom von f zu gegebener Knotenmenge $\{\xi_0, \dots, \xi_n\} \subset [a, b]$

$$\int_a^b f(t) dt \approx \int_a^b p_n(t) dt. \quad (4.4.1) \quad \text{eq:polqf}$$

Lagrange-Polynom: $L_i(t) := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - \xi_j}{\xi_i - \xi_j}, i = 0, \dots, n \triangleright p_n(t) = \sum_{j=0}^n f(\xi_j) L_j(t)$.

$$\int_a^b p_n(t) dt = \sum_{j=0}^n f(\xi_j) \int_a^b L_j(t) dt \triangleright \text{Gewichte } \omega_j := \int_a^b L_j(t) dt. \quad (4.4.2) \quad \text{eq:quadv}$$

Beispiel 205 (Newton-Cotes-Formeln).

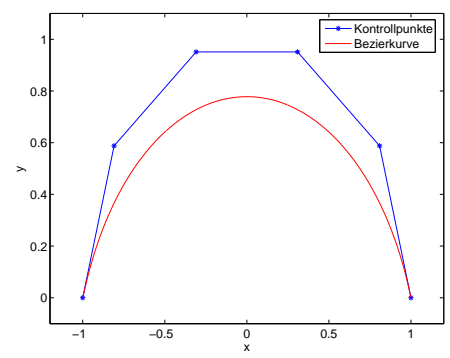
Äquidistante Knoten $v_j := a + hj, h := \frac{b-a}{n}, j = 0, \dots, n$:

4.4
p. 522

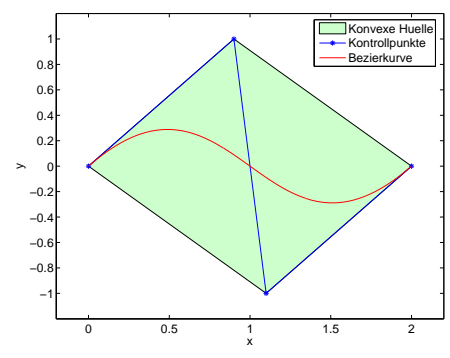
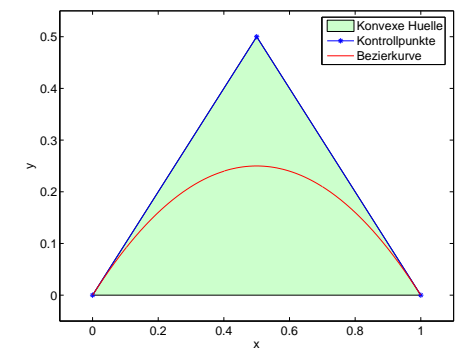
Symbolische Berechnung der Quadraturformeln mit MAPLE:

4.4
p. 524

Definition 4.3.10 (Bézier-Kurven). **Bézier-Kurve** zu **Kontrollpunkten** $\mathbf{d}_i \in \mathbb{R}^d, i = 0, \dots, n, d \in \mathbb{N}$:

$$\gamma : \begin{cases} [0, 1] \mapsto \mathbb{R}^d \\ t \mapsto \sum_{i=0}^n \mathbf{d}_i b_{i,n}(t) \end{cases}$$


- $\gamma(0) = \mathbf{d}_0, \gamma(1) = \mathbf{d}_n$ (Interpolationseigenschaft)
- $\gamma'(0) = n(\mathbf{d}_1 - \mathbf{d}_0), \gamma'(1) = n(\mathbf{d}_n - \mathbf{d}_{n-1})$ (Tangenten)



Lemma 4.3.9 (ii) $\triangleright \gamma \subset \text{convex}\{\mathbf{d}_0, \dots, \mathbf{d}_n\}$

Konvexe Hülle: $\text{convex}\{d_0, \dots, d_n\} := \{x = \sum_{i=0}^n \lambda_i \mathbf{d}_i, 0 \leq \lambda_i \leq 1, \sum_{i=0}^n \lambda_i = 1\}$.

- Praxis:
- Stückweise Bézier-Kurven \triangleright Spline-Kurven
 - X-Splines & NURBS (\rightarrow CAGD, Computergraphik)

```
> newtoncotes := n -> factor(int(interp([seq(i*h, i=0..n)],
    [seq(f(i*h), i=0..n)], z), z=0..n*h)):
```

• $n = 1$: Trapezregel

```
> trapez := newtoncotes(1);
```

$$\frac{h}{2}(f(0) + f(h)) \quad \left(= \frac{b-a}{2}(f(a) + f(b)) \right)$$

(4.4.3) [eq:trapez](#)

• $n = 2$: Simpson-Regel

```
> simpson := newtoncotes(2);
```

$$\frac{h}{3}(f(0) + 4f(h) + f(2h)) \quad \left(= \frac{b-a}{6}(f(0) + 4f(h) + f(2h)) \right)$$

(4.4.4) [eq:simps](#)

• $n = 4$: Milne-Regel

```
> milne := newtoncotes(4);
```

$$\frac{2}{45}h(7f(0) + 32f(h) + 12f(2h) + 32f(3h) + 7f(4h))$$

• $n = 6$: Weddle-Regel

```
> weddle := newtoncotes(6);
```

$$\frac{1}{140}h(41f(0) + 216f(h) + 27f(2h) + 272f(3h) + 27f(4h) + 216f(5h) + 41f(6h))$$

• $n = 8$: Quadraturformel mit *negativen* Gewichten

```
> newtoncotes(8);
```

$$\frac{4}{14175}h(989f(0) + 5888f(h) - 928f(2h) + 10496f(3h) - 4540f(4h) + 10496f(5h) - 928f(6h) + 5888f(7h) + 989f(8h))$$

Negative Quadraturgewichte beeinträchtigen numerische Stabilität !

Alternative:

Falls $\nu_j = \text{Tschebyscheff-Knoten}$ ➤ **Tschebyscheff-Quadratur**

Quadraturfehlerabschätzungen ← Direkt aus L^∞ -Interpolationsfehlerabschätzungen, Cor. 4.1.5: [cor:polintpe](#)

$$f \in C^{n+1}([a, b]) \Rightarrow \left| \int_a^b f(t) dt - Q_n(f) \right| \leq \frac{1}{(n+1)!} (b-a)^{n+1} \|f^{(n+1)}\|_{L^\infty([a, b])} \cdot (4.4.5) \quad \text{eq:polqu}$$

(Schärfere Abschätzungen für Tschebyscheff-Quadratur)

4.4.2 Gauss-Quadratur

Nach Konstruktion: Polynomiale Quadraturformeln [eq:polqf](#) (4.4.1) exakt für $f \in \mathcal{P}_n$

Mass für Qualität einer Quadraturformel Q_n :

$$\text{Ordnung}(q_n) := \max\{n \in \mathbb{N}_0: Q_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_n\} + 1$$

4.4
p. 525

4.4
p. 527

➤ Jede polynomiale Quadraturformel mit n Knoten ▷ Ordnung $\geq n$

Geht es auch besser ?

Beispiel 206 (Gauss-Legendre-Quadratur der Ordnung 4).

$$\text{Linearität} \Rightarrow Q_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_3 \Leftrightarrow Q_n(t^q) = \frac{1}{q+1}(b^{q+1} - a^{q+1}), \quad q = 0, 1, 2, 3.$$

4 Gleichungen für Gewichte ω_j , Knoten $\xi_j, j = 1, 2$ ($a = -1, b = 1$)

$$\begin{aligned} \int_{-1}^1 1 dt &= 2 = 1\omega_1 + 1\omega_2, & \int_{-1}^1 t dt &= 0 = \xi_1\omega_1 + \xi_2\omega_2 \\ \int_{-1}^1 t^2 dt &= \frac{2}{3} = \xi_1^2\omega_1 + \xi_2^2\omega_2, & \int_{-1}^1 t^3 dt &= 0 = \xi_1^3\omega_1 + \xi_2^3\omega_2. \end{aligned}$$

4.4
p. 526

4.4
p. 528


```
> 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);
```

$$\{\omega_2 = 1, \omega_1 = 1, \xi_1 = 1/3\sqrt{3}, \xi_2 = -1/3\sqrt{3}\}$$

► Quadraturformel: $\int_{-1}^1 f(x) dx \approx f\left(\frac{1}{\sqrt{3}}\right) + f\left(-\frac{1}{\sqrt{3}}\right)$

Heuristik: n Knoten + n Gewichte ► $2n$ Freiheitsgrade = $\dim \mathcal{P}_{2n-1}$

► Hoffnung: Quadraturformel der Ordnung $2n$ mit n Knoten

Herleitung: Polynomdivision mit Rest: für gegebenes $q \in \mathcal{P}_n$

$$p \in \mathcal{P}_{2n-1}: p(t) = h(t)q(t) + r(t), \quad h \in \mathcal{P}_{n-1}, r \in \mathcal{P}_{n-1}.$$

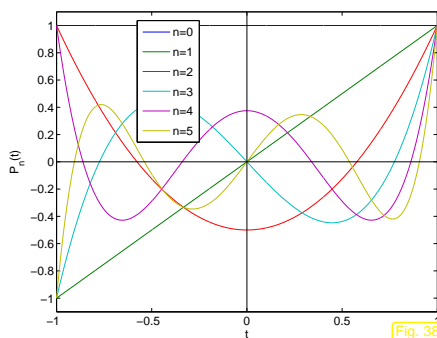
$$\int_a^b p(t) dt - \sum_{j=1}^n \omega_j p(\xi_j) = \int_a^b h(t)q(t) dt - \sum_{j=1}^n \omega_j h(\xi_j)q(\xi_j) + \int_a^b r(t) dt - \sum_{j=1}^n \omega_j r(\xi_j) \stackrel{!}{=} 0.$$

① Wähle $q \in \mathcal{P}_n$: $\int_a^b h(t)q(t) dt = 0 \quad \forall h \in \mathcal{P}_{n-1} \rightarrow \int_a^b h(t)q(t) dt = 0$

► q = Orthogonalpolynom zum Skalarprodukt $(f, g) \mapsto \int_a^b f(t)g(t) dt \rightarrow$ Abschnitt 4.1.5.2

Definition 4.4.1 (Legendre-Polynome). Das n . Legendre-Polynom P_n ist charakterisiert durch $P_n \in \mathcal{P}_n$, $\int_{-1}^1 P_n(t)q(t) dt = 0 \quad \forall q \in \mathcal{P}_{n-1}$, $P_n(1) = 1$.

Legendre-Polynome P_0, \dots, P_5



legendre

Fig. 38

Dreiterm-Rekursion analog zu Thm. 4.1.19: thm:orthpoly

$$P_{n+1}(t) = \frac{2n+1}{n+1} t P_n(t) - \frac{n}{n+1} P_{n-1}(t).$$

(4.4.6)

```
function V = legendre(N,x)
V = ones(size(x)); V = [V; x];
for n=1:N-1
    V = [V; ...
        (2*n+1)/(n+1).*x.*V(end,:)- ...
        n/(n+1)*V(end-1,:)]';
end
```

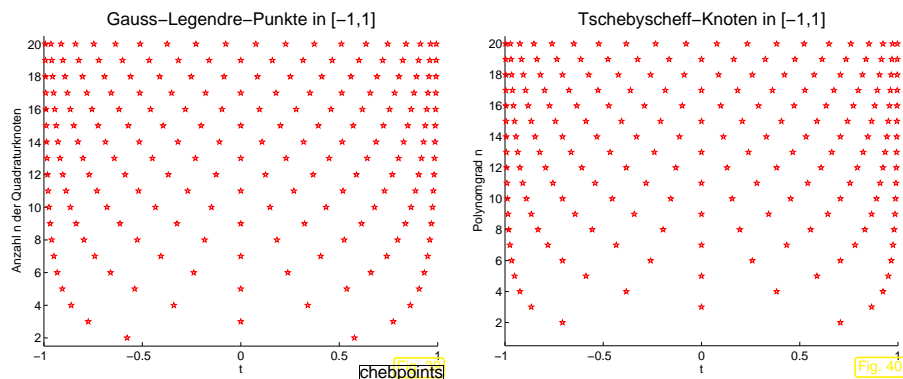
② Wähle Knoten ξ_1, \dots, ξ_n auf $[a, b]$ transformierte (vgl. (4.1.17)) Nullstellen von Legendre-Polynom P_n , $n \geq 1$ (→ Gauss-Punkte).

$$\rightarrow \sum_{j=1}^n \omega_j h(\xi_j) q(\xi_j) = 0, \quad \text{da } q = P_n!$$

Theorem 4.4.2 (Nullstellen von Orthogonalpolynomen).

P_n hat genau n einfache Nullstellen in $[-1, 1]$.

4.4
p. 529



4.4
p. 531

③ Für $[a, b] = [-1, 1]$ wähle $\omega_j \in \mathbb{R}$: $\sum_{j=1}^n \omega_j \xi_j^k = \frac{1}{k+1} (1 - (-1)^{k+1}), \quad k = 0, \dots, n-1$

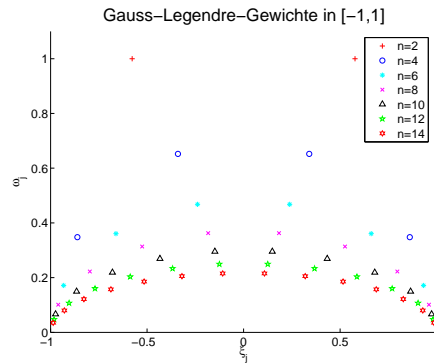
$$\rightarrow \int_a^b r(t) dt - \sum_{j=1}^n \omega_j r(\xi_j) = 0, \quad \text{da } r \in \mathcal{P}_{n-1}!$$

► Gauss-Legendre-Quadraturformeln: n Knoten & Ordnung $2n$

4.4
p. 530

4.4
p. 532

Theorem 4.4.3 (Gewichte der Gauss-Legendre-Quadraturformeln). Die Gewichte der Gauss-Legendre-Quadraturformeln sind eindeutig bestimmt und positiv.



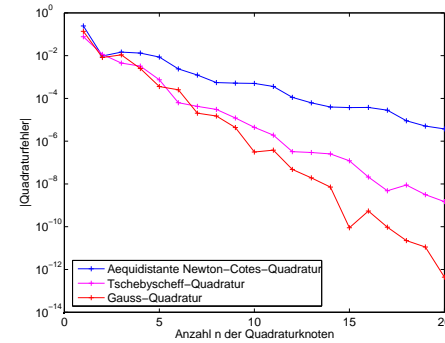
Bemerkung 207. Berechnung der Knoten/Gewichte der Gauss-Legendre-Quadraturformeln aus Eigenwertproblem!

(Golub-Welsch-Algorithmus [GGG05, Sect. 3.5.4])

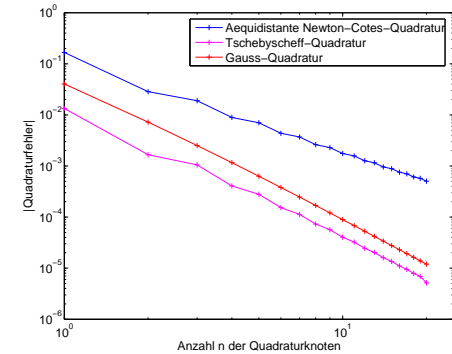
In Codes: ξ_j, w_j tabelliert

```
function [x,w]=gaussQuad(p)
b = 1:p-1;
b = b./sqrt(4*b.*b-1);
J=diag(b,-1)+diag(b,1);
[ev,ew]=eig(J);
for i=1:p
    ev(:,i)./norm(ev(:,i));
end
x=diag(ew);
w=(2*(ev(1,:).*ev(1,:)))';
```

Beispiel 208 (Asymptotisches Verhalten des Quadraturfehlers).



Quadraturfehler, $f_1(t) := \frac{1}{1+(5t)^2}$ auf $[0, 1]$



Quadraturfehler, $f_2(t) := \sqrt{t}$ auf $[0, 1]$

Asymptotik des Quadraturfehlers $\epsilon_n := \left| \int_0^1 f(t) dt - Q_n(f) \right|$ für „ $n \rightarrow \infty$ “:

- Exponentielle Konvergenz $\epsilon_n \approx O(q^n)$, $0 < q < 1$, für **analytischen Integranden** (→ Def. 4.1.14) f_1 : Newton-Cotes-Quadratur: $q \approx 0.61$, Tschebyscheff-Quadratur: $q \approx 0.40$, Gauss-Legendre-Quadratur: $q \approx 0.27$

4.4
p. 533

- Algebraische Konvergenz $\epsilon_n \approx O(n^{-\alpha})$, $\alpha > 0$, für Integranden f_2 mit **Singularität** in $t = 0$: Newton-Cotes-Quadratur: $\alpha \approx 1.8$, Tschebyscheff-Quadratur: $\alpha \approx 2.5$, Gauss-Legendre-Quadratur: $\alpha \approx 2.7$

4.4
p. 535

Geht es noch besser als Gauss-Quadratur? (etwa Ordnung $2n+1$ mit n Knoten)

Zu Knoten $a \leq \xi_1 < \xi_2 < \dots < \xi_n \leq b$ ➤ $p(t) := \prod_{j=1}^n (t - \xi_j)^2 \in \mathcal{P}_{2n}$

$$\blacktriangleright Q_n(p) = 0 \text{ aber } \int_a^b p(t) dt > 0.$$

- Ordnung $> 2n$ kann mit n Quadraturpunkten nicht erreicht werden

4.4.3 Zusammengesetzte Quadraturformeln

Analogon: Globale Polynominterpolation \longleftrightarrow stückweise Polynominterpolation
(→ Abschnitt 4.3.1)

4.4
p. 534

4.4
p. 536



- Idee: Partitionierung des Integrationsintervalls $[a, b]$ durch Gitter (\rightarrow Abschnitt 4.3) $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$ Anwendung der Quadraturformeln aus Abschnitt 4.4.1 auf Teilintervalle $I_j := [x_{j-1}, x_j], j = 1, \dots, m$, + Summation.

(Wir betrachten nur zusammengesetzte Quadraturformeln basierend auf einer einzigen lokalen Quadraturformel)

Konvergenzuntersuchung (Skalierungsargument, vgl. (4.1.16)) und (4.4.5)):

Asymptotik des Quadraturfehlers bzgl. Maschenweite h von \mathcal{M} :

Theorem 4.4.4 (Konvergenz zusammengesetzter Quadraturformeln). Für eine zusammengesetzte Quadraturformel Q basierend auf einer lokalen Quadraturformel der Ordnung m gilt

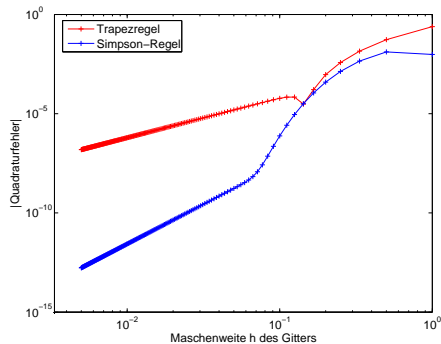
$$\exists C > 0: \left| \int_I f(t) dt - Q(f) \right| \leq C h^m \|f^{(m)}\|_{L^\infty(I)} \quad \forall f \in C^m(I), \forall \mathcal{M}.$$

Beispiel 209 (Quadraturfehler bei zusammengesetzten Regeln).

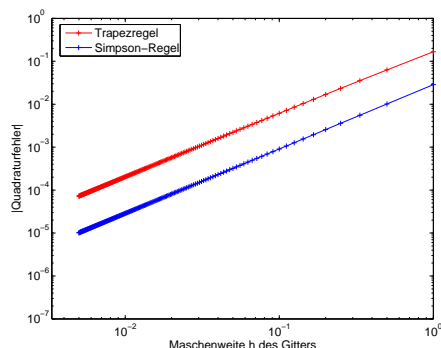
Zusammengesetzte Quadraturformeln basierend auf

- Trapezregel (4.4.3) \rightarrow Lokale Ordnung 2 (exakt für lineare Funktionen),
- Simpson-Regel (4.4.4) \rightarrow Lokale Ordnung 3 (exakt für Parabeln)

auf äquidistantem Gitter $\mathcal{M} := \{jh\}_{j=0}^n, h = 1/n, n \in \mathbb{N}$.



Quadraturfehler, $f_1(t) := \frac{1}{1+(5t)^2}$ auf $[0, 1]$



Quadraturfehler, $f_2(t) := \sqrt{t}$ auf $[0, 1]$

Asymptotik des Quadraturfehlers $\epsilon_n := \left| \int_0^1 f(t) dt - Q_n(f) \right|$ für Maschenweite „ $h \rightarrow 0$ “

Algebraische Konvergenz $\epsilon = O(h^\alpha), \alpha > 0$:

- Hinreichend glatter Integrand f_1 : α gemäss Thm. 4.4.4, Trapezregel $\rightarrow \alpha = 2$, Simpson-Regel $\rightarrow \alpha = 4$!
- Singulärer Integrand f_2 : $\alpha = 3/2$ für Trapezregel & Simpson-Regel! (Glattheit des Integranden begrenzt Konvergenz)

Simpson-Regel: Ordnung = 4? MAPLE schafft Klarheit

```
> rule := 1/3*h*(f(2*h)+4*f(h)+f(0))
> err := taylor(rule - int(f(x),x=0..2*h),h=0,6);
```

$$err := \left(\frac{1}{90} (D^{(4)})(f)(0) h^5 + O(h^6) \right), h, 6$$

- Simpson-Regel hat sogar Ordnung 4!

4.4
p. 537

4.4
p. 539

Bemerkung 210 (Auflösung einer Singularität durch Transformation).

Für $f \in C^\infty([0, b])$ approximiere $\int_0^b \sqrt{t} f(t) dt$ durch numerische Quadratur (\rightarrow Bsp. 209)

$$\text{Substitution } s = \sqrt{t}: \int_0^b \sqrt{t} f(t) dt = \int_0^{\sqrt{b}} 2s^2 f(s^2) ds.$$

Dann: Anwendung einer Quadraturformel auf glatten Integranden

Die äquidistante Trapezregel:

$\hat{=}$ stückweise lineare Approximation des Integranden auf äquidistantem Gitter:

$$\int_a^b f(t) dt \approx T_n(f) := h \left(\frac{1}{2} f(a) + \sum_{k=1}^{n-1} f(kh) + \frac{1}{2} f(b) \right), \quad h := \frac{b-a}{n}.$$

Beispiel 211 (Konvergenz für äquidistante Trapezregel).

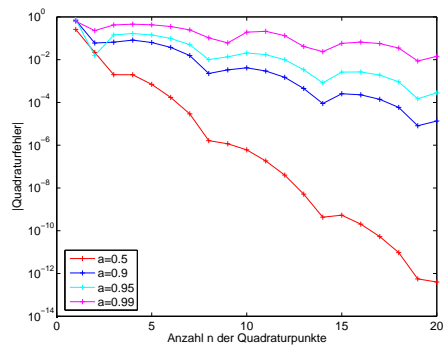
1-periodischer analytischer Integrand (\rightarrow Bsp. 166).

$$f(t) = \frac{1}{\sqrt{1 - a \sin(2\pi t - \gamma)}}, \quad 0 < a < 1, \gamma \in \mathbb{R}.$$

4.4
p. 538

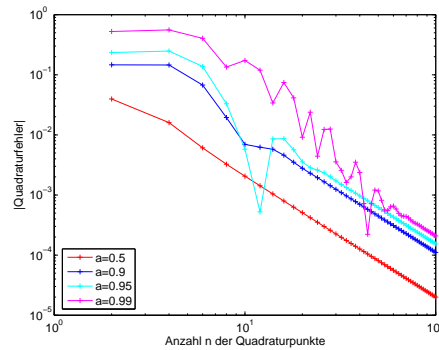
4.4
p. 540

(„Exakte Integralwerte“ aus T_{500})



Quadraturfehler für $T_n(f)$, $\gamma = 1$, auf $[0, 1]$

Exponentielle Konvergenz



Quadraturfehler für $T_n(f)$, $\gamma = 1$, auf $[0, \frac{1}{2}]$

Nur algebraische Konvergenz

◇

Erklärung:

$$f(t) = e^{2\pi i k t} \rightarrow \begin{cases} \int_0^1 f(t) dt = \begin{cases} 0, & \text{falls } k \neq 0, \\ 1 & \text{für } k = 0. \end{cases} \\ T_n(f) = \frac{1}{n} \sum_{l=0}^{n-1} e^{\frac{2\pi i l k}{n}} = \begin{cases} 0, & \text{falls } k \notin n\mathbb{Z}, \\ 1 & \text{falls } k \in n\mathbb{Z}. \end{cases} \end{cases}$$

Äquidistante Trapezregel T_n ist exakt für trigonometrische Polynome (→ Def. 4.1.8) vom Grad $< 2n$!

$T_n \hat{=}$ Approximation des Integranden durch trigonometrisches Interpolationspolynom → Abschn. 4.1.4

Glatte periodische Integranden ➤ Benutze äquidistante Trapezregel !

Übung 4.5. Was versteht man unter Tschebyscheff-Quadratur? Gib die Quadraturknoten an und bestimme numerisch eine Näherung für die Quadraturgewichte unter Verwendung der MATLAB-Funktion `intpolyval(t,y,x)` aus Abschnitt 4.1.2 und der standard MATLAB-Funktion `quad`.

Übung 4.6. Die Integralwerte

$$\int_0^1 \sqrt{1-x^2} f(x) dx \quad (4.6.1) \quad \text{eq:logqu}$$

sollen effizient, d.h. mit möglichst wenig Auswertungen von f , numerisch approximiert werden. Dabei darf $f \in C^\infty([0, 1])$ angenommen werden.

- Wie ist das Integral vor Anwendung einer Quadraturformel geeignet zu transformieren? (Hinweis: Verwende $1-x^2 = (1-x)(1+x)$)
- Schreibe eine auf Gauss-Quadratur basierende effiziente MATLAB-Funktion `logquad(f,p)` ($f \hat{=}$ handle auf Funktion f , $2p \hat{=}$ Ordnung der zugrundeliegenden Gauss-Quadratur) zur Approximation von (4.6.1). Die Routine `gaussQuad(p)` aus Bemerkung der Vorlesungsunterlagen darf verwendet werden.
- Für $f(x) = \sqrt{1+x}$ plote den Quadraturfehler in Abhängigkeit von p für $p = 1, \dots, 20$. Welche qualitative Konvergenz kann abgelesen werden?

4.4
p. 541

4.6.1 Adaptive Quadratur

Betrachte zusammengesetzte Trapezregel auf Gitter $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$

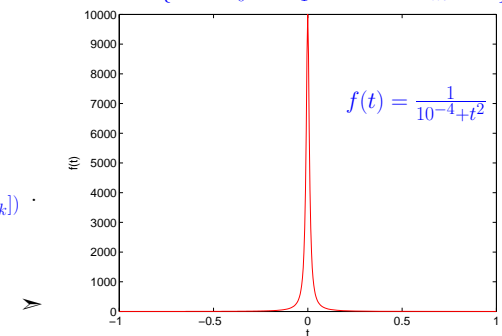
Lokaler Quadraturfehler (für $f \in C^2([a, b])$):

$$\int_{x_{k-1}}^{x_k} f(t) dt - \frac{1}{2}(f(x_{k-1}) + f(x_k)) \leq (x_k - x_{k-1})^3 \|f''\|_{L^\infty([x_{k-1}, x_k])}$$

➤ Benutze nichtäquidistante Gitter !

(fein, wo $|f''|$ gross)

Sinnvoll etwa für „Nadelfunktion“



Ziel: Gleiche Beiträge aller Gitterintervalle zum Gesamtfehler
Werkzeug: Lokale **a posteriori** Fehlerschätzung
 (Schätzen der Fehlerbeiträge von Gitterintervallen aus Zwischenergebnissen)
Technik: Lokale Gitterverfeinerung

4.6
p. 542

4.6
p. 543

4.6
p. 544

Adaptive Mehrgitterquadratur (vereinfacht nach [T2, Abschnitt 9.7])



Idee: Vergleich der Resultate zweier lokaler Quadraturformeln Q_1, Q_2 unterschiedlicher Ordnung \rightarrow lokale Fehlerschätzung

Heuristik: Fehler(Q_2) \ll Fehler(Q_1) \Rightarrow Fehler(Q_1) $\approx Q_2(f) - Q_1(f)$.

Hier: Trapezregel \leftrightarrow Simpson-Regel

Gegeben: Gitter $\mathcal{M} := \{a = x_0 < x_1 < \dots < x_m = b\}$

❶ (Fehlerschätzung) Für $I_k = [x_{k-1}, x_k]$, $k = 1, \dots, m$ ($m_k := \frac{1}{2}(x_{k-1} + x_k)$)

$$EST_k := \underbrace{\frac{h_k}{6}(f(x_{k-1}) + 4f(m_k) + f(x_k))}_{\text{Simpson-Regel}} - \underbrace{\frac{h_k}{4}(f(x_{k-1}) + 2f(m_k) + f(x_k))}_{\text{Trapezregel (auf unterteilter Gitterzelle)}}$$

❷ (Abbruch) Simpson-Regel auf $\mathcal{M} \Rightarrow$ Vorläufiger Integralwert I

Wenn $\sum_{k=1}^m EST_k \leq TOL \cdot I$ ($TOL :=$ vorgegebene Toleranz) \Rightarrow **STOP**

❸ (Lokale Gitterverfeinerung)

$$S := \{k \in \{1, \dots, m\} : EST_k \geq \eta \cdot \frac{1}{m} \sum_{j=1}^m EST_j\}, \quad \eta \approx 0.9.$$

Neues Gitter $\mathcal{M}^* := \mathcal{M} \cup \{m_k : k \in S\}$.

Dann weiter mit ❶ und $\mathcal{M} \leftarrow \mathcal{M}^*$.

Adaptive Quadratur in MATLAB (nach [GAG00, 17]):

`q = quad(fun,a,b,tol);` Adaptive Mehrgitterquadratur
(lokale Quadraturformeln niedriger Ordnung)

`q = quadl(fun,a,b,tol);` Adaptive Gauss-Quadratur
(lokale Gauss-Quadraturformeln & Ordnungserhöhung)

4.6.2 Numerische Berechnung oszillatorischer Integrale

4.7 Multiskalenbasen

[File: section-multiskalenbasen.tex, SVN: section-multiskalenbasen.tex 1010 2006-09-11 15:44:43Z hiptmair]

Ziel: Analyse von Daten (\leftrightarrow Funktionen) durch **Darstellung in geeigneter Basis**

Beispiel 212 (Nichtlokalität der DFT).

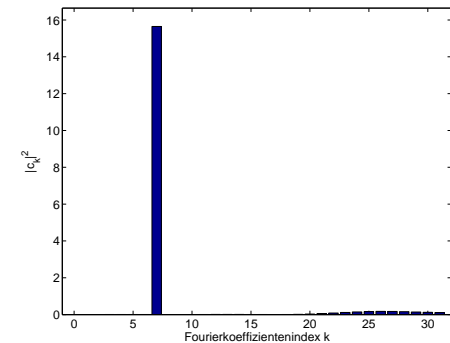
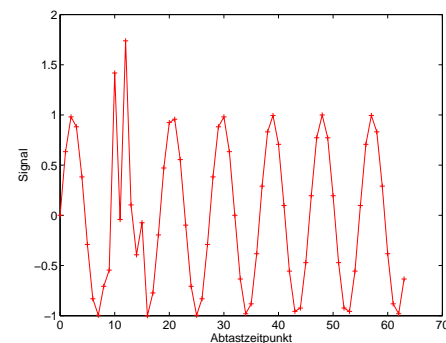
Diskrete Fouriertransformation (\rightarrow Abschn. 3.7.1, Bsp. 145) ^{sec:diskr-fourfreqfilt}

- Gute Identifikation von Frequenzanteilen
- Keine Auflösung zeitlokaler Effekte

Daten:

```
t = 0:63; y = sin(7*2*pi*t/64);
y(10:16) = y(10:16) + cos(10*t(10:16));
```

4.6
p. 545



Gegeben: Abtastzeitpunkte $t_j := j/n$, $j = 0, \dots, n-1$, $n \in \mathbb{N}$ (äquidistant), Abtastwerte y_j ,
 $j = 0, \dots, n-1 \Leftrightarrow$ Vektor $\mathbf{y} \in \mathbb{R}^n$

4.7
p. 546

4.7
p. 547

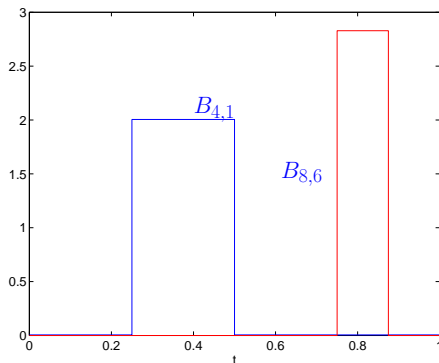
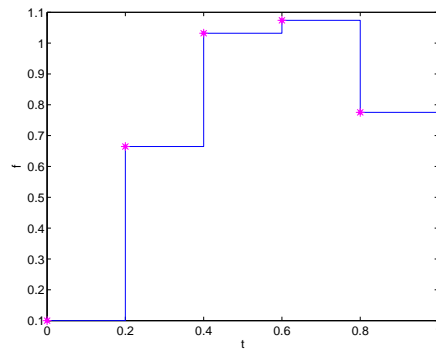
4.7
p. 548

Vektor $y \in \mathbb{R}^n$

► Assoziierte Treppenfunktion

$$f(t) = y_j \quad \text{für } t_j \leq t < t_{j+1} .$$

(mit $t_n := 1$)



►
$$f(t) = \sum_{j=0}^{n-1} \left(\frac{1}{\sqrt{n}} y_j \right) B_{n,j}(t) ,$$

$$B_{n,j}(t) := \sqrt{n} B(nt - j) ,$$

$$B(t) := \begin{cases} 1 & , \text{ falls } 0 \leq t < 1 , \\ 0 & \text{sonst.} \end{cases}$$

Beachte: $\bullet \text{supp}(B(n \cdot - j)) = [j/n, (j+1)/n]$
 $\bullet \|B_{n,j}\|_{L^2([0,1])} = 1$

$B \doteq$ „Grundbasisfunktion“ (Skalierungsfunktion)

Terminologie:

$$B(nt - j) = B\left(n\left(t - \frac{j}{n}\right)\right)$$

Dilatation Translation um j/n

► $\{B_{n,j}\}_{j=0}^{n-1} = L^2\text{-orthonormale Basis des Treppenfunktionsraums } \mathcal{C}_n \text{ auf äquidistantem Gitter}$
 $\mathcal{T} := \{j/n, j = 0, \dots, n\}$ in $[0, 1]$.

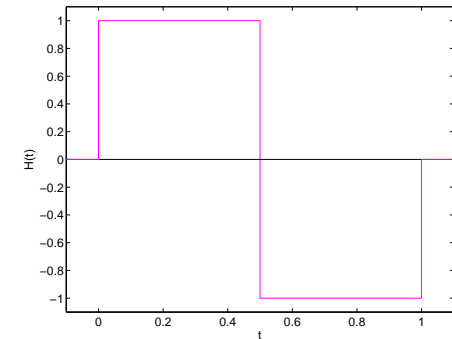
Beachte: $\mathcal{C}_n = \mathcal{S}_{0,\mathcal{T}}$ (Splinaum niedrigster Ordnung, \rightarrow Def. 4.3.5) der Spline

Eine weitere L^2 -Orthonormalbasis (für $n = 2^L, L \in \mathbb{N}$):

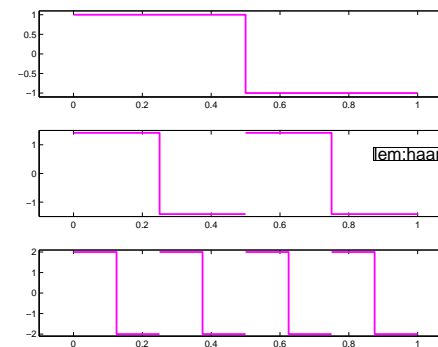
Definition 4.7.1 (Haar-Wavelet). *Haar-Wavelet* auf $[0, 1]$:

$$H(t) := \begin{cases} 1 & \text{für } 0 \leq t < \frac{1}{2} , \\ -1 & \text{für } \frac{1}{2} \leq t < 1 , \\ 0 & \text{sonst.} \end{cases}$$

$$= B(2t) - B(2t - 1) .$$



4.7
p. 549



Translierte/dilatierte Haar-Wavelets:

$$H_{l,j}(t) := 2^{l/2} H(2^l t - j) , \quad j = 0, \dots, 2^l - 1 .$$

Lemma 4.7.2.

$$\bigcup_{l=0}^{L-1} \bigcup_{j=0}^{2^l-1} \{H_{l,j}\} \cup \{B\}$$

ist eine L^2 -Orthonormalbasis von \mathcal{C}_{2^L} , die *Haar-Waveletbasis*.

► Waveletdarstellung von $f \in \mathcal{C}_{2^L}$: $f = c_0 B + \sum_{l=0}^{L-1} \sum_{j=0}^{2^l-1} c_{l,j} H_{l,j} .$

Terminologie: $c_{l,j} =$ Waveletkoeffizienten

$$\sum_{j=0}^{2^l-1} c_{l,j} H_{l,j} = \text{Level } l \text{ der Waveletzerlegung}$$

4.7
p. 550

4.7
p. 551

4.7
p. 552

Folgerung aus Orthonormalbasiseigenschaft:

$$\sum_{j=0}^{n-1} y_j B_{n,j} = c_0 B + \sum_{l=0}^{L-1} \sum_{j=0}^{2^l-1} c_{l,j} H_{l,j} \Rightarrow \sum_{j=0}^{n-1} y_j^2 = c_0^2 + \sum_{l=0}^{L-1} \sum_{j=0}^{2^l-1} c_{l,j}^2 \cdot \overset{\text{lem: dft}}{1}$$

↔ Unitarität der (skalierten) diskreten Fouriertransformation (→ Lemma 3.7.2)

Algorithmen zur Basistransformation Boxfunktionsbasis ↔ Haar-Waveletbasis:

Verfeinerungsrelationen:

$$\begin{aligned} B_{2^l,j} &= \frac{1}{2}\sqrt{2} \begin{pmatrix} B_{2^{l+1},2j} + B_{2^{l+1},2j+1} \\ H_{l,j} = \frac{1}{2}\sqrt{2} \begin{pmatrix} B_{2^{l+1},2j} - B_{2^{l+1},2j+1} \end{pmatrix} \end{pmatrix}, \Leftrightarrow \begin{pmatrix} B_{2^{l+1},2j} = \frac{1}{2}\sqrt{2} (B_{2^l,j} + H_{l,j}) \\ B_{2^{l+1},2j+1} = \frac{1}{2}\sqrt{2} (B_{2^l,j} - H_{l,j}) \end{pmatrix} \end{aligned} \quad (4.7.1) \quad \text{eq:refrel}$$

$$g = y_1 B_{2^{l+1},2j} + y_2 B_{2^{l+1},2j+1} \Rightarrow g = \frac{1}{2}\sqrt{2} \left(\underbrace{(y_1 + y_2)}_{\text{Tiefpassfilter}} B_{2^l,j} + \underbrace{(y_1 - y_2)}_{\text{Hochpassfilter}} H_{l,j} \right) \quad (4.7.2) \quad \text{eq:filter}$$

Lineare Anordnung der Waveletkoeffizienten $c_{l,j}$ für $f = c_0 B + \sum_{l=0}^{L-1} \sum_{j=0}^{2^l-1} c_{l,j} H_{l,j}$:

$$c_0, c_{0,0}, c_{1,0}, c_{1,1}, c_{2,0}, \dots, c_{2,3}, \dots, c_{l,0}, \dots, c_{l,2^l-1}, \dots, c_{L-1,0}, \dots, c_{L-1,2^{L-1}-1} \quad (4.7.3) \quad \text{eq:order}$$

Boxfunktionsbasis → Haar-Waveletbasis:

MATLAB-CODE: Basistransformation

```
function c = boxbastohwf(y,L)
c = zeros(size(y)); c(1) = y(1);
for l=2.^(L-1:-1:0)
    y = y/sqrt(2); m = 2*1;
    c(1+l:m)=y(1:2:m)-y(2:2:m);
    y(1:l) = y(1:2:m)+y(2:2:m);
end
```

Rekursiver *bottom-up* Algorithmus:

(für $\text{length}(y) = 2^L$)

(Sukzessive Berechnung der Waveletkoeffizienten auf Level l und Boxfunktionsdarstellung auf Level $l-1$)

Rechenaufwand $O(2^L)$

Haar-Waveletbasis → Boxfunktionsbasis:

MATLAB-CODE: Basistransformation

```
function y = hwftoboxbas(c,L)
y = zeros(size(c));
y(1) = c(1); ofs = 1;
for l=2.^(0:L-1)
    z = y(1:l); w=c(ofs+(1:l));
    y(1:2:2*l) = (z+w)/sqrt(2);
    y(2:2:2*l) = (z-w)/sqrt(2);
    ofs = ofs+1;
end
```

Rekursiver *top-down* Algorithmus:

(für $\text{length}(y) = 2^L$)

(Sukzessive Berechnung der Boxfunktionsdarstellung auf Level $0, 1, 2, \dots$)

Rechenaufwand $O(2^L)$

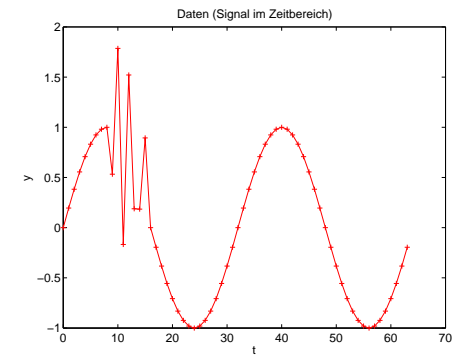
Beispiel 213 (Waveletzerlegung eines Signals). Wavelet! Zerlegung eines Signals

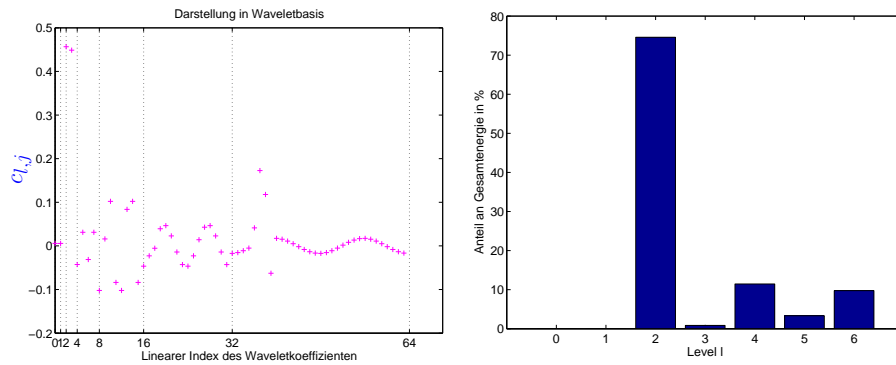
Signal (→ Bsp. 212): `lex:dftnonloc`

```
t = 0:63; y = sin(4*pi*t/64);
y(10:16) = y(10:16) + cos(10*t(10:16));
```

Signal mit

- dominierendem Grundfrequenzanteil
- zeitlokaler Störung



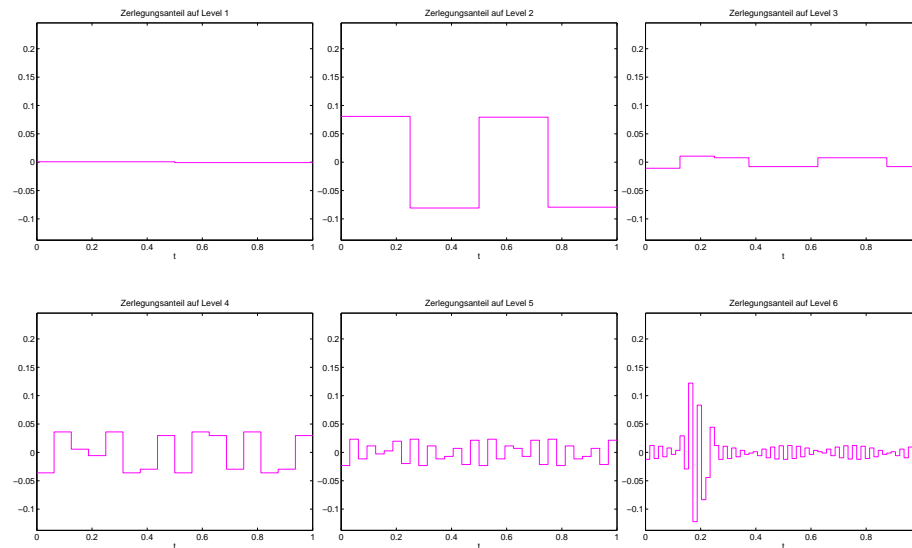


Leistung auf Level l :

$$P_l := \sum_{j=0}^{2^l-1} c_{l,j}^2 \quad \Rightarrow \quad \sum_{l=0}^L P_l = \|\mathbf{y}\|_2^2$$

➤ Auflösung der „Größenordnung“ der Grundfrequenz im Leistungsspektrum
(→ Koeffizientenindex l ↔ Frequenzzерlegung)

Wavelet-Zerlegungsanteile:



➤ Zeitliche Auflösung der Störung (→ Koeffizientenindex k)

Raum der Treppenfunktionen = Splineräum vom Grad 0/Ordnung 1

Verallgemeinerung ?

► Spline-Wavelets (7) (Komplizierte Konstruktionen)
(Herausforderung: Kleine Träger, Orthogonalität ↔ Stabilität)

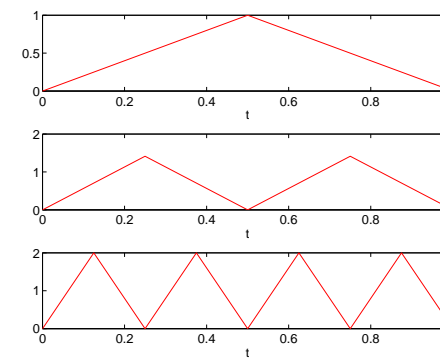
Hierarchische Basis für $\mathcal{S}_{1,\mathcal{M}}$:

Äquidistante Knoten in $[0, 1]$: $\mathcal{M} := \{t_j := j/n, j = 0, \dots, n\}, n = 2^L, L \in \mathbb{N}$

Skalierungsfunktion („Hutfunktion“) $\varphi(t) = \begin{cases} 1-t & , \text{ falls } 0 \leq t \leq 1, \\ 1+t & , \text{ falls } -1 \leq t \leq 0, \\ 0 & \text{sonst.} \end{cases} \in C^0(\mathbb{R}), \text{ stückweise linear.}$

$\varphi_{l,n}(t) := 2^{l/2} \varphi(2^l t - j) \Rightarrow \{\varphi_{l,j}\}_{j=0}^{2^l-1} = \text{Basis von } \mathcal{S}_{1,\mathcal{M}} \text{ auf } [0, 1] \text{ (Knotenbasis)}$

4.7
p. 557



Hierarchische Basis von $\mathcal{S}_{1,\mathcal{M}}$:

$$\{\varphi_{0,0}, \varphi_{0,1}\} \cup \bigcup_{l=1}^L \bigcup_{j=0}^{2^{l-1}-1} \{\varphi_{l,2j+1}\}. \quad (4.7.4)$$

Hierarchische-Basis-Zerlegung:

$$s \in \mathcal{S}_{1,\mathcal{M}} \Rightarrow s = c_0 \varphi_{0,0} + c_1 \varphi_{0,1} + \sum_{l=1}^L \sum_{j=0}^{2^{l-1}-1} c_{l,j} \varphi_{l,2j+1}.$$

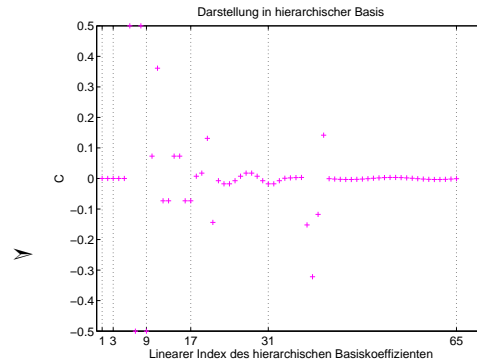
Beispiel 214 (Hierarchische-Basis-Zerlegung).

4.7
p. 558

4.7
p. 559

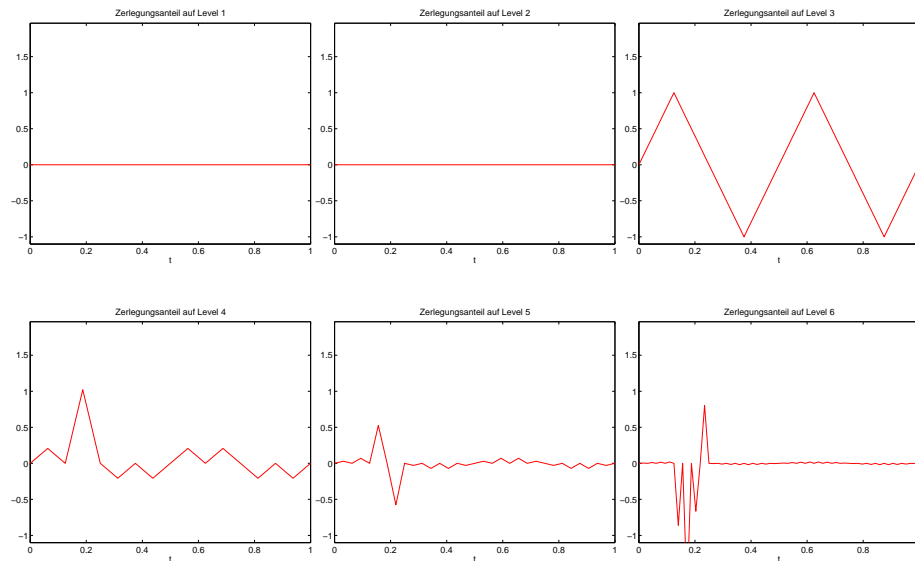
4.7
p. 560

$s \in \mathcal{S}_{1,\mathcal{M}}$, Knotenwerte $s(t_j)$ aus
 $t = 0:63$; $y = \sin(4 \cdot \pi \cdot t / 64)$;
 $y(10:16) = y(10:16) +$
 $\cos(10 \cdot t(10:16))$;
 Signal aus \rightarrow Bsp. 213
 Hierarchische-Basis-Koeffizienten $c_{l,j}$
 (Anordnung wie in (4.7.3))



➤ Auflösung von Grundschwingungen und zeitlokalen Merkmalen

Hierarchische-Basis-Zerlegungsanteile:



Beispiel 215 (L^2 -Instabilität der hierarchischen Basis).

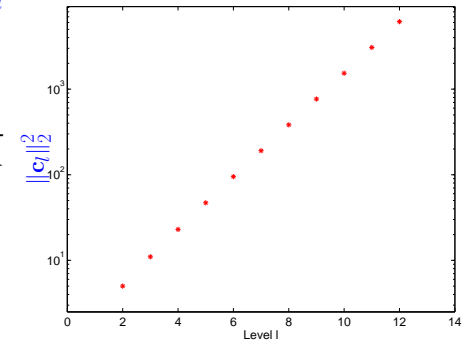
Auf Gitter $\mathcal{M}_l := \{j/n, j = 0, \dots, n\}$, $n := 2^l$
 wähle

$$s_l := \varphi_{l,2^{l-1}}.$$

Euklidische Norm der hierarchischen-Basis-
 Koeffizienten (4.7.4) von s_l

$$\|c\|_2^2 := c_0^2 + c_1^2 + \sum_{l=1}^L \sum_{j=0}^{2^{l-1}-1} c_{l,j}^2.$$

Exponentielles Wachstum von $\|c_l\|_2^2$!



➤ Hierarchische Basis *nicht* l -gleichmäßig stabil (im L^2 -Sinn)

4.7
p. 561

4.7
p. 563

5

Numerik Gewöhnlicher Differentialgleichungen

gewöhnlicher-differentialgleichungen.tex, SVN: chapter-numerik-gewoehnlicher-differentialgleichungen.tex 1010 2006-09-11 15:44:43Z hiptmair

Gegeben: • **Phasenraum** (engl. *phase space*) $\Omega \subset \mathbb{R}^d$, $d \in \mathbb{N}$ ($y \in \Omega$ = „Zustand“, engl. *state*)
 ➤ **Erweiterter Phasenraum** $\hat{\Omega} \subset \mathbb{R} \times \Omega$ „Zeitvariable“

(Annahme: Tensorproduktstruktur $\hat{\Omega} = I \times \Omega$)

• **Stetige Funktion** $f : \hat{\Omega} \mapsto \mathbb{R}^d$ („rechte Seite“)

• **Anfangsdaten** $(t_0, y_0) \in \hat{\Omega}$, (y_0 = **Anfangswert** engl. *initial value*)

Gesucht: C^1 -Lösungsfunktion $y : J \mapsto \Omega$, $J \subset \mathbb{R}$ Intervall, $t_0 \in J$, für das **Anfangswertproblem** (AWP) für die **gewöhnliche Differentialgleichung** (engl. *ordinary differential equation*, **ODE**) [1. Ordnung]:

$$y(t_0) = y_0 \quad , \quad \dot{y} = f(t, y(t)) \quad \forall t \in J.$$

Notation: Punkt $\dot{\cdot} \hat{=}$ „Zeitableitung“ $\frac{d}{dt}$

4.7
p. 562

5.0
p. 564

Mechanik
 Reaktionskinetik
 Populationsdynamik
 Schaltkreismodelle
 ...

⇒ Modellierung durch gewöhnliche Differentialgleichungen

Beispiel 216 (Mechanisches System Pendel).

Phasenraum Ω = Konfigurationsraum für **Minimal-koordinaten** (= Auslenkungswinkel)

➤ $d = 1, \Omega =] - \pi/2, \pi/2[$
 (→ zulässige Auslenkungswinkel α)

Newtonsche Bewegungsgleichungen:

$$ml \ddot{\alpha}(t) = -mg \sin \alpha(t).$$

ODE 2. Ordnung

Formale Umwandlung in gewöhnliche Differentialgleichungen 1. Ordnung (Hamiltonsche Form der Bewegungsgleichungen → Physik):

$$p := \dot{\alpha} \Rightarrow \begin{pmatrix} \dot{\alpha} \\ \dot{p} \end{pmatrix} = \begin{pmatrix} p \\ -\frac{g}{l} \sin \alpha \end{pmatrix}.$$

5.1 Theorie gewöhnlicher Differentialgleichungen

[File: section-theorie-gewöhnlicher-differentialgleichungen.tex, SVN: section-theorie-gewöhnlicher-differentialgleichungen.tex 1010 2006-09-11 15:44:43Z]

Gegeben: $f: \hat{\Omega} \subset \mathbb{R} \times \Omega \mapsto \mathbb{R}^d$ stetig, Anfangsdaten $(t_0, \mathbf{y}_0) \in \hat{\Omega}$

➤ AWP: $\mathbf{y}(t_0) = \mathbf{y}_0$, $\dot{\mathbf{y}}(t) = f(t, \mathbf{y}(t))$. (5.1.1) eq:awp

Terminologie: AWP **autonom** $\Leftrightarrow f(t, \mathbf{y}) = f(\mathbf{y})$

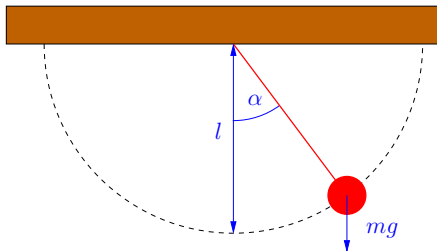
Definition 5.1.1. $J \subset \mathbb{R}$ offenes Intervall, $t_0 \in J$

$y \in C^1(J, \Omega)$ **Lösung** des AWP (5.1.1) eq:awp $\Leftrightarrow \mathbf{y}(t_0) = \mathbf{y}_0$, $\dot{\mathbf{y}}(t) = f(t, \mathbf{y}(t)) \quad \forall t \in J$.

Annahme: (**Lokale Lipschitz-Stetigkeit**)

$$\forall (t, \mathbf{y}) \in \hat{\Omega}: \exists L > 0, \delta > 0: |t - \tilde{t}| + \|\mathbf{y} - \tilde{\mathbf{y}}\| < \delta \Rightarrow \|f(t, \mathbf{y}) - f(\tilde{t}, \tilde{\mathbf{y}})\| \leq L \|\mathbf{y} - \tilde{\mathbf{y}}\|$$

Kriterium: $\frac{d}{d\mathbf{y}} f$ stetig auf $\hat{\Omega}$ (→ Kompaktheitsargument)



„Die Lösung eines Anfangswertproblems sucht sich Ihren Definitionsbereich selbst“

Lösung $y \in C^1([t_0, t_1], \Omega)$ von AWP (5.1.1) eq:awp **maximal (in die Zukunft) fortsetzbar** \Leftrightarrow

Es gibt Lösung $\tilde{y} \in C^1([t_0, t_+], \Omega)$ des AWP (5.1.1) eq:awp mit $t_+ \geq t_1$, $\mathbf{y}(t) = \tilde{\mathbf{y}}(t) \quad \forall t_0 \leq t < t_1$, wobei

$$t_+ = \infty \quad \wedge \quad \lim_{t \rightarrow t_+} \|\tilde{\mathbf{y}}(t)\| = \infty \quad \wedge \quad \lim_{t \rightarrow t_+} \text{dist}(\tilde{\mathbf{y}}(t), \partial\Omega) = 0.$$

(Analog: Maximale Fortsetzbarkeit in die Vergangenheit auf $]t_-, t_0]$) BLOWUP Blow-up

Notation: $J(t_0, \mathbf{y}_0) =]t_-, t_+[$ = maximales Existenzintervall für Lösung von AWP (5.1.1) eq:awp.

Theorem 5.1.2 (Satz von Peano & Picard-Lindelöf). Falls $f: \hat{\Omega} \mapsto \mathbb{R}^d$ lokale Lipschitz-stetig, so hat das AWP (5.1.1) eq:awp $\forall (t_0, \mathbf{y}_0) \in \hat{\Omega}$ eine eindeutige maximal fortsetzbare Lösung $y: J(t_0, \mathbf{y}_0) \mapsto \Omega$.

5.1
p. 565

Evolutionoperator $\Phi^{s,t}: \Omega \mapsto \Omega$, $\Phi^{s,t} \mathbf{y}_0 := \mathbf{y}(s)$, \mathbf{y} Lsg. von $\dot{\mathbf{y}} = f(t, \mathbf{y})$, $\mathbf{y}(t) = \mathbf{y}_0$, $s \in J(t, \mathbf{y}_0)$

$$\Phi^{t,t} = \text{Id}, \quad \Phi^{s,t} = \Phi^{s,r} \circ \Phi^{r,t}. \quad (5.1.2) \quad \text{eq:evrule}$$

Bemerkung 217. Falls AWP (5.1.1) eq:awp autonom:

$$\Phi^{s,t} = \Phi^{s-t,0} \quad \text{falls } s - t \in I.$$

\mathbf{y} Lösung mit $\mathbf{y}(t_0) = \mathbf{y}_0$ ➤ $\tilde{\mathbf{y}}(t) := \mathbf{y}(t + \tau)$ Lsg. mit $\tilde{\mathbf{y}}(t_0 - \tau) = \mathbf{y}_0$

➤ Für autonome AWP: „Kanonischer“ Anfangszeitpunkt $t_0 = 0$

Beispiel 218 (Skalare Differentialgleichungen). ➤ $d = 1$

• $f(t, y) = -\lambda y$, $\lambda \in \mathbb{R}$ ➤ Lösung des AWP $y(t) = y_0 e^{-\lambda t}$, $t \in \mathbb{R}$
 (existiert für alle Zeiten, d.h. $]t_-, t_+[= \mathbb{R}$ für jedes y_0)

• $f(t, y) = \lambda y^2$, $\lambda \in \mathbb{R}$: $\dot{y} = \lambda y^2$, $y(0) = y_0 \in \mathbb{R}$

$$\text{Lösung } y(t) = \begin{cases} -\frac{1}{y_0^{-1} - \lambda t} & , \text{ falls } y_0 \neq 0, \quad (\text{Blow-up}) \\ 0 & , \text{ falls } y_0 = 0. \end{cases}$$

$$\lambda, y_0 > 0 \Rightarrow J(0, y_0) =] - \infty, 1/\lambda y_0[.$$

5.1
p. 566

5.1
p. 567

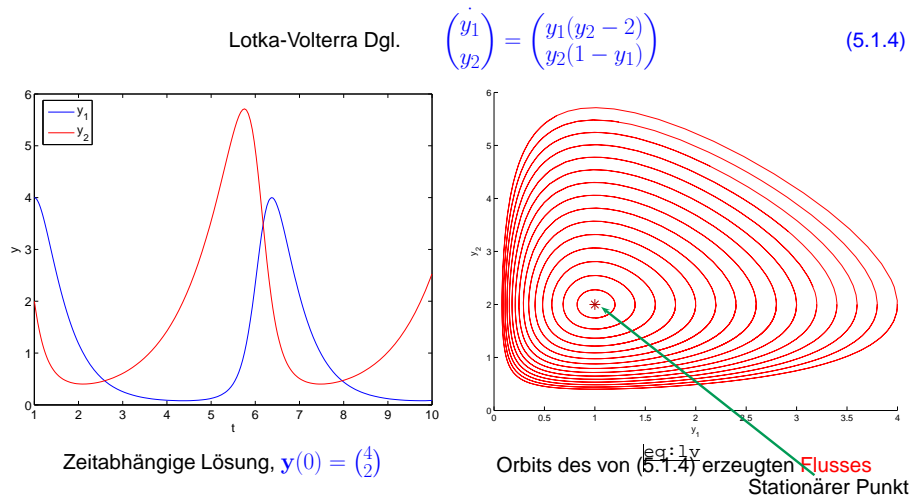
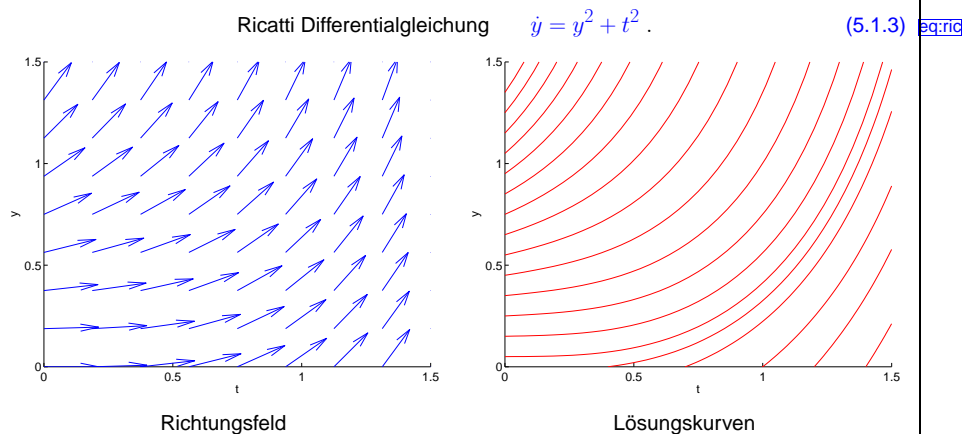
5.1
p. 568

- $f(t, y) = -\frac{1}{\sqrt{y}}$, $\Omega = \mathbb{R}^+$, Anfangswert $y(0) = 1$

$$\Rightarrow y(t) = (1 - 3t/2)^{2/3}, t_- = -\infty, t_+ = 2/3$$

(Lösung läuft zum Rand $y = 0$ des Phasenraums)

Beispiel 219 (Richtungsfeld und Orbits).



Basiswechsel im Phasenraum (kovariante Transformation): $\hat{y} = S^{-1}y$, $S \in \mathbb{R}^{d,d}$ regulär

$$y \text{ löst } \begin{cases} \dot{y} = f(t, y) \\ y(t_0) = y_0 \end{cases} \Leftrightarrow \hat{y} := S^{-1}y \text{ löst } \begin{cases} \dot{\hat{y}} = \hat{f}(t, \hat{y}) \\ \hat{y}(t_0) = S^{-1}y_0 \end{cases} \text{ mit } \hat{f}(t, y) = S^{-1}f(t, Sy).$$

(5.1.5) eq:kov

Beispiel 220 (Lineare Differentialgleichungen mit konstanten Koeffizienten).

$$\Omega = \mathbb{R}^d, \hat{\Omega} = \mathbb{R}^{d+1}, f(t, y) = Ay(t) + g(t) \text{ mit } A \in \mathbb{R}^{d,d}, g: \mathbb{R} \mapsto \mathbb{R}^d$$

Annahme: A diagonalisierbar, $\exists S \in \mathbb{R}^{d,d}$ regulär: $S^{-1}AS = \text{diag}(\lambda_1, \dots, \lambda_d)$, $\lambda_i \in \mathbb{C}$

- $g \equiv 0$: $\dot{y} = Ay$ (autonome homogene lineare Dgl.)

$$\hat{y} := S^{-1}y \text{ löst } \begin{aligned} \dot{\hat{y}}_1 &= \lambda_1 \hat{y}_1 \\ &\vdots \\ \dot{\hat{y}}_d &= \lambda_d \hat{y}_d \end{aligned} \Rightarrow \hat{y}_i(t) = (S^{-1}y_0)_i e^{\lambda_i t}, t \in \mathbb{R}.$$

$$\Rightarrow y(t) = S \underbrace{\begin{pmatrix} e^{\lambda_1 t} & & \\ & \ddots & \\ & & e^{\lambda_d t} \end{pmatrix}}_{\text{Matrixexponentialfunktion } \exp(At)} S^{-1}y_0.$$

- Inhomogener Fall $\dot{y}(t) = Ay(t) + g(t)$ partikuläre Lösung durch „Variation der Konstanten“:

Ansatz: $y(t) = \exp(At)G(t)y_0$ mit $G \in C^1(\mathbb{R}, \mathbb{R}^{d,d})$, $G(t_0) = I$

$$\dot{y}(t) = (A \exp(At)G(t) + \exp(At)\dot{G}(t))y_0 = Ay(t) = A \exp(At)G(t)y_0 + g(t)$$

$$\Rightarrow (G(t)y_0)' = g(t) \Rightarrow G(t)y_0 = y_0 + \int_{t_0}^t \exp(-A\tau)g(\tau) d\tau$$

$$\Rightarrow y(t) = \exp(At)y_0 + \int_{t_0}^t \exp(A(t-\tau))g(\tau) d\tau$$

Lsg. des homogenen Problems

Faltung mit Inhomogenität

5.2 Kondition von Anfangswertproblemen

[File: section-kondition-von-anfangswertproblemen.tex, SVN: section-kondition-von-anfangswertproblemen.tex 1010 2006-09-11 15:44:43Z hiptmail]

Wie wirken sich Störungen im Anfangswert y_0 in (5.1.1) auf die Lösung $y(t)$ aus?

Absolute Konditionszahl (\rightarrow Def. 11.6.1) durch Def:wellconditioned **differentielle Konditionsanalyse** (\rightarrow Abschnitt 11.6) sec:kondition

$$\frac{d}{ds} \Phi^{s,t} \mathbf{y} = f(s, \Phi^{s,t} \mathbf{y}) \quad \Rightarrow \quad \frac{d}{ds} \left(\frac{d}{dy} \Phi^{s,t} \mathbf{y} \right) = \frac{d}{dy} f(s, \Phi^{s,t} \mathbf{y}) = \frac{\partial f}{\partial \mathbf{y}}(s, \Phi^{s,t} \mathbf{y}) \frac{d}{dy} \Phi^{s,t} \mathbf{y}$$

Propagationsmatrix (Wronski-Matrix) $\mathbf{W}(s, t_0) := \frac{d}{dy} \Phi^{s,t_0} \mathbf{y} \in \mathbb{R}^{d,d}$ zum AWP (5.1.1) erfüllt

VARG
Variationsgleichung $\frac{d}{ds} \mathbf{W}(s, t_0) = \frac{\partial f}{\partial \mathbf{y}}(s, \Phi^{s,t_0} \mathbf{y}_0) \mathbf{W}(s, t_0) \quad , \quad \mathbf{W}(t_0, t_0) = \mathbf{I} . \quad (5.2.1)$ eq:awp eq:VARG

Beachte: Variationsgleichung = lineare Differentialgleichung mit $\Omega = \mathbb{R}^{d,d}$

$$\mathbf{y}_0 \leftarrow \mathbf{y}_0 + \delta \mathbf{y} \quad \Rightarrow \quad \delta \mathbf{y}(t) \approx \mathbf{W}(t, t_0) \delta \mathbf{y} \quad \text{für „kleine } \delta \mathbf{y} \text{“}$$

Definition 5.2.1 (Kondition des AWP). **Intervallweise Kondition** des AWP (5.1.1) bzgl. Norm $\|\cdot\|$ auf \mathbb{R}^d :

$$\kappa(t_0, t) := \max\{\|\mathbf{W}(s, t_0)\| : t_0 \leq s \leq t\} .$$

Lemma 5.2.2 (Konditionsabschätzung für AWP).

$$\left\| \frac{\partial f}{\partial \mathbf{y}}(t, \Phi^{t,t_0} \mathbf{y}_0) \right\| \leq \chi(t) \quad \Rightarrow \quad \kappa(t_0, t) \leq \exp \left(\int_{t_0}^t \chi(\tau) d\tau \right) .$$

Beweis mit *Lemma von Gronwall* (\rightarrow Analysis)

Beispiel 221 (Kondition linearer AWP mit konstanten Koeffizienten).

$$\dot{\mathbf{y}} = \lambda \mathbf{y} \quad \Rightarrow \quad \mathbf{y}(t) = e^{\lambda t} \mathbf{y}_0 \quad \Rightarrow \quad \begin{cases} \kappa(0, \infty) = \infty, \quad \kappa(-\infty, 0) = 1 & , \text{ falls } \operatorname{Re} \lambda > 0, \\ \kappa(0, \infty) = \kappa(-\infty, 0) = 1 & , \text{ falls } \operatorname{Re} \lambda = 0, \\ \kappa(0, \infty) = 1, \quad \kappa(-\infty, 0) = \infty & , \text{ falls } \operatorname{Re} \lambda < 0. \end{cases}$$

Beachte: Hier Variationsgleichung: $\dot{\mathbf{y}} = \lambda \mathbf{y}$ \triangleright Propagation von Störungen durch Dgl. selbst

lem:awpcond
 \blacktriangleright Lemma 5.2.2 liefert zu pessimistische Aussagen.

Für $\dot{\mathbf{y}} = \mathbf{A} \mathbf{y}$, $\mathbf{A} \in \mathbb{R}^{d,d}$ diagonalisierbar \rightarrow Bsp. 220: ex:linawp (Spektrum $\sigma(\mathbf{A}) \rightarrow$ Def. 3.3.1) def:ew

$$\begin{aligned} \operatorname{Re}\{\sigma(\mathbf{A})\} \subset]-\infty, 0[& \Rightarrow \kappa(0, \infty) \sim 1, \quad \kappa(-\infty, 0) = \infty, \\ \operatorname{Re}\{\sigma(\mathbf{A})\} = \{0\} & \Rightarrow \kappa(0, \infty) \sim 1, \quad \kappa(-\infty, 0) \sim 1, \\ \operatorname{Re}\{\sigma(\mathbf{A})\} \subset]0, \infty[& \Rightarrow \kappa(0, \infty) = \infty, \quad \kappa(-\infty, 0) \sim 1. \end{aligned}$$

Allgemeiner: DEB02 Kondition bzgl. Störungen in f \triangleright [11, Kap. 3]

5.3 Einschrittverfahren

[File: section-einschrittverfahren.tex, SVN: section-einschrittverfahren.tex 1010 2006-09-11 15:44:43Z hiptmail]

Gegeben: $f : \hat{\Omega} \mapsto \mathbb{R}^d$ lokale Lipschitz-stetig auf erweitertem Phasenraum $\hat{\Omega} \subset \mathbb{R} \times \Omega$
 \triangleright Definiert Familie von AWP (5.1.1) eq:awp eq:VARG

\blacktriangleright Zugehörige Evolution: $\Phi^{s,t} : \Omega \mapsto \Omega$

Gegeben: Anfangsdaten $(t_0, \mathbf{y}_0) \in \hat{\Omega}$ \triangleright Konkretes AWP

Ziel: \hookrightarrow Approximation von $\mathbf{y}(T)$ für ein $T \in J(t_0, y_0)$.

\hookrightarrow Approximation der Funktion $t \mapsto \mathbf{y}(t)$, $t \in [t_0, T]$, $T \in J(t_0, y_0)$ \triangleright $\mathbf{y}_h(t)$.

5.2
p. 573

5.3.1 Kollokation

Idee: ❶ Approximiere $\mathbf{y}(t)$, $t \in [t_0, T]$, in $s+1$ -dimensionalem **Ansatzraum** V von Funktionen $[t_0, T] \mapsto \mathbb{R}^d \triangleright \mathbf{y}_h$.

❷ Festlegung von $\mathbf{y}_h \in V$ durch **Kollokationsbedingungen**

$$\mathbf{y}_h(t_0) = \mathbf{y}_0 \quad , \quad \dot{\mathbf{y}}_h(\tau_j) = f(\tau_j, \mathbf{y}_h(\tau_j)) \quad , \quad j = 1, \dots, s, \quad (5.3.1)$$

für **Kollokationspunkte** $t_0 \leq \tau_1 < \dots < \tau_s \leq T$. eq:kollcon



„Standardoption“: **Polynomialer Ansatzraum** $V = \mathcal{P}_s$

Herleitung: Formel für $\mathbf{y}_h(T)$ ($h := T - t_0$, $\tau_j := t_0 + c_j h$, $0 \leq c_1 < c_2 < \dots < c_s \leq 1$)

$\{L_j\}_{j=1}^s \subset \mathcal{P}_{s-1} \hat{=}$ **Lagrange-Polynome** (4.1.4): eq:lagrangepol $L_j(c_i) = \delta_{ij}$, $i, j = 1, \dots, s$.

eq:kollcond (5.3.1) \blacktriangleright $\dot{\mathbf{y}}_h(t_0 + \tau h) = \sum_{j=1}^s \mathbf{k}_j L_j(\tau) \quad , \quad \mathbf{k}_j := f(t_0 + c_j h, \mathbf{y}_h(t_0 + c_j h)) .$

\blacktriangleright $\mathbf{y}_h(t_0 + \tau h) = \mathbf{y}_0 + h \sum_{j=1}^s \mathbf{k}_j \int_0^\tau L_j(\zeta) d\zeta .$

5.3
p. 574

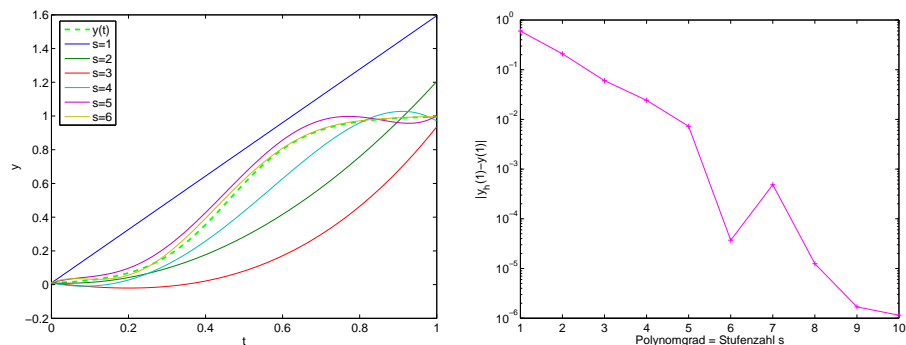
5.3
p. 575

5.3
p. 576

Logistische Differentialgleichung (→ durch Ressourcenknappheit beschränktes Wachstum)

$$\dot{y} = \lambda y(1 - y), \quad y_0 \in]0, 1[\Rightarrow y(t) = \frac{1}{1 + (y_0^{-1} - 1)e^{-\lambda t}}, \quad t \in \mathbb{R}. \quad (5.3.4)$$

Numerische Experimente mit Gauss-Kollokationsverfahren auf $[0, 1]$, $y_0 = 0.01$, $\lambda = 10$:
(Lösung der Gleichungen für Inkremente \mathbf{k}_i : MATLAB `fsolve`, Toleranz 10^{-9})



Näherungslösungen $y_h(t)$

▷ Exponentielle Konvergenz in s

! Lösbarkeit von (5.3.2) nur gesichert für „kleines h “

p. 577

Abschnitt 4.3 → Problematik globaler Polynominterpolation

▷ Stückweise Polynominterpolation

Idee: ① Wähle Gitter $\mathcal{M} := \{t_0 < t_1 < t_2 < \dots < t_n = T\}$

$$\Psi_{\mathcal{M}}^{T, t_0} := \Psi_h^{t_n, t_{n-1}} \circ \dots \circ \Psi_h^{t_2, t_1} \circ \Psi_h^{t_1, t_0}$$

② $\Psi_h^{t_k, t_{k-1}} :=$ diskrete Evolution durch Kollokation des AWP auf $[t_{k-1}, t_k]$, → Formel (5.3.2).

▷ Stückweise polynomiale Lösung $y_h \in C^0([t_0, T])$

Definition 5.3.1 (Einschrittverfahren). Ein **Einschrittverfahren** (ESV, engl. single step method) für ein AWP (5.1.1) auf einem Gitter $\mathcal{M} := \{t_0 < t_1 < t_2 < \dots < t_n = T\}$ und auf der Grundlage einer diskreten Evolution $\Psi_h^{s, t}$ erzeugt Folge von Näherungswerten

$$y_k \approx y(t_k) \quad \text{gemäss} \quad y_k := \Psi_{h_k}^{t_k, t_{k-1}} y_{k-1}, \quad k = 1, \dots, n.$$

Definition 5.3.2 (Explizite und implizite Verfahren). Ein **Einschrittverfahren** zur approximativen Lösung eines AWP heisst **explizit**, falls die zugrundeliegende diskrete Evolution durch endlich viele f -Auswertungen zu realisieren ist.

Die diskrete Evolution eines **impliziten** Einschrittverfahrens erfordert die Lösung eines Gleichungssystems.

p. 578

5.3

p. 580

$$\mathbf{k}_i = f(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^s a_{ij} \mathbf{k}_j), \quad a_{ij} = \int_0^{c_i} L_j(\tau) d\tau, \quad (5.3.2)$$

$$\mathbf{y}_h(T) = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i, \quad b_i = \int_0^1 L_i(\tau) d\tau.$$

DEVOLUTION

Diskrete Evolution $\Psi_h^{T, t_0} : \Omega \mapsto \Omega$, $\Psi_h^{T, t_0} \mathbf{y}_0 := \mathbf{y}_h(T)$

leg:rkimpl
(5.3.2)

(Nichtlineares) Gleichungssystem für Inkremente \mathbf{k}_i

Bemerkung 222. $f(t, \mathbf{y}) = f(t)$ & $\mathbf{y}_0 = 0$ ▷ Numerische Quadratur

$$\mathbf{y}(T) = \int_0^T f(t) dt \approx h \sum_{i=1}^s b_j f(t_0 + c_j h) = \text{Quadraturformel, vgl. (4.4.2)}$$

leg:gauss-quad

→ $c_1, \dots, c_s \leftrightarrow$ Knoten einer Quadraturformel (z.B. Gauss-Punkte auf $[0, 1]$ → Abschnitt 4.4.2)

• Fall $s = 1$ ▷ $c_1 = 1/2$ (↔ einfachste Gauss-Legendre-Quadraturformel)

$$L_1 \equiv 1 \Rightarrow a_{11} = 1/2, \quad b_1 = 1.$$

$$\mathbf{k}_1 = f(t_0 + 1/2h, \mathbf{y}_0 + 1/2h \mathbf{k}_1), \quad \mathbf{y}_h(T) = \mathbf{y}_0 + h \mathbf{k}_1.$$

(5.3.3) leg:rkimpl

leg:implmidrule
(5.3.3)

= implizite Mittelpunktsregel

• Fall $s = 1$ & $c_1 = 0$ (↔ linksseitige Ein-Punkt-Quadraturformel)

$$L_1 \equiv 1 \Rightarrow a_{11} = 0, \quad b_1 = 1.$$

$$\mathbf{k}_1 = f(t_0, \mathbf{y}_0), \quad \mathbf{y}_h(T) = \mathbf{y}_0 + h \mathbf{k}_1 = \mathbf{y}_0 + h f(t_0, \mathbf{y}_0).$$

leg:exEul

(5.3.1) = Explizites Eulerverfahren (kein Lösen einer Gleichung erforderlich!)

• Fall $s = 1$ & $c_1 = 1$ (↔ rechtsseitige Ein-Punkt-Quadraturformel)

$$L_1 \equiv 1 \Rightarrow a_{11} = 1, \quad b_1 = 1.$$

$$\mathbf{k}_1 = f(t_0, \mathbf{y}_0 + h \mathbf{k}_1), \quad \mathbf{y}_h(T) = \mathbf{y}_0 + h \mathbf{k}_1 = \mathbf{y}_0 + h f(t_0, \mathbf{y}_h(T)).$$

leg:implEul

(5.3.1) = implizites Eulerverfahren

Beispiel 223 (Konvergenz von Kollokationsverfahren).

Alle Kollokationsverfahren \{ Explizites Eulerverfahren \} sind implizit

5.3.2 Runge-Kutta-Verfahren

$$\text{AWP: } \dot{\mathbf{y}}(t) = f(t, \mathbf{y}(t)), \Rightarrow \mathbf{y}(t) = \mathbf{y}_0 + \int_{t_0}^t f(\tau, \mathbf{y}(\tau)) d\tau$$

$$\mathbf{y}(t_0) = \mathbf{y}_0$$

Approximation durch Quadraturformel (auf $[0, 1]$) mit s Knoten c_1, \dots, c_s :

$$\mathbf{y}(T) \approx \mathbf{y}_h(T) = \mathbf{y}_0 + h \sum_{i=1}^s b_i f(t_0 + c_i h, \mathbf{y}(t_0 + c_i h))$$

Wie bekommt man diese Werte ? \triangleright Bootstrapping

Beispiel 224 (Konstruktion einfacher Runge-Kutta-Verfahren).

Quadraturformel \rightarrow Trapezregel $\xrightarrow{\text{eq:trapezoidal}}$ & $\mathbf{y}_h(T)$ aus explizitem Eulerschritt $\xrightarrow{\text{eq:exEul}}$ (5.3.1)

$$\mathbf{k}_1 = f(t_0, \mathbf{y}_0), \quad \mathbf{k}_2 = f(t_0 + h, \mathbf{y}_0 + h\mathbf{k}_1), \quad \mathbf{y}_h(T) = \mathbf{y}_0 + \frac{h}{2}(\mathbf{k}_1 + \mathbf{k}_2). \quad (5.3.5) \quad \text{eq:exTra}$$

p. 581

(5.3.5) = **explizite Trapezregel** Quadraturformel \rightarrow Einfachste Gauss-Quadraturformel (Mittelpunktsregel) & $\mathbf{y}_h(\frac{1}{2}(T + t_0))$ aus explizitem Eulerschritt $\xrightarrow{\text{eq:exEul}}$ (5.3.1)

$$\mathbf{k}_1 = f(t_0, \mathbf{y}_0), \quad \mathbf{k}_2 = f(t_0 + \frac{h}{2}, \mathbf{y}_0 + \frac{h}{2}\mathbf{k}_1), \quad \mathbf{y}_h(T) = \mathbf{y}_0 + h\mathbf{k}_2. \quad (5.3.6) \quad \text{eq:exMP}$$

(5.3.6) = **explizite Mittelpunktsregel**

Diskrete Evolutionen der Form $\xrightarrow{\text{eq:rkiimpl}}$ (5.3.2)

Definition 5.3.3 (Runge-Kutta-Verfahren). Für $b_i, a_{ij} \in \mathbb{R}$, $i, j = 1, \dots, s$, $s \in \mathbb{N}$, $c_i = \sum_{j=1}^s a_{ij}$ \Rightarrow **s-stufiges Runge-Kutta-Einschrittverfahren (RK-ESV)** für AWP $\xrightarrow{\text{eq:awp}}$ (5.1.1)

$$\mathbf{k}_i = f(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^s a_{ij} \mathbf{k}_j), \quad i = 1, \dots, s, \quad \mathbf{y}_h(t_0 + h) = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{k}_i.$$

\mathbf{k}_i = Inkremente/Stufen

Falls $a_{ij} = 0$ für $i \leq j$ \Rightarrow **Explizites Runge-Kutta-Verfahren** \rightarrow Def. 5.3.2 $\xrightarrow{\text{def:expl}}$

Kurznotation für Runge-Kutta-Verfahren:

Butcher-Schema

$$\begin{array}{c|c} \mathbf{c} & \mathbf{A} \\ \hline \mathbf{b}^T & \end{array} := \begin{array}{c|ccc} c_1 & a_{11} & \cdots & a_{1s} \\ \vdots & \vdots & & \vdots \\ c_s & a_{s1} & \cdots & a_{ss} \\ \hline b_1 & \cdots & & b_s \end{array}. \quad (5.3.7) \quad \text{eq:BS}$$

Interpretation: Runge-Kutta-Verfahren \leftrightarrow **Polygonzugapproximation** der Lösungskurve

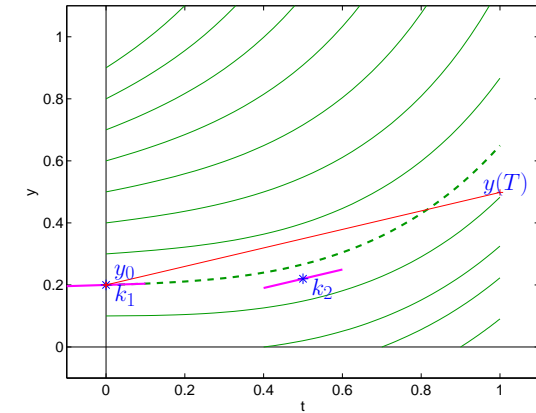
Beispiel 225 (Explizite Runge-Kutta-Schritte für Riccati-Differentialgleichung). \rightarrow Bsp 219 $\xrightarrow{\text{ex:orbits}}$

Anfangswertproblem: $\dot{y} = t^2 + y^2$, $y(0) = 0.2$.

Explizite Mittelpunktsregel:

$$\begin{array}{c|ccc} 0 & 0 & 0 & \\ \hline \frac{1}{2} & \frac{1}{2} & 0 & \\ \hline 0 & 1 & & \end{array}$$

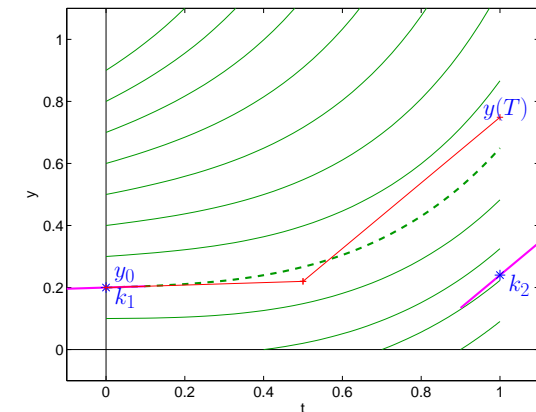
grün: Lösungskurven
magenta: Abschnittsteigungen \mathbf{k}_i
*: Punkte f -Auswertung
rot: Polygonzug



Explizite Trapezregel

$$\begin{array}{c|ccc} 0 & 0 & 0 & \\ \hline 1 & 1 & 0 & \\ \hline \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \end{array}$$

grün: Lösungskurven
magenta: Abschnittsteigungen \mathbf{k}_i
*: Punkte f -Auswertung
rot: Polygonzug



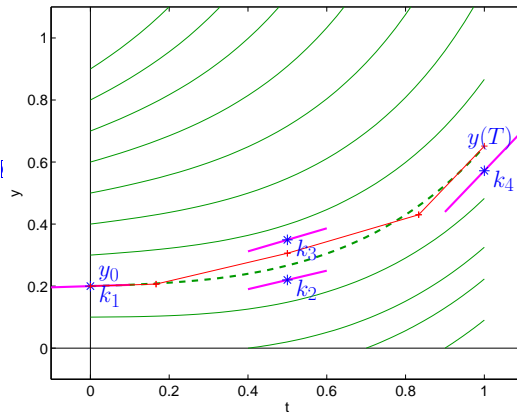
5.3
p. 583

5.3
p. 584

Klassisches Runge-Kutta-Verfahren (RK4)

$$\begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ \hline \frac{1}{4} & \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{4} & 0 & \frac{1}{2} & 0 & 0 \\ \hline \frac{1}{4} & 0 & 0 & 1 & 0 \\ \hline \frac{1}{6} & \frac{2}{6} & \frac{2}{6} & \frac{1}{6} & \frac{1}{6} \end{array} \quad (5.3.8)$$

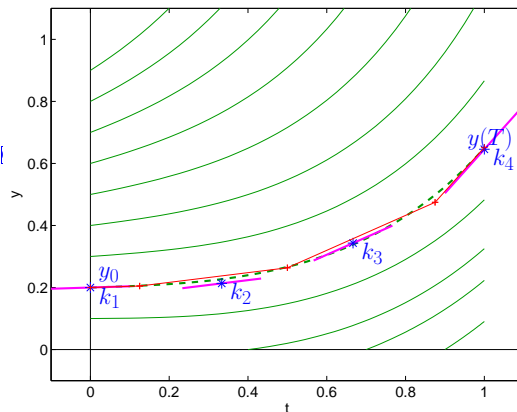
- grün Lösungskurven
- magenta Abschnittsteigungen k_i
- * Punkte f -Auswertung
- rot: Polygonzug



Kuttas 3/8-Regel

$$\begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ \hline \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & -\frac{1}{3} & 1 & 0 & 0 \\ \hline \frac{1}{8} & \frac{1}{8} & -\frac{1}{8} & 1 & 0 \\ \hline \frac{1}{8} & \frac{3}{8} & \frac{3}{8} & \frac{1}{8} & \frac{1}{8} \end{array} \quad (5.3.9)$$

- grün Lösungskurven
- magenta Abschnittsteigungen k_i
- * Punkte f -Auswertung
- rot: Polygonzug



RKKOV

Affin-Kovarianz der Runge-Kutta-Verfahren:

Für $S \in \mathbb{R}^{d,d}$ regulär, $\hat{y} := S^{-1}y$:

$\Psi_h^{s,t}$ Diskrete Evolution zu $\dot{y} = f(t, y)$,

$\hat{\Psi}_h^{s,t}$ Diskrete Evolution zu $\dot{\hat{y}} = \hat{f}(t, \hat{y}) \rightarrow$

$$S \hat{\Psi}_h^{s,t} S^{-1} y = \Psi_h^{s,t} y$$

(5.3.10)

Warum $c_i = \sum_{j=1}^s a_{ij}$ in Def. 5.3.3?

$$\text{Autonomisierung: } \begin{array}{l} \dot{y} = f(t, y), \\ y(t_0) = y_0 \end{array} \Rightarrow \begin{pmatrix} \dot{y} \\ s \end{pmatrix} = \begin{pmatrix} f(s, y) \\ 1 \end{pmatrix} =: \hat{f} \left(\begin{pmatrix} y \\ s \end{pmatrix} \right), \quad \begin{pmatrix} y(0) \\ s(0) \end{pmatrix} = \begin{pmatrix} y_0 \\ t_0 \end{pmatrix}$$

$$\begin{array}{l} \text{Evolutioner:} \\ \text{Diskrete Evt:} \end{array} \quad \begin{array}{l} \Phi^{t+h,t} \\ \Psi_h^{t+h,t} \end{array} \leftrightarrow \begin{array}{l} \hat{\Phi}^{t+h,t} \\ \hat{\Psi}_h^{t+h,t} \end{array}$$

$$\text{Wunsch: } \begin{pmatrix} \Phi^{t+h,t} y \end{pmatrix} = \hat{\Phi}^{t+h,t} \begin{pmatrix} y \\ t \end{pmatrix} \Rightarrow \begin{pmatrix} \Psi_h^{t+h,t} y \end{pmatrix} = \hat{\Psi}_h^{t+h,t} \begin{pmatrix} y \\ t \end{pmatrix} \quad (5.3.11)$$

$$\hat{\Psi}_h^{t+h,t} \begin{pmatrix} y \\ t \end{pmatrix} = \begin{pmatrix} y + h \sum_{i=1}^s b_i \hat{k}_i \\ t + h \sum_{i=1}^s b_i \hat{k}_i \end{pmatrix}, \quad \begin{pmatrix} \hat{k}_i \\ \hat{k}_i \end{pmatrix} = \begin{pmatrix} f(t + h \sum_{j=1}^s a_{ij} \hat{k}_j, y + h \sum_{j=1}^s a_{ij} \hat{k}_j) \\ 1 \end{pmatrix}$$

$$\Rightarrow \boxed{c_i = \sum_{j=1}^s a_{ij}} \geq \hat{k}_i = k_i \quad \& \quad \boxed{\sum_{i=1}^s b_i = 1} \quad (5.3.12)$$

= Hinreichende + notwendige Bedingungen für Autonomisierungsinvarianz eines RK-Verfahrens

5.3
p. 585

► Analyse von autonomisierungsinvarianten RK-Verfahren kann sich auf autonome Probleme beschränken.

5.4
p. 587

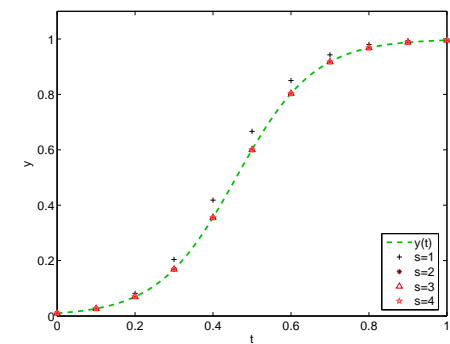
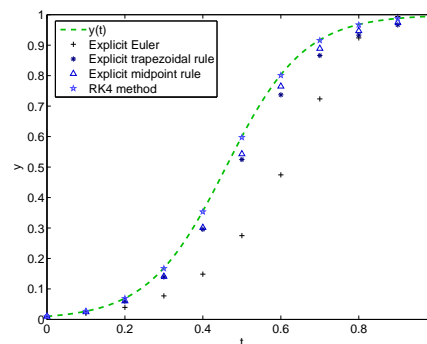
5.4 Konvergenz

[File: section-konvergenz.tex, SVN: section-konvergenz.tex 1010 2006-09-11 15:44:43Z hiptmair]

Beispiel 226 (Konvergenz von Einschrittverfahren).

Skalare logistische Differentialgleichung (5.3.4), $\lambda = 10$, $y(0) = 0.01$, $T = 1$:

Test: Explizite/implizite Runge-Kutta-Einschrittverfahren, verschiedene Schrittweiten h



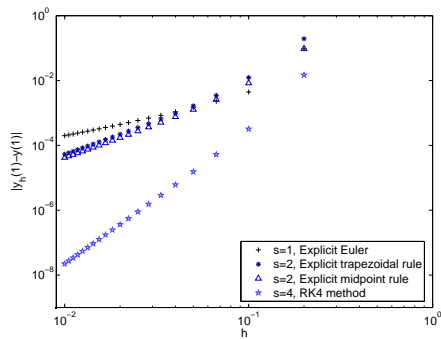
5.3

Explizite RK-Verfahren

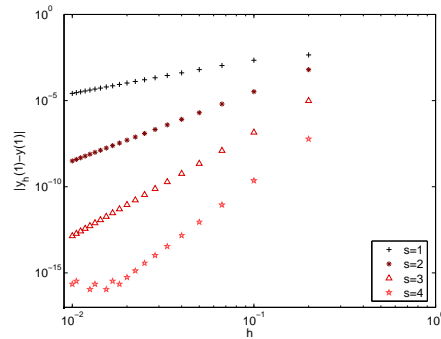
$y_h(j/10)$: Gauss-Kollokationsverfahren

5.4
p. 588

Fehler $y_n - y(1)$ zum Endzeitpunkt, äquidistantes Gitter, $n=5:5:100$



Explizite RK-Verfahren



Gauss-Kollokationsverfahren

➤ Beobachtung: Algebraische Konvergenz $|y(1) - y_n| = O(h^\alpha)$

Definition 5.4.1 (Konvergenzordnung von Einschrittverfahren). Ein ESV (\rightarrow Def. 5.3.1) für das AWP (5.1.1) auf $[t_0, T]$ hat **Konvergenzordnung** p (vgl. Thm. 4.4.4), wenn

$$\exists C > 0: \max_{j=0, \dots, n} |y_j - y(t_j)| \leq C h_{\mathcal{M}}^p \quad \forall \text{Gitter } \mathcal{M} := \{t_0 < t_1 < \dots < t_n \leq T\},$$

mit Schrittweite $h_{\mathcal{M}} := \max_{k=1, \dots, n} |t_k - t_{k-1}|$. Diskretisierungsfehler

Für ESV: Lokale Konsistenz(ordnung) \blacktriangleright (globale) Konvergenz(ordnung)

Definition 5.4.2 (Konsistenzordnung einer diskreten Evolution). Diskrete Evolution Ψ_h besitzt **Konsistenzordnung** p bzgl. einer Evolution Φ (zu einer Dgl.), falls

$$\forall (t, y \in \hat{\Omega}): \exists C > 0: \underbrace{|\Psi_h^{t+h,t} y - \Phi^{t+h,t} y|}_{\text{Konsistenzfehler}} \leq C h^{p+1} \quad \text{für kleine } h.$$

Theorem 5.4.3 (Konvergenz von Einschrittverfahren). Runge-Kutta ESVs (\rightarrow Def. 5.3.1) auf der Grundlage der diskreten Evolution $\Psi_h^{t+h,t}$ für AWP (5.1.1) erfüllen:

$$\Psi_h \text{ konsistent mit Ordnung } p \Leftrightarrow \text{ESV konvergent mit Ordnung } p$$

Merke: Konvergenzordnung erlaubt i.a. **keine** Aussage über Grösse des Diskretisierungsfehlers

Annahme: $f = f(y)$ (\rightarrow autonomes AWP) & f „hinreichend“ glatt

➤ Konsistenzanalyse von Runge-Kutta-Verfahren durch *Taylorentwicklung* (\rightarrow MAPLE)

Beispiel 227 (Konsistenzanalyse durch Computeralgebra).

Explizite Trapezregel (5.3.5)

5.4
p. 589

```

D(y) := x -> f(y(x)); y0 := y(0);
D(y) := x -> f(y(x)); y0 := y(0);
k1 := f(y0); k2 := f(y0+h*k1);
k1 := f(y0); k2 := f(y0+h*f(y0));
ym := y0+h/2*(k1+k2);
ym := y0+1/2*h*(f(y0)+f(y0+h*f(y0)));
taylor(ym-y(h), h=0, 4);
series((1/12*(D^(2))(f)(y(0))(f(y(0)))^2-1/6*(D(f)(y(0)))^2*f(y(0)))h^3+O(h^4), h, 4)

```

5.4
p. 591

5.4
p. 590

5.4
p. 592

Implizite Trapezregel: $y_1 = y_0 + \frac{h}{2}(f(y_0) + f(y_1))$

$D(y) := x \rightarrow f(y(x))$

$D(y) := x \mapsto f(y(x))$

$y0 := y(0);$

$y0 := y(0)$

$\text{solve}(y0+h/2*(f(y0) + f(yt))=yt,\{yt\});$

$\{yt = \text{RootOf}(-2y(0) - hf(y(0)) - hf(-Z) + 2-Z)\}$

$\text{assign}(\%);$

$\text{taylor}(yt-y(h),h=0,4);$

$$\text{series}\left(\left(-1/6(D^{(2)}(f)(y(0))(f(y(0)))^2 - 1/6(D(f)(y(0)))^2\right.\right. \\ \left.\left.f(y(0)) + 1/4 f(y(0))((D(f)(y(0)))^2 + f(y(0))(D^{(2)}(f)(y(0)))\right)h^3 + O(h^4), h, 4\right)$$

Für allgemeines Runge-Kutta-Verfahren \rightarrow Def. 5.3.3, $\xrightarrow{\text{def:rk}}$ autonome Dgl. $\dot{y} = f(y)$

$$\Psi_h^{h,0} y = y + h \sum_{i=1}^s b_i k_i \quad \Psi_h^{h,0} y - \Phi^{h,0} y|_{h=0} = 0, \\ \frac{d}{dh}(\Psi_h^{h,0} y - \Phi^{h,0} y)|_{h=0} = \sum_{i=1}^s b_i k_i|_{h=0} = f(y)$$

Da $k_i|_{h=0} = f(y) \quad \Rightarrow \quad \text{RK-ESV konsistent} \Leftrightarrow \sum_{i=1}^s b_i = 1$.

Analoge Herleitung von Bedingungsgleichungen für RK-Koeffizienten b_i, a_{ij} (\rightarrow Def. 5.3.3) für Konsistenzordnungen $p > 1$:

$$p \geq 2 \Leftrightarrow \text{zusätzlich} \quad \sum_{i=1}^s b_i c_i = \frac{1}{2}, \quad \sum_{i=1}^s b_i c_i^2 = \frac{1}{3}, \quad (5.4.1)$$

$$p \geq 3 \Leftrightarrow \text{zusätzlich} \quad \sum_{i,j=1}^s b_i a_{ij} c_j = \frac{1}{6}. \quad (5.4.2)$$

(Algorithmus zur Herleitung der Bedingungsgleichungen \triangleright [11, Abschn. 4.2.3], [23, Sect. III.1])

\rightarrow Konstruktion von RK-Verfahren vorgegebener Konvergenzordnung durch Lösen der (nichtlinearen) Bedingungsgleichungen (B.G.) (vom Typ (5.4.1), (5.4.2)).

p	1	2	3	4	5	6	7	8	9	10	20
#B.G.	1	2	4	8	17	37	85	200	486	1205	20247374

Einige Konvergenzordnungen \rightarrow Bsp. 226.

Explizite Verfahren		Implizite Verfahren
Explizites Eulerverfahren	$p = 1$	
Explizite Trapezregel	$p = 2$	
Explizite Mittelpunktsregel	$p = 2$	
Klassisches Runge-Kutta-V.	$p = 4$	Implizites Eulerverfahren $p = 1$
Kuttas 3/8-Regel	$p = 4$	Implizite Trapezregel $p = 2$

Viele weitere Verfahren \triangleright [25, 26]

Theorem 5.4.4 (Konvergenzordnung von Kollokationsverfahren).

Für ein Kollokations-Einschrittverfahren (5.3.2) gilt

$$\text{Konsistenzordnung } p \Leftrightarrow \text{Quadraturformel auf } [0, 1] \text{ mit Knoten } c_i \\ \text{und Gewichten } b_i, i = 1, \dots, s \text{ hat Ordnung } p \quad (\rightarrow \text{Abschnitt 4.4.2}).$$

Ordnungsschranken:

Für explizite Runge-Kutta-Verfahren $p \leq s$
 Für allgemeine Runge-Kutta-Verfahren $p \leq 2s$ (vgl. Abschnitt 4.4.2)
 \triangleright Gauss-Kollokationsverfahren realisieren maximale Ordnung

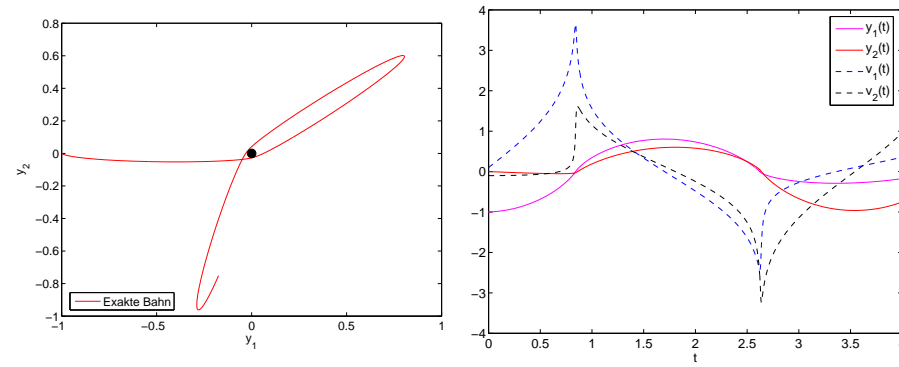
5.4.1 Schrittweitensteuerung für Einschrittverfahren

Beispiel 228 (Satellitenbahn).

5.4
p. 593

Satellitenbahn im Gravitationsfeld eines Planeten (mit Koordinaten $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$): **Keplerproblem**

$$\ddot{y} = -\frac{2y}{\|y\|^2} \Rightarrow \begin{pmatrix} \dot{y} \\ \dot{v} \end{pmatrix} = \begin{pmatrix} v \\ -\frac{2y}{\|y\|^2} \end{pmatrix}, \quad y(0) = \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \quad v(0) = \begin{pmatrix} 0.1 \\ -0.1 \end{pmatrix}.$$



Beobachtung: „Peaks“ in der Lösungskomponente $v(t)$ (= zeitlokale Charakteristika)

5.4
p. 594

5.4
p. 595

5.4
p. 596

Adaptive Wahl des Rechengitters für ESV durch **zeitlokale Fehlerschätzung**
 (→ Adaptive numerische Quadratur, Abschnitt 4.6.1)

Ziel: $\#M$ möglichst klein & $\max_{t \in [t_0, T]} \|y(t) - y_h(t)\| < \text{TOL}$, $\text{TOL} = \text{Toleranz}$
 bzw. $\|y(T) - y_h(T)\| < \text{TOL}$

Technik: Zeitlokaler Fehlerschätzer (basierend auf $y(t_k)$) → Wahl von t_{k+1}
 Heuristik: Kontrolle des (lokalen!) **Konsistenzfehlers** → Def. 5.4.2

Idee: Schätzung des Konsistenzfehlers

Vergleich zweier diskreter Evolutionen $\Psi_h^{t+h,t}, \tilde{\Psi}_h^{t+h,t}$ (für eine **aktuelle Zeitschrittweite** h) verschiedener Ordnung (→ Abschnitt 4.6.1):

Falls Ordnung($\tilde{\Psi}_h$) > Ordnung(Ψ_h)
 $\Rightarrow \underbrace{\Phi^{t+h,t} y(t_k) - \Psi_h^{t+h,t} y(t_k)}_{\text{Konsistenzfehler}} \approx \text{EST}_k := \tilde{\Psi}_h^{t+h,t} y(t_k) - \Psi_h^{t+h,t} y(t_k).$

(5.4.3)

eq:ssc1 Adaptive Integratoren für Anfangswertprobleme in MATLAB:

Vergleich $\text{EST}_k \leftrightarrow \text{TOL}$ **Absolute** Toleranz
 $\text{EST}_k \leftrightarrow \text{TOL} \|y(t_k)\|$ **Relative** Toleranz
 > Verwerfen/Akzeptieren des aktuellen Schritts

Wir wollen mehr!
 Wenn $\text{EST}_k > \text{TOL}$: **Schrittweitenkorrektur** $t_{k+1} = ?$
 Wenn $\text{EST}_k < \text{TOL}$: **Schrittweitevorschlag** $t_{k+2} = ?$

Falls Ordnung(Ψ_h) = p , Ordnung($\tilde{\Psi}_h$) > $p, p \in \mathbb{N}$,

Ziel: Effizienz

$$\begin{aligned} \Psi_h^{t+h,t} y(t_k) - \Phi^{t+h,t} y(t_k) &= ch^{p+1} + O(h^{p+2}), \quad h \ll 1 \\ \tilde{\Psi}_h^{t+h,t} y(t_k) - \Phi^{t+h,t} y(t_k) &= O(h^{p+2}) \end{aligned} \Rightarrow \text{EST}_k \approx ch^{p+1} \stackrel{!}{=} \text{TOL}.$$

„Optimale Schrittweite“:

$$h^* = h^{p+1} \sqrt{\frac{\text{TOL}}{\text{EST}_k}}$$

Korrigierte Schrittweite

Schrittweitevorschlag

Algorithmische Realisierung (ESV):

Eingebettete Runge-Kutta-Verfahren

Gleiche Inkremente k_i , verschiedene Gewichte b_i (→ Def. 5.3.3)
 realisieren RK-Evolutionen $\Phi_h, \hat{\Phi}_h$ der Ordnungen p und $p+1$.

Gebräuchlich $p=4, p=7$

$$\begin{array}{c|c} c & a \\ \hline b^T & b^T \end{array} := \begin{array}{c|ccc} c_1 & a_{11} & \cdots & a_{1s} \\ \vdots & \vdots & & \vdots \\ c_s & a_{s1} & \cdots & a_{ss} \\ \hline b_1 & \cdots & b_s \\ \hline \hat{b}_1 & \cdots & \hat{b}_s \end{array}$$

Eingebettetes RK-ESV:

Butcher-Schema

Beispiel 229 (Eingebettete Runge-Kutta-Verfahren).

0				
$\frac{1}{3}$	$\frac{1}{3}$			
$\frac{1}{3}$	$\frac{1}{6}$	$\frac{1}{6}$		
$\frac{1}{2}$	$\frac{1}{8}$	0	$\frac{3}{8}$	
1	$\frac{1}{2}$	0	$-\frac{3}{2}$	2
y_1	$\frac{1}{6}$	0	0	$\frac{2}{3}$
\hat{y}_1	$\frac{1}{10}$	0	$\frac{3}{10}$	$\frac{2}{5}$

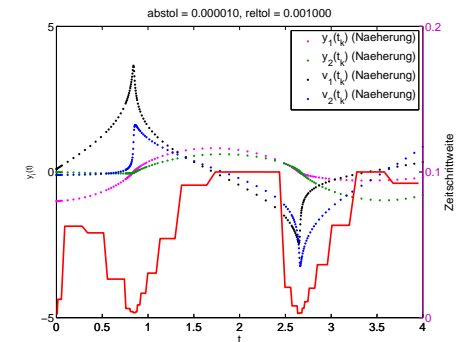
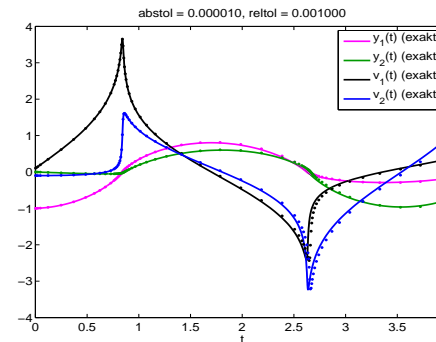
Eingebettetes RK-Verfahren der Ordnung 4(5) von Merson
 Eingebettetes RK-Verfahren der Ordnung 4(5) von Fehlberg

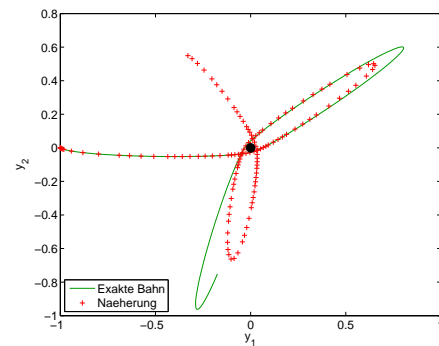
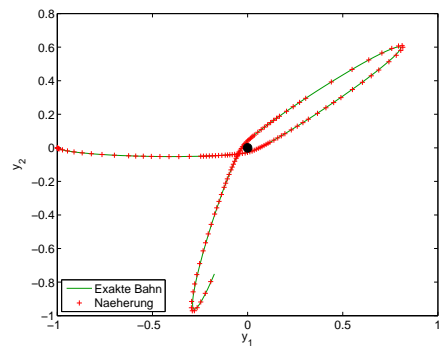
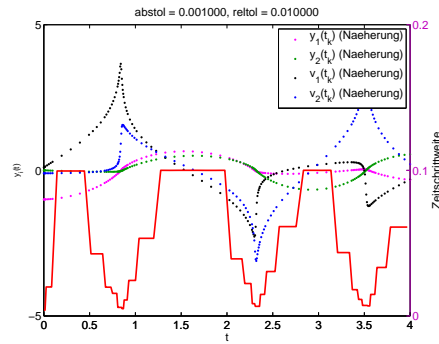
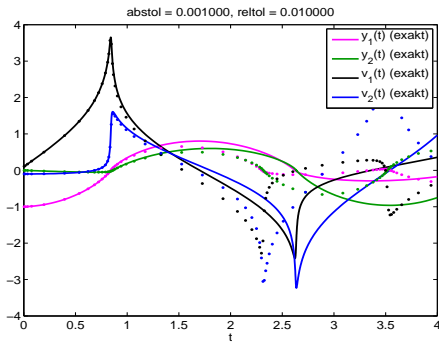
```
options = odeset('abstol',atol,'reltol',rtol,'stats','on');
[t,y] = ode45/ode23(@(t,x) f(t,x),tspan,y0,options);
(f = function handle, tspan = [t0,T], y0 = y0, t = tk, Y = yk)
```

Beispiel 230 (Adaptive RK-ESV zur Satellitenbahnberechnung). → Bsp. 228

Adaptiver Integrator: `ode45(@(t,x) satf,[0 4],[-1;0;0.1;-0.1],options))`:

- options = odeset('reltol',0.001,'abstol',1e-5);
- options = odeset('reltol',0.01,'abstol',1e-3);





reltol=0.001, abstol=1e-5

reltol=0.01, abstol=1e-3

Qualitativ falsche Lösung bei geringfügig erhöhter Toleranz !

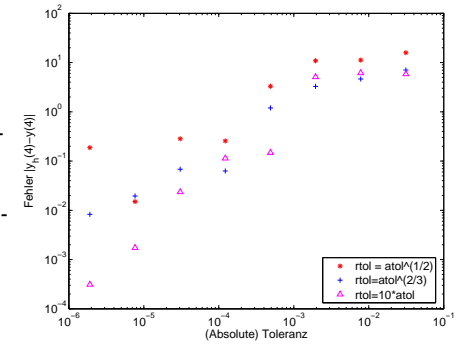
Beispiel 231 (Schrittweitensteuerung bei Satellitenbahnberechnung). → Bsp. ex:satellite

5.4
p. 601

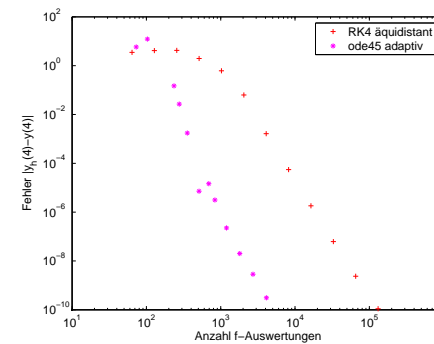
AWP aus Bsp. ex:satellite

ode45 mit verschiedenen absoluten/relativen Toleranzen

! Toleranzen sagen nichts über tatsächliche Genauigkeit



5.4
p. 603



Effizienz von Schrittweitensteuerung:

Vergleich:

- Klassisches Runge-Kutta-Verfahren ^{leg: rk4} (5.3.8)
- Eingebettetes Runge-Kutta-Verfahren mit Schrittweitensteuerung: ode45

Aufwandsmass: #f-Auswertungen

Adaptivität zahlt sich aus !

Übung 5.5. Kurve in der Ebene lässt sich als Niveaulinie einer stetigen Funktion $F : \mathbb{R}^2 \mapsto \mathbb{R}$ beschreiben:

$$\mathcal{C} := \{x \in \mathbb{R}^2 : F(x) = 0\}.$$

Falls F stetig differenzierbar ist, können zusammenhängende Teile der Kurve als Lösungskurve zu Anfangswertproblemen

$$\dot{y}(t) = \frac{\text{grad } F(y(t))^\perp}{\|\text{grad } F(y(t))\|}, \quad F(y(0)) = 0,$$

5.4
p. 602

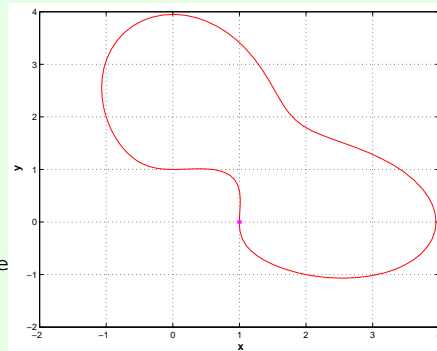
5.5
p. 604

erhalten werden. Dabei bezeichnet \perp die Rotation eines Vektor $\in \mathbb{R}^2$ um 90° gegen den Uhrzeigersinn.

Entwerfe eine MATLAB-Funktion `draweggcurve()`, die ein solches Anfangswertproblem löst, um die Kurve

$$\{x = (x, y) \in \mathbb{R}^2 : F(x) := (x^2 + y^2)^2 - x^3 - y^3 = 0\}$$

durch $(1, 0)$ zu plotten. Löse dazu das Anfangswertproblem auf $[0, 15.8]$ und erwerde dazu die MATLAB-Funktion `ode45` mit folgenden Parametern:
`options = odeset('reltol', 0.0001, 'abstol', 1e-9);`
<http://en.wikipedia.org/wiki/Crooke>

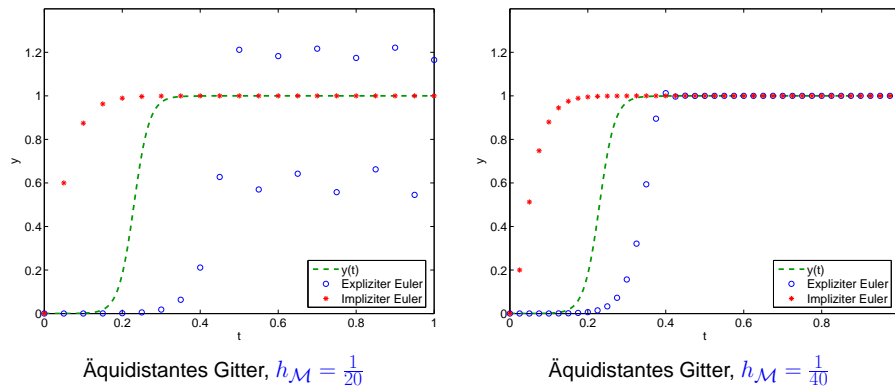


5.6 Stabilität

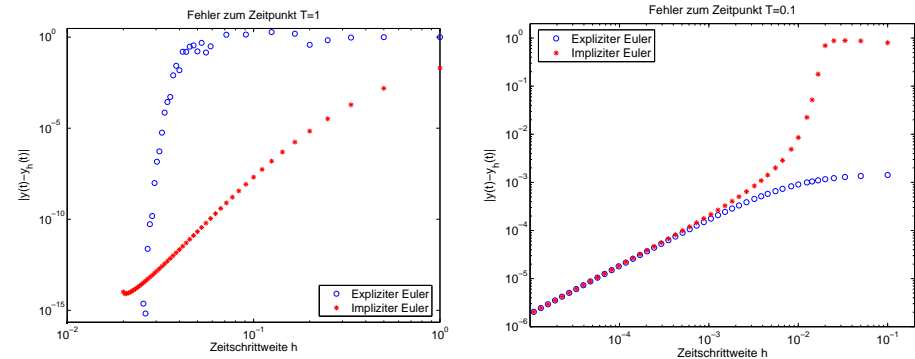
[File: section-ngdiffgleichungen-stabilitaet.tex, SVN: section-ngdiffgleichungen-stabilitaet.tex 1771 2007-08-26 02:14:34Z hiptmair]

Beispiel 232 (Unterschiedliches Verhalten expliziter und impliziter Integratoren).

Logistische Differentialgleichung $\dot{y} = f(y)$, $f(y) = \lambda y(1 - y) \rightarrow$ (5.3.4), $\lambda = 50$, Anfangswert $y_0 = 10^{-4}$, Zeitintervall $[0, 1]$:



h_M -Abhängigkeit der Fehler zum Endzeitpunkt für Integrationsintervalle $[0, T]$, $T = 0.1, 1$:



Beobachtung: $y(T) \approx 0 \Rightarrow$ Impliziter Euler *präasymptotisch* unbrauchbar
 $y(T) \approx 1 \Rightarrow$ Expliziter Euler *präasymptotisch* unbrauchbar

$y \in \{0, 1\} \Rightarrow f(y) = 0$: 0, 1 sind **stationäre Punkte** der Evolution zur Dgl.

$$\frac{df}{dy}(0) = \lambda \Rightarrow \text{für } y \approx 0 : \text{Dgl.} \approx \dot{y} = \lambda y,$$

$$\frac{df}{dy}(1) = -\lambda \Rightarrow \text{für } y \approx 1 : \text{Dgl.} \approx \dot{y} = -\lambda y.$$

5.6
p. 605

5.6.1 Modellproblemanalyse

Autonomes skalares lineares AWP: $\dot{y} = \lambda y$, $y(0) = 1$, $\text{Re } \lambda < 0$ auf $[0, \infty]$ (5.6.1) [eq:stiffm](#)

$$y(t) = e^{\lambda t} \rightarrow 0 \text{ für } t \rightarrow \infty \text{ (Asymptotische Stabilität).}$$

\leftrightarrow Diskrete Evolution von s -stufigem Runge-Kutta-Verfahren (\rightarrow Def. 5.3.3) $\models \frac{c}{b^T}$?

$$k_i = \lambda(y + h \sum_{j=1}^s a_{ij} k_j),$$

$$y_h(t+h) = y + h \sum_{i=1}^s b_i k_i$$

$$\Rightarrow \Psi_h^{h,0} y = (1 + \lambda h \mathbf{b}^T (\mathbf{I} - \lambda h \mathbf{A})^{-1} \mathbf{1}) y.$$

(Kontinuierliche) Evolution:

$$\Phi^{h,0} = e^{\lambda h}$$

Wann erbt Lösung $\{y_k\}_{k=0}^\infty$, $y_{k+1} = \Psi_h^{h,0} y_k$, eines RK-ESV auf (unendlichem) äquidistantem Gitter (Maschenweite h) asymptotische Stabilität?

$$\lim_{k \rightarrow \infty} y_k = 0 \Leftrightarrow z := \lambda h \text{ so, dass } |\Psi_h^{h,0}(z)| < 1.$$

5.6
p. 606

5.6
p. 607

5.6
p. 608

Definition 5.6.1 (Stabilitätsgebiet eines Einschrittverfahrens). Das **Stabilitätsgebiet** eines ESV für das AWP (5.6.1) auf der Grundlage der diskreten Evolution $\Psi_h^{h,0} =: S(z)$, $z := \lambda h$, $S : D_S \subset \mathbb{C} \mapsto \mathbb{C}$, ist

$$S_\Psi := \{z \in \mathbb{C} : |S(z)| < 1\} \subset \mathbb{C}.$$

Theorem 5.6.2 (Stabilitätsfunktion von Runge-Kutta-Verfahren). Ist $\Psi_h^{h,0}$ diskrete Evolution zu einem s -stufigen Runge-Kutta-Verfahren (\rightarrow Def. 5.3.3) mit Butcher-Schema $\begin{smallmatrix} \mathbf{c} & \mathbf{A} \\ \hline & \mathbf{b}^T \end{smallmatrix}$ für das AWP (5.6.1), dann

$$\Psi_h^{h,0} = \underbrace{1 + z \mathbf{b}^T (\mathbf{I} - z \mathbf{A})^{-1} \mathbf{1}}_{\text{Stabilitätsfunktion } S(z)}, \quad z := \lambda h, \quad \mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^s.$$

Übung 5.7. Das allgemeine Runge-Kutta-Verfahren aus Definition 5.3.3 werde angewandt auf das skalare Modellproblem (5.6.1), Schrittweite $h > 0$.

(i) Schreibe $\mathbf{k} \in \mathbb{R}^s$ für den Vektor $(k_1, \dots, k_s)^T$ der Inkremente und $y_1 := y_h(t_0 + h)$. Zeige, dass

$$\begin{pmatrix} \mathbf{I} - z \mathbf{A} & 0 \\ -z \mathbf{b}^T & 1 \end{pmatrix} \begin{pmatrix} \mathbf{k} \\ y_1 \end{pmatrix} = y_0 \begin{pmatrix} \mathbf{1} \\ 1 \end{pmatrix}.$$

(ii) Benutze die Cramersche Regel, um die alternative Darstellung

$$S(z) = \frac{\det(\mathbf{I} - z \mathbf{A} + z \mathbf{1} \mathbf{b}^T)}{\det(\mathbf{I} - z \mathbf{A})}$$

herzuleiten.

• Explizites Eulerverfahren: $\begin{vmatrix} 0 & 0 \\ 1 & 1 \end{vmatrix} \Rightarrow S(z) = 1 + z.$

• Implizites Eulerverfahren: $\begin{vmatrix} 1 & 1 \\ 1 & 1 \end{vmatrix} \Rightarrow S(z) = \frac{1}{1 - z}.$

• Explizite Trapezregel: $\begin{vmatrix} 0 & 0 \\ 1 & 1 \\ \hline \frac{1}{2} & \frac{1}{2} \end{vmatrix} \Rightarrow S(z) = 1 + z + \frac{1}{2}z^2.$

• Implizite Mittelpunktsregel: $\begin{vmatrix} \frac{1}{2} & \frac{1}{2} \\ \hline 1 & 1 \end{vmatrix} \Rightarrow S(z) = \frac{1 + \frac{1}{2}z}{1 - \frac{1}{2}z}.$

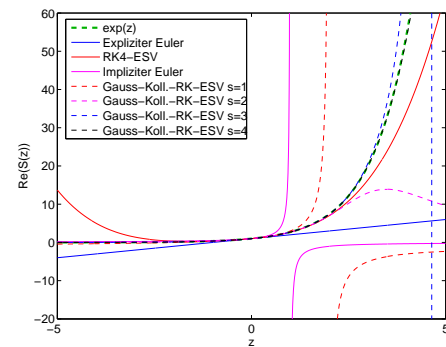
• RK4-Verfahren (5.3.8): $\begin{vmatrix} 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ \hline 1 & 0 & 0 & 1 \\ \frac{1}{6} & \frac{2}{6} & \frac{2}{6} & \frac{1}{6} \end{vmatrix} \Rightarrow S(z) = 1 + z + \frac{1}{2}z^2 + \frac{1}{6}z^3 + \frac{1}{24}z^4.$

Korollar 5.7.1.

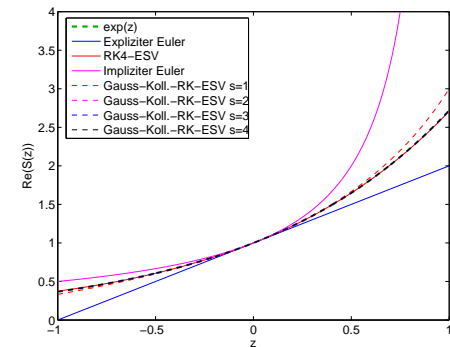
Explizite Runge-Kutta-Verfahren $\Rightarrow S(z) \in \mathcal{P}_s,$
Allgemeine Runge-Kutta-Verfahren $\Rightarrow S(z) = \frac{P(z)}{Q(z)}, P, Q \in \mathcal{P}_s.$

5.7
p. 609

EXSTABFN
Beispiele: Verhalten von Stabilitätsfunktionen (für reelles Argument z):



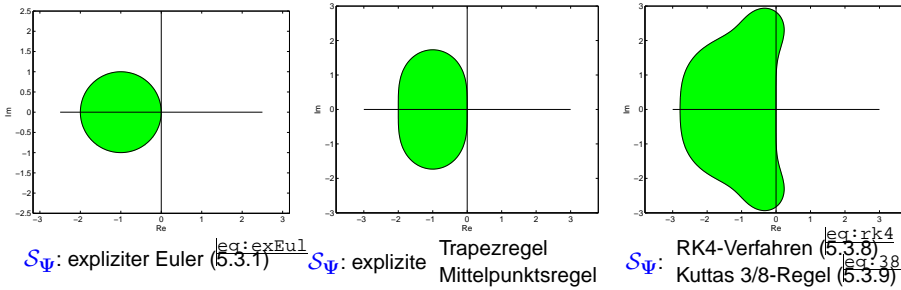
EXSTABDOM
Beispiele: Stabilitätsgebiete S expliziter RK-ESV:



5.7
p. 610

5.7
p. 611

5.7
p. 612



Das Stabilitätsgebiet expliziter RK-Verfahren ist beschränkt !

$\lim_{k \rightarrow \infty} y_k = 0$ für von explizitem RK-Verfahren für AWP (5.6.1) erzeugte Folge $\{y_k\}_{k=0}^{\infty}$ nur wenn $|\lambda h|$ hinreichend klein ($\lambda h \in S_{\Psi}$).

Nicht durch Genauigkeitsanforderungen bedingte **Schrittweitenbeschränkung** !
(für explizite (RK)-ESV)

5.7.1 Steifheit

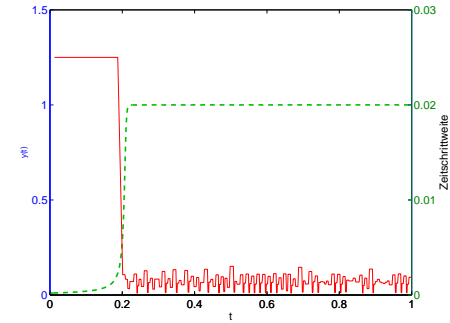
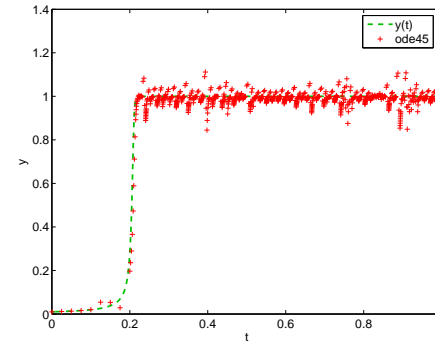
Konzept 5.7.2 (Steifes Anfangswertproblem). Ein AWP heisst **steif** (engl. stiff), falls für explizite ESV (\rightarrow Def. 5.3.2) Stabilität eine wesentlich kleinere Schrittweite verlangt als die Genauigkeitsanforderungen.

Beispiel 233 (Adaptive explizite Einschrittverfahren für steifes Problem).

$$\dot{y}(t) = \lambda y^2(1 - y), \quad \lambda = 500, \quad y(0) = \frac{1}{100}.$$

MATLAB-CODE: Adaptives ESV für steifes Problem

```
fun = @(t,x) 500*x^2*(1-x); tspan = [0 1]; y0 = 0.01;
options = odeset('reltol',0.1,'abstol',0.001,'stats','on');
[t,y] = ode45(fun,tspan,y0,options);
plot(t,y,'r+');
```



➤ Schrittweitensteuerung realisiert Schrittweitenbeschränkung !
(186 successful steps, 55 failed attempts, 1447 function evaluations)

$y = 1$ stark attraktiver Fixpunkt

➤ Extreme Schrittweitenbeschränkung für expliziten Integrator ode45

Welche Anfangswertprobleme sind steif ?

ODE-Modelle für Systeme mit schnell relaxierenden Komponenten
(mit stark unterschiedlichen Zeitkonstanten)

Beispiel 234 (Steife Schaltkreisgleichungen im Zeitbereich).

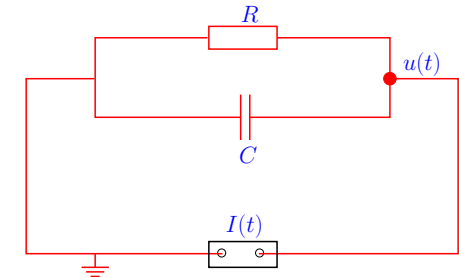
Schaltkreisanalyse im Zeitbereich:
(vgl. Frequenzbereichsanalyse Bsp. 66)

Bauelementgleichungen (Zeitbereich):

- Widerstand: $i(t) = R^{-1}u(t)$
- Kondensator: $i(t) = C\dot{u}(t)$

Dgl. aus Knotenanalyse:

$$C\dot{u}(t) = -R^{-1}u(t) + I(t).$$



Konkret: $C = 1\text{pF}$, $R = 1\text{k}\Omega$, $I(t) = \sin(2\pi \cdot 1\text{Hz}t)\text{mA}$, $u(0) = 0\text{V}$

➤ Skalierte (dimensionslose) Dgl.: $\dot{u}(t) = -10^9 u(t) + 10^9 \sin(2\pi t) \Rightarrow u(t) \approx \sin(2\pi t).$

Steife AWP ➤ (Geeignete) implizite RK-ESV (\rightarrow Def. 5.3.2) sind zu verwenden !

Übung 5.8. Das Anfangswertproblem $y(0) = 1, \dot{y}(0) = 0$ für die lineare Differentialgleichung 2. Ordnung

$$\ddot{y} + \dot{y} + y = 0 \quad (5.8.1)$$

soll mit einem 2-stufigen sogenannten SDIRK-Verfahren (singly diagonally implicit Runge-Kutta method) numerisch gelöst werden. Dies ist ein Runge-Kutta-Verfahren beschrieben durch das Butcher-Tableau

$$\begin{array}{c|cc} & \gamma & 0 \\ \hline 1-\gamma & 1-2\gamma & \gamma \\ \hline & 1/2 & 1/2 \end{array} \quad (5.8.2)$$

- Gib die Gleichungen für die Inkremente k_1 und k_2 des Runge-Kutta-Verfahrens angewandt auf ein Anfangswertproblem zur Differentialgleichung $\dot{y} = f(t, y)$ an.
- Bestimme die Stabilitätsfunktion $S(z)$ des SDIRK-Verfahrens für $\gamma = 1$. Ist das Verfahren A/L-stabil?
- Formuliere als Anfangswertproblem für ein lineares System 1. Ordnung für die Funktion $\mathbf{z}(t) = (y(t), \dot{y}(t))^T$.
- Implementiere eine MATLAB-Funktion

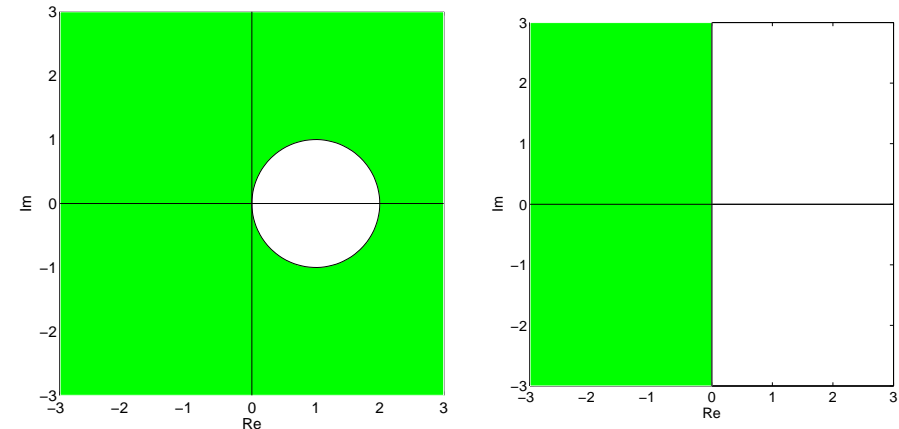
```
z = sdirkstep(z0,h),
```

die die numerische Evolution des Einschrittverfahrens für die in (ii) erhaltene Differentialgleichung realisiert.

- Führe mit Hilfe von MATLAB ein numerisches Experiment aus, das einen Hinweis auf die Ordnung des Verfahrens (mit $\gamma = \frac{3+\sqrt{3}}{6}$) für das Anfangswertproblem aus Teilaufgabe (ii) liefert.

5.8.1 Einschrittverfahren für steife AWP

Beispiele: Stabilitätsgebiete (\rightarrow Def. 5.6.1) impliziter Einschrittverfahren (\rightarrow Def. 5.3.2)



S_Ψ : implizites Eulerverfahren (5.3.1)

S_Ψ : implizite Trapezregel

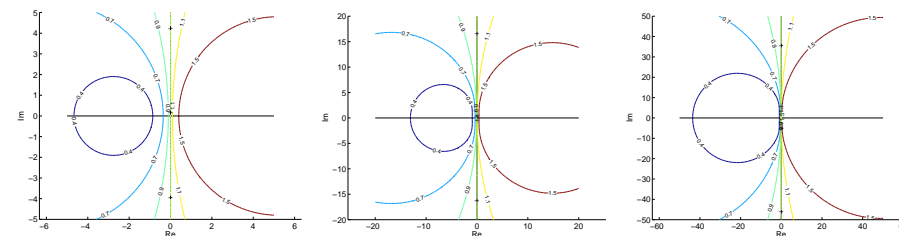
AWP & $\lambda \gg 1$: Keine stabilitätsbedingte Schrittweitenbeschränkung für beide Verfahren (vgl. Verhalten von $S(z)$, $\text{Re } z < 0$)

Definition 5.8.1 (A-stabiles Einschrittverfahren).

ESV **A-stabil** $\Leftrightarrow \{z \in \mathbb{C} : \text{Re } z < 0\} \subset S_\Psi$ (S_Ψ = Stabilitätsgebiet, Def. 5.6.1)

Autonome Dgl. $\dot{\mathbf{y}} = f(\mathbf{y})$, Stationärer Punkt $f(\mathbf{y}^*) = 0$ \wedge $\text{Re}\{\sigma(Df(\mathbf{y}^*))\} \subset \mathbb{R}^-$ ($\hat{=}$ Stabilität von \mathbf{y}^*) \Rightarrow Falls $\|\mathbf{y}_0 - \mathbf{y}^*\| < \delta$, dann $\lim_{k \rightarrow \infty} \mathbf{y}_k = 0$.

Stabilität von Gauss-Kollokations-Einschrittverfahren (\rightarrow Abschnitt 5.3.1):



Implizite Mittelpunktsregel

$s = 2$ (Ordnung 4)

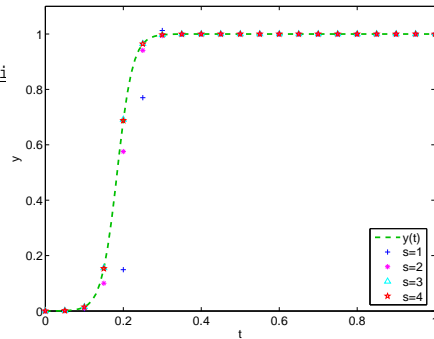
$s = 4$ (Ordnung 8)

Niveaulinien von $|S(z)|$ für Gauss-Kollokations-Einschrittverfahren

Beispiel 235 (Gauss-Kollokationsverfahren für logistische Differentialgleichung). → Bsp. 232, ^{ex:logexi}

Logistische Differentialgleichung $\dot{y} = f(y)$, $f(y) = \lambda y(1 - y) \rightarrow$ (5.3.4), $\lambda = 50$, Anfangswert $y_0 = 10^{-4}$, Zeitintervall $[0, 1]$.

Kollokations-Einschrittverfahren (→ Abschnitt 5.3.1) auf äquidistantem Gitter, $h = \frac{1}{20}$.

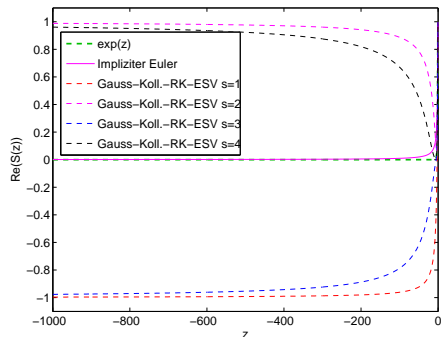


Theorem 5.8.2 (Stabilität von Gauss-Kollokations-Einschrittverfahren).
Alle Gauss-Kollokations-RK-ESV (5.3.2) sind A-stabil. ^{eg:rkimpl}

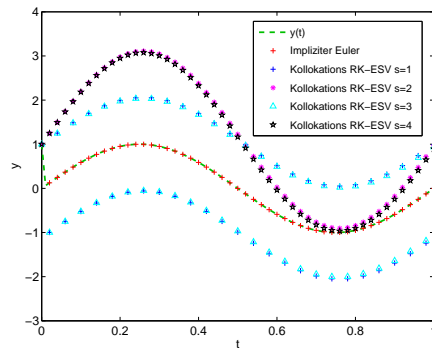
Beispiel 236 (Implizite RK-ESV bei schnellen Transienten). → Bsp. 234 ^{ex:odecircuit}

AWP: $\dot{y} = -\lambda y + \beta \sin(2\pi t)$, $\lambda = 10^6$, $\beta = 10^6$, $y(0) = 1$.

RK-ESV, äquidistantes Gitter auf $[0, 1]$, $h = \frac{1}{40}$:



Stabilitätsfunktionen für $\text{Re } z \leq 1$



Diskrete Evolutionen (Zeitverlauf)

➤ Ungenügende Dämpfung der Anfangsstörung bei Kollokations-RK-ESV!
(Oszillationen für ungerades $s \rightarrow$ vgl. Stabilitätsfunktion!)

Definition 5.8.3 (L-Stabilität).

ESV L-stabil $\Leftrightarrow \{z \in \mathbb{C} : \text{Re } z < 0\} \subset \mathcal{S}_\Psi$ & $\lim_{\text{Re } z \rightarrow -\infty} |S(z)| = 0$

Kurz:

L-stabil \Leftrightarrow A-stabil & „ $S(-\infty) = 0$ “

Wie findet man L-stabile RK-ESV?

$S(-\infty) = 1 - \mathbf{b}^T \mathbf{a}^{-1} \mathbf{1}$: Falls $\mathbf{b}^T = \mathbf{a}_{:,j}^T$ (Zeile von \mathbf{a}) $\Rightarrow S(-\infty) = 0$.



Idee: • Wähle $c_s = 1$ im Kollokations-RK-ESV (5.3.2) ^{eg:rkimpl}
• Wähle c_1, \dots, c_{s-1} als Knoten einer Quadraturformel maximaler Ordnung.
(→ Gauss-Radau-Quadratur, Ordnung $2s - 1$)

➤ Implizite s -stufige L-stabile Radau-ESV, Konvergenzordnung $2s - 1$ (→ Thm. 5.4.4) ^{thm:cvgkol}

$$\frac{1}{1} \mid \frac{1}{1}$$

$$\frac{1}{3} \mid \frac{5}{12} \mid -\frac{1}{12}$$

$$\frac{4-\sqrt{6}}{10} \mid \frac{88-7\sqrt{6}}{360} \mid \frac{296-169\sqrt{6}}{1800} \mid -\frac{2+3\sqrt{6}}{225}$$

5.8
p. 621

Impliziter Euler

Radau-ESV, Ordnung 3

Radau-ESV, Ordnung 5

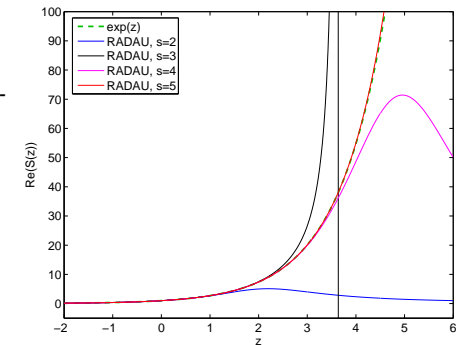
5.8
p. 623

Stabilitätsfunktion s -stufiger Radau-Kollokations-RK-ESVs:

$$S(z) = \frac{P(z)}{Q(z)}, \quad P \in \mathcal{P}_{s-1}, Q \in \mathcal{P}_s.$$

Vorsicht:

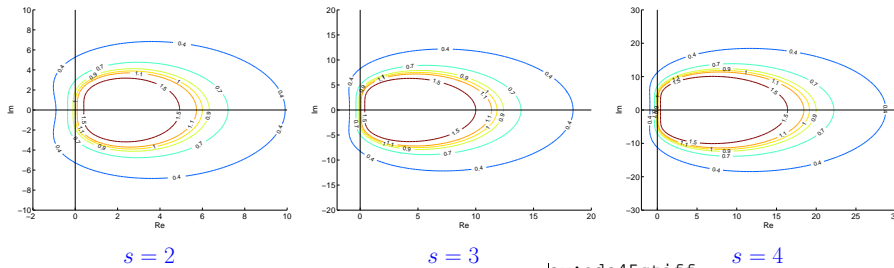
Auch „ $S(\infty) = 0$ “



Niveaus der Stabilitätsfunktionen von s -stufigen Radau-Kollokations-RK-ESVs:

5.8
p. 622

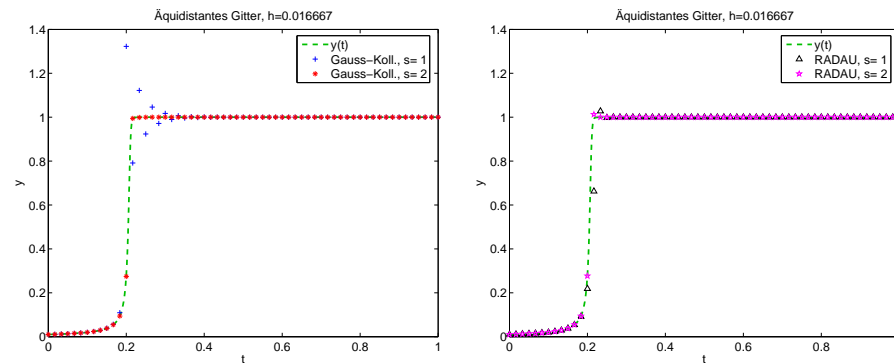
5.8
p. 624



Beispiel 237 (Implizite ESV für steifes Problem). → Bsp. 233

$$\dot{y}(t) = \lambda y^2(1-y), \quad \lambda = 500, \quad y(0) = \frac{1}{100}.$$

Qualitatives Verhalten von Gauss-Kollokations-ESV (→ Abschnitt 5.3.1) & RADAU-ESV auf äquidistantem Gitter:



➤ Falsche Oszillationen bei Gauss-Kollokations-ESV niedriger Ordnung

5.8.2 Implementierung impliziter Runge-Kutta-Einschrittverfahren

Nichtlineares Gleichungssystem für Inkremente \mathbf{k}_i (→ Def. 5.3.3):

$$\mathbf{k}_i = f(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^s a_{ij} \mathbf{k}_j) \quad \blacktriangleright \quad s \cdot d \text{ Unbekannte!}$$

$$\begin{aligned} \vec{\mathbf{g}} &= (\mathbf{g}_1, \dots, \mathbf{g}_s)^T, \\ \mathbf{g}_j &:= h \sum_{i=1}^s a_{ij} f(t_0 + c_i h, \mathbf{y}_0 + \mathbf{g}_i) \quad \blacktriangleright \quad F(\vec{\mathbf{g}}) = \vec{\mathbf{g}} - h \mathfrak{A} \begin{pmatrix} f(t_0 + c_1 h, \mathbf{y}_0 + \mathbf{g}_1) \\ \vdots \\ f(t_0 + c_s h, \mathbf{y}_0 + \mathbf{g}_s) \end{pmatrix} \stackrel{!}{=} \mathbf{0}. \end{aligned}$$

Approximative Lösung mit vereinfachtem Newton-Verfahren → Bem. 48: Startnäherung $\vec{\mathbf{g}} = \mathbf{0}$

$$\text{Jacobimatrix } DF(0) = \begin{pmatrix} \mathbf{I} - h a_{11} \frac{\partial f}{\partial \mathbf{y}}(t_0, \mathbf{y}_0) & \cdots & -h a_{1s} \frac{\partial f}{\partial \mathbf{y}}(t_0, \mathbf{y}_0) \\ \vdots & \ddots & \vdots \\ -h a_{s1} \frac{\partial f}{\partial \mathbf{y}}(t_0, \mathbf{y}_0) & \cdots & \mathbf{I} - h a_{ss} \frac{\partial f}{\partial \mathbf{y}}(t_0, \mathbf{y}_0) \end{pmatrix}.$$

➤ In jedem (vereinfachten) Newton-Schritt: Lösen eines linearen Gleichungssystems mit Matrix $DF(0) \in \mathbb{R}^{sd, sd}$.

5.8 Adaptive MATLAB-Integratoren für steife Probleme: (Schrittweitensteuerung wie in Abschnitt 5.4.1) p. 625

```
opts = odeset('abstol', atol, 'reltol', rtol, 'Jacobian', @J);
[t, y] = ode15s/ode23s(odefun, tspan, y0, opts);
```

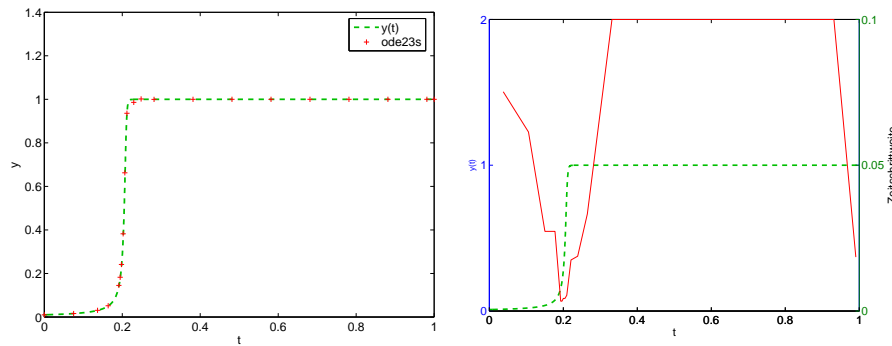
Beispiel 238 (Adaptives semi-implizites RK-ESV für steifes Problem). Bsp. 233, 237

$$\dot{y}(t) = \lambda y^2(1-y), \quad \lambda = 500, \quad y(0) = \frac{1}{100}.$$

MATLAB-CODE: Semi-Implizites ESV für steifes Problem

```
lambda = 500; tspan = [0 1]; y0 = 0.01;
fun = @(t,x) lambda*x^2*(1-x);
Jac = @(t,x) lambda*(2*x*(1-x)-x^2);
o = odeset('reltol', 0.1, 'abstol', 0.001, 'stats', 'on', 'Jacobian', Jac);
[t, y] = ode23s(fun, [0 1], y0, o);
```

Statistik: 20 successful steps 4 failed attempts 70 function evaluations



➤ Effizientes Verfahren (vgl. Bsp. 233): Keine Schrittweitenbeschränkung für $y \approx 1$

5.9 Differenziell-Algebraische Anfangswertprobleme

[File: section-differentiell-algebraische-anfangswertprobleme.tex, SVN: section-differentiell-algebraische-anfangswertprobleme.tex 1061 2006-10-20 16:58:14;

Gegeben: „Rechte Seiten“ $h : \Omega_1 \times \Omega_2 \subset \mathbb{R}^p \times \mathbb{R}^q \mapsto \mathbb{R}^p$, $g : \Omega_1 \times \Omega_2 \mapsto \mathbb{R}^q$ hinreichend glatt,
 $p, q \in \mathbb{N}$, $\mathbf{u}_0 \in \Omega_1$, $\mathbf{v}_0 \in \Omega_2$

☞ **Differenziell-algebraisches Anfangswertproblem (DAE):** Algebraische Nebenbedingung

$$\begin{aligned} \dot{\mathbf{u}} &= h(\mathbf{u}, \mathbf{v}), & \mathbf{u}(0) &= \mathbf{u}_0, & g(\mathbf{u}_0, \mathbf{v}_0) &= 0. \\ 0 &= g(\mathbf{u}, \mathbf{v}), & \mathbf{v}(0) &= \mathbf{v}_0, & & \end{aligned} \quad (5.9.1) \quad \text{eq:DAE}$$

Konsistente Anfangswerte erforderlich!

Beispiel 239 (Deskriptorform der Pendelgleichung). → Bsp. 216 lex: odemech

Pendelgleichung in nicht-minimalen Koordinaten (Ursprung $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ an Aufhängung):

$$m \begin{pmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{p}_1 \\ \dot{p}_2 \end{pmatrix} = \begin{pmatrix} p_1 \\ p_2 \\ -\lambda x_1 \\ -\lambda x_2 - g \end{pmatrix}, \quad x_1^2 + x_2^2 = l^2.$$

λ = Lagrange-Multiplikator ($\leftrightarrow \mathbf{v}$ in (5.9.1)) leg: DAE ➤ Zwangskraft Zwangsbedingung

Differenzieren der Zwangsbedingung ➤ $x_1 p_1 + x_2 p_2 = 0$, (5.9.2) eq:diffnb

Nochmaliges Differenzieren ➤ $(p_1^2 + p_2^2) - \lambda(x_1^2 + x_2^2) - g x_2 = 0$. (5.9.3) eq:diffnb

Annahme: (trifft auf Bsp. 239 nur für differenzierte Nebenbedingung (5.9.3) zu!) ex: daependulum

$\frac{\partial g}{\partial \mathbf{v}} \in \mathbb{R}^{q,q}$ regulär entlang einer Lösungskurve von (5.9.1) leg: DAE

➤ Lokale Auflösbarkeit: $\mathbf{v} = G(\mathbf{u})$: leg: DAE (5.9.1) $\Rightarrow \dot{\mathbf{u}} = f(\mathbf{u}, G(\mathbf{u})) \hat{=} \hat{\mathbf{u}} = f(\mathbf{u})$. leg: awp (5.1.1).

DAE von höherem Index:  Konzept zur Klassifikation von DAEs HAW91 [26]

Einschrittverfahren für (5.9.1)? leg: DAE



Idee: Singuläre Störungstechnik

① Betrachte AWP für Dgl. $\mathbf{M}\dot{\mathbf{y}} = f(\mathbf{y})$, $\mathbf{M} \in \mathbb{R}^{d,d}$, $d = p + q$

② Unter Annahme \mathbf{M} regulär: RK-ESV für $\dot{\mathbf{y}} = \mathbf{M}^{-1}f(\mathbf{y})$

③ Macht Verfahren noch Sinn für singuläres \mathbf{M} ? Wenn ja → 😊

➤ Inkrementgleichungen, vgl. Def. 5.3.3: $\mathbf{M}\mathbf{k}_i = f(\mathbf{y}_0 + h \sum_{j=1}^s a_{ij}\mathbf{k}_j)$. (5.9.4) eq: Minc

➤ Falls \mathbf{M} singulär: (5.9.4) nur sinnvoll für implizite RK-ESV! leg: Minc

Anwendung auf (5.9.1): eq: vns

$$\mathbf{M} = \begin{pmatrix} \mathbf{I} & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathbf{I} \in \mathbb{R}^{p,p}, \quad \mathbf{y} = \begin{pmatrix} \mathbf{u} \\ \mathbf{v} \end{pmatrix}, \quad f(\mathbf{y}) = \begin{pmatrix} h(\mathbf{u}, \mathbf{v}) \\ g(\mathbf{u}, \mathbf{v}) \end{pmatrix}.$$

Eignen sich alle impliziten ESV?

$$\begin{aligned} \text{Modellproblem: } \begin{aligned} \dot{\mathbf{u}} &= v f(\mathbf{u}), \\ 0 &= 1 - v \end{aligned} & \quad \Rightarrow \quad \begin{aligned} \dot{\mathbf{u}} &= v f(\mathbf{u}), \\ \epsilon \dot{\mathbf{v}} &= 1 - v \end{aligned}, \quad 0 < \epsilon \ll 1. \end{aligned}$$

Singulär gestörte Dgl.

Heuristik: ESV tauglich für Modellproblem \leftrightarrow ESV tauglich für singulär gestörtes Problem $\forall \epsilon \approx 0$

➤ ESV muss für AWP $\dot{v} = \epsilon^{-1}(1 - v)$, $v(0) = 1$, $0 < \epsilon \ll 1$, Folge $\{v_k\}_{k=0}^\infty$ mit $\lim_{k \rightarrow \infty} v_k = 1$ liefern!

➤ Notwendig: $S(\infty) = 0$ für Stabilitätsfunktion (\rightarrow Thm. 5.6.2) thm: SRK $S(z)$ des ESV

Kandidaten für RK-ESV für (5.9.1): leg: DAE

Radau-Verfahren

5.10 Strukturhaltung

[File: section-strukturhaltung.tex, SVN: section-strukturhaltung.tex 1010 2006-09-11 15:44:43Z hiptmair]

(Wir betrachten nur autonome Anfangswertprobleme $\dot{\mathbf{y}} = f(\mathbf{y})$, $\mathbf{y}(0) = \mathbf{y}_0$, Phasenraum $\Omega \subset \mathbb{R}^d$)

Struktur = Essentielle Qualitative Eigenschaften der Evolution Φ^t zur ODE
(Hervorragende Darstellung in Monographie [23])

Ziel: „Vererbung“ von Struktur an diskrete Evolution Ψ_h^h

Wichtig für **Langzeitintegration** zur Berechnung einer „qualitativ richtigen Lösung“
(im Sinne der Rückwärtsanalyse \rightarrow Abschnitt 11.8)

5.10.1 Nichtexpansivität

Sei $\mathbf{M} \in \mathbb{R}^{d,d}$ s.p.d. (\rightarrow Def. 3.2.8) $\xrightarrow{\text{Def: spd}}$ Norm $\|\mathbf{y}\|_M := (\mathbf{y}^T \mathbf{M} \mathbf{y})^{1/2}$ auf \mathbb{R}^d .

Definition 5.10.1 (Nichtexpansivität). Eine Evolution Φ^t zu einer autonomen Dgl. heisst **nicht-expansiv**, falls

$$\|\Phi^t \mathbf{y} - \Phi^t \tilde{\mathbf{y}}\|_M \leq \|\mathbf{y} - \tilde{\mathbf{y}}\|_M \quad \forall \mathbf{y}, \tilde{\mathbf{y}} \in \Omega.$$

Analoges Konzept für diskrete Evolutionen zu ESV [11, Abschn. 6.3.3]

Theorem 5.10.2 (B-stabile ESV). Die diskreten Evolutionen zu Gauss-Kollokations- und Radau-RK-ESV erben die Nichtexpansivität der (exakten) Evolution.

Beispiel 240 (Gradientenfluss \rightarrow „Kriechvorgänge“).

Potential $V : \mathbb{R}^d \mapsto \mathbb{R}$ strikt konvex:

$$\exists \alpha > 0: (\text{grad } V(\mathbf{x}) - \text{grad } V(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) \geq \alpha \|\mathbf{x} - \mathbf{y}\|_2^2 \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

Gradientenfluss-AWP: $\dot{\mathbf{y}}(t) = -\text{grad } V(\mathbf{y}(t))$, $\mathbf{y}(0) = \mathbf{y}_0 \in \mathbb{R}^d$.

$$\chi(t) := \|\Phi^t \mathbf{y} - \Phi^t \tilde{\mathbf{y}}\|_2^2 \Rightarrow \dot{\chi}(t) = -2 \underbrace{(\text{grad } V(\Phi^t \mathbf{y}) - \text{grad } V(\Phi^t \tilde{\mathbf{y}}))^T (\Phi^t \mathbf{y} - \Phi^t \tilde{\mathbf{y}})}_{\geq \alpha \|\Phi^t \mathbf{y} - \Phi^t \tilde{\mathbf{y}}\|_2^2} \leq 0.$$

5.10.2 Quadratische erste Integrale

Definition 5.10.3 (Quadratisches erstes Integral). Eine Funktion $I : \mathbb{R}^d \mapsto \mathbb{R}$, $I(\mathbf{x}) = \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{q}^T \mathbf{x} + q$, $\mathbf{Q} \in \mathbb{R}^{d,d}$, $\mathbf{q} \in \mathbb{R}^d$, $q \in \mathbb{R}$ ist **quadratisches erstes Integral** der Evolution Φ^t zur autonomen Dgl. $\dot{\mathbf{y}} = f(\mathbf{y})$, falls

$$I(\Phi^t \mathbf{y}) = I(\mathbf{y}) \quad \forall \mathbf{y} \in \Omega, \forall t \in J(\mathbf{y}).$$

Theorem 5.10.4 (Erhaltung quadratischer erster Integrale). Quadratische erste Integrale der (exakten) Evolution werden auch von der von Gauss-Kollokations-ESV erzeugten diskreten Evolution erhalten.

Beispiel 241 (Präzession einer Magnetnadel).

$\mathbf{y}(t)$ = Trajektorie der Spitze einer Magnetnadel (im äusseren Feld \mathbf{h} \rightarrow Bewegungsgleichung

$$\dot{\mathbf{y}} = \mathbf{y} \times \mathbf{h} \quad , \quad \text{Kreuzprodukt} \quad \mathbf{y} \times \mathbf{h} = \begin{pmatrix} y_2 h_3 - y_3 h_2 \\ y_3 h_1 - y_1 h_3 \\ y_1 h_2 - y_2 h_1 \end{pmatrix} \perp \mathbf{y}$$

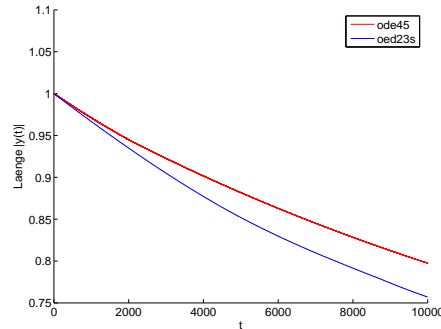
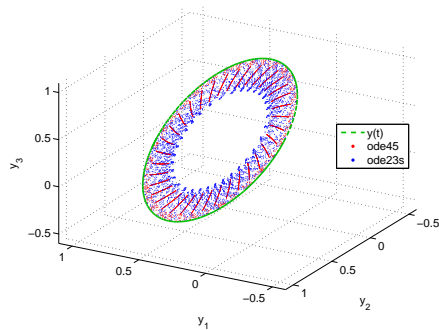
$$(\mathbf{y}(t) \times \mathbf{h}) \cdot \mathbf{y}(t) = 0 \quad \blacktriangleright \quad \frac{d}{dt} \|\mathbf{y}(t)\|_2^2 = 2\dot{\mathbf{y}} \cdot \mathbf{y}(t) = 0 \quad \Rightarrow \quad \|\mathbf{y}(t)\|_2^2 = \text{const}.$$

Anfangswert: $\mathbf{y}_0 = (\frac{1}{2}\sqrt{2}, 0, 1, \frac{1}{2}\sqrt{2})^T$

MATLAB-CODE: Berechnung des Präzession einer Magnetnadel

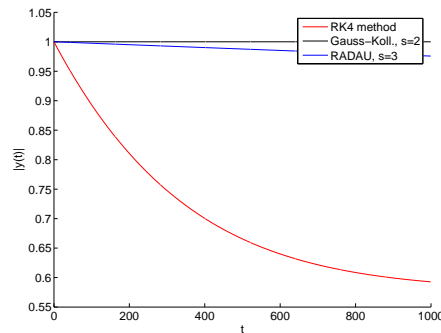
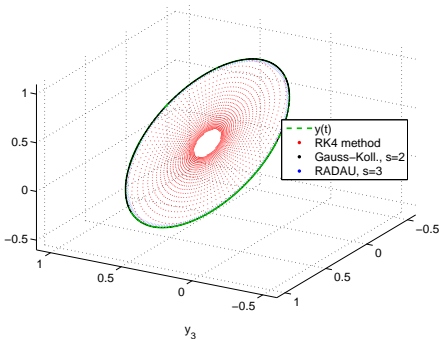
```
h = [-1;-1;-1]; tspan = [0 10000]; y0 = [0.5*sqrt(2);0;0.5*sqrt(2)];
fun = @(t,x) cross(x,h);
Jac = @(t,x) [0 h(3) -h(2); -h(3) 0 h(1); h(2) -h(1) 0];
options = odeset('reltol',0.001,'abstol',1e-4,'stats','on');
[t45,y45] = ode45(fun,tspan,y0,options);
options = odeset('reltol',0.001,'abstol',1e-4,'stats','on','Jacobian',Jac);
[t23,y23] = ode23s(fun,tspan,y0,options);
```

ode45: 24537 successful steps, 7432 failed attempts, 191815 function evaluations
ode23s: 93447 successful steps, 4632 failed attempts, 289607 function evaluations



➤ Keine Erhaltung von $\|y(t)\|_2$ über lange Zeiten

Einschrittverfahren auf äquidistanten Gittern:



5.10.3 Symplektizität

Gegeben: (Glatte) ^{HAM} Hamilton-Funktion $H : \Omega \subset \mathbb{R}^{2d} \mapsto \mathbb{R}, d \in \mathbb{N}$ (Energie)
→ Vorlesung „Theoretische Mechanik“

Oft: $H(\mathbf{p}, \mathbf{q}) = \underbrace{\frac{1}{2} \|\mathbf{p}\|_2^2}_{\text{Kinetische Energie}} + \underbrace{V(\mathbf{q})}_{\text{Potentielle Energie}}, \text{ Potential } V : \Omega_{\mathbf{q}} \mapsto \mathbb{R} \quad (5.10.1) \quad \text{Eq: Hmec}$

^{HAMGL} **Hamiltonsche (Bewegungs-)Gleichungen der klassischen Mechanik**

$$\dot{\mathbf{p}} = -\frac{\partial H}{\partial \mathbf{q}}(\mathbf{p}, \mathbf{q}), \quad \dot{\mathbf{q}} = \frac{\partial H}{\partial \mathbf{p}}(\mathbf{p}, \mathbf{q}). \quad (5.10.2) \quad \text{Eq: FAM}$$

$\mathbf{q} \hat{=}$ (verallgemeinerte) Ortskoordinate, \mathbf{p} (verallgemeinerte) Impulskoordinate

➤ Evolution Φ_H^t zu (5.10.2) ^{Eq: HAM} erfüllt $H(\Phi_H^t(\mathbf{q}, \mathbf{q})) = \text{const}$ für alle $(\mathbf{p}, \mathbf{q})^T \in \Omega$ und Zeiten t
(→ „Energieerhaltung“)

Zusätzliche Struktur: Φ_H^t ^{SYMPLECT} symplektisch → ^{HLW02} [23, Ch. VI]

Erklärung für $d = 1$:

5.10
p. 637

t fest ➤ $\Phi_H^t : \Omega \mapsto \Omega$ Diffeomorphismus des Phasenraums: $(\Phi_H^t)^{-1} = \Phi_H^{-t}$

Lokale „Flächenverzerrung“ ↔ $\det J(t; p, q), J(t; p, q) := \frac{\partial \Phi_H^t(p, q)}{\partial (p, q)}$

$t \mapsto J(t; p, q)$ erfüllt Variationsgleichung (5.2.1): ^{Eq: VARG}

$$j(t; p, q) = \begin{pmatrix} -\frac{\partial^2 H}{\partial p \partial q}(\Phi_H^t(p, q)) & -\frac{\partial^2 H}{\partial q^2}(\Phi_H^t(p, q)) \\ \frac{\partial^2 H}{\partial p^2}(\Phi_H^t(p, q)) & \frac{\partial^2 H}{\partial p \partial q}(\Phi_H^t(p, q)) \end{pmatrix} J(t; p, q).$$

$$\blacktriangleright \frac{d}{dt} \det J(t; p, q) = \frac{d}{dt} \left(\frac{\partial p(t)}{\partial p} \frac{\partial q(t)}{\partial q} - \frac{\partial p(t)}{\partial q} \frac{\partial q(t)}{\partial p} \right) = 0.$$

Ableitungen nach Anfangswerten



Evolution Φ_H^t ist flächenerhaltend

Beispiel 242 (Symplektizität des Phasenflusses für Pendelgleichung). → Bsp. 216 ^{lex: odemech}

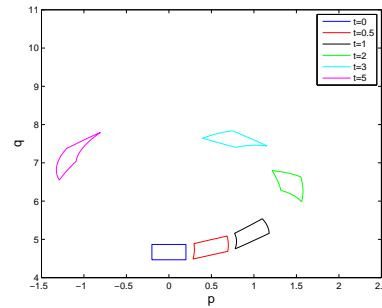
$p \leftrightarrow$ Winkelgeschwindigkeit, $q \leftrightarrow$ Winkelvariable

5.10
p. 638 $\dot{p} = -\sin q, \quad \dot{q} = p \quad \blacktriangleright \quad H(p, q) = \frac{1}{2} p^2 - \cos q$ (Gesamtenergie, vgl. (5.10.1)) ^{Eq: Hmec} ^{Eq: Hpen} (10.3)

5.10
p. 639

5.10
p. 640

Volumenerhaltung im Phasenraum:
Evolution eines quadratischen Volumens ▷



Theorem 5.10.5 (Symplektizität von Gauss-Kollokations-RK-ESV). Diskrete Evolutionen von Gauss-Kollokations-ESVs (→ Abschnitt 5.3.1) auf äquidistanten Gittern erhalten die Symplektizität einer Evolution (⇒ „symplektische Integratoren“).

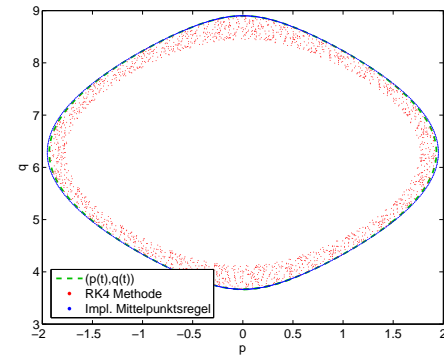
Bemerkung 243 (Warum symplektische Integratoren?).

Beobachtung: Symplektizität von Integrationsverfahren ➤ Hervorragendes qualitatives Langzeitverhalten für AWP (5.10.2)

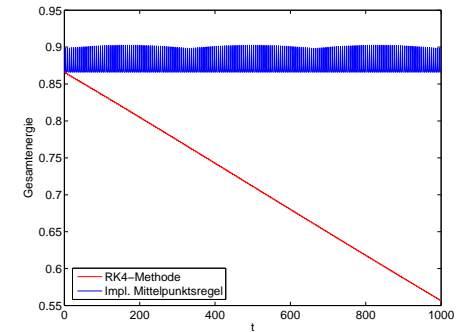
Beispiel 244 (Bewegungsgleichung des starren Pendels). → Bsp. 216

AWP für (5.10.3) auf $[0, 1000]$, $p(0) = 0$, $q(0) = 7\pi/6$.

Vergleich von klassischem Runge-Kutta-Verfahren (5.3.8) (Ordnung 4) mit 1-stufigem Gauss-Kollokations-ESV (implizite Mittelpunktsregel 5.3.3), äquidistantes Gitter, $h = \frac{1}{2}$:



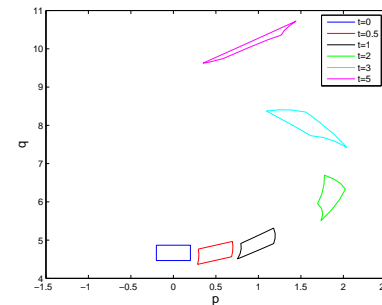
Orbits exakter/diskreter Evolutionen



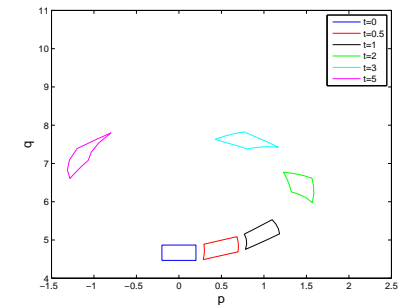
Energieerhaltung diskreter Evolutionen

➤ Keine Energiedrift bei impliziter Mittelpunktsregel (typisch für „symplektische Integratoren“)

5.10
p. 641



Evolution eines quadratischen Volumens
(Explizites Eulerverfahren)



Evolution eines quadratischen Volumens
(Implizite Mittelpunktsregel)

Beispiel 245 (Symplektisches partitioniertes Eulerverfahren). → [23, Kap. VI]

$$\begin{aligned} \dot{\mathbf{u}} &= \mathbf{f}(\mathbf{u}, \mathbf{v}), \\ \dot{\mathbf{v}} &= \mathbf{g}(\mathbf{u}, \mathbf{v}), \end{aligned} \quad \Rightarrow \quad \begin{aligned} \mathbf{u}_h(t+h) &= \mathbf{u}_h(t) + h\mathbf{f}(\mathbf{u}_h(t+h), \mathbf{v}_h(t)), \\ \mathbf{v}_h(t+h) &= \mathbf{v}_h(t) + h\mathbf{g}(\mathbf{u}_h(t+h), \mathbf{v}_h(t)), \end{aligned} \quad (5.10.4)$$

5.10
p. 642

5.10
p. 643

5.10
p. 644

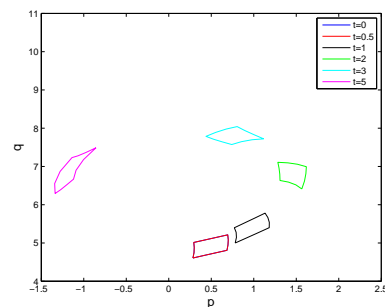
leg:HAM
(5.10.2) mit H gemäss leg:Hmech (5.10.1)

$$\begin{aligned} \mathbf{u} &:= \mathbf{q} \\ \mathbf{v} &:= \mathbf{p} \\ \Rightarrow \quad \mathbf{q}_h(t+h) &= \mathbf{q}_h(t) + h\mathbf{p}_h(t), \\ \mathbf{p}_h(t+h) &= \mathbf{p}_h(t) - h \mathbf{grad} V(\mathbf{q}_h(t+h)), \end{aligned} \quad (5.10.5) \quad \text{eq:se}$$

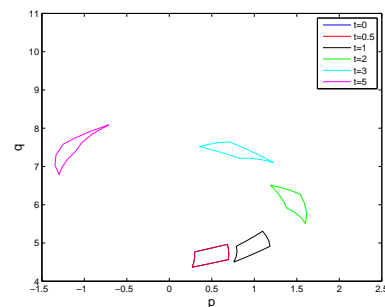
$$\begin{aligned} \mathbf{u} &:= \mathbf{p} \\ \mathbf{v} &:= \mathbf{q} \\ \Rightarrow \quad \mathbf{p}_h(t+h) &= \mathbf{p}_h(t) - h \mathbf{grad} V(\mathbf{q}_h(t)), \\ \mathbf{q}_h(t+h) &= \mathbf{q}_h(t) + h\mathbf{p}_h(t+h). \end{aligned} \quad (5.10.6) \quad \text{eq:se}$$

➤ Explizite Einschrittverfahren !

Anwendung auf Pendelgleichung leg:Hpend (5.10.3) → Bsp. 244 lex:pendsymp



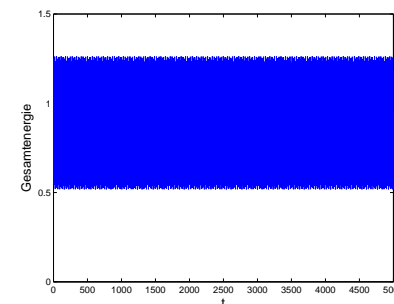
Evolution eines quadratischen Volumens
(Symplektisches partitioniertes Eulerverfahren
leg:sepend1 (5.10.5))



Evolution eines quadratischen Volumens
(Symplektisches partitioniertes Eulerverfahren
leg:sepend (5.10.6))

5.10
p. 645

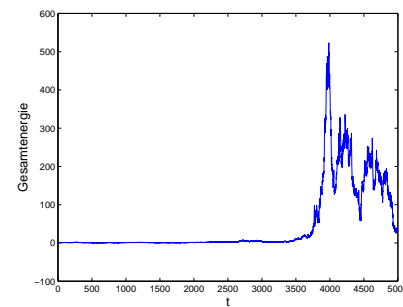
5.10
p. 646



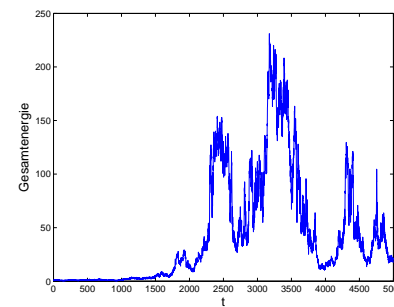
Energieerhaltung des symplektischen partitionierten Eulerverfahrens leg:sepend (5.10.6) ($p(0) = 0, q(0) = \frac{7\pi}{6}$)

Beispiel 246 (Symplektische Integratoren und variable Schrittweite). Fortsetzung Bsp. 245 lex:symlEul

Symplektisches Eulerverfahren leg:sepend (5.10.6) für leg:Hpend (5.10.3) auf $[0, T]$, $T = 5000$. Erratische variable Schrittweite $h_i = 0.5(1 + 0.5(\text{rand}() - 0.5))$, $i = 1, \dots, 10000$, $p(0) = 0, q(0) = 7\pi/6$



Energiedrift bei variabler Schrittweite
(Symplektisches partitioniertes Eulerverfahren
leg:sepend1 (5.10.5))



Energiedrift bei variabler Schrittweite
(Symplektisches partitioniertes Eulerverfahren
leg:sepend (5.10.6))

5.10
p. 647

5.10
p. 648

5.10.4 Reversibilität

$\Phi^{s,t}$ $\hat{=}$ Evolutionsoperator zum AWP (5.1.1) \rightarrow Sect. 5.1

Thm. 5.1.2 $\Rightarrow \forall (t, y) \in \hat{\Omega}: (\Phi^{t,s} \circ \Phi^{s,t})y = y \quad \forall s \in J(t, y) .$ (5.10.7)

Definition 5.10.6 (Symmetrisches Einschrittverfahren). Ein ESV für das AWP (5.1.1) mit diskreter Evolution $\Psi_h^{s,t}$ heisst **symmetrisch**, falls

$\forall (t, y) \in \hat{\Omega}: (\Psi_h^{t,t+h} \circ \Psi_h^{t+h,t})y = y \quad \text{für hinreichend kleines } h .$

Wann ist s -stufiges, $s \in \mathbb{N}$, Runge-Kutta-Einschrittverfahren (\rightarrow Def. 5.3.3)

$k_i = f(y_0 + h \sum_{j=1}^s a_{ij} k_j) \quad , \quad y_1 = y_0 + h \sum_{i=1}^s b_i k_i ,$

für autonomes AWP symmetrisch ? \Rightarrow Falls invariant gegen Vertauschung $y_0 \leftrightarrow y_1, h \leftrightarrow -h !$

Hinreichend: $b_j = a_{ij} + a_{1+s-i,1+s-j} , \quad 1 \leq i, j \leq s .$

Wichtige Konsequenz für symmetrische RK-ESV:

autonomes AWP $\dot{y} = f(y) \quad , \quad y(0) = y_0$ mit $f(Ty) = -Tf(y) \quad \forall y \in \Omega ,$ (5.10.8)

$T : \Omega \mapsto \Omega \hat{=}$ bijektive Selbstabbildung des Phasenraumes.

$\blacktriangleright \quad T \circ \Phi^t = (\Phi^t)^{-1} \circ T \quad \text{für hinreichend kleines } t .$ (5.10.9)

Lemma 5.10.7. Die diskrete Evolution Ψ_h^t für ein Runge-Kutta-Einschrittverfahren (\rightarrow Def. 5.3.3) zu autonomem AWP (5.10.8) erfüllt

$T \circ \Psi_h^h = \Psi_h^{-h} \circ T \quad \text{für kleine } h .$

Z.B. expliziter Euler (5.3.1):

$(T \circ \Phi_h^h)(y) = Ty + hTf(y) \stackrel{(5.10.8)}{=} Ty - hf(Ty) = (\Phi_h^{-h} \circ T)(y) \quad \forall y \in \Omega .$

Korollar 5.10.8 (T-Reversibilität symmetrischer RK-ESV). Die diskrete Evolution Φ_h^t für ein symmetrisches RK-ESV (\rightarrow Def. 5.10.6) zu autonomem AWP (5.10.8) erfüllt für hinreichend kleines h

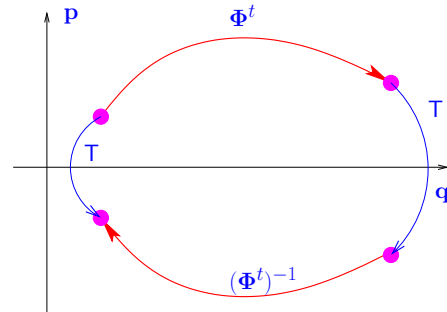
$T \circ \Psi_h^h = (\Psi_h^h)^{-1} \circ T . \quad \leftrightarrow$ (5.10.9)

Warum ist das interessant ?

Reversibilität mechanischer Systeme mit Hamilton-Funktion (5.10.1) \rightarrow Abschn. 5.10.3.

$T \begin{pmatrix} p \\ q \end{pmatrix} = \begin{pmatrix} -p \\ q \end{pmatrix} \blacktriangleright \stackrel{(5.10.2)}{=} \text{AWP der Klasse (5.10.8)}$

5.10
p. 649



Interpretation:
(5.10.7) $\hat{=}$ „Rückwärtsevolution“ nach Umkehr der Geschwindigkeitsrichtung

Verwende $\begin{pmatrix} p \\ q \end{pmatrix} \mapsto \begin{pmatrix} -p \\ q \end{pmatrix}$ -reversible Einschrittverfahren für konservative mechanische Systeme !

Kandidat: Implizite Mittelpunktsregel (5.3.3): symplektisch, symmetrisch implizit

Geht es auch explizit ? (vgl. Bsp. 245)



Idee: Hintereinanderschaltung symplektischer Eulerverfahren (5.10.5), (5.10.6) mit Schrittweiten $h \leftarrow h/2$

\blacktriangleright **Störmer-Verlet-Verfahren** [24]

5.10
p. 650

5.10
p. 651

5.10
p. 652

Für ^{eq:HAM}(5.10.2) mit $H(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \|\mathbf{p}\|_2^2 + V(\mathbf{q})$ gemäss ^{eq:Hmech}(5.10.1)

$$\begin{aligned} \mathbf{q}_h(t + \tfrac{1}{2}h) &= \mathbf{q}_h(t) + \tfrac{h}{2} \mathbf{p}_h(t), \\ \mathbf{p}_h(t+h) &= \mathbf{p}_h(t) - h \mathbf{grad} V(\mathbf{q}_h(t + \tfrac{1}{2}h)), \\ \mathbf{q}_h(t+h) &= \mathbf{q}_h(t + \tfrac{1}{2}h) + \tfrac{h}{2} \mathbf{p}_h(t+h). \end{aligned} \quad (5.10.10) \quad \text{eq:sv}$$

► Symmetrisches, symplektisches, explizites Verfahren der Ordnung 2

Beispiel 247 (Molekulardynamik). (→ Vorlesung „rechnergestützte Chemie“)

• Phasenraum für $n \in \mathbb{N}$ Atome in $d \in \mathbb{N}$ Dimensionen: $\Omega = \mathbb{R}^{2dn}$ (Positionen $\mathbf{q} = [\mathbf{q}^1; \dots; \mathbf{q}^n]^T \in \mathbb{R}^{dn}$, Impulse $\mathbf{p} = [\mathbf{p}^1; \dots; \mathbf{p}^n]^T \in \mathbb{R}^{dn}$)

• Gesamtenergie (Hamilton-Funktion): $H(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \|\mathbf{p}\|_2^2 + V(\mathbf{q})$

^{LJVP}
Lenard-Jones-Potential: $V(\mathbf{q}) = \sum_{j=1}^n \sum_{i \neq j} \mathcal{V}(\|\mathbf{q}^i - \mathbf{q}^j\|_2) \quad , \quad \mathcal{V}(\xi) = \xi^{-12} - \xi^{-6}.$

(5.10.11) eq:LJP

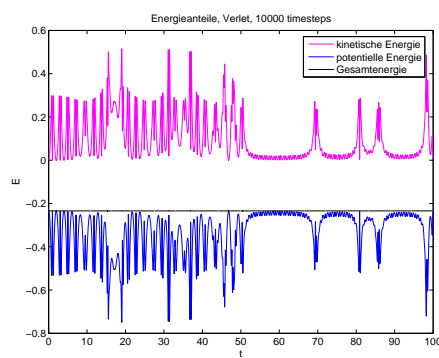
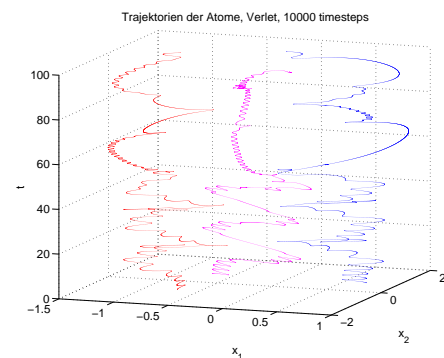
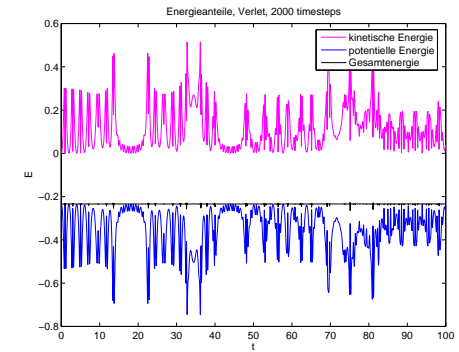
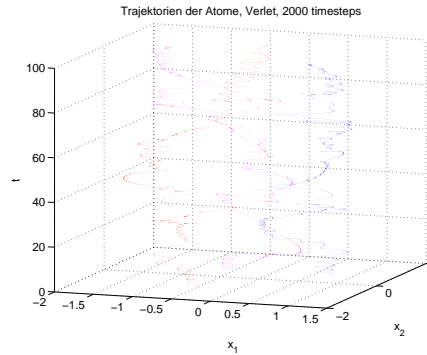
Bewegungsgleichungen ^{eq:HAM}(5.10.2):

$$\dot{\mathbf{p}}^j = - \sum_{i \neq j} \mathcal{V}'(\|\mathbf{q}^j - \mathbf{q}^i\|_2) \frac{\mathbf{q}^j - \mathbf{q}^i}{\|\mathbf{q}^j - \mathbf{q}^i\|_2} \quad , \quad \dot{\mathbf{q}}^j = \mathbf{p}^j \quad , \quad j = 1, \dots, n.$$

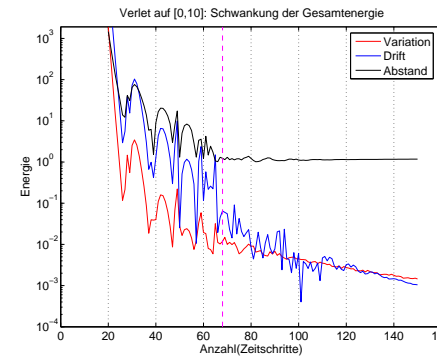
► Störmer-Verlet-Verfahren ^{eq:sv}(5.10.10):

$$\begin{aligned} \mathbf{q}_h(t + \tfrac{1}{2}h) &= \mathbf{q}_h(t) + \tfrac{h}{2} \mathbf{p}_h(t), \\ \mathbf{p}_h^j(t+h) &= \mathbf{p}_h^j(t) - h \sum_{i \neq j} \mathcal{V}'(\|\mathbf{q}_h^j(t + \tfrac{1}{2}h) - \mathbf{q}_h^i(t + \tfrac{1}{2}h)\|_2) \frac{\mathbf{q}_h^j(t + \tfrac{1}{2}h) - \mathbf{q}_h^i(t + \tfrac{1}{2}h)}{\|\mathbf{q}_h^j(t + \tfrac{1}{2}h) - \mathbf{q}_h^i(t + \tfrac{1}{2}h)\|_2}, \\ \mathbf{q}_h(t+h) &= \mathbf{q}_h(t + \tfrac{1}{2}h) + \tfrac{h}{2} \mathbf{p}_h(t+h). \end{aligned}$$

Simulation mit $d = 2, n = 3, \mathbf{q}^1(0) = \frac{1}{2}\sqrt{2} \begin{pmatrix} -1 \\ 1 \end{pmatrix}, \mathbf{q}^2(0) = \frac{1}{2}\sqrt{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \mathbf{q}^3(0) = \frac{1}{2}\sqrt{2} \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \mathbf{p}(0) = 0$, Endzeitpunkt $T = 100$



5.10
p. 653



Beispiel 248 (Harmonischer Oszillator).

Modellproblem für ^{eq:HAM}(5.10.2): $p, q \in \mathbb{R}, \Omega = \mathbb{R}^2, H(p, q) = \frac{1}{2}p^2 + \frac{1}{2}aq^2$

$$\begin{aligned} \text{eq:HAM} \quad (5.10.2) \quad \text{eq:harmosc} \quad \dot{p} &= -aq, \quad \dot{q} = p \Leftrightarrow \ddot{q} = -aq^2. \end{aligned} \quad (5.10.12) \quad \text{eq:harm}$$

Äquivalente Formulierung des Störmer-Verlet-Verfahrens ^{eq:st}(5.10.10) für ^{eq:harmosc}(5.10.12)

5.10
p. 655

5.10
p. 654

5.10
p. 656

$$\begin{aligned} q_h(t - \tfrac{1}{2}h) &= q_h(t - h) + \tfrac{1}{2}hp_h(t - h), & q_h(t + \tfrac{1}{2}h) &= q_h(t) + \tfrac{1}{2}hp_h(t), \\ p_h(h) &= p_h(t - h) - ahq_h(t - \tfrac{1}{2}h), & p_h(t + h) &= p_h(t) - ahq_h(t + \tfrac{1}{2}h), \\ q_h(h) &= q_h(t - \tfrac{1}{2}h) + \tfrac{1}{2}hp_h(t), & q_h(t + h) &= q_h(t + \tfrac{1}{2}h) + \tfrac{1}{2}hp_h(t + h), \end{aligned}$$

\Updownarrow

$$q_h(t + h) - 2q_h(t) + q_h(t - h) = -ah^2q_h(t) \quad (5.10.13) \quad \text{eq:stv3i}$$

Dreitermrekursion

Ansatz: $q_h(kh) = \xi^k \xrightarrow{\text{eq:stv3t} \text{ (5.10.13)}} \xi^2 - (2 - ah^2)\xi + 1 = 0$
CHAREQ
 Characteristische Gleichung of (5.10.13) eq:stv3t

$$\exists \xi_* \in \mathbb{C}: \xi_*^2 - (2 - ah^2)\xi_* + 1 = 0 \wedge |\xi_*| > 1 \iff ah^2 > 4.$$

\exists Lösungen von (5.10.13): $|q_h(kh)| \rightarrow \infty$ für $k \rightarrow \infty \xrightarrow{\text{eq:stv3t} \text{ (5.10.13)}} \text{Instabilität}$

5.11 Splittingverfahren

[File: section-splittingverfahren.tex, SVN: section-splittingverfahren.tex 1300 2007-02-01 15:58:07Z kalai]

Autonomes AWP mit additiv zerlegter rechter Seite:

$$\dot{y} = f(y) + g(y), \quad y(0) = y_0,$$

mit $f: \Omega \subset \mathbb{R}^d \mapsto \mathbb{R}^d, g: \Omega \subset \mathbb{R}^d \mapsto \mathbb{R}^d$ Lipschitz-stetig.

Evolutionen: $\Phi_f^t \leftrightarrow \text{AWP für } \dot{y} = g(y),$
 $\Phi_g^t \leftrightarrow \text{AWP für } \dot{y} = h(y).$

Annahme: Φ_f^t, Φ_g^t (analytisch) bekannt



Idee: Konstruiere Einschrittverfahren gemäss

$$\begin{aligned} y_h(t + h) &= (\Phi_g^h \circ \Phi_f^h) y_h(t), \\ y_h(t + h) &= (\Phi_f^{h/2} \circ \Phi_g^h \circ \Phi_f^{h/2}) y_h(t). \end{aligned}$$

Terminologie:

Splitting LIE-TROTTER
 Splitting Strang Lie-Trotter-Splitting
 Splitting (5.11.3) Strang-Splitting

Theorem 5.11.1 (Konsistenzordnung einfacher Splittingverfahren). Die ESV (5.11.2) und (5.11.3) haben die Konsistenzordnungen (\rightarrow Def. 5.4.2) 1 bzw. 2. eq:lietrotter (5.11.2) und def:ssmloccons

Beispiel 249 (Konvergenz einfacher Splittingverfahren).

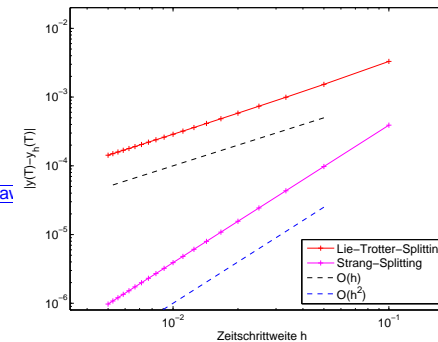
$$\dot{y} = \underbrace{\lambda y(1-y)}_{=:f(y)} + \underbrace{\sqrt{1-y^2}}_{=:g(y)}, \quad y(0) = 0.$$

► $\Phi_f^t y = \frac{1}{1 + (y^{-1} - 1)e^{-\lambda t}}, t > 0, y \in]0, 1]$ (Logistische Differentialgleichung eq:logdg1 (5.3.4))

► $\Phi_g^t y = \begin{cases} \sin(t + \arcsin(y)) & , \text{ falls } t + \arcsin(y) < \frac{\pi}{2}, \\ 1 & , \text{ sonst,} \end{cases} \quad t > 0, y \in [0, 1].$

5.11
p. 657

(5.11.1) eq:splitva



Numerisches Experiment:

$T = 1, \lambda = 1$, Vergleich von Splittingverfahren (konstante Schrittweite) mit hochgenauer numerischer Lösung erhalten durch

```
f=@(t,x) lambda*x*(1-x)+sqrt(1-x^2);
options=odeset('reltol',1.0e-10,...
               'abstol',1.0e-12);
[t,yex]=ode45(f,[0,1],y0,options);
```

◁ Fehlerverhalten zum Endzeitpunkt $T = 1$

5.11
p. 659



Idee: Ersetze

Exakte Evolutionen \rightarrow diskrete Evolutionen
 $\Phi_g^h, \Phi_f^h \rightarrow \Psi_{h,g}^h, \Psi_{h,f}^h$

(5.11.2) eq:lietrot

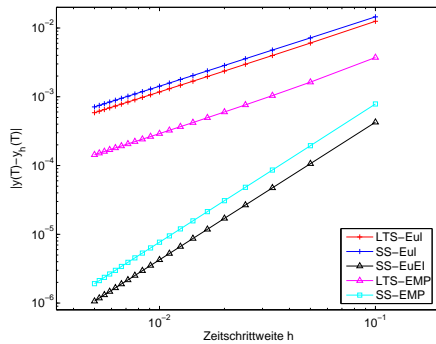
(5.11.3) eq:strang

Beispiel 250 (Inexakte Splittingverfahren). Fortsetzung Bsp. 249 ex:splitcva

AWP von Bsp. 249, inexakte Splittingverfahren auf der Grundlage verschiedener inexakter Basisverfahren: ex:splitcva

5.11
p. 658

5.11
p. 660



LTS-Eul Explizites Eulerverfahren (5.3.1) $\xrightarrow{\text{leg:exEul}}$
 $\Psi_{h,g}^h, \Psi_{h,f}^h$ + Lie-Trotter-Splitting (5.11.2) $\xrightarrow{\text{leg:lietrotter}}$
 SS-Eul Explizites Eulerverfahren (5.3.1) $\xrightarrow{\text{leg:exEul}}$
 $\Psi_{h,g}^h, \Psi_{h,f}^h$ + Strang-Splitting (5.11.3) $\xrightarrow{\text{leg:strang}}$
 SS-EuEI Strang-Splitting (5.11.3): Explizites Eulerverfahren (5.3.1) $\xrightarrow{\text{leg:exEul}}$ exakte Evolution $\xrightarrow{\text{leg:implEul}}$ implizites Eulerverfahren (5.3.1) $\xrightarrow{\text{leg:exMP}}$
 LTS-EMP Explizite Mittelpunktsregel (5.3.6) $\xrightarrow{\text{leg:exMP}}$
 $\Psi_{h,g}^h, \Psi_{h,f}^h$ + Lie-Trotter-Splitting (5.11.2) $\xrightarrow{\text{leg:lietrotter}}$
 SS-EMP Explizite Mittelpunktsregel (5.3.6) $\xrightarrow{\text{leg:exMP}}$
 $\Psi_{h,g}^h, \Psi_{h,f}^h$ + Strang-Splitting (5.11.3) $\xrightarrow{\text{leg:strang}}$

SS-EuEI = symmetrisches Einschrittverfahren \rightarrow Def. 5.10.6 $\xrightarrow{\text{def:symEul}}$

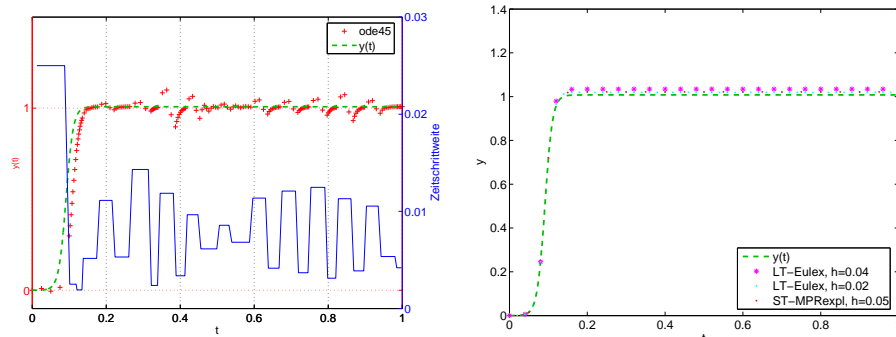
Warum das Ganze ?

$\dot{y} = f(y) + g(y)$ „schwierig“ $\xrightarrow{\text{sec:stiffheit}}$ $\dot{y} = f(y) \rightarrow$ Steifheit, aber analytisch integrierbar (z.B. steif \rightarrow Abschn. 5.7.1) $\xrightarrow{\text{sec:stiffheit}}$ $\dot{y} = g(y)$ „einfach“ (explizit) numerisch integrierbar

Beispiel 251 (Abspaltung steifer Anteile). vgl. Bsp. 233 $\xrightarrow{\text{ex:ode45stirr}}$

AWP $\dot{y} = \lambda y(1 - y) + \alpha \sin(y)$, $\lambda = 100$, $\alpha = 1$, $y(0) = 10^{-4}$.

Kleine Störung



Lösung aus ode45, siehe Bsp. 233 $\xrightarrow{\text{ex:ode45stirr}}$ inexacte Splittingverfahren: Näherungslösungen

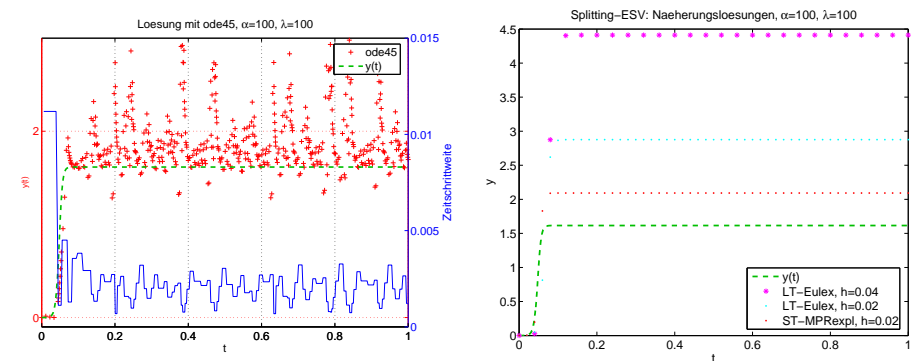
ode45: 152
 Gesamtzahl der Zeitschritte LT-Eulex, $h = 0.04$: 25
 LT-Eulex, $h = 0.02$: 50
 ST-MPRexpl, $h = 0.05$: 20

LT-Eulex: $\dot{y} = \lambda y(1 - y) \rightarrow$ exakte Evolution, $\dot{y} = \alpha \sin y \rightarrow$ expl. Euler (5.3.1) & Lie-Trotter-Splitting (5.11.2) $\xrightarrow{\text{leg:lietrotter}}$
 ST-MPRexpl: $\dot{y} = \lambda y(1 - y) \rightarrow$ exakte Evolution, $\dot{y} = \alpha \sin y \rightarrow$ expl. Mittelpunktsregel (5.3.6) & Strang-Splitting (5.11.3) $\xrightarrow{\text{leg:strang}}$

Klar: Störung $\uparrow \Rightarrow$ obige Splittingverfahren schlechter

$\lambda = 100$ und $\alpha = 100$:

5.11
 p. 661



5.12 Verfahren für oszillatorische Differentialgleichungen

[File: section-verfahren-fuer-oszillatorische-probleme.tex, SVN: section-verfahren-fuer-oszillatorische-probleme.tex 1183 2006-11-29 10:46:03Z hipmair]

5.11
 p. 662

5.11
 p. 663

5.12
 p. 664

Prototyp:

$$\ddot{y} = -\omega^2 y \quad , \quad y(0) = y_0, \dot{y}(0) = v_0$$

$$\blacktriangleright \quad y(t) = \alpha \cos(\omega t) + \beta \sin(\omega t) \quad , \quad \alpha, \beta \in \mathbb{R}$$

Verallgemeinerung (skalar):

$$\ddot{y} = -\omega^2 y + g(y) \quad , \quad y(0) = y_0, \dot{y}(0) = v_0 \quad , \quad (5.12.1)$$

mit Lipschitz-stetiger **Störung** $g : \mathbb{R} \mapsto \mathbb{R}$.

Verallgemeinerung (vektoriell)

$$\ddot{\mathbf{y}} = -\mathbf{A}\mathbf{y} + g(\mathbf{y}) \quad , \quad \mathbf{y}(0) = \mathbf{y}_0, \dot{\mathbf{y}}(0) = \mathbf{v}_0 \quad , \quad (5.12.2)$$

$\mathbf{A} \in \mathbb{R}^{d,d}$ symmetrisch positiv definit (\rightarrow Def. 3.2.8), $g : \mathbb{R}^d \mapsto \mathbb{R}^d$

Bemerkung 252.

$$\frac{d}{dt} \begin{pmatrix} y \\ v \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ -\omega^2 & 0 \end{pmatrix} \begin{pmatrix} y \\ v \end{pmatrix} + \begin{pmatrix} 0 \\ g(y) \end{pmatrix} \quad . \quad (5.12.3)$$

Lösung von (5.12.3) durch Variation der Konstanten:

$$\begin{pmatrix} y(t) \\ v(t) \end{pmatrix} = \begin{pmatrix} \cos t\omega & \omega^{-1} \sin t\omega \\ -\omega \sin t\omega & \cos t\omega \end{pmatrix} \begin{pmatrix} y_0 \\ v_0 \end{pmatrix} + \int_0^t \begin{pmatrix} \omega^{-1} \sin(t-s)\omega \\ \cos(t-s)\omega \end{pmatrix} g(y(s)) ds \quad (5.12.4)$$

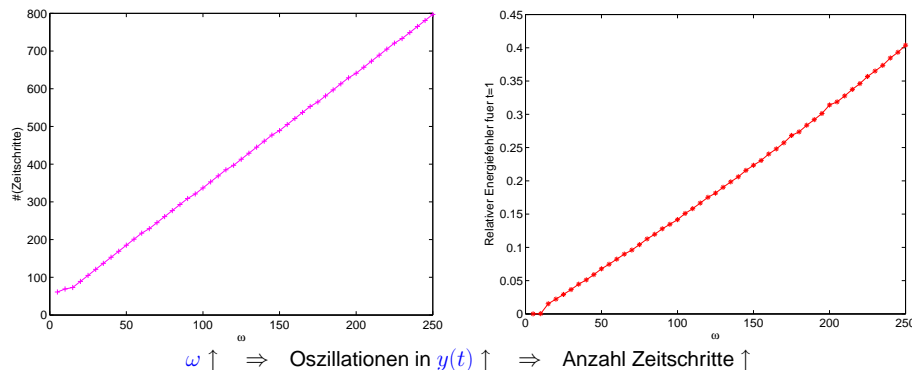
Bemerkung 253. $y(t)$ löst (5.12.1) & $G' = g \rightarrow \frac{1}{2}|\dot{y}|^2 + \frac{1}{2}\omega^2 y(t) - G(y(t)) \equiv \text{const.}$

„Energie“ für ODE (5.12.1)

Beispiel 254 (Standardintegratoren für oszillatorische Differentialgleichung).

Adaptives explizites RK-ESV (Abschn. 5.4.1) für (5.12.1):

```
y0=[1;0]; f=@(t,x) [0,1;-omega^2,0]*x + 20*[0;sin(x(1))];
options=odeset('reltol',1.0e-2,'abstol',1.0e-5);
[t,y]=ode45(f,[0,1],y0,options); („Energie“ ≙ Invariante aus Bem. 253)
```



$\omega \uparrow \Rightarrow$ Oszillationen in $y(t) \uparrow \Rightarrow$ Anzahl Zeitschritte \uparrow

Ziel: Effiziente numerische Integration von (5.12.1)/(5.12.2) auch für $\omega \gg 1$ bzw. $\lambda_{\max}(\mathbf{A}) \gg 1$



Idee: Verwende analytische Lösungsdarstellung (5.12.4) zur numerischen Integration:

$$y(t \pm h) = \cos(h\omega)y(t) \pm \frac{\sin h\omega}{\omega}\dot{y}(t) + \int_0^{\pm h} \frac{\sin(\pm h-s)\omega}{\omega} \cdot g(y(t+s)) ds \quad (5.12.5)$$

$$g \equiv \text{const.} \blacktriangleright y(t+h) - 2\cos(h\omega)y(t) + y(t-h) = h^2 \left(\frac{\sin(\frac{1}{2}h\omega)}{\frac{1}{2}h\omega} \right)^2 g \quad (5.12.6)$$

➤ **Gautschis Zweischrittverfahren** ($y_h(t+h)$ aus $y_h(t), y_h(t-h)$) für (5.12.1)

$$y_h(t+h) - 2\cos(h\omega)y_h(t) + y_h(t-h) = h^2 \left(\frac{\sin(\frac{1}{2}h\omega)}{\frac{1}{2}h\omega} \right)^2 g(y_h(t)) \quad (5.12.7)$$

Notwendig: **Startschritt** aus (5.12.4)

$$y_h(h) = \cos(h\omega)y_0 + \frac{\sin h\omega}{\omega}v_0 + \frac{1}{2}h^2 \left(\frac{\sin(\frac{1}{2}h\omega)}{\frac{1}{2}h\omega} \right)^2 g(y_0) \quad (5.12.8)$$

5.12
p. 665

5.12
p. 667

Ableitungsnaherung: Aus (5.12.5) für $g \equiv \text{const.}$:

$$\blacktriangleright y(t+h) - y(t-h) = 2h \frac{\sin h\omega}{h\omega} \dot{y}(t) \Rightarrow v_h(t) = \frac{h\omega}{\sin h\omega} \cdot \frac{y_h(t+h) - y_h(t-h)}{2h} \quad (5.12.9)$$

Bemerkung 255. Gautschi-Verfahren (5.12.7), (5.12.8) für vektorielles Problem (5.12.2)?

$$\text{Ersetze } \cos(h\omega) \mapsto \cos h\mathbf{A} \quad , \quad \left(\frac{\sin(\frac{1}{2}h\omega)}{\frac{1}{2}h\omega} \right)^2 \mapsto 4(h\mathbf{A})^{-2} \sin^2(\frac{1}{2}h\mathbf{A}) \quad .$$

Erinnerung: Funktionalkalkül für Matrizen:

$$\mathbf{S}^{-1}\mathbf{A}\mathbf{S} = \text{diag}(\lambda_1, \dots, \lambda_n) \blacktriangleright f(\mathbf{A}) = \mathbf{S} \text{diag}(f(\lambda_1), \dots, f(\lambda_n)) \mathbf{S}^{-1} \quad .$$

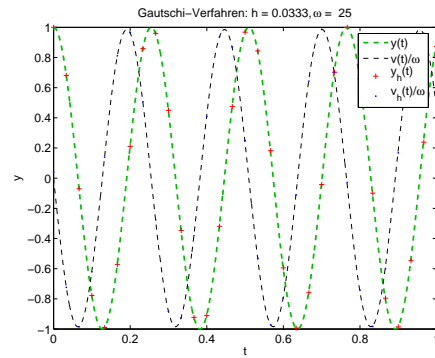
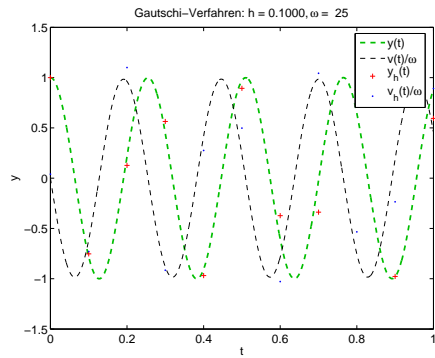
Beispiel 256 (Gautschis Zweischrittverfahren).

Anfangswertproblem vom Typ (5.12.1) auf $[0, 1]$:

$$\ddot{y} = -\omega^2 y + \sin y \quad , \quad y(0) = 1 \quad , \quad \dot{y}(0) = 0 \quad .$$

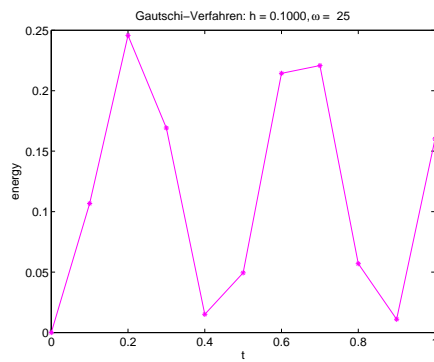
5.12
p. 666

5.12
p. 668

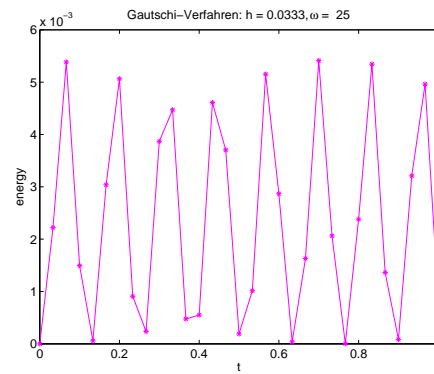


$\omega = 25$: y_h folgt (oszillatorischer Lösung), auch wenn $h \approx \frac{2\pi}{\omega}$

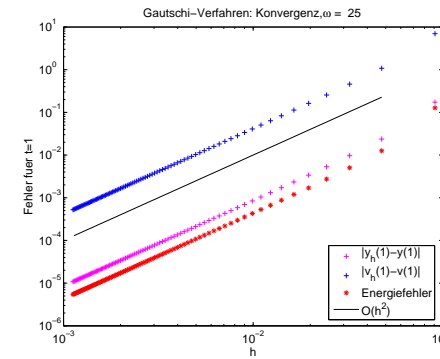
Relativer Fehler in „Energie“ (→ Bem. 253) für $t=1$: rem: Bconsosc



Zeitschritt $h = 0.1$

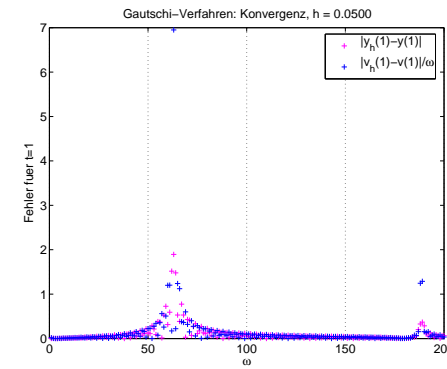


Zeitschritt $h = 0.033$



Ziel erreicht ?

Fehler für fixes h in Abhängigkeit von ω :

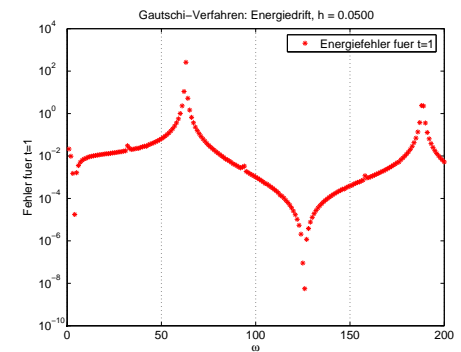


$h = 0.05$: Was passiert für $\omega \approx 61$, $\omega \approx 123$, $\omega \approx 185$?

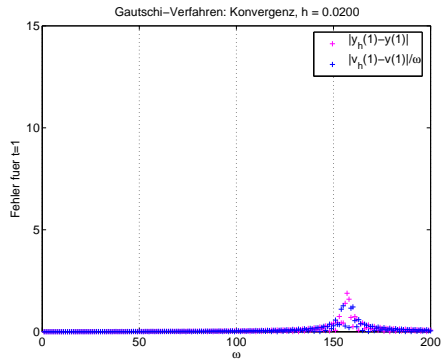
Beobachtung:

Alle Fehler $\approx O(h^2)$

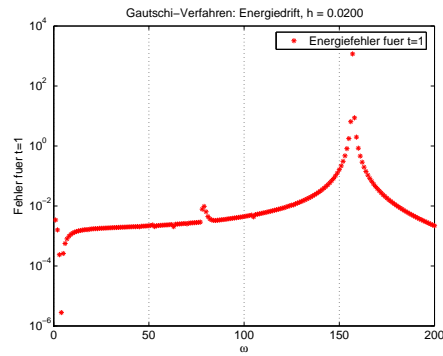
⇕
Konvergenzordnung 2



Instabilität ?



➤ h -Abhängigkeit kritischer Frequenzen

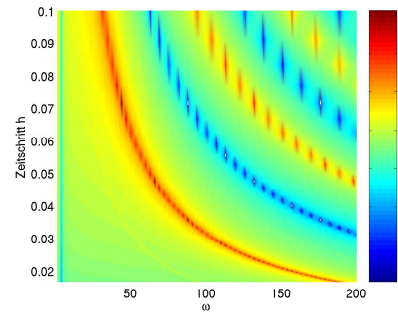


h -Abhängigkeit kritischer Frequenzen

Logarithmus \log_{10} des relativen Energiefehlers zum Endzeitpunkt $t = 0$

Beobachtung:

$h\omega$ -Abhängigkeit kritischer Frequenzen



Modellproblem:

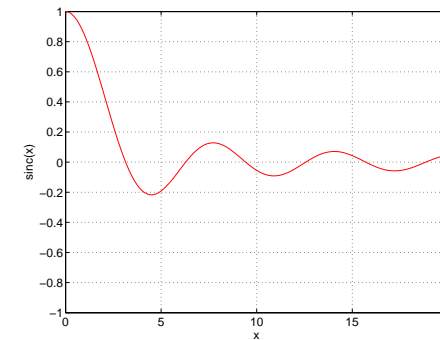
$$\ddot{y} = -\omega^2 y + \alpha y, \quad \alpha \ll \omega^2.$$

➤ Gautschi-Verfahren (5.12.7), Schrittweite h : Preitermrekursion, vgl. (5.10.13),

$$y_h(t+h) - \left\{ 2 \cos(h\omega) + h^2 \alpha \operatorname{sinc}^2\left(\frac{1}{2}h\omega\right) \right\} y_h(t) + y_h(t-h) = 0.$$

(5.12.10) [eq:osc](#)

(5.12.11) [eq:gauts](#)



Die sinc-Funktion:

$$\operatorname{sinc}(x) := \frac{\sin x}{x}$$

- Analytisch auf \mathbb{R}
- $|\operatorname{sinc}(x)| \leq 1$ mit globalem Maximum in $x = 0$

Analyse von (5.12.11): Ansatz $y_h(kh) = \xi^k \leftrightarrow$ Charakteristische (quadratische) Gleichung

$$\exists \text{ Lösungen } y_h(kh) \text{ von (5.12.11): } \lim_{k \rightarrow \infty} y_h(kh) = \pm \infty \Leftrightarrow \left| 2 \cos(h\omega) + h^2 \alpha \operatorname{sinc}^2\left(\frac{1}{2}h\omega\right) \right| > 2$$

5.12

➤ $(\cos h\omega \approx 1 \Leftrightarrow h\omega \approx 2\pi l, l \in \mathbb{Z}) \Rightarrow y_h(kh) \rightarrow \pm \infty$ auch für $h \ll 1$.

p. 673

Abhilfe [30]: „Filterung“: Dämpfung von α , falls $h\omega \approx 2\pi l$:

In (5.12.7), (5.12.8) ersetze: $g(y_h(t)) \mapsto g(\psi(h\omega)y_h(t)), \quad \psi(\xi) := \operatorname{sinc}^2 \xi (1 + \frac{1}{2}(1 - \cos \xi))$

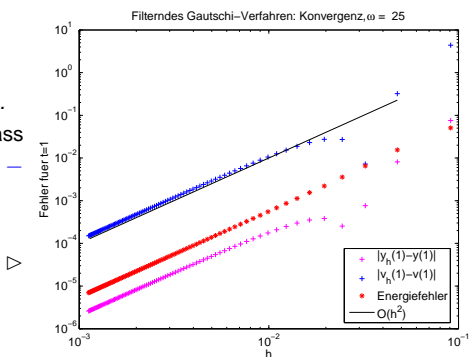
➤ Modifiziertes Gautschi-Verfahren:

$$y_h(t+h) - 2 \cos(h\omega)y_h(t) + y_h(t-h) = h^2 \left(\frac{\sin(\frac{1}{2}h\omega)}{\frac{1}{2}h\omega} \right)^2 g(\psi(h\omega)y_h(t)). \quad (5.12.12) \quad \text{eq:gaufilt}$$

Beispiel 257 (Modifiziertes Gautschi-Verfahren).

AWP aus Bsp. 5.12.7, Integration gemäss (5.12.12), Filterfunktion $\psi(\xi) := \operatorname{sinc}^2 \xi (1 + \frac{1}{2}(1 - \cos \xi))$

Konvergenzordnung 2



5.12

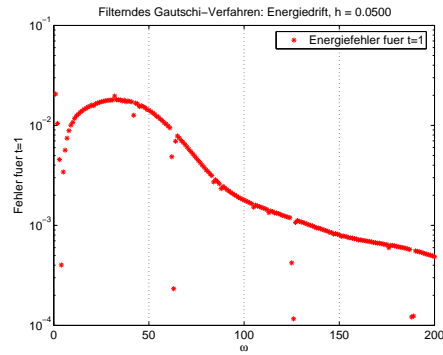
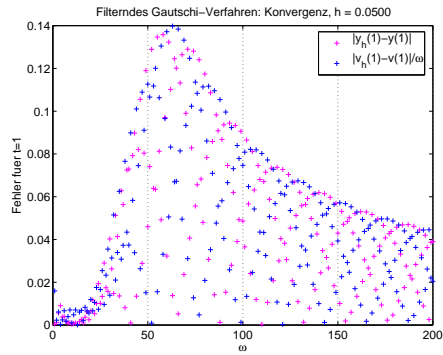
p. 674

5.12

p. 675

5.12

p. 676



Verzeichnisse

5.12
p. 677

5.12
p. 679

Index

T-Reversibilität, 629
 Ähnlichkeitstransformationen, 204
 unitäre, 209
 Überlauf, 26, 32
 Givens-Rotation, 181
 3-Term-Rekursion, 457
 für Tschebyscheff-Polynome, 377
 5-Punkt-Stern-Operator, 329
 A-Skalarprodukt, 281
 Abbruchkriterium, 62
 A posteriori, 63
 A priori, 63
 Ausnutzung von Gleitpunktarithmetik, 28
 CG-Verfahren, 291
 Maschinenunabhängig, 28, 63
 Maschinenunabhängig Quadratur des Kreises, 42
 Newton, 101
 Residuenbasiert, 62
 vorkonditioniertes CG-Verfahrens, 302
 Ableitung bei nichtlinearem Ausgleichsproblem, 117
 Adaptive Quadratur, 526
 a posteriori Fehlerabsch., 526
 Gitterverfeinerung, 526
 Mehrgitter-, 527
 Affin-Invarianz, 99
 Affin-Kovarianz
 Runge-Kutta, 568
 AGM, 61
 Aitken-Neville-Schema, 362
 Algorithmus
 Levinson, 349
 Remes, 420
 Alternanten, 419
 Tschebyscheffscher Satz, 419
 Analytizitätsgebiet, 400
 Funktionen, 400
 AnfangswertProblem, 546
 autonom, 548
 Anfangswertproblem
 Differentiell-Algebraische, 607
 Steifes, 594
 Approximation

5.12
p. 678

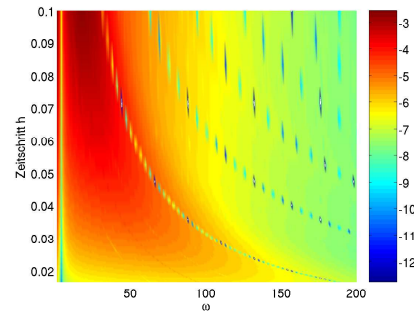
5.12
p. 680

$h\omega$ -Abhängigkeit der Energiedrift

Beobachtung:

`leg:gaufilt`
 (5.12.12): • Keine Instabilität

• Im Vergleich zu `leg:Gautschi`
 Bsp. 256) deutlich reduzierte
 Energiedrift (Skala!)



Niedrigrang, 445
 polynomiale LSG, 411
 Arnoldi-Prozess, 235
 Asymptotisches Fehlverhalten, 369
 Aufrunden, 28
 Ausgleichsproblem
 Lineares, 262
 Mit linearen Nebenbedingungen, 276
 Totales, 275
 Auslöschung, 37
 bei Householdertransformation, 180
 Definition, 38
 Exponentialfunktion, 40
 Quadratur des Kreises, 41
 Vorwärts- & Rückwärtssubstitution, 154
 Auslöschungsfrei
 Quadratur des Kreises, 42
 Axiom der Gleitpunktarithmetik, 30
 Auslöschung, 39
 Numerische Differentiation, 46
 backward error analysis, 48
 Bandbreite, 169
 obere, 169
 reduktion, 172
 untere, 169
 Zeilen-, 170
 Baryzentrische Interpolationsformel, 361
 Basis
 Hierarchisch, 541, 542
 Kosinus, 334
 Orthonormal, 184, 205, 412
 Charakteristische Gleichung, 635
 chebyscheff-Polynome
 3-Term-Rekursion, 377
 Cholesky, 165
 Cholesky-Zerlegung
 Rechenaufwand, 165
 Unvollständige, 232
 Cluster
 -ing Approximation, 425
 Clusteringapproximation
 Aufgabenstellung, 425
 Compressed-Row-Storage, 128
 Computerarithmetik, 13
 Cosinustransformation, 334
 Crout
 Algorithmus von, 149
 CRS, 128
 Cuthill & McKee, 175
 Algorithmus, 175
 Dämpfungsfaktor, 104
 Defintheit, 56
 Dezimalzahlen
 m-stellig, 23
 DFT, 312, 529
 Reelle, 320
 Diagonalisierung lokaler translationsinvarianter linearer Operatoren, 329
 Differentialgleichung
 Lineare konst. Koeffizienten, 552
 Logistische, 560
 Lotka-Volterra , 552

Orthonormal (Krylov), 235
 Sinus, 325
 Trigonometrische, 310
 Bauelementgleichungen, 135
 Bedingung
 Galerkin, 281, 283
 Bernstein-Polynome, 499
 Besetzungsmuster, 176
 Bestapproximation durch Niedrigrangmatrizen, 256
 Bewegungsgleichungen
 Molekulardynamik, 631
 Newtonsche, 547
 Newtonsche Hamiltonsche Form, 547
 Bilddatenkompression, 259
 Bisektionsverfahren, 76
 Black-Box-Methoden, 229
 BLAS-Bibliothek, 126
 Blow-up, 549, 550
 Bootstrapping, 563
 Bounding Box, 440
 Broyden
 Quasi-Newton-Verfahren, 108
 Broyden-Verfahren
 Konvergenzmonitor, 111, 112
 Broydens Rang-1-Modifikation, 108
 Butcher-Schema, 564
 CG
 Konvergenz und Spektrum, 298
 Konvergenzgeschwindigkeit, 293
 Verfahren Abbruchkriterium, 291
 CG-Verfahren, 280
 Ricatti, 551
 stationäre Punkte, Evolution, 588
 Variation der Konstanten, 554
 Differentiell-Algebraische Anfangswertprobleme, 607
 Differenzenquotient, 45
 Dilatation, 532
 Direkte Potenziteration im Unterraum, 224
 Direkte Potenzmethode, 215
 Diskrete
 Faltung, 336
 Diskrete Fouriertransformation, 310
 Diskretisierungsfehler, 572
 divide and conquer, 318
 Double Precision, 25
 Dreiecksform
 obere, 143
 Dreiecksungleichung, 56
 Dreitermrekursion, 635
 skalar, 461
 Effizienz, 92
 Eigenvektor
 Links-, 207
 Eigenwert
 -problem, 203
 -problem Verallgemeinertes, 213
 Arnoldi-Verfahren, 239
 Eingabedaten, 34
 Einpunktverfahren, 78, 80
 Einschnittverfahren
 symmetrisch, 627
 Elektrisches Netzwerk, 135

Energie
 für oszillatorische Differentialgleichung, 642
 kinetische, 617
 potentielle, 617
 Energiedrift, 621, 655
 Energieerhaltung, 617
 Energienorm, 228
 Entwickeltung
 asymptotisch, 364
 Ergebnismenge, 34
 ESV
 Radau, Ordnung 3, 602
 Radau, Ordnung 5, 602
 Euler
 Implizit, 602
 Eulersche Methode, 86
 Eulerverfahren
 explizit, 559
 explizites, Stabilitätsfunktion, 590, 591
 implizit, 560
 symplektisch, 622
 Evolution
 diskrete, 559, 562, 612
 Operator, 550
 Exponentialfunktion
 Auslöschung, 40
 Exponentialreihe, 18
 Exponentielle Konvergenz , 381
 Extrapolation, 364
 5.12
 p. 681
 Faltung
 diskrete, 336
 shandles, 99
 Skalierungs-, 532, 541
 fzero, 77
 Galerkin-Bedingung, 281, 283
 Gauss
 Kollokation, 612
 Kollokations-ESV , 613
 Punkte, 514
 Quadratur, 510
 Gauss-Radau-Quadratur, 601
 Gausselimination, 142
 Gautschi-Verfahren, 644
 Filterung, 653
 modifiziertes, 653
 Gautschis Zweischrittverfahren, 644
 Gedämpftes Newton-Verfahren, 104
 Genauigkeit
 Doppelte, 25
 Einfache, 25
 Gerschgorin-Kreise, 231
 Ghost Eigenvalues, 248
 Gibbsches Phänomen, 394
 Gitterfunktion, 330
 Gitterverfeinerung
 Lokal, adaptive Quadratur, 526
 Givens-Rotation, 180, 187, 191, 199, 202
 Gleichung
 Hamilton, 617
 nichtlineare, 53
 Gleichungssystem
 dünnbesetzt, 167

Fehler
 Relativer, 18
 Fehlergleichung, 280
 Fehlerschätzung
 linear konvergente Iteration, 64
 zeitlokal, 579
 Fehlverhalten
 Asymptotisch, 369
 Fernfeld, 435
 FFT, 315
 Fill-in, 159, 168
 Fitten
 Polynomiales, 264
 Fixpunkt, 66
 abstossend, 73
 anziehend, 73
 Fixpunktiteration, 66
 Konsistent, 66
 lokal kubische Konvergenz, 85
 Form
 erhaltende Splineinterpolation, 490
 erhaltung, 472
 Fourier
 Koeffizienten, 396
 Matrix, 311
 Reihe, 396
 Fourier-Matrix Zerlegung, 317
 Fouriertransformation
 Diskrete, 310
 Funktion
 Analytizitätsgebiet, 400
 Lineares, 135
 Gleitkommazahl, 24
 Gleitpunktarithmetik, 27
 Gradient, 118
 Gram-Schmidt
 Orthogonalisierung (CG), 285
 Orthonormalisierung, 236, 413
 Hülle, 170
 Halleysche
 Iteration, 80
 Methode, 86
 Hamilton-Funktion, 629
 Molekulardynamik, 631
 Hamiltonsche Gleichung, 617
 Handle MATLAB-Funktion bei Bisektionsverfahren, 77
 Harmonischer Oszillator, 634
 Harmonisches Mittel, 491
 Hauptkomponentenanalyse, 257
 Hermite
 Lagrange, 354
 Hermite-Interpolation
 kubisch, 357
 Hesse-Matrix
 bei Nichtlinearem Ausgleichsproblem, 116
 Definition, 118
 Homogenität, 56
 Horner Schema, 353
 Householder-Reflektion, 180
 IEEE Standard 754, 24
 Impulsordinate, 617

5.12
 p. 683

5.12
 p. 684

In situ, 149	konvergent, 54	Lagrange-Multiplikator, 277, 608	Mantisse, 24
inf, 26	Konvergenz, 54	Lanczos-Prozess, 241	Maschinengenauigkeit, 29, 30
Infinity, 26	Iterierte $x^{(k)}$, 54	Landau-Notation, 39, 125	Maschinenzahl, 24
Inkremente	Jacobi-Vorkonditionierer, 230	Langzeitverhalten, 620	Basis, 24
Kollokation, 559	Jacobimatrix, 71	Lebesgue	betragsgrösste, 24
Runge-Kutta, 564	Kardinalspline, 486	Konstante, 380	betragskleinste, 24
Instabilität, 634, 635	Keplerproblem, 578	Legendre-Polynome, 457	Exponentenbereich, 24
Gautschi-Verfahren, 649	Kernfunktion, 425, 426	Leitwertmatrix, 164	Mantisse, 24
Interpolation	separiert, 428	Lemma	Matlab, 9
Baryzentrische Formel, 361	Kettenregel, 117	Schursches, 205	Matrix
Inverse, 90	Kirchhoffsche Regel, 135	Störungs, 138	Dünnbesetzt, 127, 158
kubische Hermite-, 477	Knoten, 354	Levinson-Algorithmus, 349	Energienorm, 228
Lagrange, 354	analyse, 135	Limit, 480	Fourier, 311
Spline formerhaltend, 490	Doppelte, 354	lineare Korrelationen, 257	Fourier (Zerlegung), 317
Spline kubisch, 483	Einfache, 354	lineare Operatoren	Hermesch, 163, 205
Spline kubisch, Lokalität, 486	potentiale, 135	Diagonalisierung, 329	Hesse, 116, 118
Spline natürlich kubisch, 484	Tschebyscheff, 379, 380	Lineares Ausgleichsproblem, 262	Kollokations-, 425
Spline periodisch kubisch, 484	Knotenanalyse, 164	lineares elektrisches Netzwerk, 135	Norm, 137
Spline vollständig kubisch, 484	Kollokation, 557	Lineares Gleichungssystem, 135	normal, 205
Trigonometrische, 405	Inkremente, 559	Linearisierung	normalisierte untere Dreiecks, 147, 155
Tschebyscheff, 405	Verfahren, Konvergenz, 560	Axiom der Gleitpunktarithmetik, 39	obere Dreiecks, 147, 155, 197, 201
Inverse Iteration, 219	Kollokations	Lokale bei Newton-Iteration, 96	obere Hessenberg, 192, 236
vorkonditionierte, 225	Einschrittverfahren, 577	Linkseigenvektor, 207	Orthogonal, 178
Iteration	RK-ESV, 601	Lipschitz	Permutations, 201
-sfunktion, 66	Kollokationsmatrix, 425, 426	Stetigkeit, lokale, 548	positiv definit, 163
-sfunktion Startwert, 66	Komplexität, 125	Logistische Differentialgleichung, 637	positiv semidefinit, 163
-sverfahren, 54	Kondensator, 135	Lokalität, 470	Propagations-, 555
Inverse, 219	Kondition, 33	Lotka-Volterra Differentialgleichung, 552	Rang, 137
Inverse mit Shift, 219	analyse, differentielle, 554	LU-Zerlegung	regulär, 137
vorkonditionierte inverse, 227	sanalyse differentielle, 35	In situ, 149	Schiefhermitesches, 205
Iterationsverfahren	lokal (Newton), 102	5.12 In situ, 149	Sinus, 325
Spektrale, 229	lokal kubisch bei Fixpunktiteration, 85	p. 685 Rechenaufwand, 150	minimum degree reordering, 176
Konditionszahl	lokal quadratisch, 80	symmetrisch, 163	Mittelpunktsregel
absolut, 34, 207	lokal quadratisch bei gedämpften Newtonverfahren, 121	transponierte, 125	explizit, 564, 565
relative, 35, 137	lokal quadratisch(Newton), 100	Tridiagonal, 169, 242	implizit, 559
Konsistent	Ordnung, 60	Unitär, 178, 205	implizit, Stabilitätsfunktion, 591
Fixpunktiteration, 66	quadratisch notwendige Bedingung , 85	Vollbesetzt, 127	Modellfunktionsverfahren, 78
Konsistenz	quadratisch, Fixpunktiteration, 71	Vorkonditionierer, 229	Modellfunktion, 78
Iterationsverfahren, 55	quadratisch, Wurzeliteration, 61	Wilkinson, 155	Modellproblem
Ordnung, 577	Rate, 57	Wronski, 555	für gestörte oszillatorische Differentialgleichungen, 651
Konsistenzordnung, 636	Rate Sekantenverfahren, 89	Matrixfaktorisierung, 146	Molekulardynamik, 631
Konstante	Konvergenzanalyse	Matrixgraph, 174	Multiplikation
Lebesgue, 380	Fixpunktiteration, 69	Kantenmenge, 174	Polynom, 337
konvergent, 54	Konvergenzgeschwindigkeit	Knotengrad, 174	Multiskalenbasen, 529
Konvergenz, 54	CG, 293	Knotenmenge, 174	Nahfeld, 435
p-te Ordnung, 60	Konvergenzmonitor	planar, 174	NaN, 26
Algebraische, 369	bei Broyden-Verfahren, 111, 112	Matrixmultiplikation	Newton
Algebraische bei trig. Interp., 393	Konvergenzordnung	Blockweise, 127	Abbruchkriterium, 101
Algebraische polynom. Fitten, 417	explizit, Euler, 576	Matrixnorm, 71	Affin-Invarianz, 99
Algebraische, Quadratur, 519	explizit, Mittelpunktsregel, 576	Spaltensummen, 72, 134	Dämpfung, 104
Asymptotische, 72	explizit, Trapezregel, 576	Submultiplikativität, 134	Dämpfungsfaktor, 104
des CG-Verfahrens, 298	implizit, Euler, 576	Zeilensummen, 72, 134	Iteration, 96
Exponentielle, 369, 371	implizit, Trapezregel, 576	Matrizen	Korrektur, 96
Exponentielle , 381	Klassisches RKV, 576	Töplitz, 344	lokal quadratische Konvergenz, 100
Exponentielle bei trig. Interp., 393, 402	Kuttas 3/8-Regel, 576	Zirkulante, 336	natürlicher Monotonietest, 104
Exponentielle polynom. Fitten, 416	Kosinus	Mehrpunktverfahren, 78, 87	Numerisches Differenzieren, 99
Exponentielle, Quadratur, 518	basis, 334	Methode	vereinfachte Korrektur, 101
global, 55	transformation, 334	Eulersche, 86	vereinfachtes Verfahren, 99
Kollokationsverfahren, 560	Kosinustransformation, 334	Halleysche, 86	Newton-Verfahren
kubisch notwendige Bedingung , 85	Krylov-Unterraumverfahren, 233, 279	Quadratische inverse Interpolation, 86	1D, 78
Linear, 57	Kubische Hermite-Interpolation, 357	Metrik, 34, 35	Newtoniteration
Linear bei Gauss-Newton, 122		Milleniums-Algorithmus, 315, 456	n-dimensional, 96
		p. 686 Milne-Regel, 508	

Nichtlinear	Rundungsfehlereinfluss, 45	Gewichte, 506	LU-Zerlegung, 150
Ausgleichsproblem, 115	Numerische Rangbestimmung, 254	Knoten, 506	QR-Zerlegung, 186
Ausgleichsrechnung, 115	Numerischer Algorithmus, 47	Numerische, 506	SVD, 254
Regression, 120	Numerischer Rang, 268	Zusammengesetzte -formeln, 519	Regel
Niedrigrang-Approximation, 445	Obere Dreiecksform, 143	Quasi-Newton-Verfahren, 107, 108	Kuttas 3/8, 568
Norm, 56, 133	Obere Hessenbergmatrix, 192, 236	QZ-Algorithmus, 213	Milne, 508
L^1 , 357	Ohmscher Widerstand, 135	Rücksitution, 143	Mittelpunkt, explizit, 564, 565
L^2 , 357	ONB, 412	Rückwärtsfehleranalyse, 48	Mittelpunkt, implizit, 559
Äquivalenz, 57	Operationen	Rückwärtsrekursion, 464	Mittelpunkt, implizit, Stabilitätsfunktion, 591
1-, 56, 133	Elementare arithmetische, 30	Radau-ESV, Ordnung 3, 602	Simpson, 508
Energie-, 281	LEVEL I, 126	Radau-ESV, Ordnung 5, 602	Trapez, 508
Euklidisch, 56, 133	LEVEL II, 126	Radau-RK-ESV, 612	Trapez äquidistante, 523
Matrix, 71, 134	LEVEL III, 126	Radau-Verfahren, 610	Trapez, explizit, 564, 566, 573
Maximum, 56, 133	Ordnung	Rang	Trapez, explizit, Stabilitätsfunktion, 591
Semi, 374	Gauss-Quadratur, 510	einer Matrix, 137	Trapez, implizit, 575
Sobolev-Semi, 374	Ordnungsschranken, 577	Numerischer, 268	Weddle, 509
Supremum, 357	Orthogonalisierung	Spaltenrang, 137	Regression
Normalengleichung, 272	Gram-Schmidt (CG), 285	Zeilenrang, 137	nichtlinear, 120
Erweiterte, 277	Orthogonalitätsverlust bei Lanczos, 245	Rang-1-Modifikation	Reihe
Normalengleichungen, 272	Orthogonalpolynome, 457, 513	Broyden, 108	Fourier, 396
Normierungsbedingung	Orthonormalbasis, 184, 205	Rang-1-Modifikationen, 189	Harmonische, 21
Besselfunktionen, 465	Orthonormalisierung	Rangbestimmung	Rekursion
Not a Number, 26	Gram-Schmidt, 413	numerisch, 254	3-Term, 377, 457
Nullabfrage	Ortskoordinate, 617	Raum	Fläche reguläres n -Eck, 15
Numerische, 31	oszillatorische Differentialgleichungen, 641	Phasen-, 546	Relation
Nullstellen, 49	PCG	Phasen-, erweiterter, 546	Verfeinerungs-, 535
Bestimmen von, 49	vorkonditioniert, 301	Rayleigh-Quotienten-Iteration, 222	Remes-Algorithmus, 420
Nullstellenbestimmung	Pendelgleichung, 618	Rechenaufwand, 92	Residuenminimierende Verfahren, 304
Modellfunktionsverfahren, 78	Permutationsmatrix, 197	$\text{eig}()$, 211	Residuum, 226, 230
Nullstellenformel, 50	Pfeilmatrix, 158	Cholesky-Zerlegung, 165	Resonanzen elektrischer Netzwerke, 203
Numerische Differentiation	power method, 215	Gausselimination, 146	Reversibilität
Phänomen	Problem	mechanisches System, 629	Komplement, 152
Gibbsches, 394	Anfangswert, 546	Ricatti Differentialgleichung, 551	Schursches Lemma, 205
Phasenraum, 546	Anfangswert, autonom, 548	Ritz-Projektion, 224, 238	Sekantenbedingung, 107
erweiterter, 546	Eigenwert (quadratisch), 203	RK4, 567	Sekantenverfahren, 87, 107
Molekulardynamik, 631	Gut konditioniertes, 34	Verfahren, Stabilitätsfunktion, 591	Seminorm, 374
Pivot	Mathematisch, 33	Rundung, 28	Sherman-Morrison-Woodbury-Formel, 189
-strategie, 152	Nichtlineares Ausgleichs-, 115	Rundungsfehler	Shift (Inverse Iteration), 219
-suche, 152, 155	Sattelpunkt, 277	Kapitel über, 27	Simpson-Regel, 508
-wahl, 154	steifes Anfangswert, 594	Relativer, 29	sinc-Funktion, 652
Element, 143, 144	Produktregel, 117	Unvermeidlich, 28	Single Precision, 25
Zeile, 143, 144	Projektion	Runge-Kutta	Singulärwertzerlegung, 250
Pivotstrategie, 152	Ritz, 224, 234, 238	3/8-Regel, 568	sparsame, 253
Pivotsuche, 152	Propagationsmatrix, 555	Affin-Kovarianz, 568	Sinus
Spalten-, 154	Prozess	Autonomisierung, 569	basis, 325
Polynom, 352	Arnoldi, 235	Autonomisierungsinvarianz, 569	matrix, 325
-iale LSQ-Approximation, 411	Lanczos, 241	Eingebettet, 580	transformation, 325
Bernstein, 499	Punkt	Inkrement, 564	Sinustransformation, 324
charakteristisches, 204	stationär, 552, 598	klassisch, 567	Skalarprodukt
Lagrange, 355	QR-Algorithmus (mit Shift), 209	Runge-Kutta-Verfahren, 563	A-, 281
monomiale Darstellung, 352	Quadratisch inverse Interpolation, 86	Sattelpunktproblem, 277	Euklidisches, 126
Multiplikation, 337	Mehrpunktverfahren, 91	Satz	Skalierungsfunktion, 532
Orthogonal, 513	Quadratur	Peano & Picard-Lindelöf, 549	Spaltenpivotsuche, 154
Probleme bei Auswertung, 21	Adaptive, 526	SAXPY, 126, 289	Sparsame Singulärwertzerlegung, 253
trigonometrisch, 386	Adaptive, a posteriori Fehlerabsch., 526	Schema	Speicherung
Polynomiales Fitten, 264	Adaptive, Gitterverfeinerung, 526	Aitken-Neville, 362	Hüllenorientiert, 171
Polynominterpolation	Adaptive, Mehrgitter-, 527	Horn, 353	Spline
allgemein, 354	formel, 506	Schnitt von Geraden, 34	funktionen, 482
Potenziteration	Gauss, 510	Schrittweiten	Interpolation formerhaltend, 490
im Unterraum, 224	Gauss Ordnung, 510	beschränkung, 593	Interpolation kubisch, 483
Potenzmethode	Gauss-Radau, 601	steuerung, ESV, 577	Interpolation kubisch, Lokalität, 486
direkte, 215		Schur	Interpolation natürlich kubisch, 484

5.12
p. 689

5.12
p. 691

5.12
p. 690

5.12
p. 692

Interpolation periodisch kubisch, 484
 Interpolation vollständig kubisch, 484
 Kardinal, 486
 Splitting
 Lie-Trotter, 636
 Strang, 636
 Splittingverfahren, 635
 inexakte, 638
 Spule, 135
 Störmer-Verlet-Verfahren, 630
 Molekulardynamik, 632
 Störungsanfälligkeit von Eigenwerten, 207
 Stabil
 Algorithmus, 47
 numerisch, 47
 Stabilität, 47
 -sfunktion, 602
 -sfunktion von Runge-Kutta-Verfahren, 590
 -sgebiet, 598
 Gauss-Kollokations-ESV, 600
 Startschritt, 644
 Startwert $x^{(0)}$, 54
 stationäre Punkte
 Evolution Dgl., 588
 steif
 Differentialgleichung, 639
 Stetigkeit
 lokal Lipschitz, 548
 Strahlungstransport, 338
 Struktursymmetrie, 169
 Substitution
 äquidistante, 523
 Trapezregel, 508
 explizit, 564, 566, 573
 explizit, Stabilitätsfunktion, 591
 Implizit, 575
 Tridiagonalmatrix, 169
 Trigonometrische Basis, 310
 Trigonometrische Interpolation, 405
 trigonometrisches Polynom, 386
 Tschebyscheff-Interpolation, 405
 Tschebyscheff-Knoten, 379, 380
 Tschebyscheff-Polynome, 457
 Tschebyscheffscher Alternantensatz, 419

 Unitäre Ähnlichkeitstransformationen, 209
 Unnormalisierte Zahlen, 26
 Unterlauf, 26, 32
 Unterraumkorrektur, 283

 Variationsgleichung, 555
 Variation der Konstanten, 554, 642
 Variationsgleichung, 618
 Verallgemeinertes Eigenwertproblem, 213
 Verfahren
 Arnoldi, 239
 Bisektions-, 76
 Einschritt, implizit, 562
 Einschritt, Schrittweitensteuerung, 577
 Euler, explizit, 559
 Euler, explizit, Stabilitätsfunktion, 590
 Euler, implizit, 560
 Euler, implizit, Stabilitätsfunktion, 591

 Rück, 150
 Rück (bei Cholesky), 165
 Vorwärts, 150
 Vorwärts (bei Cholesky), 165
 Supremumsnorm, 357
 SVD, 267
 Ausgleichsrechnung, 278
 Symmetrischer Gauss-Seidel Vorkonditionierer, 230
 symplektisch, 617
 symplektisches Integrationsverfahren, 619
 Symplektizität, 616

 Töplitz-Matrizen, 344
 Taylorformel, 70
 mehrdimensional, 71
 Teile-und-herrsche, 318
 Tensorprodukt
 Interpolationspolynom, 430
 Tschebyscheff-Polynominterpolation, 431, 445
 tic,toc, 131
 Toleranz
 Absolut, 579
 Relativ, 579
 Totales Ausgleichsproblem, 275
 Transformation
 Ähnlichkeits, 204
 Kosinus, 334
 kovariant, 552
 schnelle Fourier, 315
 Sinus, 325
 Translation, 532
 Trapez-Regel
 Gauss-Newton, 118
 gedämpftes Newton-, 104
 Iterationen-, 54
 Kollokation, Konvergenz, 560
 Kollokations-Einschritt, 577
 Krylov-Unterraum, 233
 kurze Rekursionen, 306
 Mehrpunkt, 87
 Newton, 95
 Quasi-Newton-, 107
 Radau, 610
 Residuenminimierend, 304
 RK4, Stabilitätsfunktion, 591
 Runge-Kutta, 563
 Runge-Kutta, Eingebettet, 580
 Runge-Kutta, klassisch, 567
 Sekanten, 87
 Störmer-Verlet, 630
 Trust-Region, 119
 vorkonditioniertes CG, 301
 Verfeinerungsrelationen, 535
 Vorkonditionierer, 229
 Gauss-Seidel, 230
 Jacobi, 230
 Vorkonditionierte inverse Iteration, 225, 227
 Vorkonditioniertes CG-Verfahren (PCG), 301
 Vorkonditionierung, 299, 300

 Wavelet
 Koeffizienten, 534
 Zerlegung, 534
 Weddle-Regel, 509

Wurzeliteration, 28
 Zahldarstellung, 23
 Zeilen
 permutation, 155, 201
 umformungen, 142
 Zerlegung
 Cholesky, 165, 194
 der Eins, 503
 Fouriermatrix, 317
 LU, 150
 QR, 178, 182, 197, 201, 266
 QR (Aufwand), 186
 QR (Rücksubstitution), 186
 unvollständige Cholesky-, 232
 Zirkulante Matrizen, 336
 Zulässigkeitsbedingung, 439
 Zweidimensionale diskrete Fouriertransformation, 320
 Zweischrittverfahren
 Gautschis, 644

5.12
 p. 693

Beispiele und Bemerkungen

L^2 -Instabilität der hierarchischen Basis, 562
 m -stellige Dezimalzahlen, 23
 Über- und Unterlauf, 33
 $B = B^H$ s.p.d. mit Cholesky-Zerlegung, 223
 fft
 Effizienz, 325

 Abbruch des vorkonditionierten CG-Verfahrens, 312
 Abbruchkriterium
 Potenziteration, 229
 Abbruchkriterium für nicht-vorkonditioniertes CG-Verfahren, 301
 Abspaltung steifer Anteile, 662
 Adaptive explizite Einschrittverfahren für steifes Problem, 614
 Affin-Invarianz, 102
 Akkumulation von unitären Transformationen, 194
 Algorithmus
 Golub-Welsch, 534
 Remes Konvergenz, 438
 Aliasing, 413

5.12
 p. 694

Analytizitätsgebiete von Funktionen, 416
 Anfangswert
 Problem, autonom, 568
 Approximation
 Bernstein, 516
 Best- vs. Tschebyscheff-Interp., 439
 Fitpolynome, 431
 Konvergenz Clustering- mit Kollokationsmatrix, 468
 Arnoldi-Verfahren zur Eigenwertberechnung, 249
 Auflösung Singularität durch Transformation, 540
 Aufspüren fast singulärer Matrizen, 198
 Auslöschung 2-stellige Dezimalarithmetik, 38
 Auswertung der Exponentialfunktion, 18, 40
 AWP
 linear, Kondition, konst. Koeff., 574

 Bandbreite, 177
 Basis
 L^2 -Instabilität, hierarchischen, 562
 Hierarchische Zerlegung, 560
 Berechnung der Eulerschen Zahl, 14

5.12
 p. 695

5.12
 p. 696

<p>Bernstein-Approximation, 516</p> <p>Besetzungsmuster der LU-Faktoren, 155</p> <p>Bilddatenkompression, 269</p> <p>Bisektion in MATLAB: fzero, 80</p> <p>BLAS-Bibliothek, 130</p> <p>Blockweise Matrixmultiplikation:, 131</p> <p>Broyden-Verfahren für grosses nichtlineares Gleichungssystem, 117</p> <p>Broydens Quasi-Newton-Verfahren: Konvergenz, 113</p> <p>CG</p> <ul style="list-style-type: none"> direkter Löser, 292 <p>CG-Konvergenz und Spektrum, 308</p> <p>CG-Verfahren und quadratische Minimierung, 293</p> <p>Charakteristische Grössen der IEEE Gleitpunktarithmetik, 26</p> <p>Clusterbaum, 458</p> <p>CRS-Matrixspeicherformat, 133</p> <p>DFT</p> <ul style="list-style-type: none"> Frequenzanalyse, 323 Nichtlokalität, 547 <p>Differentialgleichungen</p> <ul style="list-style-type: none"> Skalare, 568 <p>Differentiation, 45</p> <p>Differenzieren Wiederholung, 121</p> <p>Direkte Potenziteration im Unterraum, 235</p> <p>Effiziente Initialisierung von Sparse-Matrizen, 134</p> <p>Effizienz FFT-basierter Gleichungslöser, 345</p> <p>Effizienz von Iterationsverfahren, 97</p> <p>Effizienz von fft, 325</p> <p>Eigenschaften von Leitwertmatrizen, 169</p> <p>Globale separierbare Approximation bei glatter Kernfunktion, 447</p> <p>Globale separierbare Approximation bei nichtglatter Kernfunktion, 448</p> <p>GMRES und Arnoldi-Prozess, 317</p> <p>Golub-Welsch Algorithmus, 534</p> <p>Gradientenfluss, 634</p> <p>Gravitationskräfte in Galaxien, 442</p> <p>Hülle einer Matrix, 179</p> <p>Hüllenorientierte Speicherung, 180</p> <p>Halleysche Iteration, 83</p> <p>Harmonischer Oszillator, 656</p> <p>Hauptkomponentenanalyse, 267</p> <p>Hermite-Interpolation, 496</p> <ul style="list-style-type: none"> Theorem, 386 <p>Hierarchische Basis-Zerlegung, 560</p> <p>IEEE Standard</p> <ul style="list-style-type: none"> Sonderfälle, 26 <p>IEEE Standard 754 für Maschinenzahlen, 24</p> <p>Inexakte Splittingverfahren, 660</p> <p>Interaktionsberechnung Vielteilchensysteme, 441</p> <p>Interne Binärdarstellung einer Gleitkommazahl doppelter Genauigkeit (MATLAB):, 25</p> <p>Interpolation</p> <ul style="list-style-type: none"> Hermite Theorem, 386 kubische Hermite-, 496 kubische Spline- Lokalität, 502 kubische Spline- vollständig, 502 Polynom Versagen, 484 Spline quadratisch formerhaltend, 513 	<p>Eigenwert</p> <ul style="list-style-type: none"> Arnoldi-Verfahren, 249 erw. Krylov-Unterraum-Verfahren, 259 Lanczos/Arnoldi, 257 <p>Eingabefehler und Rundungsfehler, 28</p> <p>Eingebettete Runge-Kutta-Verfahren, 599</p> <p>Endlichkeit von \mathbb{M}, 24</p> <p>Erweiterung Krylov-Unterraum-Verfahren, 259</p> <p>ESV</p> <ul style="list-style-type: none"> implizit, steifes Problem, 625 <p>Exponentialfunktion</p> <ul style="list-style-type: none"> Auswertung, 18, 40 <p>Exponentielle Konvergenz der trigonometrischen Interpolation, 420</p> <p>Falls $DF(x)$ nicht verfügbar, 102</p> <p>Fehlerschätzung für linear konvergente Iteration, 66</p> <p>FFT</p> <ul style="list-style-type: none"> Effizienz Gleichungslöser, 345 Primzahl, 354 <p>Fill-in-Minimierung durch Umordnung, 186</p> <p>Fitten</p> <ul style="list-style-type: none"> Hyperebene, 279 Linear, 273 <p>Fixpunktiteration, 69</p> <p>Frequenzanalyse mit DFT, 323</p> <p>Frequenzfilterung diskreter Signale, 333</p> <p>Gauss-Kollokationsverfahren, logistische Dgl., 621</p> <p>Gautschis Zweischrittverfahren, 668</p> <p>Gedämpftes Newtonverfahren, 109</p> <p>Spline quadratisch formerhaltend, Hilfskonstruktion, 510</p> <ul style="list-style-type: none"> stückweise, 490 stückweise kubisch monotoniererhaltend, 497 stückweise Polynom- von Messpunkten, 486 <p>Trigonometrische, 408</p> <p>Trigonometrische analytische Funktionen, 410</p> <p>Tschebyscheff vs. Bestapproximation, 439</p> <p>Tschebyscheff, 396</p> <p>Interpolationsfehler, 387</p> <p>Iteration</p> <ul style="list-style-type: none"> Rayleigh-Quotienten, 232 <p>Knotenanalyse eines linearen elektrischen Netzwerks, 139</p> <p>Kollokation</p> <ul style="list-style-type: none"> Gauss-, logistische Dgl., 621 <p>Kompression</p> <ul style="list-style-type: none"> Bilddaten, 269 <p>Kondition</p> <ul style="list-style-type: none"> erweitertes System, 283 Nullstellen, 36 <p>Konsistenzanalyse durch Computeralgebra, 591</p> <p>Kontraktionszahl, 310</p> <p>Konvergenz</p> <ul style="list-style-type: none"> äquidistante Trapezregel, 540 bereich, Newton-Verfahrens, 106 CG-Verfahren, 304 Clusteringapproximation mit Kollokationsmatrix, 468 Einschrittverfahren, 588 Exponentielle trigonometrische Interpolation, 420 Krylovverfahren, nichtsymmetrische Matrix, 319 rate, CG-Verfahren, 307 Remes-Algorithmus, 438 	<p>theorie, PCG, 311</p> <ul style="list-style-type: none"> vorkonditionierte inverse Iteration, 243 <p>Konvergenz der Clusteringapproximation, 467</p> <p>Konvergenz einfacher Splittingverfahren, 659</p> <p>Lösen eines Rang-1-modifizierten LGS, 199</p> <p>Lanczos-Prozess, 253</p> <ul style="list-style-type: none"> CG, 297 <p>Lanczos/Arnoldi, 257</p> <p>Lebesgue-Konstanten, 374</p> <p>Linear Konvergente Iteration, 58</p> <p>Lineare Regression bei stationären Markov-Ketten, 360</p> <p>Lineare zeitinvariante Systeme, 358</p> <p>Logistische Differentialgleichung, 579</p> <p>Lokale Konvergenz des Newton-Verfahren, 106</p> <p>Lokale Konvergenz, Sekantenverfahren, 93</p> <p>LU-Faktorisierung dünnbesetzter Matrizen, 176</p> <p>Magnetnadel</p> <ul style="list-style-type: none"> Präzession, 636 <p>Maschinengenauigkeit von MATLAB, 29</p> <p>Massgeschneidertes Newton-Verfahren, 85</p> <p>Matrix</p> <ul style="list-style-type: none"> Wilkinson, 160 <p>Mechanisches System</p> <ul style="list-style-type: none"> Pendel, 565 <p>Mehrdimensionale Fixpunktiteration, 77</p> <p>Minimierung von C^2-Funktion, 168</p> <p>Modifiziertes Gautschi-Verfahren, 676</p> <p>Molekulardynamik, 653</p> <p>Multiplikation</p> <ul style="list-style-type: none"> Polynom, 349 monomiale Darstellung, 367 Multiplikation, 349 trigonometrisch Auswertung, 469 <p>Potenzmethode</p> <ul style="list-style-type: none"> Direkte, 227 <p>Präzession</p> <ul style="list-style-type: none"> Magnetnadel, 636 <p>Primzahl-FFT, 354</p> <p>Problem</p> <ul style="list-style-type: none"> Anfangswert, autonom, 568 steifes, Adaptives semi-implizites RK-ESV, 628 steifes, implizites ESV, 625 <p>Probleme bei Polynomauswertung, 21</p> <p>Prozess</p> <ul style="list-style-type: none"> Lanczos mit CG, 297 <p>QR</p> <ul style="list-style-type: none"> basiertes Lösen eines tridiagonalen Gleichungssystems, 197 Orthogonalisierung, 193 <p>Quadratische inverse Interpolation, 95</p> <p>Quadratisches Eigenwertproblem, 224</p> <p>Quadrat</p> <ul style="list-style-type: none"> -Fehler zusammengesetzte Regeln, 538 -Fehler, Asympt. Verhalten, 534 des Kreises, 15, 41 Gauss-Legendre Ordnung 4, 528 <p>Qualitätsmass für Kernapproximation, 446</p> <p>Rayleigh-Quotienten-Iteration, 232</p> <p>Reelle Nullstellen eines quadratischen Polynoms, 49</p> <p>Regression</p>	<p>Neustart von GMRES, 317</p> <p>Newton</p> <ul style="list-style-type: none"> vereinfachtes Verfahren, 102 Verfahren in 1D, 82 Verfahren in 2D, 100 <p>Newton-Cotes-Formeln, 524</p> <p>Newton-Verfahren, modifiziert, 89</p> <p>Nichtassoziativität der Maschinenaddition, 31</p> <p>Nichtlineare Regression, 124</p> <p>Normalengleichung vs. Orthogonaltransformationsmethode, 284</p> <p>Nullabfrage</p> <ul style="list-style-type: none"> Numerische, 31 <p>Numerische Differentiation durch Extrapolation, 381</p> <p>Numerische Quadratur, 577</p> <p>Operationen $\hat{+}, \hat{\cdot}$ nicht assoziativ, 31</p> <p>Optimierung der Anordnung \rightarrow Fill-in, 186</p> <p>Orbits, 569</p> <p>Orthogonalitätsverlust der Residuen, 301</p> <p>Oszillatorische Interpolationspolynome, 376</p> <p>Pendel</p> <ul style="list-style-type: none"> Mechanisches System, 565 <p>Pendelgleichung, 642</p> <ul style="list-style-type: none"> Deskriptorform, 630 <p>Pivotelement = 0, 231</p> <p>Pivotstrategie und Rundungsfehler, 157</p> <p>Polynom</p> <ul style="list-style-type: none"> Interpolation, 490 Interpolation stückweise von Messpunkten, 486 Interpolation Versagen, 484 Linear (Fitten), 273 <p>Remes-Algorithmus</p> <ul style="list-style-type: none"> Konvergenz, 438 <p>Residuum</p> <ul style="list-style-type: none"> Grundlage Krylov-Raum, 294 <p>Richtungsfeld, 569</p> <p>Rundungsfehler</p> <ul style="list-style-type: none"> Lanczos-Prozess, 253 <p>Runge-Kutta</p> <ul style="list-style-type: none"> Adaptives semi-implizites RK-ESV, steifes Problem, 628 explizite Schritte, Ricatti Dgl., 583 Implizite ESV, schnelle Transienten, 621 Konstruktion, 581 <p>Runges Beispiel, 388</p> <p>Satellitenbahn, 595</p> <ul style="list-style-type: none"> Adaptive RK-ESV, 600 Schrittweitensteuerung, 602 <p>Schaltkreis</p> <ul style="list-style-type: none"> steife -gleichung Zeitbereich, 615 <p>Schlecht konditioniertes lineares Gleichungssystem, 144</p> <p>Schnitt von Geraden, 35</p> <p>Sekantenverfahren, 91</p> <p>Skalare Differentialgleichungen, 568</p> <p>Spline</p> <ul style="list-style-type: none"> interpolanten, Approx. vollst. kubisch, 504 Interpolation natürlich kubisch Lokalität, 502 Interpolation quadratisch formerhaltend, 513 Interpolation quadratisch formerhaltend, Hilfskonstruktion, 510 Interpolation vollständig kubisch, 502 <p>Stückweise</p>	<p>5.12 p. 697</p> <p>5.12 p. 700</p>
---	---	--	---	---

kubische Hermite-Interpolation, 496
 Polynominterpolation, 490
 Stückweise Polynominterpolation von Messpunkten, 486
 Stabilität
 LU-Zerlegung, 162
 von Algorithmus A aus Bsp. 24, 48
 Stabilitätsfunktionen, 611
 Standardintegratoren für oszillatorische Differentialgleichung, 666
 Steife Schaltkreisgleichungen im Zeitbereich, 615
 Strahlungstransport, 350
 Summation der Harmonischen Reihe, 21
 Symplektische Integratoren und variable Schrittweite, 647
 Symplektisches Eulerverfahren, 644
 Symplektizität der Pendelgleichung, 640
 Tensorprodukt-Tschebyscheff-Interpolation auf vriablen Rechtecken, 454
 Tensorprodukt-Tschebyscheff-Interpolation auf zulässigen Rechtecken, 453
 Theorem für Hermite-Interpolation, 386
 Toeplitzlöser, Superschnell, 363
 Totalpivotsuche, 162
 Trigonometrische Interpolation, 408
 analytische Funktionen, 410
 Tschebyscheff-Interpolation, 396
 Unitäre Ähnlichkeitstransformation auf Tridiagonalgestalt, 219
 Unterschiedliches Verhalten expliziter und impliziter Integratoren, 606
 Verfahren

Adaptives semi-implizites RK-ESV, steifes Problem, 628
 Arnoldi, 249
 CG (direkter Löser), 292
 CG (Konvergenz), 304
 CG (Konvergenzrate), 307
 Einschnitt, Konvergenz, 588
 Erweiterung Krylov-Unterraum, 259
 ESV adaptiv, explizit, steifes Problem, 614
 ESV implizit, steifes Problem, 625
 EW/EV, 219
 Runge-Kutta ESV, schnelle Transienten, 621
 Runge-Kutta, Konstruktion, 581
 Versagen des gedämpften Newton-Verfahrens, 110
 Versagen von Krylov-Raum basierten iterativen Lösern, 318
 Vorkommen von Clustern in Partitionsrechecken, 465
 Vorkonditionierer
 Gauss-Seidel, 241
 Jacobi-, 240
 Vorkonditionierung, 310
 Wurzeliteration, 28, 62
 Zahldarstellung durch Exponent und Mantisse, 23
 Zerlegung
 Block-LU-, 156
 Hierarchische Basis-, 560
 Teil-LU, 156
 Zerlegung der Fourier-Matrix, 328
 Zweidimensionale diskrete Fouriertransformation, 331

Definitionen und Konzepte

A-Stabilität, 620
 Analytische Funktion, 416
 Anfangswert
 Problem, Lösung, 566
 Auslöschung, 38
 AWP
 Einschrittverfahren, 580
 Kondition, 573
 Verfahren explizit, 580
 Verfahren implizit, 580
 Bézier-Kurven, 521
 Bernstein-Polynome, 519
 Cluster
 Baum, 456
 Daten
 konkav, 488
 konvex, 488
 monoton, 487
 Diagonaldominanz, 169

Differenzenquotient, 45
 Diskrete
 Faltung, 348
 Fouriertransformation, 323
 Eigenraum, 214
 Eigenvektor, 214
 Eigenwert, 214
 Evolution
 Diskrete, Konsistenzordnung, 590
 Nichtexpansivität, 634
 Faltung
 diskrete, 348
 Fehler
 Relativer, 18
 Fill-In, 177
 Fixpunkt
 abstossend, 76
 anziehend, 76
 Fouriertransformation, 323
 Frobeniusnorm, 266

Funktion
 Analytisch, 416
 konkav, 489
 konvex, 489
 Haar-Wavelet, 551
 Integral
 quadratisches erstes, 635
 Inverse
 Pseudo-, 273
 Kondition
 AWP, 573
 Eigenwert, 218
 Spektrale, 239
 Konkav
 Daten, 488
 Funktion, 489
 Konsistent
 Fixpunktiteration, 68
 Konsistenz
 Iterationsverfahren, 55
 Ordnung diskrete Evolution, 590
 Konvergenz, 54
 p -te Ordnung, 61
 global, 55
 Linear, 57
 lokal, 55
 Ordnung, 61
 Ordnung von ESV, 590
 Konvex
 Daten, 488
 Evolution, 634
 Norm, 56, 137
 Energie-, 291
 Frobeniusnorm, 266
 Numerischer Algorithmus, 47
 Pfeilmatrix, 163
 Polynom
 Bernstein, 519
 Interpolation Tensorprodukt, 446
 Lagrange verallg., 372
 Laurent, 403
 Legendre, 530
 Tschebyscheff, 391
 Problem
 Anfangswert, Lösung, 566
 Gut konditioniertes, 34
 Mathematisch, 34
 Pseudoinverse, 273
 Quadratisches erstes Integral, 635
 Raum
 Spline-, 499
 Rayleigh-Quotient, 226
 Rundungsfehler bei Potenzmethoden, 229
 Runge-Kutta
 ES-Verfahren, 582
 Verfahren, 582
 Singulärwert, 262
 zerlegung, 262
 Spektrum, 214

Funktion, 489
 Krylov-Unterraum, 245
 L-Stabilität, 622
 Landau-Notation, 39
 Laurent-Polynome, 403
 Legendre-Polynom, 530
 Maschinengenauigkeit, 29
 Maschinenzahl, 24
 Matrix
 Dünnbesetzt, 131
 Diagonal, 152
 Diagonaldominant, 169
 Hülle, 178
 Kondition, 144
 Kondition allg., 274
 normalisiert, 152
 obere Dreiecks, 152
 Orthogonal, 188
 Pseudoinverse, 273
 Rang, 141
 s.p.d., 168
 symmetrisch positiv definit, 168
 Toeplitz, 358
 Unitär, 188
 untere Dreiecks, 152
 Zirkulant, 349
 Matrixnorm, 74, 138
 Monoton
 Daten, 487
 Nichtexpansivität
 Spline
 Raum, 499
 Stabiler Algorithmus, 47
 Stabilität
 A-, 620
 Gebiet bei ESV, 609
 L-, 622
 Struktursymmetrie, 178
 SVD, 262
 Symmetrisches Einschrittverfahren, 649
 Tensorprodukt
 Polynominterpolation, 446
 Toeplitzmatrix, 358
 Transformation
 Fourier, 323
 Trigonometrisches Polynom, 402
 Tschebyscheff-Polynome, 391
 Verfahren
 AWP, explizit, 580
 AWP, implizit, 580
 Einschritt, AWP, 580
 ESV, Runge-Kutta, 582
 Runge-Kutta, 582
 Vielfachheit
 geometrische, 214
 Wavelet
 Haar, 551
 Zerlegung
 Singulärwert, 262

5.12
 p. 701

5.12
 p. 703

5.12
 p. 702

5.12
 p. 704

MATLAB-CODE-Fragmente

bisect , 80 sinetrans, 339 symamd, 186 symrcm, 186 bisect Kapitel.2.NichtlineareGleichungen/ Iterationsverfahren/ bisect.m, 80 blockgs Kapitel.3.NumerischeLineareAlgebra/ 3.2.NumerischeLoesung_linearer_Gleichungssysteme/ blockgs.m, 149 broyden Kapitel.2.NichtlineareGleichungen/ Newton-Verfahren/Broyden-Verfahren/ fastbroyd.m , 116 dampnewton Kapitel.2.NichtlineareGleichungen/ Newton-Verfahren/ dampnewton.m, 109 gn Kapitel.3.NumerischeLineareAlgebra/ 3.1.Grundbegriffe.und.operationen/ spinit.m, 135 Adaptives ESV, steifes Problem, 615 arnoldidev, 249 arnoldi, 247 bernstein, 520 bicgstab, 317 blockgs, 149 boxbastohwf, 554 broyden, 116 cg, 301 cholinc, 243 cholupdate, 206 chol, 172, 186 cond, 146 conv, 354 costrans, 347 ct.rect, 457 dampnewton, 109 divide, 459 dorthp, 430 eigs, 260 eig, 221, 223 evaliptrig, 407 evlortho, 431 expeval, 18 Kapitel.1.ComputerArithmetics/ ExponentialFunction/ expeval.simple.m, 18 fftsolve, 345	Kapitel.2.NichtlineareGleichungen/ NichtlineareAusgleichsprobleme/ gaussnewton.m , 123 lupattern Kapitel.3.NumerischeLineareAlgebra/ 3.2.NumerischeLoesung_linearer_Gleichungssysteme/ lupattern.m, 155 lurec Kapitel.3.NumerischeLineareAlgebra/ 3.2.NumerischeLoesung_linearer_Gleichungssysteme/ lurec.m, 154 newton Kapitel.2.NichtlineareGleichungen/ Newton-Verfahren/ newton.m, 102 secant Kapitel.2.NichtlineareGleichungen/ Iterationsverfahren/ secant.m, 91 spinit fft, 323 fzero, 80 gallery, 235, 238, 243, 250, 254 gaussQuad, 534 gmres, 316 gn, 123 hwftoboxbas, 555 icostrans, 347 ifft, 323 ipoleval, 380 lanczos, 253 legendre, 531 lsgsvd, 278 lsgtotal, 286 lurec, 154 lu, 160, 177 minres, 316 newton, 102 ode15s, 628 ode23s, 628 ode23, 599 ode45, 599, 615, 660 odeset, 599, 615, 628, 660 partition, 459 pcg, 300, 313 pchip, 497 planerot, 191 polyfit, 274, 380 polyval, 380 qmr, 318 grupdate, 203	qr, 193 quadl, 546 quadzero, 51 Kapitel.1.ComputerArithmetics/ QuadZero/ quadzero.m, 51 quad, 546 qzfool, 51 Kapitel.1.ComputerArithmetics/ QuadZero/ qzfool.m, 51 qztest Kapitel.1.ComputerArithmetics/ QuadZero/ qztest.m, 51 remes, 438 roudchol, 206 rqui, 232 sa1, 165 sa2, 165 sa3, 167 schur, 221 secant, 91 sinetrans, 341 sinft2d, 344 smw, 200 sparse, 134 spdiags, 134 speye, 134 spline, 502 spones, 134 Näherung für e, 14 Kapitel.1.ComputerArithmetics/ Calculate.EulerNumber/ Calculate.EulerNumber.m, 14 Niedrigrang-Bestapproximation, 271 Numerische Differentiation, 45 Kapitel.1.ComputerArithmetics/ NumericDifferentiation/ NumericDifferentiation.m, 45 Parameter der Gleitkommaarithmetik, 27 Plotten von Funktionensystemen, 476 Polynomauswertung, 22 Kapitel.1.ComputerArithmetics/ EvaluationOfPolynom/ pvunst.m, 22 Rückwärtsrekursion für Besselfunktion, 482 Rekursive Gausselimination, 149 Rekursive Gausselimination mit Pivotsuche, 159 Rekursive LU-Zerlegung, 154 Summation rückwärts Kapitel.1.ComputerArithmetics/ SummationHarmonic/ SummationHarmonic.m, 21 Unterlauf, 33	spowit, 234 spy, 186 squareroot, 28 Kapitel.1.ComputerArithmetics/ SquareRoot/ squareroot.m, 28 svds, 263 svd, 263 symamd, 186 symrcm, 186 tic,toc, 135, 345 toeplitz, 361 tril, 240 triu, 240 Abfrage der Maschinengenauigkeit, 30 Adaptive Rückwärtsrekursion für Besselfunktion 484 Berechnung von pi Instabile, 16 Kapitel.1.ComputerArithmetics/ CircleQuadrature/ Approx.PI_instable.m, 16 Stabile, 43 Kapitel.1.ComputerArithmetics/ CircleQuadrature/ Approx.PI_stable.m, 43 Demonstration von Rundungsfehlern, 29 Givens-Rotation, 191 linsolve, 174	5.12 p. 705 5.12 p. 707 5.12 p. 706 5.12 p. 707
--	---	---	--	--

Symbolverzeichnis

$O(\epsilon)$, 39	$\ \cdot\ $, 56, 137
$O(n)$, 129	$\rho_A(\mathbf{u})$, 226
$\mathcal{R}_k(m, n)$, 266	\star , 29
eps, 29	$\ \mathbf{x}\ _1$, 56
l_T , 371	$\ \mathbf{x}\ _2$, 56
$\mathcal{K}_l(\mathbf{A}, \mathbf{z})$, 245	$\ \mathbf{x}\ _\infty$, 56
$\ \mathbf{A}\ _F^2$, 266	
$\ \mathbf{x}\ _A$, 291	
$\ f\ _{L^\infty(I)}$, 373	
$\ f\ _{L^1(I)}$, 373	
$\ f\ _{L^2(I)}^2$, 373	
\mathcal{P}_k , 367	
\mathbf{T}_n , 404	
\mathbf{A}^+ , 273	
\mathbb{M} , 24	
ϵ_F , 18	
$\kappa(\mathbf{A})$, 239	
κ_{abs} , 34	
κ_{rel} , 35	
rd, 28	

Literaturverzeichnis

[1] M. ABRAMOWITZ UND I. STEGUN, *Handbook of Mathematical Functions*, Dover Publications, New York, 1970.

[2] P. AMESTOY, T. DAVIS UND I. DUFF, *An approximate minimum degree ordering algorithm*, SIAM Journal on Matrix Analysis and Applications, 17 (1996), S. 886–905.

[3] R. BARRETT, M. BERRY, T. F. CHAN, J. DEMMEL, J. DONATO, J. DONGARRA, V. EIJKHOUT, R. POZO, C. ROMINE UND H. V. DER VORST, *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*, SIAM, Philadelphia, PA, 2nd ed., 1994.

[4] S. BÖRM, L. GRASEDYCK UND W. HACKBUSCH, *Introduction to hierarchical matrices with applications*, Engineering Analysis with Boundary Elements, 27 (2003), S. 405–422.

[5] R. BRENT, C. PERCIVAL UND P. ZIMMERMANN, *Error bounds on complex floating-point multiplication*, Math. Comp., (2007).

[6] Q. CHEN UND 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), S. 405–417.

[7] C. CHUI, *Wavelets and spline interpolation*, in Advances in numerical analysis. Vol. 2: Wavelets, subdivision algorithms, and radial basis functions, Proc. 4th Summer Sch., Lancaster/UK 1990, Clarendon Press, oxford ed., 1992, S. 1–35.

[8] D. COPPERSMITH UND T. RIVLIN, *The growth of polynomials bounded at equally spaced points*, SIAM J. Math. Anal., 23 (1992), S. 970–983.

[9] P. DAVIS, *Interpolation and Approximation*, Dover, New York, 1975.

[10] P. DEUFLHARD, *Newton Methods for Nonlinear Problems*, Band 35 in Springer Series in Computational Mathematics, Springer, Berlin, 2004.

[11] P. DEUFLHARD UND F. BORNEMANN, *Numerische Mathematik II*, DeGruyter, Berlin, 2 ed., 2002.

[12] P. DEUFLHARD UND A. HOHMANN, *Numerische Mathematik I*, DeGruyter, Berlin, 3 ed., 2002.

[13] P. DUHAMEL UND M. VETTERLI, *Fast fourier transforms: a tutorial review and a state of the art*, Signal Processing, 19 (1990), S. 259–299.

[14] A. DUTT UND V. ROKHLIN, *Fast Fourier transforms for non-equispaced data II*, Appl. Comput. Harmon. Anal., 2 (1995), S. 85–100.

[15] F. FRITSCH UND R. CARLSON, *Monotone piecewise cubic interpolation*, SIAM J. Numer. Anal., 17 (1980), S. 238–246.

[16] M. GANDER, W. GANDER, G. GOLUB UND D. GRUNTZ, *Scientific Computing: An introduction using MATLAB*, Springer, 2005. In Vorbereitung.

[17] W. GANDER UND W. GAUTSCHI, *Adaptive quadrature - revisited*, BIT, 40 (2000), S. 84–101.

[18] J. GILBERT, C. MOLER UND R. SCHREIBER, *Sparse matrices in MATLAB: Design and implementation*, SIAM Journal on Matrix Analysis and Applications, 13 (1992), S. 333–356.

[19] G. GOLUB UND C. VAN LOAN, *Matrix computations*, John Hopkins University Press, Baltimore, London, 2nd ed., 1989.

[20] L. GREENGARD UND V. ROKHLIN, *A new version of the fast multipole method for the Laplace equation in three dimensions*, Acta Numerica, (1997), S. 229–269.

[21] W. HACKBUSCH, *Iterative Lösung großer linearer Gleichungssysteme*, B.G. Teubner-Verlag, Stuttgart, 1991.

[22] W. HACKBUSCH UND S. BÖRM, *Data-sparse approximation by adaptive \mathcal{H}^2 -matrices*, Computing, 69 (2002), S. 1–35.

[23] E. HAIRER, C. LUBICH UND G. WANNER, *Geometric numerical integration*, Band 31 in Springer Series in Computational Mathematics, Springer, Heidelberg, 2002.

[24] —, *Geometric numerical integration illustrated by the Störmer-Verlet method*, Acta Numerica, 12 (2003), S. 399–450.

[25] E. HAIRER, S. NORSETT UND G. WANNER, *Solving Ordinary Differential Equations I. Nonstiff Problems*, Springer-Verlag, Berlin, Heidelberg, New York, 2 ed., 1993.

- [26] E. HAIRER UND G. WANNER, *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic Problems*, Springer-Verlag, Berlin, Heidelberg, New York, 1991.
- [27] C. HALL UND W. MEYER, *Optimal error bounds for cubic spline interpolation*, J. Approx. Theory, 16 (1976), S. 105–122.
- [28] M. HANKE-BOURGEOIS, *Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens*, Mathematische Leitfäden, B.G. Teubner, Stuttgart, 2002.
- [29] N. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, PA, 2 ed., 2002.
- [30] M. HOCHBRUCK UND C. LUBICH, *A Gautschi-type method for oscillatory second-order differential equations*, Numer. Math., 83 (1999), S. 403–426.
- [31] D. MCALLISTER UND J. ROULIER, *An algorithm for computing a shape-preserving osculatory quadratic spline*, ACM Trans. Math. Software, 7 (1981), S. 331–347.
- [32] J. MELENK, *Vorlesungsaufzeichnungen numerik*, skriptum, Universität Regensburg, 2003.
- [33] K. NEYMEYR, *A geometric theory for preconditioned inverse iteration applied to a subspace*, Tech. Report 130, SFB 382, Universität Tübingen, Tübingen, Germany, November 1999. Submitted to Math. Comp.
- [34] —, *A geometric theory for preconditioned inverse iteration: III. Sharp convergence estimates*, Tech. Report 130, SFB 382, Universität Tübingen, Tübingen, Germany, November 1999.
- [35] M. OVERTON, *Numerical Computing with IEEE Floating Point Arithmetic*, SIAM, Philadelphia, PA, 2001.
- [36] A. D. H.-D. QI, L.-Q. QI UND H.-X. YIN, *Convergence of Newton's method for convex best interpolation*, Numer. Math., 87 (2001), S. 435–456.
- [37] C. RADER, *Discrete Fourier transforms when the number of data samples is prime*, Proceedings of the IEEE, 56 (1968), S. 1107–1108.
- [38] R. RANNACHER, *Einführung in die numerische mathematik*. Vorlesungsskriptum Universität Heidelberg, 2000. <http://gaia.iwr.uni-heidelberg.de/>.
- [39] J. RIVLIN, *The Chebyshev Polynomials*, Wiley-Interscience, 1984.
- [40] V. ROKHLIN, *Rapid solution of integral equations of classical potential theory*, J. Comp. Phys., 60 (1985), S. 187–207.
- [41] H. SCHWARZ, *Methode der finiten Elemente*, Band 47 in Leitfäden der angewandten Mathematik und Mechanik, Teubner-Verlag, Stuttgart, 3rd ed., 1991.
- [42] M. STEWART, *A superfast toeplitz solver with improved numerical stability*, SIAM J. Matrix Analysis Appl., 25 (2003), S. 669–693.
- [43] J. STOER, *Einführung in die Numerische Mathematik*, Heidelberger Taschenbücher, Springer, 4 ed., 1983.
- [44] F. TISSEUR UND K. MEERBERGEN, *The quadratic eigenvalue problem*, SIAM Review, 43 (2001), S. 235–286.
- [45] H. WERNER UND R. SCHABACK, *Praktische Mathematik II*, Hochschulktext, Springer, Berlin, 1972.

5.12
p. 713

5.12
p. 714