

EINFÜHRUNG IN DIE NUMERIK

Vorlesungsmanuskript, Stand 29.01.2010

Wintersemester 2009/2010

Peter Bastian
Interdisziplinäreszentrum für Wissenschaftliches Rechnen
Universität Heidelberg, Im Neuenheimer Feld 368, D-69120 Heidelberg
Peter.Bastian@iwr.uni-heidelberg.de

INHALTSVERZEICHNIS

1	Grundbegriffe der Numerik	4
1.1	Stabilität und Konditionszahl	4
1.2	Numerische Auswertung von Funktionsvorschriften	9
1.3	Lösen von Gleichungen	12
2	Fließkommazahlen	17
2.1	Zahldarstellung	17
2.2	Runden und Rundungsfehler	22
2.3	Fließkommaarithmetik	25
2.4	Fehleranalyse	28
2.5	Auslöschung	33
3	Motivation linearer Gleichungssysteme	34
3.1	Strömung in Rohrleitungsnetzen	34
3.2	Radiosity-Methode in der Computergrafik	38
4	Konditionierung der Lösung linearer Gleichungssysteme	41
4.1	Lösbarkeit	41
4.2	Vektornormen	42
4.3	Matrixnormen	44
4.4	Eigenwerte und Eigenvektoren	46
4.5	Die Spektralnorm	47
4.6	Positiv definite Matrizen	49
4.7	Störungstheorie	51
5	Eliminationsverfahren zur Lösung linearer Gleichungssysteme	57
5.1	Dreieckssysteme	57
5.2	Gauß Elimination	58
5.3	LR-Zerlegung	63
5.4	Rundungsfehleranalyse der LR-Zerlegung	71
5.5	Pivotisierung	81
5.6	Spezielle Systeme	86
6	Interpolation und Approximation	96
6.1	Einführung	96
6.2	Polynominterpolation	98
6.3	Anwendungen der Polynominterpolation	106
6.4	Bernstein-Polynome und Kurveninterpolation	112
6.5	Splines	118
6.6	Trigonometrische Interpolation	125
6.7	Approximation von Funktionen	134
7	Numerische Integration	146
7.1	Newton-Cotes Formeln	146
7.2	Summierte Quadraturformeln	150
7.3	Quadraturen höherer Ordnung	152
7.4	Ausblick	154
8	Iterative Lösung von Gleichungssystemen	156
8.1	Newton-Verfahren	157
8.2	Sukzessive Approximation	164
8.3	Iterationsverfahren zur Lösung linearer Gleichungssysteme	166

1.1. Stabilität und Konditionszahl

Betrachte das Auswerten einer Funktionsvorschrift

$\text{Gegeben } x, \text{ berechne } y = F(x).$

(1.1)

Wobei $F: X \rightarrow Y$ mit X, Y normierte Räume.

(Es geht darum F später näherungsweise im Rechner auszuwerten).

Frage: Wie wirkt sich Änderung in x im Ergebnis $F(x)$ aus?

M.a.W. welche Stetigkeit erfüllt F ?



Definition 1.1 (Stabilität) Wir nennen die Auswertung der Funktion $F: X \rightarrow Y$

a) stabil, wenn F stetig ist.

d.h. $\forall \epsilon > 0 \exists \delta = \delta(\epsilon) : \forall x, x' \in X : \|x - x'\|_X \leq \delta \Rightarrow \|F(x) - F(x')\|_Y \leq \epsilon$

b) lokal L -stabil, wenn F lokal Lipschitz-stetig ist

d.h. $\forall x \in X \exists \delta(x) > 0 \exists K_0(x) > 0 : \forall x' \in X : \|x - x'\|_X \leq \delta \Rightarrow \|F(x) - F(x')\|_Y \leq K_0(x) \|x - x'\|_X$

\neq heißt L -stabil, falls $K_0(x) = K_0$ unabhängig von x gilt. L -stabil ist „stärker“ als stabil, d.h. L -stabil \Rightarrow stabil aber nicht umgekehrt.

Die Definition 1.1 b ist in der Praxis nützlicher, da sie eine Quantifizierung der Abhängigkeit $y = F(x)$ erlaubt:

$$\frac{\| \delta y \|_Y}{\| \delta x \|_X} = \frac{\| F(\overset{x'}{x + \delta x}) - F(x) \|_Y}{\| \delta x \|_X} \leq K_0(x)$$

(Betrag der Ableitung kann nicht unendlich sein).

Definition 1.2 (Konditionszahlen)

Die absolute Kondition einer Abbildung $F: X \rightarrow Y$ im Punkt x ist

$$K_{abs}(x) = \sup \left\{ \frac{\|F(x+\delta x) - F(x)\|_Y}{\delta \cdot \|\delta x\|} \mid \delta x \neq 0, x + \delta x \in X \right\} \quad (1.2)$$

Entsprechend lautet die relative Kondition im Punkt x :

$$K(x) = \sup \left\{ \frac{\|F(x+\delta x) - F(x)\| / \|F(x)\|}{\|\delta x\| / \|x\|} \mid \delta x \neq 0, x + \delta x \in X \right\} \quad (1.3)$$

← relative Abweichung.

wobei $x \neq 0$ und $F(x) \neq 0$ angenommen sei

$F: X \rightarrow Y$ heißt

- gut konditioniert (in x) falls $K(x)$ „klein“.
- schlecht " " " " " „groß“.

Erwartung :- Gut konditionierte Vorschriften erlauben eine Berechnung auf dem Computer.

- Bei schlecht konditionierten Vorschriften ist mit „Problemen“ zu rechnen. Und zwar unabhängig vom genauen numerischen Verfahren !

Beispiel 1.3 (Kondition der Addition)

Betrachte $y = F(x_1, x_2) = x_1 + x_2$.

Mit $x = (x_1, x_2)^T$ und $\delta x = (\delta x_1, \delta x_2)^T$ lautet die Taylorreihe

$$F(x + \delta x) = F(x) + \nabla F(x) \cdot \delta x + O(\|\delta x\|^2)$$

- Restglied
- Landausche Symbole

$$\Rightarrow |F(x + \delta x) - F(x)| \leq |\nabla F(x) \cdot \delta x| + O(\|\delta x\|^2)$$

Cauchy-Schwarz $\Rightarrow \leq \|\nabla F(x)\| \|\delta x\| + O(\|\delta x\|^2)$

$$\Leftrightarrow \frac{|F(x + \delta x) - F(x)|}{\|\delta x\|} \leq \|\nabla F(x)\| + O(\|\delta x\|)$$

Nun ist $\nabla F(x) = (1, 1)^T$ und somit

$$K_{\text{abs}}(x) \stackrel{\substack{\uparrow \\ \text{„in erster Naherung“}}}{=} \|(1, 1)^T\| = \sqrt{2} \quad (\text{fur } \|\cdot\| = \|\cdot\|_2 \text{ euklidische Norm})$$

gilt auch fur andere Normen

Fur die relative Kondition gilt

$$K(x) \stackrel{\substack{\uparrow \\ \text{„in erster Naherung“}}}{=} \|(1, 1)^T\| \frac{\|x\|}{|x_1 + x_2|}$$

Fur $x_1 \approx -x_2$ gilt $|x_1 + x_2| \ll \|x\| = \sqrt{x_1^2 + x_2^2} \rightarrow \sqrt{2} |x_1|$

Somit ist die Addition schlecht konditioniert fur $x_1 \approx -x_2$.

Beispiel 1.4

Wir betrachten die Lösung der quadratischen Gleichung

$$x^2 - 2px + 1 = 0.$$

Quadratische Ergänzung liefert

$$x^2 - 2px + p^2 = p^2 - 1$$

$$\Leftrightarrow (x-p)^2 = p^2 - 1$$

$$\Leftrightarrow x_{1,2} = p \pm \sqrt{p^2 - 1}$$

Uns interessiert die Vorschrift

$$F(p) = \begin{pmatrix} p + \sqrt{p^2 - 1} \\ p - \sqrt{p^2 - 1} \end{pmatrix} \text{ mit } p \in [1, \infty)$$

Taylor:

$$F(p + \delta p) = F(p) + \frac{dF}{dp}(p) \delta p + O(|\delta p|^2)$$

$$\|F(p + \delta p) - F(p)\| \leq \left\| \frac{dF}{dp}(p) \right\| |\delta p| + O(|\delta p|^2)$$

$$\Leftrightarrow \frac{\|F(p + \delta p) - F(p)\|}{|\delta p|} \leq \left\| \frac{dF}{dp}(p) \right\| + O(|\delta p|) \quad \text{euklidische Norm:}$$

$$\begin{aligned} \left\| \frac{dF}{dp} \right\|_2 &= \sqrt{\left(1 + \frac{p}{\sqrt{p^2 - 1}}\right)^2 + \left(1 - \frac{p}{\sqrt{p^2 - 1}}\right)^2} \\ &= \sqrt{1 + \frac{p^2}{p^2 - 1} + 1 - \frac{p^2}{p^2 - 1}} \\ &= \sqrt{2 \frac{p^2 - 1 + p^2}{p^2 - 1}} = \sqrt{2} \sqrt{\frac{2p^2 - 1}{p^2 - 1}} \end{aligned}$$

$$\frac{dF}{dp} = \left(1 + \frac{p}{\sqrt{p^2 - 1}}, 1 - \frac{p}{\sqrt{p^2 - 1}}\right)^T$$

|| euklidische Norm:

$$K_{abs}(p) = \sqrt{2} \sqrt{\frac{2p^2 - 1}{p^2 - 1}} \text{ schlecht konditioniert f\u00fcr } p=1$$

$$K(p) = \left\| \frac{dF}{dp}(p) \right\| \frac{|p|}{\|F(p)\|} = \sqrt{2} \sqrt{\frac{2p^2 - 1}{p^2 - 1}} \frac{|p|}{\sqrt{\frac{(p + \sqrt{p^2 - 1})^2}{\rightarrow 1} + \frac{(p - \sqrt{p^2 - 1})^2}{\rightarrow 1}}}$$

ebenfalls schlecht konditioniert f\u00fcr $p \rightarrow 1$.

\"Ubung: Kondition der L\u00f6sungsformel von $x^2 - \frac{1+t^2}{t}x + 1$ $t \in [1, \infty]$
 $\frac{1+t^2}{t}$ bildet $[1, \infty]$ bijektiv auf sich selbst ab.

1.2 Numerische Auswertung von Funktionsvorschriften

Im Rechner kann die Auswertung $y = F(x)$ nicht exakt realisiert werden.

Stattdessen realisieren wir dort für $k \in \mathbb{N}$

$$y^{(k)} = F^{(k)}(x^{(k)}) \quad \text{mit} \quad F^{(k)} : X^{(k)} \rightarrow Y^{(k)}, \quad (1.4)$$

mit normierten Räumen $X^{(k)}, Y^{(k)}$.

Beispiel 1.5

(a)
$$F^{(k)} : (F(10, k, 3))^2 \rightarrow F(10, k, 2)$$

$$F^{(k)}(x_1^{(k)}, x_2^{(k)}) = x_1^{(k)} \oplus x_2^{(k)}$$

$$F(x_1, x_2) = x_1 + x_2$$

$F(10, k, 3)$: Zahlen der Form
 $\pm 0, \underbrace{m_1 \dots m_k}_{0 \leq m_i < 10} \cdot 10^{\pm e}$
 Exponent mit 2 Stellen
 Ziffern der Mantisse

$k \rightarrow \infty$ nutzt immer mehr Stellen in Dezimaldarstellung.

(b) $F^{(k)} : \mathbb{R} \rightarrow \mathbb{R}$

$$F^{(k)}(x) = 1 + \sum_{i=1}^k \frac{x^i}{i!}$$
 abgebr. Reihe für e-Funktion.

$$F(x) = \exp(x)$$

Hier gilt $X^{(k)} \subseteq X, Y^{(k)} \subseteq Y$. Das muss aber nicht so sein.

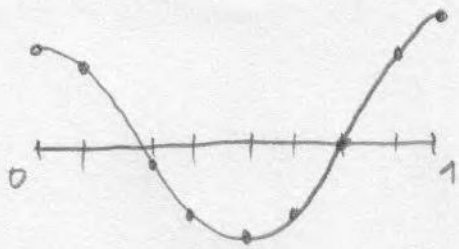
Im folgenden gehen wir davon aus, dass $X^{(k)} \subseteq X$ aber $Y^{(k)} \not\subseteq Y$.

Es gebe aber eine Abbildung $R^{(k)} : Y \rightarrow Y^{(k)}$ mit

$R^{(k)}$ linear und lokal L-stabil

$$R^{(k)}(x+x') = R^{(k)}(x) + R^{(k)}(x')$$

$$R^{(k)}(\alpha x) = \alpha R^{(k)}(x)$$



$$(R^{(k)} y)_i = y\left(\frac{i}{k}\right) \quad 0 \leq i \leq k.$$

Frage: $x^{(k)} \rightarrow x$, geht dann auch $F^{(k)}(x^{(k)}) \rightarrow F(x)$?

Mathematisch: Konvergenz der Folge $\{F^{(k)}(x^{(k)})\}$

Definition 1.7 (Konvergenz) Die numerische Auswertung (1.4)

heißt konvergent, genau dann wenn für jedes $x \in X$ gilt:

$$\forall \epsilon > 0 \exists k_0(\epsilon) \in \mathbb{N} \exists \delta(\epsilon, k_0) > 0:$$

$$\forall k > k_0, \forall x' \in X^{(k)}: \|x - x'\| \leq \delta \Rightarrow \|R^{(k)}F(x) - F^{(k)}(x')\| \leq \epsilon$$

Aus Def. 1.7 folgt auch die Konvergenz

dann ist auch $F(x) - F^{(k)}(x)$ klein

$$\|R^{(k)}F(x') - F^{(k)}(x')\| \rightarrow 0 \text{ für } k \rightarrow \infty \text{ und } x' \in X^{(k)}. \quad (1.5)$$

↑
gleiches Argument!

Denn sei $x' \in X^{(k)}$ $x \in X$ beliebig

$$\|R^{(k)}F(x') - F^{(k)}(x')\| = \|R^{(k)}F(x') - R^{(k)}F(x) + R^{(k)}F(x) - F^{(k)}(x')\|$$

Dreiecks-
ungl. \leq

$$\underbrace{\|R^{(k)}F(x') - R^{(k)}F(x)\|}_{\text{nur } F!} + \underbrace{\|R^{(k)}F(x) - F^{(k)}(x')\|}_{\leq \frac{\epsilon}{2} \text{ wenn } k > k_0(\frac{\epsilon}{2}), \|x - x'\| \leq \delta(\frac{\epsilon}{2}, k_0)}$$

$\leq R_0 K_0(x) \|x - x'\|$ wenn F lokal L -stabil

also $\leq \frac{\epsilon}{2}$ falls $\|x - x'\| \leq \frac{\epsilon}{2R_0 K_0}$

Gilt umgekehrt die Konvergenz für $\|R^{(k)}F(x') - F^{(k)}(x')\|$ für $x' \in X^{(k)}$ und $k \rightarrow \infty$, dann

$$\|R^{(k)}F(x) - F^{(k)}(x')\| \leq \underbrace{\|R^{(k)}F(x) - R^{(k)}F(x')\|}_{\text{wieder Stabilität von } F} + \underbrace{\|R^{(k)}F(x') - F^{(k)}(x')\|}_{\text{Konvergenz zweiter Art}}$$

↑
Bed. in 1.7

Somit sind die beiden Bedingungen aus Def. 1.7 und Gl. (1.5) äquivalent.

Für die Praxis ist diese Aufspaltung praktisch, da

- $\|R^{(k)}F(x) - R^{(k)}F(x')\| \leq \|R^{(k)}\| \|F(x) - F(x')\|$
nur die Stabilität der Auswertung F ist
- $\|R^{(k)}F(x') - F^{(k)}(x')\|$ den Rundungs- oder Abruchfehler darstellt (siehe etwa Beispiel 1.5).

Bisher haben wir die Stabilität der $F^{(k)}$ nicht gefordert. Es zeigt sich, dass diese eine notwendige Voraussetzung für Konvergenz ist.

Wir zeigen: Konvergenz $F^{(k)} \rightarrow F \Rightarrow$ Stabilität der $F^{(k)}$

damit gilt dann: \neg Stabilität $\Rightarrow \neg$ Konvergenz.

Satz 1.8 Es sei F stabil und $F^{(k)}$ konvergent nach Definition 1.7.

Dann sind die $F^{(k)}$ stabil. (nicht lokal L -stabil!)

Beweis: Sei $x, x' \in X^{(k)}$

$$\|F^{(k)}(x) - F^{(k)}(x')\| = \underbrace{\|F^{(k)}(x) - R^{(k)}F(x)\|}_{\text{Konvergenz nach Gl. (1.5)}} + \underbrace{\|R^{(k)}F(x) - F^{(k)}(x')\|}_{\text{Konvergenz wie in Def. 1.7}}$$

$$\leq \|R^{(k)}F(x) - F^{(k)}(x)\| + \|R^{(k)}F(x) - F^{(k)}(x')\|$$

Konvergenz nach Gl. (1.5)

Konvergenz wie in Def. 1.7.

Ist nun $k > k_0(\frac{\epsilon}{2})$ und $\|x - x'\| \leq \delta(\frac{\epsilon}{2}, k)$ (solche k_0 und δ existieren nach Def. 1.7)

so sind beide Terme jeweils kleiner als $\frac{\epsilon}{2}$ und $F^{(k)}$, $k > k_0$, ist stabil.

1.3 Lösen von Gleichungen

24.09.09

Bisher haben wir nur die Auswertung von Funktionsvorschriften betrachtet.

Oft ist die unbekannte Größe y nur implizit durch eine Gleichung bestimmt:

Gegeben $x \in X$ ("Daten"), finde $y \in Y$ ("Lösung") so dass

$$G(y, x) = 0$$

(1.6)

Definition 1.10 (Sachgemäß gestellt).

Das Problem (1.6) heißt sachgemäß gestellt falls gilt:

(a) zu jedem $x \in X$ gibt es genau eine Lösung $y \in Y$.

(b) Die Lösung y hängt stetig von den Daten x ab.

Die Funktion $F: X \rightarrow Y$ mit

$$\forall x \in X : G(F(x), x) = 0$$

heißt Lösungsoperator zur Gleichung (1.6).

Sachgemäß gestellt bedeutet also, dass genau ein solches F existiert und F stabil ist.

\neg (sachgemäß gestellt) = schlecht gestellt.

Beispiel 1.11 Sei $x \in \mathbb{R}^n$ und $A \in \mathbb{R}^{n \times n}$ invertierbar.

Setze $G(y, x) = Ay - x$, dann ist $F(x) = A^{-1}x$.

9
24.09.09

Im Rechner können wir (1.6) nur näherungsweise zu
genäherten Daten lösen. Für $k \in \mathbb{N}$ betrachte die Probleme

$$\text{Gegeben } x^{(k)} \in X^{(k)}, \text{ finde } y^{(k)} \in Y^{(k)} \text{ sodass } G^{(k)}(y^{(k)}, x^{(k)}) = 0 \quad (1.7)$$

Zu jedem $G^{(k)}$ gebe es einen numerischen Lösungsoperator
 $F^{(k)}: X^{(k)} \rightarrow Y^{(k)}$:

$$\forall x^{(k)} \in X^{(k)}: G^{(k)}(F^{(k)}(x^{(k)}), x^{(k)}) = 0$$

Wie oben sei $X^{(k)} \subseteq X$ und $R^{(k)}: Y \rightarrow Y^{(k)}$ da $Y^{(k)} \not\subseteq Y$ zulässig ist.

Wie oben gezeigt, ist die Stabilität der $F^{(k)}$ eine notwendige
Voraussetzung für die Konvergenz. $R^k y - y^{(k)} \rightarrow 0$.

Wir skizzieren nun eine einfache Konvergenztheorie, die
in der Praxis häufig bei gewöhnlichen und partiellen DGL Einsatz findet

Zentral ist dabei

Definition 1.12 (Konsistenz)

Das numerische Verfahren (1.7) heißt konsistent zu (1.6)
falls gilt

$$G^{(k)}(R^{(k)} y, x + \delta x^{(k)}) \rightarrow 0 \text{ für } k \rightarrow \infty, \delta x^{(k)} \rightarrow 0. \quad \square$$

Bei vielen praktischen Problemen lässt sich diese
Eigenschaft leicht nachweisen. Siehe Beispiel.

10
24.09.09

Die Funktion $G^{(k)}: Y^{(k)} \times X^{(k)} \rightarrow Z^{(k)}$ sei invertierbar
 bezüglich dem ersten Argument. D.h. es existiere $H^{(k)}: Z^{(k)} \times X^{(k)} \rightarrow Y^{(k)}$
 so dass

$$\forall x \in X^{(k)}, \forall y \in Y^{(k)}: H^{(k)}(\underbrace{G^{(k)}(y, x)}_z, x) = y \quad (1.8)$$

Definition 1.13

Das numerische Verfahren ^(1.7) heißt stabil falls $H^{(k)}$ stabil bezüglich dem ersten Argument ist.

Damit zeigt man den zentralen

Satz 1.14 Das numerische Verfahren (1.7) sei konsistent und stabil. Dann ist es auch konvergent.

Beweis: Zu $x \in X$ sei $G(y, x) = 0$, $G^{(k)}(y^{(k)}, x + \delta x^{(k)}) = 0$, $\|\delta x^{(k)}\| \rightarrow 0$ für $k \rightarrow \infty$.

$$\|R^{(k)}y - y^{(k)}\| = \left\| H^{(k)}\left(\underbrace{G^{(k)}(R^{(k)}y, x + \delta x^{(k)})}_{= R^{(k)}y \text{ nach } 1.8}, x + \delta x^{(k)}\right) - \underbrace{H^{(k)}\left(G^{(k)}(y^{(k)}, x + \delta x^{(k)}), x + \delta x^{(k)}\right)}_{= y^{(k)}} \right\|$$

$$\leq k_0 \left\| \underbrace{G^{(k)}(R^{(k)}y, x + \delta x^{(k)}) - G^{(k)}(y^{(k)}, x + \delta x^{(k)})}_{= 0 \text{ n. Vor.}} \right\|$$

↑
Stabilität von $H^{(k)}$

$$= k_0 \left\| \underbrace{G^{(k)}(R^{(k)}y, x + \delta x^{(k)})}_{\rightarrow 0 \text{ wg. Konsistenz.}} \right\|$$

Bemerkung 1.14

Der Begriff „Stabilität“ tritt in der Numerik häufig auf. Er meint immer einen Zusammenhang der Form $\|f(x) - f(x')\| \leq K(x)\|x - x'\|$. Was f genau ist variiert von Problem zu Problem.

Beispiel 1.15 (Anfangswertproblem) Wir zeigen, dass die abstrakte Definitionen eine Anwendung haben! 11
24.09.09

Betrachte die DGL

$$y'(t) = f(t, y(t)) \quad \text{für } t \in I = [0, T],$$

$$y(0) = x \quad (\text{Anfangswert}). \quad (1.9)$$

Hier ist $Y = C^1(I)$, $X = \mathbb{R}$, $Z = \mathbb{R}$ und eine mögliche Wahl für G :

$$G: Y \times X \rightarrow Z, \quad G(y, x) = |y(0) - x| + \sup_{t \in I} |y'(t) - f(t, y(t))|$$

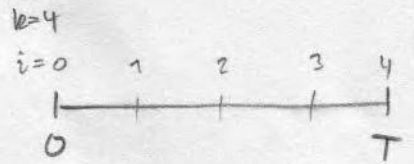
Das schreibt man in der Praxis nicht so hin.
 ↳ alternative Normen möglich
 ↳ bietet sich an wegen $y \in C^1(I)$

Lösungsoperator $F: X \rightarrow Y$ liefert zu gegebenem Anfangswert x die Lösung.

Wir setzen voraus, dass f in (1.9) darstellt, dass (1.9) nachgewiesen gestellt.

Numerik: „Expliziter Euler“. Sei $k \in \mathbb{N}$.

Setze $\Delta t^{(k)} = \frac{T}{k}$ und $t_i^{(k)} = i \cdot \Delta t^{(k)}$, $0 \leq i \leq k$.



Taylorentwicklung der Lösung von (1.9): Für $i > 0$

$$y(t_i^{(k)}) = y(t_{i-1}^{(k)} + \Delta t^{(k)}) = y(t_{i-1}^{(k)}) + \Delta t^{(k)} y'(t_{i-1}^{(k)}) + \frac{(\Delta t^{(k)})^2}{2} y''(\xi_i^{(k)}), \quad \xi_i^{(k)} \in [t_{i-1}^{(k)}, t_i^{(k)}]$$

(1.9) Einsetzen \downarrow

$$= y(t_{i-1}^{(k)}) + \Delta t^{(k)} f(t_{i-1}^{(k)}, y(t_{i-1}^{(k)})) + \frac{(\Delta t^{(k)})^2}{2} \underbrace{[2f + (2yf)' f]}_{\text{Kann man ausrechnen}}(\xi_i^{(k)}, y(\xi_i^{(k)}))$$

Umstellen: Für die Lösung $y(t)$ von (1.9) gilt für $0 < i \leq k$

$$\frac{1}{\Delta t^{(k)}} (y(t_i^{(k)}) - y(t_{i-1}^{(k)})) - f(t_{i-1}^{(k)}, y(t_{i-1}^{(k)})) = O(\Delta t^{(k)}) \quad (1.10)$$

$\rightarrow 0$ für $k \rightarrow \infty$

„konsistent von erster Ordnung“

Idee: Betrachte (1.10) mit rechter Seite 0 als Definitionsgleichung für $y_i^{(k)} \approx y(t_i^{(k)})$

Formal: $Y^{(k)} = \mathbb{R}^{k+1}$, $X^{(k)} = \mathbb{R}$, $Z^{(k)} = \mathbb{R}^{k+1}$

$G^{(k)}: Y^{(k)} \times X^{(k)} \rightarrow Z^{(k)}$ mit den einzelnen Komponenten:

$$G_i^{(k)}(Y_0^{(k)}, Y_1^{(k)}, \dots, Y_k^{(k)}, X) = \begin{cases} Y_0^{(k)} - X & i=0 \\ \frac{1}{\Delta t^{(k)}}(Y_i^{(k)} - Y_{i-1}^{(k)}) - f(t_{i-1}^{(k)}, Y_{i-1}^{(k)}) & 0 < i \leq k \end{cases} \quad (1.11)$$

Der Lösungsoperator $F^{(k)}: X^{(k)} \rightarrow Y^{(k)}$ lautet ("vorwärts einsetzen")

$$Y_i^{(k)} = F_i^{(k)}(X) = \begin{cases} X & i=0 \\ F_{i-1}^{(k)}(X) + \Delta t^{(k)} f(t_{i-1}^{(k)}, F_{i-1}^{(k)}(X)) & i>0 \end{cases} \quad \text{rekursive Definition}$$

Zur Konsistenz: $R^{(k)}: Y \rightarrow Y^{(k)}$ also $C^1([0, T]) \rightarrow \mathbb{R}^{k+1}$

mit $(R^{(k)} Y)_i = Y(t_i^{(k)})$ $0 \leq i \leq k$

Somit ist $G^{(k)}(R^{(k)} Y, X) = O(\Delta t^{(k)})$ wegen (1.10).

Damit ist das Verfahren konsistent nach Definition (1.12).

Zur Stabilität: Was ist H aus (1.8)? $H(z^{(k)}, X) = ?$

Für gegebenes $z \in \mathbb{R}^{k+1}$ löse $G^{(k)}(Y^{(k)}, X) = z^{(k)}$ nach $Y^{(k)}$ auf.

Nach (1.11)

$$\left. \begin{aligned} Y_0^{(k)} - X &= z_0^{(k)} \\ \frac{1}{\Delta t^{(k)}}(Y_i^{(k)} - Y_{i-1}^{(k)}) - f(t_{i-1}^{(k)}, Y_{i-1}^{(k)}) &= z_i^{(k)} \end{aligned} \right\} \Rightarrow z_i^{(k)} = H_i^{(k)}(z^{(k)}, X) = \begin{cases} X + z_0^{(k)} \\ H_{i-1}^{(k)}(z^{(k)}, X) \\ + \Delta t^{(k)} f(t_{i-1}^{(k)}, H_{i-1}^{(k)}(z^{(k)}, X)) \\ + \Delta t^{(k)} z_i^{(k)} \end{cases}$$

Man zeigt: $f(x, y)$ Lipschitzstetig in y ,
dann ist H stabil.

Und damit hat man Konvergenz. Weiteres in Numerik 1.

Kapitel 2: Fließkommazahlen

1
04.10.09

Im Computer gibt es verschiedene Typen zur Repräsentation von Zahlen. Etwa in C/C++:

unsigned int	\mathbb{N}_0
int	\mathbb{Z}
float	\mathbb{R}
double	\mathbb{R}
complex<double>	\mathbb{C}

int: exakt, aber endlicher Bereich

float, double, ...: Approximation, endlicher Bereich

Was hat das für Folgen?

Beispiel 2.1 (Potenzreihe für e^x). Auf Folie

□

2.1 Zahlendarstellung

Zellenwertsystem.

$$x = (\pm) \dots m_n \beta^n + \dots + m_1 \beta^1 + m_0 + m_{-1} \beta^{-1} + \dots + m_{-k} \beta^{-k} + \dots$$

\pm ist das Vorzeichen (ein Bit genügt).

$\beta \in \mathbb{N}, \beta \geq 2$ heißt Basis.

$m_i \in \{1, 2, \dots, \beta-1\}$ heißen Ziffern.

Jede reelle Zahl $x \in \mathbb{R}$ kann mit unendlich vielen Ziffern dargestellt werden.

Geschichte: Siehe [Knuth, Band 2, p. 194]

Babylonier 1750 v. Chr.: $\beta = 60$

Basis 10 in Europa ab ca 1585

Blaise Pascal: Jedes $\beta \geq 2$ möglich

Festkommazahlen:

2
04.10.09

Wähle maximalen und minimalen Exponenten:

$$x = (-1)^s \sum_{i=-k}^n m_i \beta^i$$

Problem: Wissenschaftl. Anwendungen brauchen Zahlen sehr unterschiedlicher Größe

Plancksches Wirkungsquantum: $6.6260693 \cdot 10^{-34} \text{ Js}$

Ruhemasse Elektron: $9.11 \cdot 10^{-28} \text{ g}$

Avogadrokonstante: $6.021415 \cdot 10^{23} \frac{1}{\text{mol}}$

$$\Rightarrow \beta = 2 \quad 2^{n+k} \approx 10^{23+34} \rightarrow n+k = \frac{57}{\log_2} \approx \underline{\underline{190 \text{ Stellen!}}}$$

Fließkommazahlen erlauben effizientere Darstellung unterschiedlich großer Zahlen.

Definition 2.2 (normierte Fließkommazahlen)

Sei $\beta, r, s \in \mathbb{N}$ und $\beta \geq 2$. $\mathbb{F}(\beta, r, s) \subset \mathbb{R}$ besteht aus den Zahlen mit folgenden Eigenschaften:

a) $\forall x \in \mathbb{F}(\beta, r, s)$ gilt $x = m(x) \beta^{e(x)}$ mit

$$m(x) = \pm \sum_{i=1}^r m_i \beta^{-i}, \quad e(x) = \pm \sum_{j=0}^{s-1} e_j \beta^j$$

m heißt Mantisse, e Exponent.

b) $\forall x \in \mathbb{F}(\beta, r, s)$ gilt $x = 0 \vee m_1 \neq 0$.

$m_1 \neq 0$ heißt Normierung. Bedingung b) macht die Fließkommadarstellung eindeutig. □

Ist $x \in \mathbb{F}(\beta, r, s)$ und $x \neq 0$ dann gilt wegen b):

$$\beta^{-1} \leq |m(x)| < 1 \quad \text{und damit} \quad \beta^{e(x)-1} \leq |x| < \beta^{e(x)}. \quad (2.1)$$

Beispiel 2.3

a) $F(10, 3, 1)$ besteht aus Zahlen der Form:

$$x = \pm (m_1 \cdot 0.1 + m_2 \cdot 0.01 + m_3 \cdot 0.001) \cdot 10^{\pm e_0}$$

mit $m_i \neq 0 \vee (m_1 = m_2 = m_3 = 0)$

z.B. $0.999 \cdot 10^7, 0.123 \cdot 10^{-1}, 0$, aber

$0,000\,000\,000\,000\,014 = 0.140 \cdot 10^{-10} \notin F(10, 3, 1)$ da Exponent zu klein.

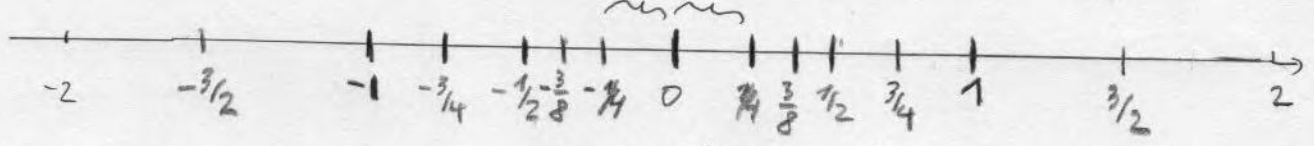
b) $F(2, 2, 1)$ besteht aus Zahlen der Form

$$x = \pm \left(m_1 \cdot \frac{1}{2} + m_2 \cdot \frac{1}{4} \right) \cdot 2^{\pm e_0}$$

$$m_i = 1 \quad \left\{ 0, \frac{1}{2}, \frac{3}{4} \right\} \quad \{2^{-1}, 2^0, 2^1\} = \left\{ \frac{1}{2}, 1, 2 \right\}$$

$$\Rightarrow F(2, 2, 1) = \left\{ \underbrace{-\frac{3}{2}, -1, -\frac{3}{4}, -\frac{1}{2}}_{\left\{ -\frac{3}{4}, -\frac{1}{2} \right\} \cdot 2}, \underbrace{-\frac{3}{8}, -\frac{1}{4}, 0, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}}_{\left\{ -\frac{3}{4}, -\frac{1}{2} \right\} \cdot 1}, \underbrace{\frac{3}{4}, 1, 2}_{\left\{ -\frac{3}{4}, -\frac{1}{2} \right\} \cdot \frac{1}{2}} \right\}$$

keine wg. Normalisierung.



Mögliche Lösung: Nehme für $|x| < \beta^{-1}$ die nicht normalisierten Zahlen hinzu.

Zahlenbereich:

Größte / kleinste darstellbare Zahl in $F(\beta, r, s)$

$$X_{+/-} = \pm \underbrace{(\beta-1) \{ \beta^{-1} + \dots + \beta^{-r} \}}_{= 1 - \beta^{-r}} \cdot \beta^{\underbrace{(\beta-1) \{ \beta^{s-1} + \dots + \beta^0 \}}_{= \beta^s - 1}}$$

← Wertigkeit der letzten Stelle

$$= \pm (1 - \beta^{-r}) \beta^{\beta^s - 1}$$

Kleinste positive / größte negative Zahl in $F(\beta, r, s)$:

$$x_{+/-} = \pm \underbrace{\beta^{-1}}_{\text{kleinste Mantisse wg. Normierung}} \cdot \beta^{\underbrace{-1}_{\text{negativ!}}} \cdot \underbrace{(\beta-1) \{ \beta^{s-1} + \dots + \beta^0 \}}_{\beta^s - 1} = \pm \beta^{\beta^s}$$

Damit

$$F(\beta, r, s) \subset D(\beta, r, s) = [X_-, x_-] \cup \{0\} \cup [x_+, X_+] \subset \mathbb{R}.$$

Abstände zwischen Fließkommazahlen

04.10.09

Festkommazahlen: Abstand zwischen Zahlen konstant β^{-k}

Fließkommazahlen: Abstand von $x \in \mathbb{F}$ zum nächsten $x' \in \mathbb{F}$ hängt von x ab.

(Vorzeichen spielt keine Rolle)

$$\beta^{e-1} = 0.\overset{m_1}{1}0\dots\overset{m_r}{0} \cdot \beta^e$$

$$\text{nächste Zahl: } 0.10\dots1 \cdot \beta^e \quad \left. \vphantom{\beta^{e-1}} \right\} \text{Abstand } \beta^{-r} \cdot \beta^e = \underline{\beta^{e-r}}$$

$$\text{nächste Zahl } +\beta^{e-r} \quad 0.\beta\beta\dots\beta \cdot \beta^e \quad \left. \vphantom{\beta^{e-1}} \right\} \text{Abstand } \beta^{e-r} \text{ (da Exponent immer noch gleich)}$$

$$\beta^e = 0.10\dots0 \cdot \beta^{e+1} \quad \left. \vphantom{\beta^{e-1}} \right\} \text{Abstand } \beta^{e-r} \quad \downarrow \text{Abstand springt!}$$

$$\text{nächste Zahl: } 0.10\dots1 \cdot \beta^{e+1} \quad \left. \vphantom{\beta^{e-1}} \right\} \text{Abstand } \beta^{-r} \cdot \beta^{e+1} = \underline{\beta^{e+1-r}}$$

\Rightarrow Im Intervall $[\beta^{e-1}, \beta^e]$ beträgt der absolute Abstand zwischen zwei Zahlen β^{e-r} .

Sei $x' \in \mathbb{F}$ die nächstliegende Fließkommazahl zu $x \in \mathbb{F}$. Dann gilt

$$|x-x'| = \beta^{e(x)-r} = \frac{|m(x)| \beta^{e(x)}}{|m(x)|} \cdot \beta^{-r} = \frac{|x|}{|m(x)|} \beta^{-r}$$

\uparrow
 $x \neq 0$

Da $\beta^{-1} \leq |m(x)| < 1$ folgt β^{-1} groß, wenn $|m(x)|$ klein, d.h. β^{-1}

$$|x| \beta^{-1} \beta^{1-r} < |x-x'| \leq |x| \beta^{1-r}$$

Für den relativen Abstand gilt:

$$\beta^{-1} \beta^{1-r} < \frac{|x-x'|}{|x|} \leq \beta^{1-r}$$

Er schwankt um den Faktor β^{-1} je nach Größe von $|x|$. Dieser Faktor heißt wobble.

Konsequenz: Kleine β sind besser, β^{-1} dann kleiner ist!

Übung:

β^{1-r} : Die kleinste Zahl, $x \in \mathbb{F}$ das $1+x \neq 1$.

Schreibe Programm, das diese x bestimmt.

Teste das für float, double.

Ziel: Portabilität von Programmen mit Fließkommazahlenarithmetik.
Verabschiedet 1985.

$\beta = 2$ mit vier Genauigkeitsstufen und normierter Darstellung:

	Format			
	single	single-act	double	double-ext
e_{\min}	+127	≥ 1024	1023	≥ 16384
e_{\max}	-126	≤ -1021	-1022	≤ -16381
Bits Expon.	8	≤ 11	11	15
Bits total	32	≥ 43	64	≥ 79

Betrachte double genauer:

- total 64 bit
- 11 bit für Exponent. Dieser wird vorzeichenlos als Zahl $c \in [1, 2046]$ gespeichert.

$$\text{Setze } e = c - 1023 \Rightarrow e \in \left[\underset{2046-1023}{-1022}, \underset{1-1023}{1023} \right]$$

Die Werte $c \in \{0, 2047\}$ werden anderweitig genutzt:

$c = 0 \wedge m = 0$ kodiert die Null

$c = 2047 \wedge m \neq 0$ kodiert NaN = "not a number"

$c = 2047 \wedge m = 0$ kodiert ∞ (Überlauf). (z.B. Division durch Null)

- $64 - 11 = 53$ Bit Mantisse, 1 Bit für Vorzeichen, bleiben 52 Nachkommastellen.
Wg $\beta = 2$ und Normierung ist immer $m_1 = 1$ und diese Ziffer wird nicht gespeichert.
Heißt hidden bit.

Somit gilt $r = 53$.

2.2 Runden und Rundungsfehler

6
04.10.09

Um $x \in \mathbb{R}$ im $\mathbb{F}(\beta, r, s)$ zu approximieren brauchen wir

$$\text{rd} : \mathcal{D}(\beta, r, s) \rightarrow \mathbb{F}(\beta, r, s) \quad (2.2)$$

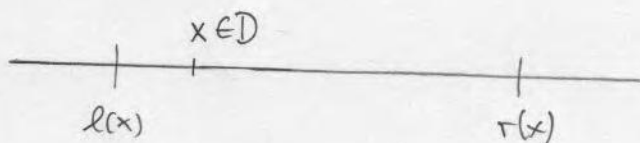
Achtung: rd setzt voraus, dass x im darstellbaren Bereich liegt!

Im Falle eines Über-/Unterlaufes ist r, s zu ändern.

Sinnvollerweise soll für rd gelten:

$$|x - \text{rd}(x)| \leq \min_{y \in \mathbb{F}} |x - y| \quad \forall x \in \mathcal{D} \quad (\text{Bestapproximation}).$$

Mit



$$l(x) = \max\{y \in \mathbb{F} \mid y \leq x\} \quad r(x) = \min\{y \in \mathbb{F} \mid y \geq x\}$$

gilt dann

$$\text{rd}(x) = \begin{cases} l(x) & |x - l(x)| < |x - r(x)| \\ r(x) & |x - l(x)| > |x - r(x)| \\ ? & |x - l(x)| = |x - r(x)| \end{cases}$$

Im letzten Fall ist eine Rundung erforderlich. Dafür gibt es verschiedene Möglichkeiten.

Sei $x = \text{sign}(x) \cdot \left(\sum_{i=1}^{\infty} m_i \beta^{-i} \right) \beta^e$ die normierte Darstellung von $x \in \mathcal{D}$.

Natürliche Rundung (das was jeder kennt)

$$\text{rd}(x) = \begin{cases} l(x) = \text{sign}(x) \left(\sum_{i=1}^r m_i \beta^{-i} \right) \beta^e \\ r(x) = l(x) + \beta^{e-r} \end{cases}$$

Wichtigkeit
der letzten Stelle

falls $0 \leq m_{r+1} < \beta/2$ nächste Stelle

falls $\beta/2 \leq m_{r+1} < \beta$

Gerade Rundung (β sei gerade)

04.10.09

$$rd(x) = \begin{cases} l(x) & (|x-l(x)| < |x-r(x)|) \vee \\ & (|x-l(x)| = |x-r(x)| \wedge m_r \text{ gerade}) \\ r(x) = l(x) + \beta^{e-r} & \text{sonst.} \end{cases}$$

Damit ist $m_r rd(x)$ immer gerade, wenn gerundet wurde (d.h. $\frac{|x-l(x)|}{|x-r(x)|} =$

- ist $rd(x) = l(x)$ so ist das per Def so

- sonst ist $rd(x) = \underbrace{l(x)}_{m_r \text{ unger.}} + \underbrace{\beta^{e-r}}_{+1 \text{ in der letzten Stelle.}} \Rightarrow m_r \text{ gerade.}$

Diese Wahl vermeidet eine mögliche Drift beim Aufrunden.

Definition 2.4 (absoluter und relativer Fehler)

Sei $x' \in \mathbb{R}$ eine Näherung von $x \in \mathbb{R}$. Dann heißt

$$\Delta x = x' - x \quad \text{absoluter Fehler}$$

und für $x \neq 0$

$$\epsilon_{x'} = \frac{\Delta x}{x} \quad \text{relativer Fehler.} \quad \square$$

Umformen liefert:

$$x' = x + \Delta x = x \left(1 + \frac{\Delta x}{x}\right) = x (1 + \epsilon_{x'})$$

\uparrow
abs. Fehler

Motivation:

Es sei $\Delta x = x' - x = 100 \text{ km}$.

Für $x = \text{Entfernung Erde-Sonne} \approx 1.5 \cdot 10^8 \text{ km}$ ist

$$\epsilon_{x'} = \frac{10^2 \text{ km}}{1.5 \cdot 10^8 \text{ km}} \approx 6.6 \cdot 10^{-7}$$

relativ klein. Für $x = \text{Entfernung Heidelberg-Paris} \approx 500 \text{ km}$ ist

$$\epsilon_{x'} = \frac{100 \text{ km}}{500 \text{ km}} = 0.2$$

dagegen relativ groß.

Lemma 2.5 (Rundungsfehler)

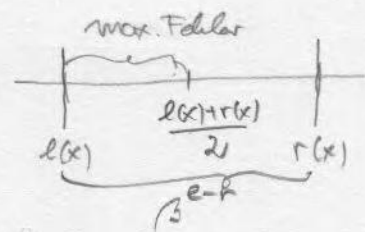
04.10.09

Bei der Rundung in $\mathbb{F}(\beta, r, s)$ gilt für den absoluten Fehler

$$|x - rd(x)| \leq \frac{1}{2} \beta^{e(x)-r} \quad (2.2)$$

und für den relativen Fehler ($x \neq 0$)

$$\frac{|x - rd(x)|}{|x|} \leq \frac{1}{2} \beta^{1-r}$$



Beweis: Es ist $x = m(x) \beta^{e(x)}$ die normierte Darstellung von x .

Wie oben gezeigt gilt

$$|r(x) - l(x)| = \beta^{e(x)-r}$$

zwei aufeinanderfolgende Zahlen in \mathbb{F} , $l(x) \in [\beta^{e(x)-1}, \beta^{e(x)}]$

Der maximale Fehler ergibt sich für $x = \frac{l(x)+r(x)}{2}$, also

$$|x - rd(x)| \leq \left| \frac{l(x)+r(x)}{2} - l(x) \right| = \frac{1}{2} |r(x) - l(x)| = \frac{1}{2} \beta^{e(x)-r}$$

↑
egal

Für den relativen Fehler ($x \neq 0$) gilt

$$\frac{|x - rd(x)|}{|x|} \leq \frac{\frac{1}{2} \beta^{e(x)-r}}{|m(x)| \beta^{e(x)}} = \frac{1}{2} \frac{1}{|m(x)|} \beta^{-r} \leq \frac{1}{2} \beta^{1-r}$$

$|m(x)| \geq \beta^{-1}$

Die Zahl $\text{eps} := \frac{1}{2} \beta^{1-r}$ heißt Maschinengenauigkeit. (Nimmt Rannacher auch so).
Wird oft auch mit ϵ (z.B. in MATLAB) abgekürzt.

Vorricht: Die Bezeichnungen gehen in der Literatur durcheinander.

[Goldberg] nennt $\frac{1}{2} \beta^{1-r}$ machine epsilon.

[Quarteroni] nennt β^{1-r} machine epsilon und

$\frac{1}{2} \beta^{1-r}$ machine precision.

2.3 Fließkommaarithmetik

05.10.09

Wir benötigen eine Arithmetik auf \mathbb{F} :

$$\otimes : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F} \quad \text{mit } \otimes \in \{\oplus, \ominus, \odot, \oslash\},$$

die den bekannten Operationen $*$ \in $\{+, -, \cdot, /$ auf \mathbb{R} entsprechen.

Problem: $x, y \in \mathbb{F} \Rightarrow x \otimes y \notin \mathbb{F}$ in der Regel!

Somit ist das Ergebnis $x \otimes y$ wieder zu runden, d. h. wir definieren

$$x \otimes y = \text{rd}(x * y) \quad \forall x, y \in \mathbb{F}. \quad (2.3)$$

Man sagt \otimes ist „exakt gerundet“. Dies ist nicht trivial!
(implizite Annahme: $x * y \in \mathbb{D}$!)

Beispiel 2.6 (Guard digit)

Sei $\mathbb{F} = \mathbb{F}(10, 3, 1)$, $x = 0.215 \cdot 10^8$, $y = 0.125 \cdot 10^{-5}$. Wir betrachten die Subtraktion $x \ominus y = \text{rd}(x - y)$.

① $x - y$:

$x =$	0.215	000	000	000	000	0	$\cdot 10^8$
$y =$	0.000	000	000	000	0	125	$\cdot 10^8$ ←
			13 Nullen!			1 1 1	← Schiebe y auf gleichen Exponenten.
$x - y =$	0.214	999	999	999	9	875	$\cdot 10^8$

$y = 0.125 \cdot 10^{-5} = 0.125 \cdot 10^{-13} \cdot 10^8 = 10^{-5}$

② Runden

$$x \ominus y = \text{rd}(0.2149 \dots \cdot 10^8) = 0.215$$

Problem: Schritt 1 erfordert extrem hohe Stellenzahl $\mathcal{O}(\beta^5)$!

Geht es einfacher? Z. B. Runde y schon nach dem Schieben.

Dies liefert im obigen Fall das gleiche Ergebnis.

Ü: Tabellenmachardilemma. Tabellieren von transzendenten Funktionen.

$\mathbb{F} = \mathbb{F}(10, 4, 1)$, erstelle Liste $y = \exp(x) \forall x \in \mathbb{F}$ mit $y = \text{rd}(\exp(x))$

Problem: $\exp(1.626) = \underbrace{5.0835}_{5 \text{ Stellen}}$

ich verstehe es nicht,
wann 5 eine gültige Ziffer ist,
dann muss man aufrunden

Im allgemeinen ist das gefährlich:

$$\begin{array}{lcl} x = 0.101 \cdot 10^1 & \longrightarrow & 0.101 \cdot 10^1 \\ y = 0.993 \cdot 10^0 & \xrightarrow{\text{schreiben}} & \frac{0.099 \cdot 10^1}{0.002 \cdot 10^1} \end{array} \quad \begin{array}{l} \\ \\ \text{Schreiben und} \\ \text{Runden} \end{array}$$

relativer Fehler im Ergebnis:

$$\frac{(x \ominus y) - (x - y)}{x - y} = \frac{0.02 - 0.017}{0.017} \approx 0.176 \approx 35 \text{ eps}$$

$$\text{mit } \text{eps} = \frac{1}{2} 10^{1-3} = 0.005.$$

→ Fehler ist 35 mal größer als erwartet.

○ Eine Stelle mehr im Addierer (also $r+1$) liefert das exakte Ergebnis!

Mit einer zusätzlichen Stelle erreicht man

$$\frac{(x \ominus y) - (x - y)}{x - y} \leq 2 \text{ eps.}$$

Mit zwei Stellen erreicht man exakte Rundung!

Die zusätzlichen Stellen nennt man „guard digits“.

Zusätzliche Probleme bei der Arithmetik $\oplus: \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$:

- - Assoziativ und Distributivgesetz gelten nicht, es kommt also auf die Reihenfolge der Operationen an!
- $\exists y \in \mathbb{F}$ so dass $x \oplus y = x$
- Allerdings gilt das Kommutativgesetz!
- Es gelten auch folg. einfache Regeln:
 $(-x) \oplus y = -(x \oplus y), \quad 1 \oplus x = x, \quad x \oplus y = 0 \Rightarrow x = 0 \vee y = 0,$
 $x \oplus z \leq y \oplus z \quad \text{falls } x \leq y \text{ und } z > 0.$

2.4 Fehleranalyse

Fortpflanzung von Rundungsfehlern in Rechnungen.

Wie in Kap. 1 betrachten wir die Funktionsauswertung.

- Sei $F: \mathbb{R}^m \rightarrow \mathbb{R}^n$, in Komponenten

$$F(x) = \begin{pmatrix} F_1(x_1, \dots, x_m) \\ \vdots \\ F_n(x_1, \dots, x_m) \end{pmatrix}$$

- Zur Berechnung von F im Rechner nutze numerische Realisierung

$F': \mathbb{F}^m \rightarrow \mathbb{F}^n$, F' wird durch einen Algorithmus realisiert, d.h. aus

-- endlich vielen (= Terminierung)

-- elementaren (= bekannte) Rechenoperationen

Zusammengesetzt: $F'(x) = \varphi_2(\dots \varphi_2(\varphi_1(x)) \dots)$.

Wichtig: i) Zu einem F gibt es i.d.R. viele Realisierungen im Sinne unterschiedlicher Reihenfolgen

$$a + b + c \approx (a \oplus b) \oplus c \neq a \oplus (b \oplus c)!$$

ii) Jedes φ_i steuert einen (unbekannten) Fehler bei

iii) Im Prinzip kann die Rechengenauigkeit beliebig gesteigert werden, d.h. eigentlich Folge $F^{(k)}, (F^{(k)})^m \rightarrow (F^{(k)})^n$.

Das machen wir aber nicht so formal.

Wie in Kap. 1 nutze Aufspaltung:

$$F(x) - F'(rd(x)) = \underbrace{F(x) - F(rd(x))}_{\text{① Konditionsanalyse von } F} + \underbrace{F(rd(x)) - F'(rd(x))}_{\text{② Rundungsfehleranalyse: Unterschied } F, F' \text{ bei gleicher Eingabe (aus } F \text{)}} \quad (2.4)$$

$\underbrace{F(x)}_{\text{exaktes Ergebnis}} \xrightarrow{F \in \mathbb{R}} \underbrace{F'(rd(x))}_{\substack{\text{Eingabe runden} \\ \text{numerische Auswertung}}}$

von hier aus:

- Analyse „in erster Näherung“
- absolute / relative Fehler.
- Normen bilden, das lassen wir aber i.d.R. weg.

① Differentielle Konditionsanalyse

Wir nehmen an, dass $F : \mathbb{R}^m \rightarrow \mathbb{R}^n$ zweimal stetig differenzierbar.

Nach dem Satz von Taylor gilt für die F_i :

$$F_i(x + \Delta x) = F_i(x) + \sum_{j=1}^m \frac{\partial F_i}{\partial x_j}(x) \Delta x_j + R_i^F(x; \Delta x) \quad i=1, \dots, n.$$

Für das Restglied gilt (unter diesen Voraussetzungen)

$$R_i^F(x; \Delta x) = O(\|\Delta x\|^2).$$

Definition 2.7 (Landausche Symbole)

Man schreibt:

$$g(t) = O(h(t)) \quad (t \rightarrow 0)$$

falls es $t_0 > 0$ und $c_0 > 0$ gibt so dass für alle $t \in (0, t_0]$ die Abschätzung

$$|g(t)| \leq c |h(t)|$$

gilt. Sprechweise: „ $g(t)$ geht wie $h(t)$ gegen 0“. Man will also quantifizieren „wie schnell“ eine Funktion (mindestens) gegen 0 geht.

Weiter bedeutet

$$g(t) = o(h(t)) \quad (t \rightarrow 0),$$

dass es $t_0 > 0$ und eine Funktion $c(t)$, $\lim_{t \rightarrow 0} c(t) = 0$ gibt, sodass für alle $t \in (0, t_0]$ gilt

$$|g(t)| \leq c(t) |h(t)|.$$

Bedeutung: „ $g(t)$ geht schneller als $h(t)$ gegen Null“ (falls $h(t) \rightarrow 0$). ■

Somit können wir die Taylorformel umformen:

$$F_i(x + \Delta x) - F_i(x) = \underbrace{\sum_{j=1}^m \frac{\partial F_i}{\partial x_j}(x) \Delta x_j}_{\text{"führende (erste) Ordnung"}} + \underbrace{O(\|\Delta x\|^2)}_{\text{höhere Ordnung}}$$

Oft lässt man die Terme höherer Ordnung weg und schreibt " \doteq " statt " $=$ ". Sprechweise: "ist in erster Näherung gleich".

Nun gehen wir zu relativen Fehlern über. Sei $F_i(x) \neq 0$ und $x_j \neq 0$

$$\begin{aligned} \frac{F_i(x + \Delta x) - F_i(x)}{F_i(x)} &\doteq \sum_{j=1}^m \frac{\partial F_i}{\partial x_j}(x) \frac{\Delta x_j}{F_i(x)} \\ &\doteq \sum_{j=1}^m \underbrace{\left(\frac{\partial F_i}{\partial x_j}(x) \frac{x_j}{F_i(x)} \right)}_{\text{Verstärkungsfaktor } \bar{k}_{ij}(x)} \left(\frac{\Delta x_j}{x_j} \right) \end{aligned} \quad (2.5)$$

↑ relativer Eingabefehler.

Fasst man die Verstärkungsfaktoren in einer Matrix zusammen

$$\left(\bar{K}(x) \right)_{ij} = \bar{k}_{ij}(x)$$

Zusammen so kann man (für geeignete Normen) zeigen

$$\|\bar{K}(x)\| \leq K(x)$$

mit $K(x)$ der relativen Konditionszahl aus Kapitel 1.

U: Setze für $x = (x_1, \dots, x_m)^T$, $x_i \neq 0$ " $\frac{\Delta x}{x}$ " = $\left(\frac{\Delta x_1}{x_1}, \dots, \frac{\Delta x_m}{x_m} \right)^T$

$$\frac{\|\Delta x\|_\infty}{\|x\|_\infty} \leq \left\| \frac{\Delta x}{x} \right\|_\infty \leq \frac{\|\Delta x\|_\infty}{\min_i x_i}$$

$$\Rightarrow \frac{\|F(x + \Delta x) - F(x)\|_\infty}{\|F(x)\|_\infty} \leq \frac{\|\Delta x\|_\infty}{\|x\|_\infty} = K(x)$$

Definition 2.8 Wir nennen die Auswertung $y = F(x)$ „schlecht konditioniert“,
 im Punkt x , falls $|k_{ij}(x)| \gg 1$, andernfalls „gut konditioniert“.
 $|k_{ij}(x)| < 1$ heißt Fehlerdämpfung, $|k_{ij}(x)| > 1$ Fehlerverstärkung. \square

Warum relative Kondition?

Wegen Lemma 2.5 gilt

$$\left| \frac{x - rd(x)}{x} \right| \leq \epsilon_{ps} = \frac{1}{2} \beta^{1-r}$$

d.h. es gibt $\epsilon \in \mathbb{R}$, $|\epsilon| \leq \epsilon_{ps}$, sodass

$$\frac{x - rd(x)}{x} = \epsilon \Leftrightarrow x - rd(x) = \epsilon x \Leftrightarrow rd(x) = x + \underbrace{\epsilon x}_{=: \Delta x}$$

d.h. für die relativen Eingabefehler in (2.5) gilt gerade

$$F(x) - F(rd(x)) \frac{\Delta x_j}{x_j} = \frac{\epsilon_j \cdot x_j}{x_j} = \epsilon_j.$$

Beispiel 2.9

Ü: Kondition von $F(x_1, x_2) = x \cdot y$, x/y , $F(x) = \sqrt{x}$

a) Addition. $F(x_1, x_2) = x_1 + x_2$ $\frac{\partial F}{\partial x_1} = 1$, $\frac{\partial F}{\partial x_2} = 1$. Nach obiger Formel (2.5):

$$\frac{F(x_1 + \Delta x_1, x_2 + \Delta x_2) - F(x_1, x_2)}{F(x_1, x_2)} = 1 \cdot \frac{x_1}{x_1 + x_2} \underbrace{\frac{\Delta x_1}{x_1}}_{1.1 \leq \epsilon_{ps}} + 1 \cdot \frac{x_2}{x_1 + x_2} \underbrace{\frac{\Delta x_2}{x_2}}_{1.1 \leq \epsilon_{ps}}$$

\Rightarrow Für $x_1 \rightarrow -x_2$ gehen beide Verstärkungsfaktoren gegen ∞ (sofern $|x_1|, |x_2| > \delta$).
 Schlecht konditioniert!

b) $F(x_1, x_2) = x_1^2 - x_2^2$, $\frac{\partial F}{\partial x_1} = 2x_1$, $\frac{\partial F}{\partial x_2} = -2x_2$.

$$\frac{F(x_1 + \Delta x_1, x_2 + \Delta x_2) - F(x_1, x_2)}{F(x_1, x_2)} = \underbrace{2x_1 \cdot \frac{x_1}{x_1^2 - x_2^2}}_{k_1 = 2 \frac{x_1^2}{x_1^2 - x_2^2}} \underbrace{\frac{\Delta x_1}{x_1}}_{1.1 \leq \epsilon_{ps}} - \underbrace{2x_2 \cdot \frac{x_2}{x_1^2 - x_2^2}}_{k_2 = 2 \frac{x_2^2}{x_2^2 - x_1^2}} \underbrace{\frac{\Delta x_2}{x_2}}_{1.1 \leq \epsilon_{ps}}$$

Schlecht konditioniert für $|x_1| \approx |x_2|$.

② Rundungsfehleranalyse

d.h. Nach (2.4): $F(x) - F'(x)$ mit $x \in \mathbb{F}^m$ Maschinenzahl.

F' „zusammengesetzt“ aus Einzeloperationen $\otimes \in \{+, \ominus, \odot, \oslash\}$.

Wegen (2.3) (exakt gerundete Arithmetik) und Lemma 2.5 gilt

$$\frac{(x \otimes y) - (x * y)}{(x * y)} = \epsilon \quad \text{mit } |\epsilon| \leq \text{eps.}$$

Vorsicht: ϵ ist abhängig von x und y , d.h. für jede Operation verschieden!

Und damit

$$x \otimes y = (x * y) (1 + \epsilon) \quad \text{für ein } |\epsilon(x, y)| \leq \text{eps.}$$

Analyse, in erster Näherung

Beispiel 2.10 $F(x_1, x_2) = x_1^2 - x_2^2$ mit zwei Realisierungen

$$F_a(x_1, x_2) = (x_1 \odot x_2) \ominus (x_2 \odot x_2)$$

$$F_b(x_1, x_2) = (x_1 \ominus x_2) \odot (x_1 \oplus x_2)$$

a) $u = x_1 \odot x_1 = (x_1 \cdot x_1) (1 + \epsilon_1)$
 $v = x_2 \odot x_2 = (x_2 \cdot x_2) (1 + \epsilon_2)$ $\epsilon_1 \neq \epsilon_2$ aber $|\epsilon_i| \leq \text{eps}!$

$$F_a(x_1, x_2) = u \ominus v = (u - v) (1 + \epsilon_3)$$

$$= (x_1^2 (1 + \epsilon_1) - x_2^2 (1 + \epsilon_2)) (1 + \epsilon_3)$$

$$= (x_1^2 + \epsilon_1 x_1^2 - x_2^2 - \epsilon_2 x_2^2) (1 + \epsilon_3)$$

$$= \underbrace{x_1^2 - x_2^2}_{= F(x_1, x_2)} + \underbrace{\epsilon_1 x_1^2 - \epsilon_2 x_2^2 + \epsilon_3 x_1^2 - \epsilon_3 x_2^2}_{\text{erste Ordnung}} + \underbrace{\epsilon_1 \epsilon_3 x_1^2 - \epsilon_2 \epsilon_3 x_2^2}_{\text{zweite Ordnung}}$$

relativer Fehler

$$\frac{F_a(x_1, x_2) - F(x_1, x_2)}{F(x_1, x_2)} \stackrel{\text{in erster Näherung}}{=} \frac{x_1^2}{x_1^2 - x_2^2} (\epsilon_1 + \epsilon_3) + \frac{x_2^2}{x_2^2 - x_1^2} (\epsilon_2 + \epsilon_3)$$

Bemerkung! gleiche Verstärkungsfaktoren wie in Bsp. 2.9!

b) $u = x_1 \ominus x_2 = (x_1 - x_2)(1 + \epsilon_1)$

$v = x_1 \oplus x_2 = (x_1 + x_2)(1 + \epsilon_2)$

$$\begin{aligned}
F(x_1, x_2) &= u \odot v = (u \cdot v)(1 + \epsilon_3) \\
&= \left((x_1 - x_2)(1 + \epsilon_1) \cdot (x_1 + x_2)(1 + \epsilon_2) \right) (1 + \epsilon_3) \\
&= \underbrace{(x_1 - x_2)(x_1 + x_2)} \cdot \underbrace{(1 + \epsilon_1)(1 + \epsilon_2)(1 + \epsilon_3)} \\
&= (x_1^2 - x_2^2) \left(1 + \epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_1\epsilon_2 + \dots + \epsilon_1\epsilon_2\epsilon_3 \right)
\end{aligned}$$

$$\frac{F_e(x_1, x_2) - F(x_1, x_2)}{F(x_1, x_2)} \stackrel{!}{=} \frac{x_1^2 - x_2^2}{x_1^2 - x_2^2} (\epsilon_1 + \epsilon_2 + \epsilon_3) = \epsilon_1 + \epsilon_2 + \epsilon_3$$

⇒ Verstärkungsfaktor 1, besser als vorher.

Definition 2.11

Wir nennen einen numerischen Algorithmus „numerisch stabil“, wenn die im Lauf der Rechnung akkumulierten Rundungsfehler aus ② den unvermeidbaren Problemfehler aus der Konditionsanalyse ① nicht übersteigen.

N.a.w. Verstärkungsfaktoren aus Rundungsfehleranalyse ≤ denen aus Konditionsanalyse ⇒ „numerisch stabil“

Beide Realisierungen a, b aus Bsp. 2.10 sind numerisch stabil.

2.5 Auslöschung

17
06.10.09

Obrige Beispiele 2.9, 2.10 enthalten das Phänomen der Auslöschung.

Dies tritt auf bei

- Addition $x_1 + x_2$ mit $x_1 \approx -x_2$,
- Subtraktion $x_1 - x_2$ mit $x_1 \approx x_2$.

Bemerkung 2.12 Bei der Auslöschung werden vor der antipredanda Addition bzw. Subtraktion eingeführte Fehler extra verstärkt. \square

Sind $x_1, x_2 \in \mathbb{F}$ Maximanzahlen, so gilt wie oben gezeigt:

$$\left| \frac{(x_1 \oplus x_2) - (x_1 - x_2)}{x_1 - x_2} \right| \leq \text{eps.}$$

Also kein Problem. Das Problem tritt erst ein, wenn x_1, x_2 selbst schon mit Fehlern behaftet sind.

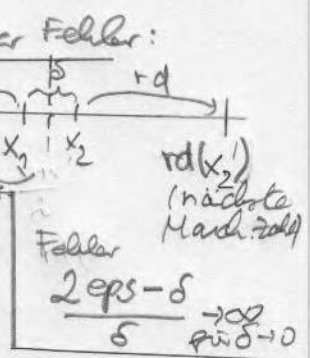
Beispiel 2.13 $\mathbb{F} = \mathbb{F}(10, 4, 1)$

$$x_1 = 0.11258762 \cdot 10^2 \rightarrow \text{rd}(x_1) = 0.1126 \cdot 10^2$$

$$x_2 = 0.11244891 \cdot 10^2 \rightarrow \text{rd}(x_2) = 0.1124 \cdot 10^2$$

$$\begin{aligned} x_1 \ominus x_2 &= 0.00013871 \cdot 10^2 \\ &= 0.13871 \cdot 10^{-1} \end{aligned}$$

$$\begin{aligned} \text{rd}(x_1) - \text{rd}(x_2) &= 0.0002 \cdot 10^2 \\ &= 0.2 \cdot 10^{-1} \end{aligned}$$



relativer Fehler:

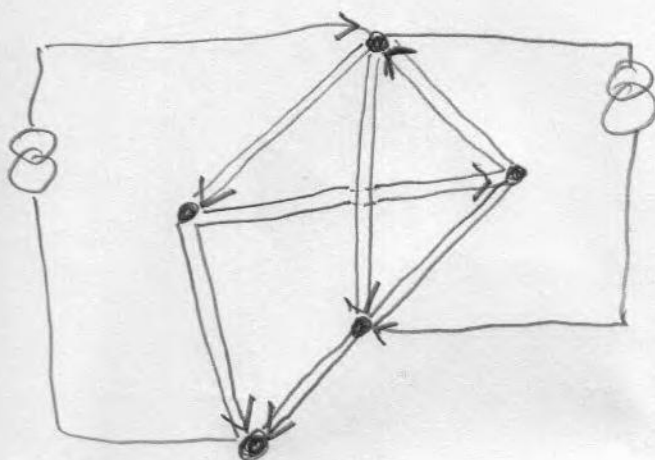
$$\frac{0.2 \cdot 10^{-1} - 0.13871 \cdot 10^{-1}}{0.13871 \cdot 10^{-1}} \approx 0.44 \approx 44\% \cdot \frac{1}{2} 10^{-3} = \text{eps} ! \quad \square$$

Hier: Rundung der Eingangsgrößen.

Ursprung der Fehler ist egal, tritt ebenso auf, wenn x_1, x_2 mit Fehlern vorhergehender Rechenschritte behaftet sind.

Regel 2.14 Setze potentiell gefährliche Operationen möglichst früh im Algorithmus ein. (Siehe Beispiel 2.10). \square

3.1 Strömung in Rohrleitungsketzten



1) Netzwerk von Röhren beschrieben durch gerichteten Graphen.

- Knotenmenge $V = \{v_1, \dots, v_n\}$, $|V| = n$.

- Kantenmenge $E = \{e_1, \dots, e_M\}$, $|E| = M$,

$E \subseteq V \times V$, mit: $(v, w) \in E \Rightarrow (w, v) \notin E$

- $E = E_R \cup E_P$ „Röhren und Pumpen“, $|E_R| = m$,

$E_R = \{e_1, \dots, e_m\}$, $E_P = \{e_{m+1}, \dots, e_M\}$.

← eigentlich hier!

2) Gesetz von Hagen-Poiseuille.

Röhre $e = (v, w)$ (von Knoten v nach Knoten w)

$$(3.1) \quad q_e = \frac{\pi r_e^4}{8\eta d_e} \Delta p_e$$

$= i_e$ „Leitfähigkeit“

r_e : Radius der Röhre e [m]

d_e : Länge der Röhre e [m]

η : dyn. Viskosität Flüssigk. [Pas]

q_e : Volumenstrom [m³/s]

Orientierung: $q_e > 0$ falls Fluss von Knoten v nach w .

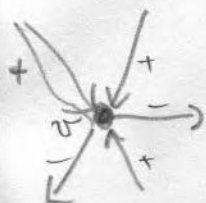
$q_e < 0$ Fluss von w nach v

Δp_e : gerichtete Druckdifferenz über Rohr e [Pa] = [N/m²]

$\Delta p_e > 0 \rightarrow q_e > 0$

3) Knotenregel (1. Kirchhoffsches Gesetz)

07.10.09



$$E_v^+ = \{ (u, w) \in E \mid u = v \} \text{ "ausgehend"}$$

$$E_v^- = \{ (u, w) \in E \mid w = v \} \text{ "eingehend"}$$

$$\sum_{e \in E_v^+} q_e - \sum_{e \in E_v^-} q_e = 0 \quad \forall v \in V. \quad (3.2)$$

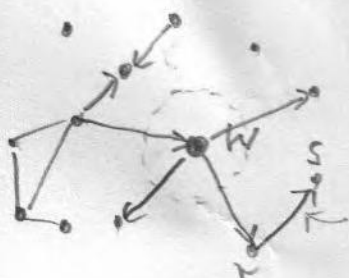
Grund: Massenerhaltung, Knoten speichert keine Flüssigkeit.

Nur $n-1$ der Beziehungen (3.2) sind linear unabhängig:

Wähle $w \in V$.

$$\sum_{v \in V - \{w\}} \left(\sum_{e \in E_v^+} q_e - \sum_{e \in E_v^-} q_e \right) = \sum_{v \in V - \{w\}} \sum_{e \in E_v^+} q_e - \sum_{v \in V - \{w\}} \sum_{e \in E_v^-} q_e$$

$= 0!$



$$e' = (r, s) \quad r \neq w, s \neq w$$

Kommt genau zweimal

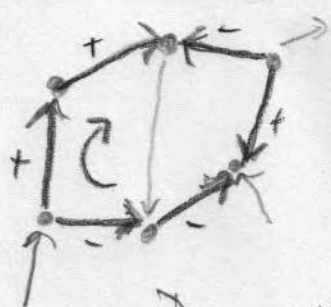
vor: $+q_{e'}$ für Knoten s (eingehend)

$-q_{e'}$ für Knoten r (ausgehend)

$$= \sum_{e \in E_w^-} q_e - \sum_{e \in E_w^+} q_e = 0$$

Knotenregel für $v \in V - \{w\}$ erfüllt \Rightarrow Knotenregel für w erfüllt.

4) Maschenregel (2. Kirchhoffsches Gesetz)



- Betrachte $C \subseteq E$, C ^{beliebigen} geschlossener Pfad

$$C^+ = \{ e \in C \mid e \text{ genauso wie } C \text{ orientiert} \}$$

$$C^- = \{ e \in C \mid e \text{ entgegen } C \text{ orientiert} \}$$

Dann gilt

Grund:

$$\sum_{e \in C^+} \Delta p_e - \sum_{e \in C^-} \Delta p_e = 0 \quad \forall \text{ Pfad}$$



$$\int_v^w \Delta p_e \cdot A \cdot ds$$

ist Energie die notwendig ist um Probenkörper von v nach w zu bringen

- Betr. zweier bel. Knoten r, s . Energie im K. von r nach s zu bringen unabhängig von Weg.
- Lineare Abhängigkeit!

Folge der Maschenregel:

(„Potentiale“)

- Man darf $n-1$ Knotendrücker p_v als Unbekannte wählen
- $e=(v,w) : \Delta p_e = p_v - p_w$
- Druck in einem (Referenz-) Knoten r kann willkürlich festgelegt werden. z.B. $p_r = 0$. ($p_v, v \in V$ Lösung $\Rightarrow p_v + const$ auch Lösung)
- Maschenregel ist dann für alle Pfade erfüllt.

a) $\forall e \in E_R$ schreibe Druckdifferenzen:

$$e=(v,w) \in E_R : \Delta p_e = \begin{cases} p_v - p_w & v \neq r \wedge w \neq r \\ p_v & w = r \\ p_w & v = r \end{cases} \quad r: \text{Referenzknoten}$$

Wähle „Anordnung“ $e_k \leftrightarrow k \quad v_i \leftrightarrow i$ Wahl

$$e_k=(v_i, v_j) : k \quad \begin{bmatrix} \Delta p \end{bmatrix} = \begin{bmatrix} B \end{bmatrix} \begin{bmatrix} p \end{bmatrix} \quad \begin{matrix} \tau = v_n \\ \Rightarrow p = (p_{v_1}, \dots, p_{v_{n-1}})^T \\ \Delta p = (\Delta p_1, \dots, \Delta p_m) \end{matrix}$$

$m \times (n-1)$ Matrix
 $|E_R| \quad |V|-1$

b) Leitfähigkeiten:

$$e \in E_R : q_e = L_e \Delta p_e$$

als Matrix:

$$\begin{bmatrix} q \end{bmatrix} = \begin{bmatrix} L \end{bmatrix} \begin{bmatrix} \Delta p \end{bmatrix}$$

$m \times m$ Diagonalmatrix $m = |E_R|$

c) Knotenregeln: $n-1$ Stück für Knoten v_1, \dots, v_{n-1} (exklusive Referenzknoten)

$$v \in V - \{v_n\} : \sum_{e \in E_v^+} q_e - \sum_{e \in E_v^-} q_e = 0$$

Pumpenströme auf rechte Seite \Leftrightarrow

$$\sum_{e \in E_v^+ \cap E_R} q_e - \sum_{e \in E_v^- \cap E_R} q_e = \sum_{e \in E_v^- \cap E_p} q_e - \sum_{e \in E_v^+ \cap E_p} q_e$$

$$\Leftrightarrow \begin{bmatrix} B^T \end{bmatrix} \begin{bmatrix} q \end{bmatrix} = \begin{bmatrix} b \end{bmatrix} \quad b = (b_1, \dots, b_{n-1})$$

$(n-1) \times m$ Matrix
enthält Pumpenströme

Alles zusammen:

$$\underbrace{B^T L B}_{=: A} p = b$$

$$\Leftrightarrow A p = b$$

- A ist $(n-1) \times (n-1)$ Matrix

- A ist symmetrisch und positiv definit, d.h. $x^T A x > 0 \forall x \neq 0$

- damit insbesondere invertierbar

- A ist dünn besetzt:

$$a_{ij} \neq 0 \Leftrightarrow (v_i, v_j) \in E \vee (v_j, v_i) \in E.$$

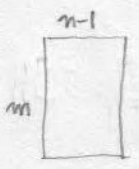
sehr viele Einträge sind Null

$$|\{(i,j) \mid a_{ij} \neq 0\}| = O(n) \text{ (statt } n^2).$$

Bew. 1) $B \in \mathbb{R}^{m \times (n-1)}$ hat max. Rang, also $n-1$ da $m \geq n-1$, $n-1$ linear unabh. Spalten!

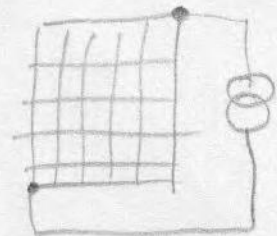
2) L diagonal mit $l_{ii} > 0$

3) Sei $x \neq 0$, $y := Bx$ dann ist $x^T A x = \sum_i l_{ii} y_i^2 > 0$ da $\dim(\ker(B)) + \text{rang}(B) = n-1$ und $\text{rang}(B) = n-1 \Rightarrow \dim(\ker(B)) = 0$ also $y=0$ nur für $x=0$.



ü1: ein kleines Netzwerk

ü2: New York City



- völlig analog lassen sich elektrische Netzwerke behandeln (aus Widerständen)

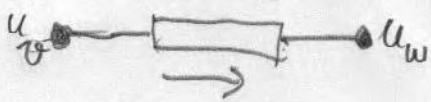
Ohmsches Gesetz

$$i_e = \frac{1}{R_e} (u_v - u_w)$$

i_e : Strom durch Widerstand R_e

u_v : Knotenpotentiale

R_e : Widerstand



- RLC Netzwerke mit harmonischer Anregung

→ komplexe Ströme und Spannungen

$$A x = b \text{ mit } x, b \in \mathbb{C}^n, A \in \mathbb{C}^{n \times n}$$

L: Spulen

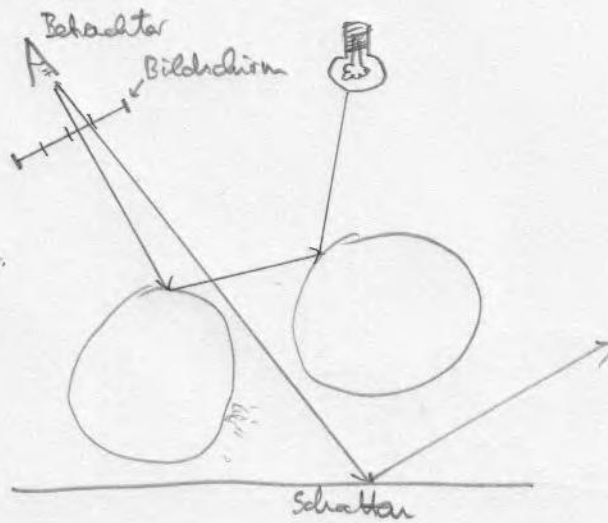
C: Kapazitäten

3.2 Radiosity-Methode in der Computergrafik

Bleuchtung einer "Szene".

Ray-Tracing:

Nachteil: starke Schatten.



Radiosity-Methode:

$$S = \{ x \in \mathbb{R}^3 \mid x \text{ auf Oberfläche eines Objekts} \}$$

Bestimme $B : S \rightarrow \mathbb{R}$ "Energiedichte"

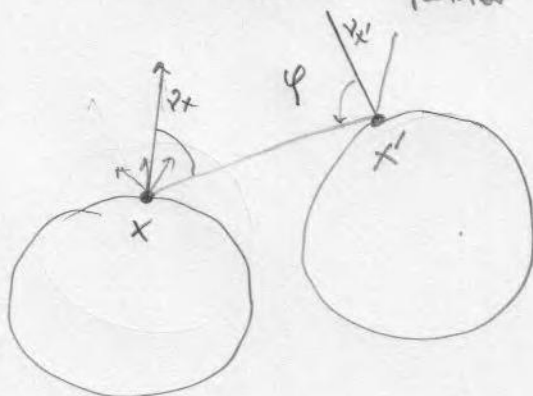
$$\int_{\omega} B \, dA = \text{von } \omega \text{ abgestrahlte Energie}$$

Energie wird von einem Punkt $x \in S$ in alle Richtungen abgestrahlt.

Bestimmungsgleichung für $B(x)$: reflektiertes Licht

$$B(x) = E(x) + \rho(x) \int_S B(x') \underbrace{\frac{\cos \varphi_{x,x'} \cos \varphi_{x',x}}{\pi \|x-x'\|^2} V(x,x')}_{=: \lambda(x,x') \text{ "Kern" }} \, dA'$$

\uparrow $x \in S$ \uparrow Eigenstrahlung (Lichtquelle) \uparrow Reflexionsfaktor \uparrow Licht von x'



$$\varphi_{a,b} = \angle(v_a, b-a)$$

$$\cos \varphi_{a,b} = \begin{cases} 1 & \text{Licht trifft aus Richtung } v_a \text{ ein} \\ 0 & \text{Licht trifft tangential zur Oberfläche ein} \end{cases}$$

Sichtbarkeit (Visibility)

$$V(x,x') = \begin{cases} 1 & x' \text{ von } x \text{ aus sichtbar} \\ 0 & \text{sonst} \end{cases}$$

(Eigenschaft der Szene)

Integralgleichung für $B(x): S \rightarrow \mathbb{R}$:

$$B(x) - s(x) \int_S B(x') \lambda(x, x') dA' = E(x) \quad \forall x \in S.$$

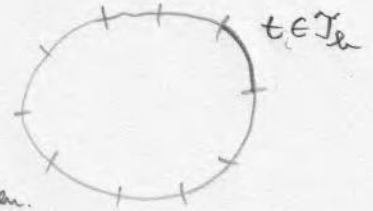
Numerische Lösung mit "Kollokationsmethode".

a) Zerlegung der Oberfläche: $\mathcal{T}_n = \{t_1, \dots, t_n\}$

$$t_i \subset S, \quad t_i \cap t_j = \emptyset \quad \forall i \neq j, \quad \bigcup_{i=1}^n \bar{t}_i = S$$

t_i : offene Gebiete.

Könnte man abschwächen.



Zu $t \in \mathcal{T}_n$ wähle $x_t =$ Mittelpunkt von t .

Dieser Prozess heißt auch Diskretisierung. Üblich bei Differential- und Integralgleichungen.

b) Approximiere $B: S \rightarrow \mathbb{R}$ durch diskrete Funktion $B_n: S \rightarrow \mathbb{R}$.

$$B_n(x) = \sum_{j=1}^n z_j \varphi_j(x)$$

$z_j \in \mathbb{R}$ Koeffizient

$\varphi_j: S \rightarrow \mathbb{R}$ "Basisfunktion"

d.h. $\varphi_1, \dots, \varphi_n$ linear unabhängig.

stückweise konstante Funktionen:

$$\varphi_j(x) = \begin{cases} 1 & \text{falls } x \in t_j \\ 0 & \text{sonst.} \end{cases}$$

c) Erfülle Integralgleichung nur für $x \in X_n = \{x_{t_1}, \dots, x_{t_n}\}$.

$$B_n(x_i) - s(x_i) \int_S B_n(x') \lambda(x_i, x') dA' = E(x_i) \quad i=1, \dots, n$$

$$\Leftrightarrow \sum_{j=1}^n z_j \varphi_j(x_i) - s(x_i) \int_S \sum_{j=1}^n z_j \varphi_j(x') \lambda(x_i, x') dA' = E(x_i) \quad i=1, \dots, n$$

$$\Leftrightarrow \sum_{j=1}^n z_j \left\{ \varphi_j(x_i) - s(x_i) \int_S \varphi_j(x') \lambda(x_i, x') dA' \right\} = \underbrace{E(x_i)}_{b_i} \quad i=1, \dots, n$$

$$\Leftrightarrow \boxed{Az = b}.$$

- Integral in a_{ij} wird i. d. R. auch numerisch berechnet.
 \Rightarrow Methoden später in der Vorlesung.

- Man begeht einen Diskretisierungsfehler

$$\|B - B_h\| = O(h^\alpha) \quad (\text{Konvergenz})$$

mit $h = \max_{t \in T_h} \text{diam}(t)$, α : "Konvergenzordnung".

d.h. je feiner die Unterteilung ($n \rightarrow \infty$), desto besser approximiert B_h die gesuchte Funktion B .

\Rightarrow Gleichungssysteme können beliebig groß werden.
 (im Gegensatz zum Röhrennetzwerk).

- A ist in diesem Fall nicht dünn besetzt sondern voll besetzt!

4: Konditionen der Lösung linearer Gleichungssysteme

4.1 Lösbarkeit

Gegeben Matrix A , Vektor b :

$$A = (a_{ij})_{i,j=1}^{m,n} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}, \quad b = (b_i)_{i=1}^m = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$$

Gesucht ist $x = (x_j)_{j=1}^n = (x_1, \dots, x_n)^T$ so dass

$$\forall i: \sum_{j=1}^n a_{ij} x_j = b_i,$$

bzw.

$$Ax = b.$$

(4.1)

Die Zahlen a_{ij}, b_i, x_j können aus \mathbb{R} oder \mathbb{C} sein.

Das Gleichungssystem (4.1) heißt

- unterbestimmt falls $m < n$
- quadratisch falls $m = n$
- überbestimmt falls $m > n$

$$\begin{array}{l} \square \quad | \quad = \quad | \\ \square \quad | \quad = \quad | \\ \square \quad | \quad = \quad | \end{array}$$

bei maximalem Rang von A :

i. allg. viele Lösungen

genau eine Lösung

$b \in \text{Bild}(A)$.

Es gibt mindestens eine Lösung falls

$$\text{Rang}(A) = \text{Rang}([A|b]) = \text{Rang} \left(\begin{bmatrix} a_{11} & \dots & a_{1n} & b_1 \\ \vdots & & \vdots & \vdots \\ a_{m1} & \dots & a_{mn} & b_m \end{bmatrix} \right)$$

Für quadratische Matrizen sind folg. Aussagen äquivalent:

- $Ax = b$ ist für jedes b eindeutig lösbar.
- $\text{Rang}(A) = n$.
- $\det(A) \neq 0$.
- A hat keinen Eigenwert Null.

Quantifizieren von Fehlern erfordert Normen.

Im folgenden sei $K = \mathbb{R}$ oder $K = \mathbb{C}$ der zugrunde liegende Körper.

Definition 4.1 (Vektornorm)

Eine Abbildung $\|\cdot\| : K^n \rightarrow \mathbb{R}_+$ heißt Norm falls gilt

$$(N1) \quad \|x\| > 0 \quad x \neq 0 \quad (\text{Definitheit})$$

$$(N2) \quad \|\alpha \cdot x\| = |\alpha| \|x\| \quad x \in K^n, \alpha \in K \quad (\text{positive Homogenität})$$

$$(N3) \quad \|x+y\| \leq \|x\| + \|y\| \quad x, y \in K^n \quad (\text{Subadditivität, Dreiecksungl.})$$

Beispiel 4.2 Häufig verwendete Normen sind:

$$\|x\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2} \quad \text{Euklidische Norm, } \ell_2\text{-Norm}$$

$$\|x\|_\infty = \max_{i=1, \dots, n} |x_i| \quad \text{Maximumnorm, } \ell_\infty\text{-Norm}$$

$$\|x\|_1 = \sum_{i=1}^n |x_i| \quad \ell_1\text{-Norm}$$

Man zeigt über die Konvergenz von Folgen von Vektoren:

$$\forall i=1, \dots, n: x_i^{(t)} \rightarrow x_i \quad (t \rightarrow \infty) \Leftrightarrow \|x^{(t)} - x\| \rightarrow 0 \quad (t \rightarrow \infty)$$

(Konvergenz der Normen äquivalent zu komponentenweiser Konvergenz).
→ siehe nächster Satz

Eine wichtige Folgerung aus (N3) ist

$$\|x-y\| \geq \left| \|x\| - \|y\| \right| \quad \forall x, y \in K^n. \quad (4.2)$$

Beweis:

$$\|x\| = \|x-y+y\| \stackrel{(N3)}{\leq} \|x-y\| + \|y\| \Rightarrow \|x-y\| \geq \|x\| - \|y\|$$

$$\|y\| = \|y-x+x\| \leq \|y-x\| + \|x\| \Rightarrow \|x-y\| \geq \|y\| - \|x\| = -(\|x\| - \|y\|)$$

Aus (4.2) folgt:

$$\|x^{(t)} - x\| \rightarrow 0 \quad (t \rightarrow \infty) \Rightarrow \|x^{(t)}\| \rightarrow \|x\| \quad (t \rightarrow \infty)$$

Satz 4.3 (Äquivalenz aller Normen)

Auf \mathbb{K}^n sind alle Normen im folgenden Sinne äquivalent:

Zu $\|\cdot\|, \|\cdot\|'$ gibt es Zahlen $m, M > 0$ aus \mathbb{R} so dass gilt

$$m \|x\|' \leq \|x\| \leq M \|x\|' \quad \forall x \in \mathbb{K}^n.$$

Beweis: Es genügt dies zu zeigen für $\|\cdot\|' = \|\cdot\|_\infty$.

Seien e_1, \dots, e_n die kanonischen Einheitsvektoren. Jedes $x \in \mathbb{K}^n$ hat die Darstellung $x = \sum_{i=1}^n x_i e_i$, also $\|x\| = \left\| \sum_{i=1}^n x_i e_i \right\| \leq \sum_{i=1}^n |x_i| \|e_i\| \leq \left(\max_{1 \leq i \leq n} |x_i| \right) \sum_{i=1}^n \|e_i\| = \|x\|_\infty \sum_{i=1}^n \|e_i\| = \delta \|x\|_\infty$

wg. $\|x\| \leq \delta \|x\|_\infty$ ist $\|\cdot\|$ stetig bezügl. Komponentenweiser Konvergenz.

Punktmenge $S = \{x \in \mathbb{K}^n : \|x\|_\infty = 1\} \subset \mathbb{K}^n$ ist - beschränkt (bezügl. $\|\cdot\|$ da $\|x\| \leq \delta$)
 - abgeschlossen, d.h. Grenzwerte von Folgen aus S sind in S (da \mathbb{K}^n metrischer Raum).
 \Rightarrow kompakt.

Die Funktion $\|\cdot\|: S \rightarrow \mathbb{R}_+$ nimmt auf S ihr Minimum und Maximum an; d.h.:

Es gibt $\underline{x}, \bar{x} \in S$:

$$0 < \underbrace{\| \underline{x} \|}_{m} \leq \|x\| \leq \| \bar{x} \| = \delta \underbrace{\| \bar{x} \|_\infty}_{=1} = \delta \quad \forall x \in S$$

da $0 \notin S, \underline{x} \in S$ und $N1$.

$$\| \bar{x} \| = \left\| \sum_{k=1}^n \bar{x}_k e_k \right\| = \sum_{k=1}^n |\bar{x}_k| \|e_k\| \leq \left(\max_{k=1, \dots, n} |\bar{x}_k| \right) \sum_{k=1}^n \|e_k\| = \|\bar{x}\|_\infty \sum_{k=1}^n \|e_k\| = \|\bar{x}\|_\infty \delta$$

da $\bar{x} \in S$ und $\|\bar{x}\| = \delta$ folgt $\delta = \|\bar{x}\|_\infty \delta \Rightarrow \|\bar{x}\|_\infty = 1$.

Nun sei $y \in \mathbb{K}^n - \{0\}$ beliebig.

Dann ist $y / \|y\|_\infty \in S$ (da Norm 1) und somit

$$\|z\| \leq \|y / \|y\|_\infty\| = \frac{1}{\|y\|_\infty} \|y\| \leq \| \bar{x} \|$$

$$\Leftrightarrow \underbrace{\|z\|}_{=m} \|y\|_\infty \leq \|y\| \leq \underbrace{\| \bar{x} \|}_{=M} \|y\|_\infty$$

Vorsicht: m, M hängen i.d.R. von der Dimension n ab.

Auffassen wenn man etwas für $n \rightarrow \infty$ beweist.

Nachmal zur Konvergenz:

$$x_i^{(t)} \rightarrow x_i \quad (t \rightarrow \infty) \Leftrightarrow \|x^{(t)} - x\|_\infty \rightarrow 0 \quad (t \rightarrow \infty) \text{ überlegt man leicht}$$

Satz 4.3 zeigt dann, dass man auch eine beliebige Norm nehmen darf.

4.3 Matrixnormen

Der $K^{m \times n}$ stellt auch einen Vektorraum dar und kann mit dem K^{mn} identifiziert werden.

Damit definiert jede Norm auf $K^{m \times n}$ auch eine Norm auf $A \in K^{m \times n}$.

$A^{(t)} \rightarrow A$ ($t \rightarrow \infty$) meint dann $a_{ij}^{(t)} \rightarrow a_{ij}$ ($t \rightarrow \infty$) $\forall ij \in \{1, \dots, m\} \times \{1, \dots, n\}$

Es zeigt sich, dass folgende Eigenschaften hilfreich sind:

Definition 4.4 Eine Norm $\|\cdot\|$ auf $K^{n \times n}$ (man würde zwei V-Normen erfordern) heißt verträglich mit einer Vektornorm $\|\cdot\|$ auf K^n falls gilt

$$\|Ax\| \leq \|A\| \|x\| \quad x \in K^n, A \in K^{n \times n}$$

Sie heißt Matrixnorm, wenn sie submultiplikativ ist:

$$\|AB\| \leq \|A\| \|B\| \quad A, B \in K^{n \times n}$$

Beispiel 4.5 Die Frobenius-Norm

$$\|A\|_{Fr} := \left(\sum_{i,j=1}^n |a_{ij}|^2 \right)^{1/2}$$

(heißt matrixnorm)

ist eine Matrixnorm, die verträglich mit der euklidischen Norm ist.

Beweis: (als Ü?) c.s. $\langle x, y \rangle \leq \|x\| \|y\|$

$$\|Ax\|_2^2 = \sum_{i=1}^n \left(\sum_{j=1}^n a_{ij} x_j \right)^2 \leq \sum_{i=1}^n \left(\sum_{j=1}^n |a_{ij}|^2 \cdot \sum_{j=1}^n |x_j|^2 \right) = \|x\|_2^2 \sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2 = \|A\|_{Fr}^2 \|x\|_2^2$$

Definition 4.6 (Zugeordnete Matrixnorm)

Es sei $\|\cdot\|$ eine beliebige (Vektor-)Norm auf K^n . Dann heißt

$$\|A\| := \sup_{x \in K^n - \{0\}} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\|$$

die $\|\cdot\|$ zugeordnete (oder natürliche) Matrixnorm, d.h. $\|\cdot\|$ ist verträglich und submultiplikativ.

Übung: Beweis dazu $y \neq 0$ beliebig $\frac{\|Ay\|}{\|y\|} \leq \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \|A\| \Rightarrow \|Ay\| \leq \|A\| \|y\|$

$$\|AB\| = \sup_{x \neq 0} \frac{\|ABx\|}{\|x\|} \leq \sup_{x \neq 0} \frac{\|A\| \|Bx\|}{\|x\|} = \|A\| \sup_{x \neq 0} \frac{\|Bx\|}{\|x\|} = \|A\| \|B\|$$

Hilfssatz 4.7

22.10.09

Die zugeordneten Matrixnormen zu $\|\cdot\|_\infty$ und $\|\cdot\|_1$ sind

$$\|A\|_\infty := \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| \quad (\text{Zeilensummennorm}),$$

$$\|A\|_1 := \max_{j=1, \dots, n} \sum_{i=1}^n |a_{ij}| \quad (\text{Spaltensummennorm}).$$

Beweis: siehe Rannacher, Numerik S. 104. Teile hieraus:

$$\begin{aligned} \text{a) } \|Ax\|_\infty &= \max_{i=1, \dots, n} \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| |x_j| \\ &\leq \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| \underbrace{\left(\max_{k=1, \dots, n} |x_k| \right)}_{= \|x\|_\infty} \\ &\leq \|x\|_\infty \underbrace{\max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}|}_{= \|A\|_\infty} = \|A\|_\infty \|x\|_\infty \end{aligned}$$

Dies zeigt die Verträglichkeit

b) Somit

$$\sup_{\|x\|_\infty=1} \|Ax\|_\infty \stackrel{(a)}{\leq} \sup_{\|x\|_\infty=1} \|A\|_\infty \underbrace{\|x\|_\infty}_{=1} = \|A\|_\infty$$

c) zu zeigen ist nun

$$\sup_{\|x\|_\infty=1} \|Ax\|_\infty \geq \|A\|_\infty$$

Dies geschieht konstruktiv indem man so einen Vektor angibt $\rightarrow Ra$.

Sei $A \in \mathbb{K}^{n \times n}$, $\lambda \in \mathbb{K}$ und $e \in \mathbb{K}^n$ so dass

$$Ae = \lambda e,$$

dann heißt λ Eigenwert von A und e zugehöriger Eigenvektor.
 (λ, e) heißt auch Eigenpaar.

$(A - \lambda I)e = 0$ ist für nichttriviales $e \neq 0$ nur für $\det(A - \lambda I) = 0$

$$p(\lambda) = \det(A - \lambda I) = 0$$

erfüllbar.

Das charakteristische Polynom $p(\lambda)$ hat genau n Nullstellen in \mathbb{C} (Vielfachheiten mitgezählt).

Zu jedem Eigenwert λ gibt es mindestens einen Eigenvektor.

Mit den Normregeln folgt für ein Eigenpaar (λ, e) , $\|e\| = 1$:

$$|\lambda| \underset{\|e\|=1}{=} |\lambda| \|e\| = \|\lambda e\| = \|Ae\| \leq \|A\| \underbrace{\|e\|}_{=1} = \|A\| \quad (4.3)$$

Also alle Eigenwerte $\lambda \in \mathbb{C}$ liegen innerhalb eines Kreises um 0 mit Radius $\|A\|$.

Z.B. mit $\|A\|_\infty$ erhält man so eine konkrete Abschätzung für den größten Eigenwert.

Definition 4.8 (Spezielle Matrizen)

Zu $A \in \mathbb{K}^{m \times n}$ heißt $A^T \in \mathbb{K}^{n \times m}$ transponierte Matrix und es ist $(A^T)_{ij} = (A)_{ji}$

Zu $A \in \mathbb{K}^{m \times n}$ ist $\bar{A} \in \mathbb{K}^{m \times n}$ gegeben durch $(\bar{A})_{ij} = \overline{(A)_{ij}}$ konjugiert komplex

a) $A \in \mathbb{K}^{n \times n}$ heißt hermitesch falls $A = \bar{A}^T$ d.h. $a_{ij} = \overline{a_{ji}}$
 Manche Autoren schreiben auch A^H für \bar{A}^T .

Reelle hermitesche Matrizen heißen symmetrisch (dann ist $A = A^T$)

b) Für komplexe Matrizen $A \in \mathbb{C}^{n \times n}$ heißt normal also Arcell dann ist hermitisch (\Leftrightarrow symmetrisch) in \mathbb{C} nicht!

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

$$A \bar{A}^T = \bar{A}^T A = I, \text{ also } A^{-1} = \bar{A}^T$$

unitär

c) Für reelle Matrizen $A \in \mathbb{R}^{n \times n}$ heißt

$$A A^T = A^T A, \text{ also } A^{-1} = A^T \quad \text{orthogonal.}$$

Definition 4.9 (Skalarprodukt)

22.10.09

Eine Abbildung $(\cdot, \cdot) : \mathbb{K}^n \times \mathbb{K}^n \rightarrow \mathbb{K}$ heißt Skalarprodukt, falls gilt:

$$(S1) \quad (x, y) = \overline{(y, x)} \quad x, y \in \mathbb{K}^n \quad (\text{Symmetrie})$$

$$(S2) \quad (\alpha x + \beta y, z) = \alpha (x, z) + \beta (y, z) \quad x, y, z \in \mathbb{K}^n, \alpha, \beta \in \mathbb{K}$$

$$(S3) \quad (x, x) > 0 \quad x \in \mathbb{K}^n - \{0\} \quad \begin{array}{l} \text{(Linearität)} \\ \text{(Definitheit)} \end{array}$$

aus (S1) folgt $(x, x) = \overline{(x, x)}$ also $\text{Im}(x, x) = 0$

Ein Skalarprodukt erzeugt immer eine zugehörige Norm

$$\|x\| := \sqrt{(x, x)}, \quad x \in \mathbb{K}^n.$$

Es gilt die Cauchy-Schwarzsche Ungleichung

$$|(x, y)| \leq \|x\| \|y\|$$

Das Euklidische Skalarprodukt (wird im folg. oft verwendet) lautet:

$$(x, y)_2 = \sum_{i=1}^n x_i \overline{y_i} = x^T \overline{y}$$

Damit gilt (andere Schreibweise für hermitesch):

$$A = \overline{A}^T \iff (Ax, y)_2 = (x, Ay)_2 \quad \forall x, y \in \mathbb{K}^n.$$

4.5 Die Spektralnorm

Die der Euklidischen Norm zugeordnete Matrixnorm $\|A\|_2$ heißt Spektralnorm.

Hilfssatz 4.10 Für die Spektralnorm hermitescher Matrizen gilt

$$\|A\|_2 = \max \{ |\lambda| : \lambda \text{ Eigenwert von } A \}.$$

Für jede Matrix $A \in \mathbb{K}^{n \times n}$ gilt

$$\|A\|_2 = \max \{ |\lambda|^{1/2} : \lambda \text{ ist Eigenwert von } \overline{A}^T A \}$$

Beweis: a) Sei A hermitesch. Dann hat A n reelle

23.10.09

Eigenwerte und einen vollständigen Satz von orthonormalen Eigenvektoren:

$$\{w^1, \dots, w^n\} \subset \mathbb{K}^n : Aw^i = \lambda_i w^i, \quad (w^i, w^j)_2 = \delta_{ij} \quad i, j = 1, \dots, n$$

↑ orthogonales SP!

Jedes $x \in \mathbb{K}^n$ lässt sich in der Basis $\{w^1, \dots, w^n\}$ darstellen:

$$x = \sum_{i=1}^n \alpha_i w^i \quad \text{mit} \quad \alpha_i = (x, w^i)_2 \in \mathbb{K} \quad (\text{eindeutig!})$$

Und es gilt:

$$\begin{aligned} \|x\|_2^2 &= (x, x)_2 = \left(\sum_{i=1}^n \alpha_i w^i, \sum_{j=1}^n \alpha_j w^j \right)_2 \\ &= \sum_{i,j=1}^n \alpha_i \overline{\alpha_j} \underbrace{(w^i, w^j)_2}_{=\delta_{ij}} = \sum_{i=1}^n |\alpha_i|^2 \end{aligned}$$

$$\begin{aligned} \overline{(a+ib)(c+id)} &= \overline{(a-ib)(c+id)} \\ &= (a+ib)(c-id) \\ &= (ac-bd) + i(ad+bc) \end{aligned}$$

$$\begin{aligned} \|Ax\|_2^2 &= (Ax, Ax)_2 = \left(\sum_i \underbrace{A\alpha_i}_{\alpha_i \lambda_i} w^i, \sum_j \underbrace{A\alpha_j}_{\alpha_j \lambda_j} w^j \right)_2 \\ &= \sum_{i,j=1}^n \alpha_i \lambda_i \overline{\alpha_j \lambda_j} \underbrace{(w^i, w^j)_2}_{=\delta_{ij}} = \sum_{i=1}^n \lambda_i^2 |\alpha_i|^2 \end{aligned}$$

↑ Betrag, da $\alpha \in \mathbb{K}$

Also

$$\|A\|_2^2 = \sup_{x \in \mathbb{K}^n \setminus \{0\}} \frac{\|Ax\|_2^2}{\|x\|_2^2} = \sup_{x \in \mathbb{K}^n \setminus \{0\}} \frac{\sum_{i=1}^n \lambda_i^2 |\alpha_i|^2}{\sum_{i=1}^n |\alpha_i|^2} \leq \max_{1 \leq i \leq n} |\lambda_i|^2$$

Andererseits zeigt (4.3) $\|A\| \geq |\lambda|$ für jede Norm und jeden EW, also muss $\|A\|_2 = \max_{1 \leq i \leq n} |\lambda_i|$ sein.

b) Sei $A \in \mathbb{K}^n$ eine beliebige Matrix. Dann gilt Rechenregeln für konjug. komplex

$$\|Ax\|_2^2 = (Ax, Ax)_2 = (Ax)^T (\overline{Ax}) = x^T A^T (\overline{Ax}) = x^T (\overline{A^T A} x) = (x, \overline{A^T A} x)_2$$

$\overline{A^T A}$ ist hermitesch und habe EW λ_i sowie EV w^i orthonormal wie oben

also

$$\|Ax\|_2^2 = \left(\sum_{i=1}^n \alpha_i w^i, \sum_{j=1}^n \overline{A^T A} \alpha_j w^j \right)_2 = \sum_{i,j=1}^n \alpha_i \underbrace{\lambda_j}_{\text{reell}} \overline{\alpha_j} \underbrace{(w^i, w^j)_2}_{=\delta_{ij}} = \sum_{i=1}^n \lambda_i |\alpha_i|^2$$

und damit

$$\|A\|_2^2 = \sup_{x \in \mathbb{K}^n \setminus \{0\}} \frac{\|Ax\|_2^2}{\|x\|_2^2} = \sup_{x \in \mathbb{K}^n \setminus \{0\}} \frac{\sum_{i=1}^n \lambda_i |\alpha_i|^2}{\sum_{i=1}^n |\alpha_i|^2} \leq \max_{1 \leq i \leq n} |\lambda_i|$$



4.6 Positiv definite Matrizen

24.10.09

Wichtige Klasse von Matrizen mit vorteilhaften Eigenschaften.

in \mathbb{K} !

Definition 4.11 Eine Matrix $A \in \mathbb{K}^{n \times n}$ heißt positiv definit wenn

a) $(Ax, x)_2 \in \mathbb{R} \quad \forall x \in \mathbb{K}^n \setminus \{0\} \quad (\mathbb{K} = \mathbb{C})$

b) $(Ax, x)_2 > 0 \quad \forall x \in \mathbb{K}^n \setminus \{0\}$ (eigentliche Bedingung)

Ziel: Charakterisierung positiv definiten Matrizen.

Im Fall $\mathbb{K} = \mathbb{C}$ bedeutet die Bedingung a) eine Einschränkung an die Matrix A .

Eigenschaft 4.12 Für $A \in \mathbb{C}^{n \times n}$ gilt A hermitesch genau dann wenn $(Ax, x)_2 \in \mathbb{R} \quad \forall x \in \mathbb{C}^n \setminus \{0\}$.

Beweis: " \Rightarrow " (als Ü-Aufgabe!)
 $(Ax, x)_2 \in \mathbb{R} \Leftrightarrow \underbrace{(Ax, x)_2}_{a+ib} = \overline{\underbrace{(Ax, x)_2}_{a-ib}}$ geht nur für $b=0$

also $(Ax, x)_2 \stackrel{\uparrow}{=} (x, Ax) \stackrel{(5.1)}{=} \overline{(Ax, x)}$ damit ist \Rightarrow gezeigt.
hermitesch 5.7

" \Leftarrow " Auftrennen in Real- und Imaginärteil:

$\mathbb{C}^n \ni x = a+ib \quad a, b \in \mathbb{R}^n \quad A = V+iW, \quad V, W \in \mathbb{R}^{n \times n}$

und ausrechnen:

$$(Ax, x)_2 = ((V+iW)(a+ib), a+ib)_2 = ((Va-Wb)+i(Vb+Wa), a+ib)_2$$
$$\stackrel{(x,y)_2 = \overline{(y,x)}_2!}{=} \underbrace{(Va-Wb, a)_2 + (Vb+Wa, b)_2}_{\text{Re}} + i \underbrace{[(Vb+Wa, a)_2 - (Va-Wb, b)_2]}_{\text{Im}}$$

also $\text{Im}((Ax, x)_2) = (Vb, a)_2 + (Wa, a)_2 - (Va, b)_2 + (Wb, b)_2 \stackrel{!}{=} 0 \quad (*)$

Fall I) $W=0$ d.h. $A=V$ (reell), aber $x=a+ib$ ist komplex. Dann reduziert sich (*) auf

$(Vb, a)_2 - (Va, b)_2 = 0 \quad \forall a, b \in \mathbb{R}^n \Leftrightarrow (Vb, a)_2 = (V^T b, a)_2 \quad \forall a, b \in \mathbb{R}^n$
alles reell

damit muss $V=V^T$ sein.

Fall II) Setze $b=0$, dann reduziert sich (*) auf

10
25.10.09

$$(W a, a)_2 = 0 \quad \forall a \in \mathbb{C}^n \setminus \{0\}$$

IIa: Setze $a = e^i \Rightarrow w_{ii} = 0$. Hier: $e_j^i = \delta_{ij}$ Kartesische Einheitsv.

IIb: Sei also $w_{ii} = 0 \forall i$ und setze $a = e^i + e^j$ ($j \neq i$)

$$\text{dann folgt } w_{ij} + w_{ji} = 0 \quad \forall i, j$$

$$\Leftrightarrow w_{ij} = -w_{ji}$$

Damit ist $A = V + iW$ mit $V = V^T$ und $W = -W^T$ also A hermitesch. \blacksquare

Positiv definite Matrizen in $\mathbb{C}^{n \times n}$ sind also immer hermitesch.

Positiv definite Matrizen in $\mathbb{R}^{n \times n}$ sind nicht notwendigerweise symmetrisch.

Charakterisierung über Eigenwerte:

Lemma 4.13 Eine hermitesche Matrix A ist genau dann positiv definit, wenn alle ihre (reellen) Eigenwerte positiv sind. Alle Hauptdiagonalelemente sind (reell und) positiv.

Beweis:

- A sei hermitesch mit lauter positiven Eigenwerte.

$x \in \mathbb{C}^n$ hat Darstellung $x = \sum \alpha_i w^i$ w^i : Eigenvektor zum EW λ_i .

$$(Ax, x)_2 = \sum_{i,j=1}^n \lambda_i \alpha_i \bar{\alpha}_j \underbrace{(w^i, w^j)_2}_{=\delta_{ij}} = \sum_{i=1}^n \lambda_i |\alpha_i|^2 > 0$$

- A sei hermitesch und positiv definit. A hat also Eigenpaare (λ_i, w^i) mit $\lambda_i \in \mathbb{R}$. Zu zeigen ist λ_i positiv

$$0 < (A w^i, w^i)_2 = (\lambda_i w^i, w^i)_2 = \lambda_i \underbrace{(w^i, w^i)_2}_{=1} = \lambda_i$$

- Setze $e^i \in \mathbb{R}^n$, $e_j^i = \delta_{ij}$ kartesische Einheitsvektoren.

$$0 < (A w_i, w_i)_2 = a_{ii} \in \mathbb{R}. \quad \blacksquare$$

Die Aussage gilt natürlich auch für reelle symmetrische und positiv definite Matrizen, da auch diese hermitesch sind.

Speziell für reelle Matrizen gilt:

11
25.10.09

Lemma 4.14 Sei $A \in \mathbb{R}^{n \times n}$ symmetrisch und positiv definit.

Dann liegt das betragsmäßig größte Element auf der Hauptdiagonalen.

Beweis: e^i sei wieder der i -te kanonische Einheitsvektor und a_{ij} , $j \neq i$, das betragsmäßig größte Element. Dann zeigt man den Widerspruch.

$$0 < (A(e^i - \text{sign}(a_{ij})e^j), e^i - \text{sign}(a_{ij})e^j)_2$$

$$= (Ae^i, e^i)_2 - 2 \text{sign}(a_{ij}) \underbrace{(Ae^i, e^j)}_{a_{ij}} + \text{sign}(a_{ij})^2 (Ae^j, e^j)$$

$$\stackrel{A=A^T}{=} a_{ii} - 2|a_{ij}| + a_{jj} \leq 0 \quad \downarrow$$

da $|a_{ij}| \geq a_{ii}$, $|a_{ij}| \geq a_{jj}$ ($a_{ii}, a_{jj} > 0$!) ■

- Manche Autoren (auch Rannacher) verlangen, dass reelle positiv definite Matrizen symmetrisch sind. Wir fordern das extra

Ü $A \in \mathbb{R}^{n \times n}$ A ist positiv definit genau dann wenn $A_S = \frac{1}{2}(A + A^T)$ positiv definit ist.

A pos. definit, dann sind alle Hauptuntermatrizen (definiert) positiv definit

4.7 Störungstheorie

12
27.10.09

Sei $A \in \mathbb{K}^{n \times n}$ invertierbar und $x, b \in \mathbb{K}^n$.

$F(A, b) = A^{-1}b$ ist der Lösungsoperator zur Gleichung $G(x) = Ax - b = 0$.

Betrachte die relative Kondition von F , wobei wir zunächst nur Änderungen in b zulassen wollen:

$$\frac{\|F(A, b + \delta b) - F(A, b)\|}{\|\delta b\|} \frac{\|b\|}{\|F(A, b)\|} = \frac{\|A^{-1}(b + \delta b) - A^{-1}b\|}{\|\delta b\|} \frac{\|b\|}{\|A^{-1}b\|}$$

$$\leq \frac{\|A^{-1}\| \|\delta b\|}{\|\delta b\|} \frac{\|A\| \|A^{-1}b\|}{\|A^{-1}b\|} = \|A^{-1}\| \|A\|$$

$\| \cdot \|$ sei verträglich
und submultiplikativ

Definition 4.15 Die Zahl

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

für irgendeine verträgliche und submultiplikative Matrixnorm heißt
Konditionszahl von A . ▣

Nun wollen wir auch Änderungen in A selbst zulassen. Sei $A \in \mathbb{K}^{n \times n}$
regulär. Frage: Wann ist $A + \delta A$ regulär?

Hilfssatz 4.16 $B \in \mathbb{K}^{n \times n}$ habe Norm $\|B\| < 1$. Dann ist $I + B$
regulär und es gilt

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

Beweis: i) Für alle $x \in \mathbb{K}^n$ gilt

$$\|x\| = \|x + Bx - Bx\| \leq \|(I+B)x\| + \|Bx\|$$

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

\uparrow
 $\|Bx\| \leq \|B\| \|x\|$

Damit gilt für $x \neq 0$, dass $(I+B)x \neq 0$ also $I+B$ regulär.

ii)

$$1 = \|I\| = \|(I+B)(I+B)^{-1}\| \stackrel{\text{erste Klammer auflösen}}{=} \|(I+B)^{-1} + B(I+B)^{-1}\|$$

wie oben: $\geq \|(I+B)^{-1}\| - \|B\| \|(I+B)^{-1}\| = \|(I+B)^{-1}\| (1 - \|B\|) > 0$

$$\|V+W-W\| \leq \|V+W\| + \|W\|$$

$$\Leftrightarrow \|V+W\| \geq \|V\| - \|W\|$$

Daraus folgt die Behauptung. da $\|B\| < 1$

Damit gilt das folgende

Satz 4.17 (Störungssatz) $A \in \mathbb{K}^{n \times n}$ sei regulär und $\|\delta A\| < \frac{1}{\|A^{-1}\|}$.

Dann ist $\tilde{A} = A + \delta A$ ebenfalls regulär und es gilt für den relativen Fehler des gestörten Systems $(A + \delta A)(x + \delta x) = b + \delta b$:

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

Beweis:

i) $A + \delta A = A(I + \underbrace{A^{-1}\delta A}_{\text{„B“}})$ und $\|A^{-1}\delta A\| \leq \|A^{-1}\| \|\delta A\| < \frac{\|A^{-1}\|}{\|A^{-1}\|} = 1$

Vor
 \downarrow

nach Hilfssatz 4.16 ist $A + \delta A$ regulär.

ii) $(A + \delta A)(x + \delta x) = b + \delta b$

$$\Leftrightarrow Ax + \delta Ax + (A + \delta A)\delta x = b + \delta b$$

\uparrow
 Ax Lösung von $Ax = b$

$$\Leftrightarrow (A + \delta A)\delta x = \delta b - \delta Ax$$

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

$$\|\delta x\| \leq \| (A+\delta A)^{-1} \| \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

Norm-
regeln

$$= \| [A(I+A^{-1}\delta A)]^{-1} \| \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$= \| (I+A^{-1}\delta A)^{-1} A^{-1} \| \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$\leq \| (I+A^{-1}\delta A)^{-1} \| \|A^{-1}\| \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

Hilfssatz 4.16

$$\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\delta A\|} \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\delta A\|} \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$\|A^{-1}\delta A\| \leq 1$$

$\|A^{-1}\| \|\delta A\| < 1$
nach Vor.!

mehr abziehen macht
1-x kleiner, also
den Bruch größer!

ausgeklammert

$$= \frac{\|A^{-1}\| \|A\| \|x\|}{1 - \|A^{-1}\| \|\delta A\| \|A\| \|A\|^{-1}} \left\{ \frac{\|\delta b\|}{\|A\| \|x\|} + \frac{\|\delta A\|}{\|A\|} \right\}$$

= 1
hinzugefügt

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

$$\|b\| = \|Ax\| \leq \|A\| \|x\|$$

Nenner verkleinern
macht Bruch größer



Beispiel 4.18

Es sei $\frac{\|\delta A\|}{\|A\|} \approx 10^{-k}$ und $\frac{\|\delta b\|}{\|b\|} \approx 10^{-k}$ sowie $\text{cond}(A) \approx 10^s$.

Weiter nehmen wir an, dass $10^s \cdot 10^{-k} \ll 1$ also etwa $s-k \leq -3$.

Dann gilt

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{10^s}{1 - \underbrace{10^s 10^{-k}}_{\ll 1}} \cdot 2 \cdot 10^{-k} \approx 10^{s-k} = 10^{-(k-s)}$$

Nenner ≈ 1

\Rightarrow man verliert s Stellen an Genauigkeit!

Eingabefehler ist in der k -ten Nachkommastelle, Fehler im Ergebnis ist in der $k+s$ -ten Stelle. ▣

Man kann zeigen, dass diese Abschätzung im wesentlichen "scharf" ist [Ra].
S. 117.

Beispiel 4.19 (Kondition + Determinante) Evtl als Übung.

(a) Betrachte die 2×2 Matrix

$$A = \begin{bmatrix} -1 & 1 \\ 1+\epsilon & -1 \end{bmatrix}, \quad A^{-1} = \frac{1}{\epsilon} \begin{bmatrix} -1 & 1 \\ -(1+\epsilon) & -1 \end{bmatrix}, \quad \epsilon = \det(A)$$

Es gilt

$$\|A\|_{\infty} = \max(2, 2+\epsilon), \quad \|A^{-1}\|_{\infty} = \frac{1}{\epsilon} \max(2, 2+\epsilon), \quad \text{cond}_{\infty}(A) = \frac{(2+\epsilon)^2}{\epsilon}$$

Also $\text{cond}_{\infty}(A) = O\left(\frac{1}{\det(A)}\right)$.

Dies ist aber nicht immer so:

(b)

$$B = \begin{bmatrix} 10^{-10} & 0 \\ 0 & 10^{10} \end{bmatrix}, \quad B^{-1} = \begin{bmatrix} 10^{10} & 0 \\ 0 & 10^{-10} \end{bmatrix}$$

$$\text{cond}_{\infty}(B) = \|B\|_{\infty} \|B^{-1}\|_{\infty} = 10^{-10} \cdot 10^{10} = 1 \text{ aber}$$

$$\det(B) = 10^{-20}!$$



4.6 Positiv definite Matrizen

24.10.09

Wichtige Klasse von Matrizen mit vorteilhaften Eigenschaften.

in \mathbb{K} !

Definition 4.11 Eine Matrix $A \in \mathbb{K}^{n \times n}$ heißt positiv definit wenn

a) $(Ax, x)_2 \in \mathbb{R} \quad \forall x \in \mathbb{K}^n \setminus \{0\} \quad (\mathbb{K} = \mathbb{C})$

b) $(Ax, x)_2 > 0 \quad \forall x \in \mathbb{K}^n \setminus \{0\}$ (eigentliche Bedingung)

Ziel: Charakterisierung positiv definiten Matrizen.

Im Fall $\mathbb{K} = \mathbb{C}$ bedeutet die Bedingung a) eine Einschränkung an die Matrix A .

Eigenschaft 4.12 Für $A \in \mathbb{C}^{n \times n}$ gilt A hermitesch genau dann wenn $(Ax, x)_2 \in \mathbb{R} \quad \forall x \in \mathbb{C}^n \setminus \{0\}$.

Beweis: " \Rightarrow " (als Ü-Aufgabe!)
 $(Ax, x)_2 \in \mathbb{R} \Leftrightarrow \underbrace{(Ax, x)_2}_{a+ib} = \overline{\underbrace{(Ax, x)_2}_{a-ib}}$ geht nur für $b=0$

also $(Ax, x)_2 \stackrel{\uparrow}{=} (x, Ax) \stackrel{(5.1)}{=} \overline{(Ax, x)}$ damit ist \Rightarrow gezeigt.
hermitesch 5.7

" \Leftarrow " Auftrennen in Real- und Imaginärteil:

$\mathbb{C}^n \ni x = a+ib \quad a, b \in \mathbb{R}^n \quad A = V+iW, \quad V, W \in \mathbb{R}^{n \times n}$

und ausrechnen:

$$(Ax, x)_2 = ((V+iW)(a+ib), a+ib)_2 = ((Va-Wb) + i(Vb+Wa), a+ib)_2$$
$$\stackrel{(x,y)_2 = \overline{(y,x)}_2!}{=} \underbrace{(Va-Wb, a)_2 + (Vb+Wa, b)_2}_{\text{Re}} + i \underbrace{[(Vb+Wa, a)_2 - (Va-Wb, b)_2]}_{\text{Im}}$$

also $\text{Im}((Ax, x)_2) = (Vb, a)_2 + (Wa, a)_2 - (Va, b)_2 + (Wb, b)_2 \stackrel{!}{=} 0 \quad (*)$

Fall I) $W=0$ d.h. $A=V$ (reell), aber $x=a+ib$ ist komplex. Dann reduziert sich (*) auf

$(Vb, a)_2 - (Va, b)_2 = 0 \quad \forall a, b \in \mathbb{R}^n \Leftrightarrow (Vb, a)_2 = (V^T b, a)_2 \quad \forall a, b \in \mathbb{R}^n$
alles reell

damit muss $V=V^T$ sein.

Fall II) Setze $b=0$, dann reduziert sich (*) auf

10
25.10.09

$$(W a, a)_2 = 0 \quad \forall a \in \mathbb{C}^n \setminus \{0\}$$

IIa: Setze $a = e^i \Rightarrow w_{ii} = 0$. Hier: $e_j^i = \delta_{ij}$ Kartesische Einheitsv.

IIb: Sei also $w_{ii} = 0 \forall i$ und setze $a = e^i + e^j$ ($j \neq i$)

$$\text{dann folgt } w_{ij} + w_{ji} = 0 \quad \forall i, j$$

$$\Leftrightarrow w_{ij} = -w_{ji}$$

Damit ist $A = V + iW$ mit $V = V^T$ und $W = -W^T$ also A hermitesch. \blacksquare

Positiv definite Matrizen in $\mathbb{C}^{n \times n}$ sind also immer hermitesch.

Positiv definite Matrizen in $\mathbb{R}^{n \times n}$ sind nicht notwendigerweise symmetrisch.

Charakterisierung über Eigenwerte:

Lemma 4.13 Eine hermitesche Matrix A ist genau dann positiv definit, wenn alle ihre (reellen) Eigenwerte positiv sind. Alle Hauptdiagonalelemente sind (reell und) positiv.

Beweis:

- A sei hermitesch mit lauter positiven Eigenwerte.

$x \in \mathbb{C}^n$ hat Darstellung $x = \sum \alpha_i w^i$ w^i : Eigenvektor zum EW λ_i .

$$(Ax, x)_2 = \sum_{i,j=1}^n \lambda_i \alpha_i \bar{\alpha}_j \underbrace{(w^i, w^j)_2}_{=\delta_{ij}} = \sum_{i=1}^n \lambda_i |\alpha_i|^2 > 0$$

- A sei hermitesch und positiv definit. A hat also Eigenpaare (λ_i, w^i) mit $\lambda_i \in \mathbb{R}$. Zu zeigen ist λ_i positiv

$$0 < (A w^i, w^i)_2 = (\lambda_i w^i, w^i)_2 = \lambda_i \underbrace{(w^i, w^i)_2}_{=1} = \lambda_i$$

- Setze $e^i \in \mathbb{R}^n$, $e_j^i = \delta_{ij}$ kartesische Einheitsvektoren.

$$0 < (A w_i, w_i)_2 = a_{ii} \in \mathbb{R}. \quad \blacksquare$$

Die Aussage gilt natürlich auch für reelle symmetrische und positiv definite Matrizen, da auch diese hermitesch sind.

Speziell für reelle Matrizen gilt:

11
25.10.09

Lemma 4.14 Sei $A \in \mathbb{R}^{n \times n}$ symmetrisch und positiv definit.

Dann liegt das betragsmäßig größte Element auf der Hauptdiagonalen.

Beweis: e^i sei wieder der i -te kanonische Einheitsvektor und a_{ij} , $j \neq i$, das betragsmäßig größte Element. Dann zeigt man den Widerspruch.

$$0 < (A(e^i - \text{sign}(a_{ij})e^j), e^i - \text{sign}(a_{ij})e^j)_2$$

$$= (Ae^i, e^i)_2 - 2 \text{sign}(a_{ij}) \underbrace{(Ae^i, e^j)}_{a_{ij}} + \text{sign}(a_{ij})^2 (Ae^j, e^j)$$

$$\stackrel{A=A^T}{=} a_{ii} - 2|a_{ij}| + a_{jj} \leq 0 \quad \downarrow$$

da $|a_{ij}| \geq a_{ii}$, $|a_{ij}| \geq a_{jj}$ ($a_{ii}, a_{jj} > 0$!) ■

- Manche Autoren (auch Rannacher) verlangen, dass reelle positiv definite Matrizen symmetrisch sind. Wir fordern das extra

Ü $A \in \mathbb{R}^{n \times n}$ A ist positiv definit genau dann wenn $A_S = \frac{1}{2}(A + A^T)$ positiv definit ist.

A pos. definit, dann sind alle Hauptuntermatrizen (definiert) positiv definit

4.7 Störungstheorie

12
27.10.09

Sei $A \in \mathbb{K}^{n \times n}$ invertierbar und $x, b \in \mathbb{K}^n$.

$F(A, b) = A^{-1}b$ ist der Lösungsoperator zur Gleichung $G(x) = Ax - b = 0$.

Betrachte die relative Kondition von F , wobei wir zunächst nur Änderungen in b zulassen wollen:

$$\frac{\|F(A, b + \delta b) - F(A, b)\|}{\|\delta b\|} \frac{\|b\|}{\|F(A, b)\|} = \frac{\|A^{-1}(b + \delta b) - A^{-1}b\|}{\|\delta b\|} \frac{\|b\|}{\|A^{-1}b\|}$$

$$\leq \frac{\|A^{-1}\| \|\delta b\|}{\|\delta b\|} \frac{\|A\| \|A^{-1}b\|}{\|A^{-1}b\|} = \|A^{-1}\| \|A\|$$

$\| \cdot \|$ sei verträglich
und submultiplikativ

Definition 4.15 Die Zahl

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

für irgendeine verträgliche und submultiplikative Matrixnorm heißt
Konditionszahl von A . ▣

Nun wollen wir auch Änderungen in A selbst zulassen. Sei $A \in \mathbb{K}^{n \times n}$
regulär. Frage: Wann ist $A + \delta A$ regulär?

Hilfssatz 4.16 $B \in \mathbb{K}^{n \times n}$ habe Norm $\|B\| < 1$. Dann ist $I + B$
regulär und es gilt

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

Beweis: i) Für alle $x \in \mathbb{K}^n$ gilt

$$\|x\| = \|x + Bx - Bx\| \leq \|(I+B)x\| + \|Bx\|$$

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

\uparrow
 $\|Bx\| \leq \|B\| \|x\|$

Damit gilt für $x \neq 0$, dass $(I+B)x \neq 0$ also $I+B$ regulär.

ii)

$$1 = \|I\| = \|(I+B)(I+B)^{-1}\| \stackrel{\text{erste Klammer auflösen}}{=} \|(I+B)^{-1} + B(I+B)^{-1}\|$$

wie oben: $\geq \|(I+B)^{-1}\| - \|B\| \|(I+B)^{-1}\| = \|(I+B)^{-1}\| (1 - \|B\|) > 0$

$$\|V+W-W\| \leq \|V+W\| + \|W\|$$

$$\Leftrightarrow \|V+W\| \geq \|V\| - \|W\|$$

Daraus folgt die Behauptung. da $\|B\| < 1$

Damit gilt das folgende

Satz 4.17 (Störungssatz) $A \in \mathbb{K}^{n \times n}$ sei regulär und $\|\delta A\| < \frac{1}{\|A^{-1}\|}$.

Dann ist $\tilde{A} = A + \delta A$ ebenfalls regulär und es gilt für den relativen Fehler des gestörten Systems $(A + \delta A)(x + \delta x) = b + \delta b$:

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

Beweis:

i) $A + \delta A = A(I + \underbrace{A^{-1}\delta A}_{\text{„B“}})$ und $\|A^{-1}\delta A\| \leq \|A^{-1}\| \|\delta A\| < \frac{\|A^{-1}\|}{\|A^{-1}\|} = 1$

Vor
 \downarrow

nach Hilfssatz 4.16 ist $A + \delta A$ regulär.

ii) $(A + \delta A)(x + \delta x) = b + \delta b$

$$\Leftrightarrow Ax + \delta Ax + (A + \delta A)\delta x = b + \delta b$$

\uparrow
 Ax Lösung von $Ax = b$

$$\Leftrightarrow (A + \delta A)\delta x = \delta b - \delta Ax$$

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

$$\|\delta x\| \leq \| (A+\delta A)^{-1} \| \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

Norm-
regeln

$$= \| [A(I+A^{-1}\delta A)]^{-1} \| \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$= \| (I+A^{-1}\delta A)^{-1} A^{-1} \| \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$\leq \underbrace{\| (I+A^{-1}\delta A)^{-1} \|}_{\text{Hilfssatz 4.16}} \|A^{-1}\| \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\delta A\|} \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\delta A\|} \{ \|\delta b\| + \|\delta A\| \|x\| \}$$

$$\|A^{-1}\delta A\| \leq 1$$

$\|A^{-1}\| \|\delta A\| < 1$
nach Vor.!

mehr abziehen macht
1-x kleiner, also
den Bruch größer!

ausgeklammert

$$= \frac{\|A^{-1}\| \|A\| \|x\|}{1 - \|A^{-1}\| \|\delta A\| \|A\| \|A\|^{-1}} \left\{ \frac{\|\delta b\|}{\|A\| \|x\|} + \frac{\|\delta A\|}{\|A\|} \right\}$$

= 1
hinzugefügt

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

$$\|b\| = \|Ax\| \leq \|A\| \|x\|$$

Nenner verkleinern
macht Bruch größer



Beispiel 4.18

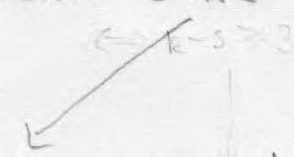
Es sei $\frac{\|\delta A\|}{\|A\|} \approx 10^{-k}$ und $\frac{\|\delta b\|}{\|b\|} \approx 10^{-k}$ sowie $\text{cond}(A) \approx 10^s$ ($s, k > 0$)

Weiter nehmen wir an, dass $10^s \cdot 10^{-k} \ll 1$ also etwa $s - k \leq -3$.

Dann gilt

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{10^s}{1 - \underbrace{10^s 10^{-k}}_{\ll 1}} \cdot 2 \cdot 10^{-k} \approx 10^{s-k} = 10^{-(k-s)}$$

Nenner ≈ 1



\Rightarrow man verliert s Stellen an Genauigkeit!

Eingabefehler ist in der k -ten Nachkommastelle, Fehler im Ergebnis ist in der $(k-s)$ -ten Stelle. ▣

Man kann zeigen, dass diese Abschätzung im wesentlichen "scharf" ist [Ra].
S. 117.

Beispiel 4.19 (Kondition + Determinante) Evtl als Übung.

(a) Betrachte die 2×2 Matrix

$$A = \begin{bmatrix} -1 & 1 \\ 1-\epsilon & -1 \end{bmatrix}, \quad A^{-1} = \frac{1}{\epsilon} \begin{bmatrix} -1 & -1 \\ -(1-\epsilon) & -1 \end{bmatrix}, \quad \epsilon = \det(A), \quad |\epsilon| < 1$$

Es gilt

$$\|A\|_{\infty} = \max(2, 2-\epsilon), \quad \|A^{-1}\|_{\infty} = \frac{1}{\epsilon} \max(2, 2-\epsilon), \quad \text{cond}_{\infty}(A) = \frac{(2+\epsilon)^2}{\epsilon}$$

Also $\text{cond}_{\infty}(A) = O\left(\frac{1}{\det(A)}\right)$.

Dies ist aber nicht immer so:

(b)

$$B = \begin{bmatrix} 10^{-10} & 0 \\ 0 & 10^{-10} \end{bmatrix}, \quad B^{-1} = \begin{bmatrix} 10^{10} & 0 \\ 0 & 10^{10} \end{bmatrix}$$

$$\text{cond}_{\infty}(B) = \|B\|_{\infty} \|B^{-1}\|_{\infty} = 10^{-10} \cdot 10^{10} = 1 \quad \text{aber}$$

$$\det(B) = 10^{-20}!$$

(c) $\text{cond}(\epsilon A) = \text{cond}(A) \quad \forall \epsilon \neq 0$, aber $\det(\epsilon A) = \epsilon^n \det(A)$ ▣

Beispiel 4.20 (oder Übung)

16
18.11.09

$$A = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}, \quad A^{-1} = \begin{bmatrix} d_1^{-1} & 0 \\ 0 & d_2^{-1} \end{bmatrix}, \quad d_1, d_2 \neq 0$$

$$\text{cond}_\infty(A) = \|A\|_\infty \|A^{-1}\|_\infty = \frac{\max(|d_1|, |d_2|)}{\min(|d_1|, |d_2|)}$$

$$\text{d.h. } \text{cond} \left(\begin{bmatrix} 10^{10} & 0 \\ 0 & 1 \end{bmatrix} \right) = 10^{10} !$$

Andererseits sind alle Gleichungen in $Ax = b$ unabhängig und

○ $x_i = b_i/d_i$ ist gut konditioniert.

Problem liegt hier in der Verwendung von Normen, die Abschätzung für einzelne Komponenten ist relativ schlecht.

Man lernt: schlechte Kondition muss nicht unbedingt Probleme beim Lösen des LGS bewirken, es ist jedoch „wahrscheinlich“.

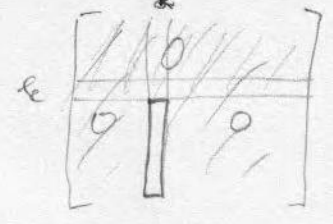
Wir definieren ^{für Schritt (2)} sog. Frobeniusmatrizen welche folgende Struktur haben:

$$G_k \in \mathbb{R}^{n \times n}, \quad 1 \leq k < n$$

$$G_k = I + G'_k \quad \text{mit} \quad (G'_k)_{ij} = 0 \quad \text{für} \quad (j \neq k) \vee (i \leq k)$$

d.h.

$$G_k = \begin{pmatrix} 1 & & & & & & \\ & \ddots & & & & & \\ & & 1 & & & & \\ & & & \ddots & & & \\ & & & & g_{k+1,k} & & \\ & & & & \vdots & & \\ & & & & & & g_{n,k} & \\ & & & & & & & 1 \end{pmatrix}$$



Hilfssatz 5.5 Die Frobeniusmatrizen haben folgende Eigenschaften:

i) Für $\tilde{A} = G_k A$ gilt

$$\tilde{a}_{ij} = \begin{cases} a_{ij} & i \leq k \\ a_{ij} + g_{ik} a_{kj} & i > k \end{cases}$$

(Dies ist gerade der Eliminationsschritt in der Gauß-Elimination wenn $g_{ik} := -a_{ik}$)

ii) $G'_k G'_k = 0$

iii) $G_k^{-1} = (I - G'_k)$

iv) $G_1 G_2 \dots G_k = I + \sum_{j=1}^k G'_j$, analog: $G_1^{-1} \dots G_k^{-1} = I - \sum_{j=1}^k G'_j$

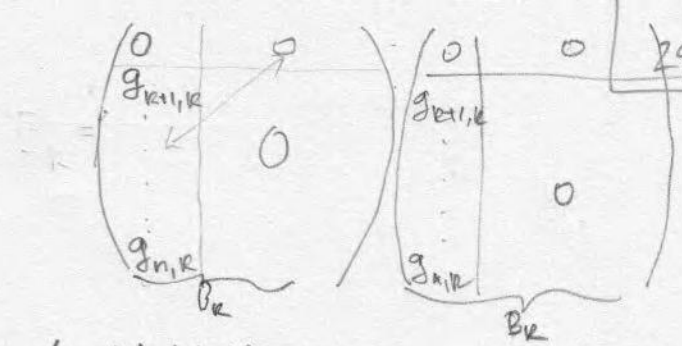
v) $\forall k < \alpha \leq \beta$: $P_{\alpha\beta} G_k = (I + P_{\alpha\beta} G'_k) P_{\alpha\beta}$ \rightarrow $I + P_{\alpha\beta} G'_k$ ist eine geänderte Frobeniusmatrix!

Beweis: Wird alles mit der Nullstruktur von G begründet:

i) $\tilde{a}_{ij} = \sum_{s=1}^n (G_k)_{is} a_{sj} = \begin{cases} a_{ij} \\ \sum_{s=k} g_{ik} a_{kj} + a_{ij} \end{cases}$

$i \leq k \rightarrow (G_k)_{ii} = 1$
da $s = k$
sonst:
 $(G_k)_{ik} = g_{ik}$
 $(G_k)_{ii} = 1$

ii) $G'_k = \begin{pmatrix} 0 & 0 \\ 0 & B_k \end{pmatrix}$
 Dimensions: $\begin{matrix} k-1 \\ n-k+1 \end{matrix}$



Es genügt $B_k B_k$ zu betrachten, das ist Null

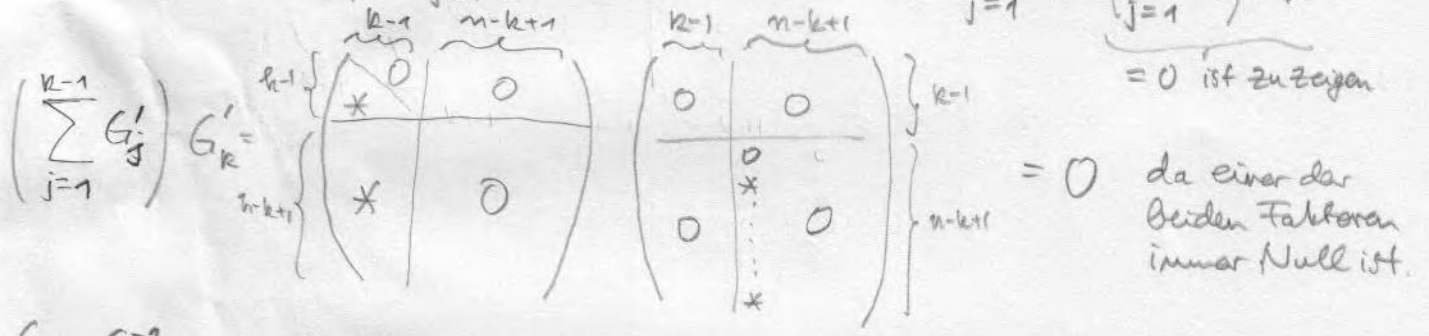
iii) $(I + G'_k)(I - G'_k) = I - G'_k + G'_k - \underbrace{G'_k G'_k}_{=0 \text{ wg ii}} = I$

iv) Per Induktion.

$k=1$: $G_1 = I + G'_1$ hat die geforderte Gestalt

$k-1 \rightarrow k$: Es sei also $G_1 \dots G_{k-1} = I + \sum_{j=1}^{k-1} G'_j$ bereits gezeigt

$G_1 \dots G_k = \left(I + \sum_{j=1}^{k-1} G'_j \right) (I + G'_k) = I + G'_k + \sum_{j=1}^{k-1} G'_j + \underbrace{\left(\sum_{j=1}^{k-1} G'_j \right) G'_k}_{=0 \text{ ist zu zeigen}}$

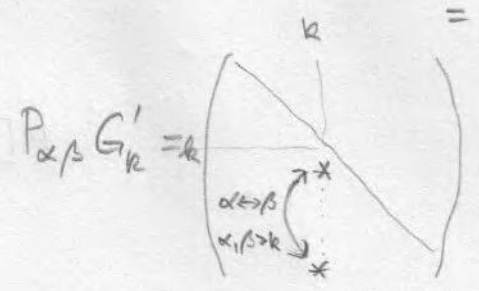


G_k, G_k^{-1} unterscheiden sich nur durch Vorzeichen Beweis hängt nur an Nullstruktur

v) $\boxed{k < \alpha \leq \beta}$: $P_{\alpha\beta} G_k = P_{\alpha\beta} (I + G'_k) = P_{\alpha\beta} + P_{\alpha\beta} G'_k$

$\begin{matrix} \alpha & \beta \\ \alpha & \beta \end{matrix} = P_{\alpha\beta} + \underbrace{P_{\alpha\beta} G'_k P_{\alpha\beta}}_{\substack{\text{nur Spalte tauscht zwei Spalten } \alpha, \beta > k, \\ k \text{ besetzt. Diese sind Null!}}}$

$= (I + P_{\alpha\beta} G'_k) P_{\alpha\beta}$



Somit ist $(I + P_{\alpha\beta} G'_k)$ wieder Frob. Matrix



Wir definieren das Produkt von Matrizen:

$$\prod_{i=a}^b B_i = B_b \cdots B_{a+1} B_a.$$

Beachte die Reihenfolge, da die Matrixmultiplikation nicht kommutativ ist.

Satz 5.6 (LR-Zerlegung)

Sei $A \in \mathbb{R}^{n \times n}$ regulär, dann gibt es eine Zerlegung

$$PA = LR$$

wobei

$$L = \begin{pmatrix} 1 & & & \\ l_{21} & & & \\ & \ddots & & \\ l_{n1} & \dots & l_{n,n-1} & 1 \end{pmatrix}, \quad R = \begin{pmatrix} r_{11} & & & \\ & r_{22} & & \\ & & \ddots & \\ 0 & & & r_{nn} \end{pmatrix},$$

und $P = \prod_{k=1}^{n-1} P_{k, \tau_k}$ ein Produkt von Permutationsmatrizen mit $\tau_k \geq k$.

Im Fall $P = I$ ist die Zerlegung eindeutig.

Beweis a) Wir behandeln zunächst den Fall $P = I$, d.h. ohne Zeilertausch.

Das Gaußsche Eliminationsverfahren lässt sich schreiben als

$$\begin{aligned} Ax &= b \\ \text{Schritt 1: } G_1 Ax &= G_1 b \quad \text{mit Frobeniusmatrix} \\ \text{Schritt 2: } G_2 G_1 Ax &= G_2 G_1 b \quad (G_1)_{i1} = -q_{i1} = -\frac{a_{i1}}{a_{11}} \end{aligned}$$

nach $n-1$ Schritten $\underbrace{G_{n-1} \cdots G_1}_R Ax = G_{n-1} \cdots G_1 b$ mit $(G_k)_{ik} = -q_{ik}$ aus GEM

Ergebnis der GEM: $G_{n-1} \cdots G_1 A = R$ (rechte obere Dreiecksmatrix).

Nun nutze Hilfssatz 5.5 i)

29.10.09

$$A = G_1^{-1} G_2^{-1} \dots G_{n-1}^{-1} R$$

$$\stackrel{\text{(iii)}}{=} (I - G'_1) \cdot (I - G'_2) \dots (I - G'_{n-1}) R$$

$$\stackrel{\text{(iv)}}{=} \underbrace{\left(I - \sum_{j=1}^{n-1} G'_j \right)}_{=: L} R$$

L hat die geforderte Gestalt.

b) Nun mit den Zeilenaustauschen. GEM liefert

$$G_{n-1} P_{n-1, \tau_{n-1}} \dots P_{2, \tau_2} G_1 P_{1, \tau_1} A = R$$

$\alpha < \beta$

τ_k ist der in Schritt k bestimmte Tauschungsindex $\tau_k \geq k$.

Nende sukzessive HS 5.5 v) an

$$G_{n-1} P_{n-1, \tau_{n-1}} \dots P_{2, \tau_2} G_1 P_{1, \tau_1} A$$

$$= G_{n-1} P_{n-1, \tau_{n-1}} P_{3, \tau_3} G_2 \underbrace{\left(I + P_{2, \tau_2} G'_1 \right)}_{\text{wieder Frobenius m.}} P_{2, \tau_2} P_{1, \tau_1} A$$

$$= G_{n-1} P_{n-1, \tau_{n-1}} P_{4, \tau_4} G_3 \left(I + P_{3, \tau_3} G'_2 \right) \left(I + P_{3, \tau_3} P_{2, \tau_2} G'_1 \right) P_{3, \tau_3} P_{2, \tau_2} P_{1, \tau_1} A$$

$$= \prod_{k=1}^{n-1} \left(I + \underbrace{\left(\prod_{\alpha=k+1}^{n-1} P_{\alpha, \tau_\alpha} \right)}_{\substack{\text{alle sp\u00e4teren} \\ \text{Tauschoperationen}}} G'_k \right) \underbrace{\prod_{k=1}^{n-1} P_{k, \tau_k}}_{=: P} A = R$$

nach rechts bringen

$$\Leftrightarrow PA = \underbrace{\prod_{k=n-1}^1 \left(I - \left(\prod_{\alpha=k+1}^{n-1} P_{\alpha, \tau_\alpha} \right) G'_k \right)}_{\substack{\text{wg} \\ \text{Inverse}}} R$$

$=: L$

Beobachtung: Transformation auf obere Dreiecksgestalt mittels GEM
äquivalent zu LR-Zerlegung berechnen + $Ly = Pb$ lösen.

Somit: Der Aufwand zur Berechnung der LR-Zerlegung ist

$$N_{LR}(n) = \frac{2}{3} n^3 + O(n^2).$$

LR-Zerlegung ist insbesondere interessant, wenn $Ax = b_i$ zu mehreren
rechten Seiten b_i zu lösen ist.

Zur Permutation

$$P = P_{n-1, r_{n-1}} \cdots P_{2, r_2} P_{1, r_1} \quad \text{mit } r_k \geq k$$

speichert man nicht als Matrix sondern man speichert nur
die Zahlen r_1, \dots, r_{n-1} in einem Vektor (Feld).

Algorithmus zur LR-Zerlegung

Input: $A \in \mathbb{R}^{n \times n}$ (wird überschrieben)

Output: $L \in \mathbb{R}^{n \times n}$ in $a_{ij}, j < i$, $l_{ii} = 1$ implizit

$R \in \mathbb{R}^{n \times n}$ in $a_{ij}, j \geq i$

$p: \{0, \dots, n-1\} \rightarrow \{0, \dots, n-1\}$

for ($k=1; k \leq n; k=k+1$) {

Finde $r \in \{k, \dots, n\}$ sodass $a_{rk} \neq 0$; sonst Fehler;

if ($r \neq k$) // tausche Zeilen

for ($j=1; j \leq n; j=j+1$) {

$t = a_{kj}; a_{kj} = a_{rj}; a_{rj} = t;$

}

$p[k] = r$; // merke Permutation

for ($i=k+1; i \leq n; i=i+1$) {

$a_{ik} = a_{ik} / a_{kk}$

for ($j=k+1; j \leq n; j=j+1$)

$a_{ij} = a_{ij} - a_{ik} a_{rj}$

}

}

5.4 Rundungsfehleranalyse der Gauß-Elimination (LR-Zerlegung)

10.11.09

Absolutwertnotation

Für $A \in \mathbb{R}^{m \times n}$ ist

$$B = |A| \Rightarrow b_{ij} = |a_{ij}| \quad 1 \leq i \leq m, 1 \leq j \leq n$$

Die Abbildung $rd: \mathbb{R} \rightarrow \mathbb{F}$ erweitern wir entsprechend auf $\mathbb{R}^n, \mathbb{R}^{m \times n}$.

Dann gilt

$$rd(A) = A + A' \quad \text{mit} \quad |A'| \leq |A| \epsilon_{ps}$$

dies sind $m \cdot n$ Ungleichungen

Wdh. aus Analyse von rd :

$$rd(a_{ij}) = a_{ij}(1 + \epsilon_{ij}) = a_{ij} + \underbrace{a_{ij} \epsilon_{ij}}_{a'_{ij}}$$

$$|a'_{ij}| \leq |a_{ij}| \epsilon_{ps}$$

fl-Notation → nächstes Mal früher einführen.
Stör: gl

Sei E eine Formel, dann bezeichne $fl(E)$ eine Berechnung der Formel E in

Fließkommaarithmetik (dabei wird die Ausführungsreihenfolge meist offensichtlich sein, sonst ist sie angegeben).
(links nach rechts, Klammern)

Somit ist zum Beispiel für

$$fl(A+B) = \underbrace{(A+B)}_{\text{exakte Rech. in } \mathbb{R}} + H \quad \text{mit} \quad |H| \leq \epsilon_{ps} (|A+B|)$$

$A, B \in \mathbb{F}^{m \times n}$

(folgt aus)

$$fl(a_j + b_j) = a_j \oplus b_j$$

$$= (a_j + b_j)(1 + \epsilon_{ij}) \quad |\epsilon_{ij}| \leq \epsilon_{ps}$$

$$= (a_j + b_j) + \epsilon_{ij}(a_j + b_j)$$

Man kann auch

$$fl(\sqrt{x})$$

oder $fl(\sin(x))$

schreiben.

Rückwärtsanalyse

10.11.09

Bisher haben ^{wir} die „Vorwärtsanalyse“ der Rundungsfehler betrieben.
Z.B. gilt für das Skalarprodukt $x^T y := x^T y$

$$|fl(x^T y) - x^T y| \leq n \text{ eps} |x|^T |y| + O(\text{eps}^2) \quad \text{absolut}$$

oder

$$\frac{|fl(x^T y) - x^T y|}{|x^T y|} \leq n \text{ eps} \frac{|x|^T |y|}{|x^T y|} + O(\text{eps}^2) \quad \text{relativ.}$$

Eine Alternative ist die sog. „Rückwärtsanalyse“.

Dort versucht man das Fließkommaergebnis als exakter Ergebnis eines modifizierten Ausdrucks zu schreiben.

Beispiel 5.7 Betrachte Lösung des LGS $Ax = b$. Die GEM berechnet die numerische Lösung \hat{x} .

Vorwärtsanalyse: $\|\hat{x} - x\| \leq F(\text{eps}, n, A, b)$

Rückwärtsanalyse: $(A+E)\hat{x} = b$ mit $\|E\| \leq F'(\text{eps}, n, A)$

Mit dem Störungssatz 4.17 und $\|E\|_\infty = \| |E| \|_\infty$ folgt dann für d. Rundungsfehler:

$$\frac{\|\hat{x} - x\|_\infty}{\|x\|_\infty} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\| |E| \|_\infty}{\|A\|_\infty}} \underbrace{\frac{\| |E| \|_\infty}{\|A\|_\infty}}$$

Rounding: $rd(A) = A + \delta A$
 $|\delta A| \leq |A| \cdot \text{eps}$
Störungssatz
 $\frac{\|\delta A\|_\infty}{\|A\|_\infty} = \frac{\| |\delta A| \|_\infty}{\|A\|_\infty} \leq \frac{\|A\|_\infty \text{eps}}{\|A\|_\infty} = \text{eps}$

ist noch entsprechend abzuschätzen, ideal wäre $\| |E| \|_\infty \leq n^2 \text{eps} \|A\|_\infty$

Somit ist ein direkter Vergleich mit der Konditionsanalyse möglich.

Rückwärtsanalyse des Skalarproduktes.

Hilfssatz 5.8 Es gilt ^{für} $x, y \in \mathbb{F}^n$:

3
10.11.09

$$\hat{s} = fl(x^T y) = (x+f)^T y, \quad |f| \leq n \text{ eps } |x| + O(\text{eps}^2)$$

Beweis: Induktion über n .

$n=1$: $\hat{s}_1 = fl(x_1 y_1) = x_1 y_1 (1 + \delta_1) = (x_1 + \underbrace{x_1 \delta_1}_{=: f_1}) y_1$

also $|f_1| = |\delta_1| |x_1| \leq \text{eps } |x_1|$.

$n \geq 2$: $\hat{s}_n = fl(x^T y) = fl(\underbrace{fl(\tilde{x}^T \tilde{y})}_{x_n y_n} + x_n y_n)$

$x = \begin{pmatrix} \tilde{x} \\ x_n \end{pmatrix}, y = \begin{pmatrix} \tilde{y} \\ y_n \end{pmatrix}$

$$= (fl(\tilde{x}^T \tilde{y}) + fl(x_n y_n)) (1 + \epsilon_n)$$

$$= ((\tilde{x} + \tilde{f})^T \tilde{y} + (x_n y_n (1 + \delta_n))) (1 + \epsilon_n)$$

$$= (\tilde{x} + \tilde{f})^T \tilde{y} (1 + \epsilon_n) + x_n y_n (1 + (\delta_n + \epsilon_n) + \delta_n \epsilon_n)$$

$$= (\tilde{x} + \tilde{f} + \underbrace{\epsilon_n \tilde{x} + \epsilon_n \tilde{f}}_{=: f})^T \tilde{y} + (x_n + (\delta_n + \epsilon_n)x_n + \delta_n \epsilon_n x_n) y_n$$

$$= \left[\underbrace{\begin{pmatrix} \tilde{x} \\ x_n \end{pmatrix}}_{=: x} + \underbrace{\begin{pmatrix} \tilde{f} + \epsilon_n \tilde{x} + \epsilon_n \tilde{f} \\ (\delta_n + \epsilon_n)x_n + \delta_n \epsilon_n x_n \end{pmatrix}}_{=: f} \right] \underbrace{\begin{pmatrix} \tilde{y} \\ y_n \end{pmatrix}}_{=: y}$$

mit $|\tilde{f} + \epsilon_n \tilde{x} + \epsilon_n \tilde{f}| \leq (n-1) \text{ eps } |\tilde{x}| + \text{eps } |\tilde{x}| + O(\text{eps}^2) \leq n \text{ eps } |\tilde{x}| + O(\text{eps}^2)$

und $|(\delta_n + \epsilon_n)x_n + \delta_n \epsilon_n x_n| \leq 2 \text{ eps } |x_n| + O(\text{eps}^2)$

Wg $n \geq 2$ gilt also $|f| \leq n \text{ eps } |x| + O(\text{eps}^2)$ ▣

Bemerkung 5.9 $n \text{ eps}$ in 5.8 ist der schlechteste Fall und sehr pessimistisch. Rundungsfehler in einer Operation hängt von den Argumenten ab und sie können sich auch wegheben (positiv/negativ!). Besser wäre eine statistische Betrachtung. ▣

Satz 5.10 (Lösung von Dreieckssystemen)

10.11.09

Es seien \hat{x} bzw \hat{y} die numerischen Lösungen des unteren bzw oberen Dreieckssysteme $Lx=b$ und $Ry=c$. Dann gilt

$$(L+F)\hat{x} = b \quad |F| \leq n \text{ eps } |L| + O(\text{eps}^2)$$

$$(R+G)\hat{y} = c \quad |G| \leq n \text{ eps } |R| + O(\text{eps}^2)$$

Beweis: (evtl als Übung). Induktion über n .

$n=1$ $l_{11}x_1 = b_1 \rightarrow \hat{x}_1 = fl(b_1/l_{11}) = (b_1/l_{11})(1+\epsilon_1)$

$$\Leftrightarrow \frac{1}{1+\epsilon_1} \hat{x}_1 = \frac{b_1}{l_{11}} \Leftrightarrow \frac{l_{11}}{1+\epsilon_1} \hat{x}_1 = b_1$$

$$\frac{1}{1+\epsilon_1} = 1 +$$

$$\frac{\partial}{\partial \epsilon} \frac{1}{1+\epsilon} = -\frac{1}{(1+\epsilon)^2}$$

$$\Downarrow \Leftrightarrow (1 - \epsilon_1 + \tilde{R}(\epsilon^2)) l_{11} \hat{x}_1 = b_1$$

$$\Leftrightarrow (l_{11} \underbrace{-\epsilon_1 l_{11} + l_{11} \tilde{R}(\epsilon^2)}_{=f_1}) \hat{x}_1 = b_1$$

mit $|f_1| \leq \text{eps } |l_{11}| + O(\text{eps}^2)$

$n \geq 2$ $Lx=b \Leftrightarrow \begin{matrix} n-1 & 1 \\ L_1 & 0 \\ 1 & v^T \end{matrix} \begin{pmatrix} x_1 \\ w \end{pmatrix} = \begin{pmatrix} b_1 \\ \beta \end{pmatrix}$

Berechne Lösung \hat{x}_1 von $L_1x_1=b_1$ und rechne

$$\hat{w} = fl((\beta - v^T \hat{x}_1)/\alpha)$$

↙ Fehler in α

$$= fl((\beta - v^T \hat{x}_1)/\alpha) (1+\epsilon_n)$$

↙ Skalarprodukt

$$= ((\beta - fl(v^T \hat{x}_1))(1+\delta_n)/\alpha) (1+\epsilon_n) = ((\beta - fl(v^T \hat{x}_1))/\alpha) (1+\delta_n)(1+\epsilon_n)$$

HS. 5.8

$$\Downarrow = ((\beta - (v+f)^T \hat{x}_1)/\alpha) (1 + (\epsilon_n + \delta_n) + \epsilon_n \delta_n)$$

Mit $\frac{1}{1 + (\varepsilon_n + \delta_n) + \varepsilon_n \delta_n} = 1 - (\varepsilon_n + \delta_n) + R(\varepsilon_n^2 + \delta_n^2 + \varepsilon_n \delta_n)$ gilt

10.11.09

$$(1 - (\varepsilon_n + \delta_n) + R(\dots)) \hat{w} = (\beta - (v + f)^T \hat{x}_1) / \alpha$$

$$\Leftrightarrow (v + f)^T \hat{x}_1 + (1 - (\varepsilon_n + \delta_n) + R(\dots)) \alpha \hat{w} = \beta$$

Mit der Induktionsvoraussetzung $(L_1 + F_1) \hat{x}_1 = b_1$ folgt also für \hat{x} :

$$\begin{pmatrix} L_1 + F_1 & 0 \\ (v + f)^T & (1 - (\varepsilon_n + \delta_n) + R(\dots)) \alpha \end{pmatrix} \begin{pmatrix} \hat{x}_1 \\ \hat{w} \end{pmatrix} = \begin{pmatrix} b_1 \\ \beta \end{pmatrix}$$

$$\Leftrightarrow \left[\underbrace{\begin{pmatrix} L_1 & 0 \\ v^T & \alpha \end{pmatrix}}_{= L} + \underbrace{\begin{pmatrix} F_1 & 0 \\ f^T & (-(\varepsilon_n + \delta_n) + R(\dots)) \alpha \end{pmatrix}}_{= F} \right] \begin{pmatrix} \hat{x}_1 \\ \hat{w} \end{pmatrix} = \begin{pmatrix} b_1 \\ \beta \end{pmatrix}$$

Wegen $|F_1| \leq (n-1) \text{eps} |L| + O(\text{eps}^2)$

$|f^T| \leq (n-1) \text{eps} |v^T| + O(\text{eps}^2)$

$|(-(\varepsilon_n + \delta_n) + R(\dots)) \alpha| \leq 2 \text{eps} |\alpha| + O(\text{eps}^2)$

hier kommt das n-1 her!
nicht durch die Rekursion.
es geht also nicht besser

gilt für $n \geq 2$ dann

$$|F| \leq n \text{eps} |L| + O(\text{eps}^2)$$

(möglich wäre auch $\max(n-1, 2)$, ist aber nicht wirklich besser).

Analog für $Ry = c$.



Satz 5.11 (Rückwärtsanalyse der LU-Zerlegung)

11.11.09

Sei $A \in \mathbb{F}^{n \times n}$. Es werde die LR-Zerlegung ohne Pivotierung berechnet. Dann gilt für die numerisch berechneten \hat{L}, \hat{R} :

$$\hat{L} \hat{R} = A + H \quad \text{mit} \quad |H| \leq 3(n-1)\epsilon_{ps}(|A| + |\hat{L}| |\hat{R}|) + O(\epsilon_{ps}^2)$$

Beweis: [Golub/van Loan, THM 3.3.1] Induktion über n .

$n=1$: Es ist $\hat{L}_{11} = 1$ und $\hat{R}_{11} = a_{11}$ und damit $H = 0$.

$n \geq 2$: Schreibe

$$A = \begin{bmatrix} \alpha & w^T \\ v & B \end{bmatrix} \begin{matrix} 1 \\ n-1 \end{matrix}$$

Dann macht die LR-Zerlegung:

a) $\hat{z} = fl(v/\alpha)$

b) $\hat{A}_1 = fl(B - \hat{z}w^T)$

c) Berechne LR-Zerlegung von \hat{A}_1

Fehler in \hat{z} :

(P₁) $\hat{z} = \frac{v}{\alpha} + f$ mit $|f| \leq \epsilon_{ps} \frac{|v|}{|\alpha|}$

Fehler in \hat{A}_1 :

$\hat{A}_1 = fl(B - \hat{z}w^T)$

← exakt gerundet

Bistuff $\rightarrow B - fl(\hat{z}w^T) + G$

$|G| \leq \epsilon_{ps} |B - fl(\hat{z}w^T)|$

$= B - (\hat{z}w^T + G') + G$
 $=: F$

$|G'| \leq \epsilon_{ps} |\hat{z}w^T| \leq \epsilon_{ps} |\hat{z}| |w|^T$

← outer product! genau eine Operation pro matrixeintrag.

(P₂) $\hat{A}_1 = B - \hat{z}w^T + F$

mit $|F| = |-G' + G| \leq |G'| + |G|$

$\leq \frac{\epsilon_{ps} |\hat{z}| |w|^T}{|G'|} + \epsilon_{ps} |B - \hat{z}w^T - G'|$

$\leq \epsilon_{ps} |\hat{z}| |w|^T + \epsilon_{ps} (|B| + |\hat{z}| |w|^T) + \epsilon_{ps}^2 |\hat{z}| |w|^T$

also $|F| \leq 2\epsilon_{ps} (|B| + |\hat{z}| |w|^T) + O(\epsilon_{ps}^2)$

Nun wird A_1 LR-zersetzt und es gilt die Induktionsannahme 10.11.09

$$(D_3) \quad \hat{L}_1 \hat{R}_1 = \hat{A}_1 + H_1 \quad \text{mit} \quad |H_1| \leq 3(n-2)\epsilon \left(|\hat{A}_1| + |\hat{L}_1| |\hat{R}_1| \right) + O(\epsilon^2)$$

also ist die Blockform der LR-Zerlegung von A :

$$\hat{L} \hat{R} = \begin{bmatrix} 1 & 0 \\ \hat{z} & \hat{L}_1 \end{bmatrix} \begin{bmatrix} \alpha & w^T \\ 0 & \hat{R}_1 \end{bmatrix} = \begin{bmatrix} \alpha & w^T \\ \hat{z}\alpha & \hat{z}w^T + \hat{L}_1 \hat{R}_1 \end{bmatrix}$$

$= \hat{A}_1 + H_1$

Darstellung
von
oben
linker
 $(D_1), (D_2), (D_3)$

$$= \begin{bmatrix} \alpha & w^T \\ (\frac{\sigma}{\alpha} + f)\alpha & \hat{z}w^T + \underbrace{(B - \hat{z}w^T + F)}_{=\hat{A}_1} + H_1 \end{bmatrix}$$

$$= \underbrace{\begin{bmatrix} \alpha & w^T \\ v & B \end{bmatrix}}_{=A} + \underbrace{\begin{bmatrix} 0 & 0 \\ \alpha f & H_1 + F \end{bmatrix}}_{=:H}$$

Damit hat man schon mal die Form $\hat{L}\hat{U} = A + H$. Bleibt noch H abzuschätzen.

In H steckt $H_1 + F$, in $|H_1|$ steckt $|\hat{A}_1|$ also:

$$|\hat{A}_1| = |B - \hat{z}w^T + F| \leq |B| + |\hat{z}| |w|^T + |F|$$

$$\text{Absh. für } |F| \rightarrow \leq |B| + |\hat{z}| |w|^T + 2\epsilon \left(|B| + |\hat{z}| |w|^T \right) + O(\epsilon^2)$$

$$\leq (1 + 2\epsilon) \left(|B| + |\hat{z}| |w|^T \right) + O(\epsilon^2)$$

$$|H_1 + F| \leq |H_1| + |F|$$

$$\text{Indukt. annahme} \rightarrow \leq 3(n-2)\epsilon \left(|\hat{A}_1| + |\hat{L}_1| |\hat{R}_1| \right) + 2\epsilon \left(|B| + |\hat{z}| |w|^T \right) + O(\epsilon^2)$$

$$\text{Absh. für } |\hat{A}_1| \text{ ein} \rightarrow \leq 3(n-2)\epsilon \left[(1 + 2\epsilon) \left(|B| + |\hat{z}| |w|^T \right) + |\hat{L}_1| |\hat{R}_1| \right] + 2\epsilon \left(|B| + |\hat{z}| |w|^T \right) + O(\epsilon^2)$$

$$\leq 3(n-1)\epsilon \left[|B| + |\hat{z}| |w|^T + |\hat{L}_1| |\hat{R}_1| \right] + O(\epsilon^2)$$

Und damit

17.11.09

$$|H| = \begin{bmatrix} 0 & 0 \\ |\alpha f| & |H_1 + F| \end{bmatrix} \leq \begin{bmatrix} 0 & 0 \\ \epsilon |\nu| & 3(n-1)\epsilon (|B| + |\hat{z}| |w|^T + |\hat{L}_1| + |\hat{R}_1|) \end{bmatrix} + O(\epsilon^2)$$

(0,1)

$$\leq 3(n-1)\epsilon \begin{bmatrix} 0 & 0 \\ |\nu| & |B| + |\hat{z}| |w|^T + |\hat{L}_1| + |\hat{R}_1| \end{bmatrix} + O(\epsilon^2)$$

$$\leq 3(n-1)\epsilon \left(\underbrace{\begin{bmatrix} |\alpha| & |w|^T \\ |\nu| & |B| \end{bmatrix}}_{|A|} + \underbrace{\begin{bmatrix} 1 & 0 \\ |\hat{z}| & |\hat{L}_1| \end{bmatrix}}_{|\hat{L}|} \underbrace{\begin{bmatrix} |\alpha| & |w|^T \\ 0 & |\hat{R}_1| \end{bmatrix}}_{|\hat{R}|} \right)$$

○ Terme sind hinzugefügt.
0 ≤ (1.1)

Nun sind noch die Dreieckssysteme aufzulösen.

Seien \hat{L} und \hat{R} die numerisch berechnete LR-Zerlegung von $A \in \mathbb{F}^{n \times n}$ aus Satz 5.11. Sei weiter $\hat{y} \in \mathbb{F}^n$ die numerische Lösung von $\hat{L}y = b$ und schließlich $\hat{x} \in \mathbb{F}^n$ die numerische Lösung von $\hat{R}x = \hat{y}$. Dann gilt für \hat{x} die Beziehung

$$(A+E)\hat{x} = b$$

mit

$$|E| \leq n \text{ eps} (3|A| + 5|\hat{L}||\hat{R}|) + O(\text{eps}^2).$$

Beweis:

Wg Satz 5.10 gilt

$$\begin{aligned} (\hat{L}+F)\hat{y} &= b & |F| &\leq n \text{ eps} |\hat{L}| + O(\text{eps}^2) \\ (\hat{R}+G)\hat{x} &= \hat{y} & |G| &\leq n \text{ eps} |\hat{R}| + O(\text{eps}^2) \end{aligned}$$

Einsetzen liefert:

$$(\hat{L}+F)\hat{y} = (\hat{L}+F)(\hat{R}+G)\hat{x} = (\underbrace{\hat{L}\hat{R}}_{=A+H} + F\hat{R} + \hat{L}G + FG)\hat{x} = b$$

Wg Satz 5.11 gilt

$$\hat{L}\hat{R} = A+H \quad |H| \leq 3(n-1) \text{ eps} (|A| + |\hat{L}||\hat{R}|) + O(\text{eps}^2)$$

also

$$(A+E)\hat{x} = b \quad \text{mit} \quad E = H + F\hat{R} + \hat{L}G + FG$$

und

$$|E| \leq |H| + |F||\hat{R}| + |\hat{L}||G| + O(\text{eps}^2)$$

$$\leq \underbrace{3(n-1) \text{ eps} (|A| + |\hat{L}||\hat{R}|)}_{|H|} + \underbrace{n \text{ eps} |\hat{L}'||\hat{R}|}_{|F|} + \underbrace{n \text{ eps} |\hat{L}'||\hat{R}|}_{|G|} + O(\text{eps}^2)$$

$$\leq 3n \text{ eps} |A| + 5n \text{ eps} |\hat{L}'||\hat{R}| + O(\text{eps}^2). \quad \square$$

Folgerung

Nach Satz 5.12 ist \hat{x} exakte Lösung des modifizierten Systems

$$(A+E)\hat{x} = b.$$

Mit dem Störungsatz gilt dann in $\|\cdot\|_\infty$ -Norm (Beachte: $\|B\|_\infty = \| |B| \|_\infty$)

$$\frac{\|\hat{x} - x\|_\infty}{\|x\|_\infty} \leq$$



$$\text{Cond}(A) \cdot \left\{ \underbrace{-3n \text{ eps} \frac{\|A\|_\infty}{\|A\|_\infty}}_{=1} + 5n \text{ eps} \frac{\|\hat{L}\| \|\hat{R}\|}{\|A\|_\infty} + O(\text{eps}^2) \right\}$$

Vergleichbar mit Rundungsfehler $\| \text{rd}(A) \| / \|A\|$
O.K.

Problem!

Nehmen an, dass $\|E\|_\infty / \|A\|_\infty \ll 1$ (brauchen ohnehin $\|E\| < \frac{1}{\|A^{-1}\|}$)

\rightarrow Numer im Vorfaktor ist ≈ 1

- Erster Term vergleichbar mit dem aus der Konditionsanalyse
- Zweiter Term ist möglicherweise problematisch.

\hat{L}^1 enthält Einträge der Form $\frac{\tilde{a}_{ij}}{\tilde{a}_{ii}}$, also $\frac{1}{\text{Pivotelement}}$.

|Pivotelement| klein \Rightarrow $|\hat{L}^1|$ groß \Rightarrow großer Rundungsfehler!

- Dies kann trotz guter Kondition von A passieren!

Beispiel: $A = \begin{bmatrix} \epsilon & 1 \\ 1 & 0 \end{bmatrix}$

\Rightarrow Gauß-Elimination (LR-Zerlegung) ist in dieser Form nicht numerisch stabil!

5.5 Pivotalisierung

16.11.09

Die Rundungsfehleranalyse in Satz 5.12 führt auf den unvorteilhaften Term $\|\hat{L}\|_{\infty}$.

Mit der Wahl von r im Gauß-Algorithmus so dass

$$|a_{rk}^{(k)}| \geq |a_{ik}^{(k)}| \quad \forall k \leq i \leq n$$

gilt dann

$$|\hat{l}_{ij}| \leq 1 \quad \text{und damit} \quad \|L\|_{\infty} \leq n.$$

Diese Wahl nennt man „Spaltenpivotalisierung“.

Beispiel 5.13 Aus [GVL]

$$\begin{bmatrix} -10^{-5} & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

In exakter Arithmetik führt das Gaußsche Verfahren nach Elimination von a_{21} auf

$$\begin{bmatrix} -10^{-5} & 1 \\ 0 & 1 + 2 \cdot 10^5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \cdot 10^5 \end{bmatrix}$$

mit der Lösung

$$x_1 = -0.4999975, \quad x_2 = 0.999995.$$

Nun führen wir das Verfahren in $\mathbb{F}(10, 4, 1)$ durch. Beim Multiplikator

$$q_{21} = (0.2 \cdot 10^1) \ominus (-0.1 \cdot 10^{-4}) = -0.2 \cdot 10^6$$

ergibt sich kein Rundungsfehler.

Für das neue a_{22} ergibt sich

$$\begin{aligned} a_{22}^{(1)} &= 0.1 \cdot 10^1 \ominus (-0.2 \cdot 10^6) \ominus (0.1 \cdot 10^1) \\ &= 0.1 \cdot 10^1 \oplus 0.2 \cdot 10^6 = \boxed{0.2 \cdot 10^6}. \end{aligned}$$

Hier wurde auf vier Stellen gerundet.

Damit ergibt sich (ohne Fehler)

$$b_2^{(1)} = -(-0.2 \cdot 10^6) \ominus (0.1 \cdot 10^1) = 0.2 \cdot 10^6$$

und

$$\begin{aligned} x_2 &= b_2^{(1)} \ominus a_{22}^{(1)} = 0.2 \cdot 10^6 \ominus 0.2 \cdot 10^6 = \boxed{1}, \\ x_1 &= (0.1 \cdot 10^1 \ominus 0.1 \cdot 10^1 \ominus 1) \ominus (-0.1 \cdot 10^{-4}) = \boxed{0}. \end{aligned}$$

Es ist also *keine* Stelle im Ergebnis korrekt obwohl nur an einer *einzigen* Stelle (in der Berechnung von $a_{22}^{(1)}$) ein Rundungsfehler eingeführt wurde.

16.11.09

Darüberhinaus überprüfe man, dass für die Kondition von A gilt:

$$\kappa(A) = 3 .$$

Demnach ist das System gut konditioniert! Der Algorithmus, so wie er ist, ist numerisch nicht stabil.

Das Problem ist offensichtlich der große Multiplikator q_{21} der aus dem sehr kleinen a_{11} resultiert und der dafür sorgt, dass das ursprüngliche a_{22} in $a_{22}^{(1)}$ vollkommen ignoriert wird.

Im Prinzip haben wir in Fließkommaarithmetik das System

$$\begin{bmatrix} -10^{-5} & 1 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

exakt gelöst, was eine völlig andere Lösung hat als das ursprüngliche (11.1) (Rückwärtsanalyse).

Der große Multiplikator kann ganz einfach vermieden werden indem man eine Zeilenvertauschung durchführt, d. h. wir lösen

$$\begin{bmatrix} 2 & 1 \\ -10^{-5} & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} . \quad (11.2)$$

Nun erhält man

$$q_{21} = -0.1 \cdot 10^{-4} \oslash 0.2 \cdot 10^1 = -0.5 \cdot 10^{-5},$$

$$a_{22}^{(1)} = 0.1 \cdot 10^1 \ominus (-0.5 \cdot 10^{-5}) \odot 0.1 \cdot 10^1 = 0.1 \cdot 10^1 \oplus 0.5 \cdot 10^{-5} = 0.1 \cdot 10^1,$$

$$b_2^{(1)} = 0.1 \cdot 10^1 \ominus 0.5 \cdot 10^{-5} \odot 0 = 0.1 \cdot 10^1,$$

$$x_2 = 0.1 \cdot 10^1 \oslash 0.1 \cdot 10^1 = \boxed{1},$$

$$x_1 = (0 \ominus 0.1 \cdot 10^1 \odot 0.1 \cdot 10^1) \oslash 0.2 \cdot 10^1 = \boxed{-0.5},$$

was in $\mathbb{F}(10, 4, 1)$ völlig in Ordnung ist. □

Spaltenpivotisierung ist nicht ausreichend wie folgendes Beispiel zeigt.

Fort. von Beispiel 5.13

(ebenfalls aus [GO96]). Wir betrachten das 2×2 System

$$\begin{bmatrix} 10 & -10^6 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -10^6 \\ 0 \end{bmatrix} .$$

welches aus (11.1) durch Multiplikation der ersten Zeile mit -10^6 entsteht.

Die Spaltenpivotisierung erfordert keine Vertauschung. Allerdings entsteht für $a_{22}^{(1)} = 1 + 2 \cdot 10^5$ genau dasselbe Problem wie oben! □

Diese Probleme kann man durch eine Skalierung des Gleichungssystems vermindern:

16.11.09

$$Ax = b \rightarrow D^{-1}Ax = D^{-1}b \quad \text{mit } d_{ii} = \sum_{j=1}^n |a_{ij}|$$
$$\Leftrightarrow \tilde{A}x = \tilde{b} \quad (\text{also } \|\tilde{A}\|_{\infty} = 1)$$

Rundungsfehleranalyse bei Pivotisierung

Analog zu Satz 5.12 zeigt man, dass für die Lösung \hat{x} bei Spaltenpivotisierung gilt:

$$(A+E)\hat{x} = b \quad \text{mit } \|E\| \leq n \text{ eps } (3\|A\| + 5 P^T \|L\| \|\hat{U}\|) + O(\text{eps}^2).$$

Nach Konstruktion ist $\|L\|_{\infty} \leq n$ und mit der Definition

$$g = \max_{i,j,k} \frac{|a_{ij}^{(k)}|}{\|A\|_{\infty}} \quad \text{„Wachstumsfaktor“}$$

Zeigt man [GVL] also auch $\|E\|_{\infty} \leq g n^3 \|A\|_{\infty} \text{ eps} + O(\text{eps}^2)$.

$$\|E\|_{\infty} \leq g n^3 \|A\|_{\infty} \text{ eps} + O(\text{eps}^2).$$

In der Praxis ist $g \approx 10$

schlechtester Fall bei Spaltenpivotisierung ist $g = 2^{n-1}$.

Mit totaler Pivotisierung erreicht man

$$|a_{ij}^{(k)}| \leq k^{1/2} (2 \cdot 3^{1/2} \cdots k^{1/(k-1)})^{1/2} \max |a_{ij}|$$

also deutlich kleineres Wachstum.

Totale Pivotisierung

4
16.11.09

Wähle $r, s \in \{1, \dots, n\}$ so dass

$$|a_{rs}^{(k)}| \geq |a_{ij}^{(k)}| \quad \forall k \leq i, j \leq n$$

und erreiche durch Zeilen- und Spaltenvertauschung, dass

$$\tilde{a}_{kk}^{(k)} = a_{rs}^{(k)}$$

In Matrixform:

Schritt 1: $G_1 P_{r_1} A P_{s_1} P_{s_1}^T x = G_1 P_{r_1} b$

Schritt 2: $G_2 P_{r_2} G_1 P_{r_1} A P_{s_1} P_{s_2} P_{s_2}^T P_{s_1}^T x = G_2 P_{r_2} G_1 P_{r_1} b$

⋮

Schritt n und Umformulierung

$$\underbrace{G'_n \dots G'_1 P_{r_n} \dots P_{r_1}}_P \underbrace{A P_{s_1} \dots P_{s_n}}_{Q^T} \underbrace{P_{s_n} \dots P_{s_1}}_Q x = G'_n \dots G'_1 P_{r_n} \dots P_{r_1} b$$

$= R$

und damit

$$PAQ^T z = Pb$$

$$PAQ^T = LR \quad z = Qx$$

Lösen des LGS gelingt dann mit

$$Lb' = Pb'$$

$$Ly = b'$$

$$Rz = y$$

$$x = Q^T z$$

Aufwand der Pivotisierung:

5
16.11.09

- $n^2/2$ Vergleiche bei Spaltenpivotisierung
- $n^3/3$ Vergleiche bei totaler Pivotisierung

Da Speicherzugriffe ^{im Grunde} teurer als eigentliche Rechung \Rightarrow Verdopplung des Zeitbedarfs bei totaler Pivotisierung.

Praktische Erfahrung zeigt keine Vorteile für Rundungsfehler bei totaler Pivotisierung

\Rightarrow Spaltenpivotisierung mit Zeilenskalierung ist effizient und numerisch stabil.

Symmetrisch positiv definite Matrizen

Satz 5.14 Eine symmetrisch positiv definite Matrix $A \in \mathbb{R}^{n \times n}$ ist stets ohne Pivotisierung LR-zerlegbar. Für die Diagonalelemente der im Eliminationsprozess auftretenden Matrizen gilt

$$a_{ii}^{(k)} \geq \lambda_{\min}(A), \quad k \leq i \leq n.$$

Beweis: Betrachte einen Schritt in der LR-Zerlegung

c).

$$A = \begin{bmatrix} \alpha & v^T \\ v & B \end{bmatrix} \begin{matrix} 1 \\ n-1 \end{matrix}$$

Elimination der Spalte v liefert

$$\begin{bmatrix} 1 & 0 \\ -v/\alpha & I \end{bmatrix} \underbrace{\begin{bmatrix} \alpha & v^T \\ v & B \end{bmatrix}}_A = \begin{bmatrix} \alpha & v^T \\ 0 & B - \frac{1}{\alpha} v v^T \end{bmatrix} = \begin{bmatrix} \alpha & 0 \\ 0 & B - \frac{1}{\alpha} v v^T \end{bmatrix} \begin{bmatrix} 1 & v^T/\alpha \\ 0 & I \end{bmatrix}$$

Mit $\begin{bmatrix} 1 & v^T/\alpha \\ 0 & I \end{bmatrix}^{-1} = \begin{bmatrix} 1 & -v^T/\alpha \\ 0 & I \end{bmatrix}$ folgt:

$$\underbrace{\begin{bmatrix} 1 & 0 \\ -v/\alpha & I \end{bmatrix}}_{X^T} \underbrace{\begin{bmatrix} \alpha & v^T \\ v & B \end{bmatrix}}_A \underbrace{\begin{bmatrix} 1 & -v^T/\alpha \\ 0 & I \end{bmatrix}}_X = \underbrace{\begin{bmatrix} \alpha & 0 \\ 0 & B - \frac{1}{\alpha} v v^T \end{bmatrix}}_{\tilde{A}}$$

X hat vollen Rang, A ist s.p.d. nach Vor. $\Rightarrow X^T A X = \tilde{A}$ ist s.p.d.

$B - \frac{1}{\alpha} v v^T$ ist Hauptuntermatrix von \tilde{A} und somit auch sym.-pos. definit.

b) Es gilt (Rayleigh-Quotient):

16-11.09

$$(Ax, x)_2 \geq \lambda_{\min}(A) (x, x)_2$$

und damit für $x = e^{(i)}$ (kanonischer Einheitsvektor)

$$a_{ii} = (Ae^{(i)}, e^{(i)})_2 \geq \lambda_{\min}(A) \underbrace{(e^{(i)}, e^{(i)})_2}_{=1} = \lambda_{\min}(A).$$

Für die Diagonalelemente von $B - \frac{1}{\alpha} vv^T$ gilt mit $\tilde{e}^{(i)} = \begin{pmatrix} 0 \\ e^{(i)} \end{pmatrix}$

$$\begin{aligned} (B - \frac{1}{\alpha} vv^T)_{ii} &= (\tilde{A} \tilde{e}^{(i)}, \tilde{e}^{(i)})_2 = (X^T A X \tilde{e}^{(i)}, \tilde{e}^{(i)})_2 \\ &= (A X \tilde{e}^{(i)}, X \tilde{e}^{(i)})_2 \\ &\geq \lambda_{\min}(A) (X \tilde{e}^{(i)}, X \tilde{e}^{(i)})_2 = \lambda_{\min}(A) \left(1 + \frac{v_i^2}{\alpha^2}\right) \end{aligned}$$

$$\underbrace{\begin{bmatrix} 1 & -v^T/\alpha \\ 0 & I \end{bmatrix}}_X \underbrace{\begin{bmatrix} 0 \\ e^{(i)} \end{bmatrix}}_{\tilde{e}^{(i)}} = \begin{bmatrix} -v_i/\alpha \\ e^{(i)} \end{bmatrix} \geq \lambda_{\min}(A).$$

□

Dies zeigt, dass die Pivotelemente bei der Gaußelimination echt nach unten beschränkt bleiben.

ü) Rundungsfehleranalyse. Kann man $|\tilde{C}|, |\tilde{U}|$ abschätzen?
→ Ranacler Lemma 4.1.

Die Pivotelemente die bei der LR-Zerlegung auftreten sind nach 5.14 stets positiv.

Mit $D = \text{diag}(R)$ gilt

$$A = L D \underbrace{D^{-1} R}_U = L D U$$

U ist obere Dreiecksmatrix mit $u_{ii} = 1$. Wegen der Symmetrie gilt $U = L^T$, also

$$A = L D L^T.$$

Da $d_{ii} > 0$ ist die Matrix " $D^{1/2}$ "

$$(D^{1/2})_{ii} = \sqrt{d_{ii}}$$

wohldefiniert und es gilt

$$A = \underbrace{L}_{=: \tilde{L}} D^{1/2} D^{1/2} L^T = \tilde{L} \tilde{L}^T.$$

Dies ist die sog. Cholesky-Zerlegung einer symmetrisch positiv definiten Matrix.

Ausnutzung der Symmetrie erlaubt die Berechnung der Cholesky-Zerlegung in $\frac{n^3}{3} + O(n^2)$ Operationen (Faktor 2 schneller).

Diagonaldominante Matrizen

16.11.09

Definition 5.15 Eine Matrix $A \in \mathbb{R}^{n \times n}$ heißt diagonaldominant falls

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

Satz 5.16 Diagonaldominante, reguläre Matrizen erlauben eine LR-Zerlegung ohne Pivotisierung.

Beweis. $A = \begin{bmatrix} \alpha & w^T \\ v & B \end{bmatrix}$, ein Schritt der GEM liefert

$$\begin{bmatrix} 1 & 0 \\ -v/\alpha & I \end{bmatrix} \begin{bmatrix} \alpha & w^T \\ v & B \end{bmatrix} = \begin{bmatrix} \alpha & w^T \\ 0 & B - \frac{1}{\alpha} v w^T \end{bmatrix}$$

Da $|\alpha| \neq 0$ wg. Diagonaldominanz ist dies wohldefiniert.

Aus der Diagonaldominanz von A ergibt sich:

$$\text{(Zeile 1)} \quad \sum_{j=1}^{n-1} |w_j| \leq |\alpha| \Leftrightarrow \sum_{j=1}^{n-1} \frac{|w_j|}{|\alpha|} \leq 1$$

$$\text{(Zeile 2...n)} \quad |v_i| + \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij}| \leq |b_{ii}|$$

Für $B - \frac{1}{\alpha} v w^T$ rechnen wir dann nach:

$$\begin{aligned} \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij} - \frac{1}{\alpha} v_i w_j| + \left| \frac{v_i w_i}{\alpha} \right| &\leq \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij}| + |v_i| \sum_{\substack{j=1 \\ j \neq i}}^{n-1} \frac{|w_j|}{|\alpha|} + |v_i| \frac{|w_i|}{|\alpha|} \\ &= \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij}| + |v_i| \underbrace{\sum_{j=1}^{n-1} \frac{|w_j|}{|\alpha|}}_{\leq 1} \leq |v_i| + \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij}| \leq |b_{ii}| \end{aligned}$$

$$\Leftrightarrow \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij} - \frac{1}{\alpha} v_i w_j| \leq |b_{ii}| - \left| \frac{v_i w_i}{\alpha} \right| \leq |b_{ii} - \frac{v_i w_i}{\alpha}| = |(B - \frac{1}{\alpha} v w^T)_{ii}|$$

Somit ist $B - \frac{1}{\alpha} v w^T$ wieder diagonaldominant.

$$\begin{aligned} |x| &= |x - y + y| \leq |x - y| + |y| \\ \Leftrightarrow |x| - |y| &\leq |x - y| \end{aligned}$$

Rangbestimmung

10
16.11.09

Sei $A \in \mathbb{K}^{n \times n}$. Gilt nach k -Schritten der Gauß-Elimination

$$A^{(k)} = \begin{array}{c|c} \overbrace{\quad}^k & \overbrace{\quad}^{n-k} \\ \hline U^{(k)} & * \\ \hline 0 & 0 \end{array} \begin{array}{l} k \\ n-k \end{array}$$

mit $u_{ii}^{(k)} \neq 0$ so ist $\text{Rang}(A) = k$. Somit

- Kann die GE zu Ende geführt werden so gilt $k=n$, mithin $\text{Rang}(A)=n$.
- Rangbestimmung erfordert totales Pivoting, da die erste Spalte der $n-k \times n-k$ Restmatrix 0 sein kann, was aber nicht $\text{Rang}(A)=k$ bedeutet.

Inversenberechnung

Zur Berechnung der Inversen geht man so vor:

- Berechne LR-Zerlegung von A . Aufwand $\frac{2}{3}n^3 + O(n^2)$
bei Spaltenpivoting

- Für $i=1, \dots, n$ löse $Ax^{(i)} = e^{(i)}$.

Aufwand $n \cdot 2n^2$ für Lösen der Dreieckssysteme

- $A^{-1} = [x^{(1)}, \dots, x^{(n)}]$ besteht spaltenweise aus den $x^{(i)}$.

- Gesamtaufwand ist $\frac{8}{3}n^3 + O(n^2)$.

Tridiagonalsysteme

11
16.11.09

Definition 5.17

$A \in \mathbb{R}^{n \times n}$ heißt Tridiagonalmatrix falls

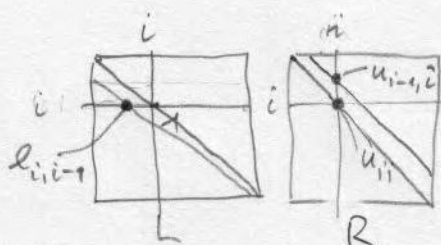
$$a_{ij} = 0 \text{ für } |i-j| > 1, \quad 1 \leq i \leq n$$

Eine Tridiagonalmatrix ist Spezialfall einer Bandmatrix:

$$a_{ij} = 0 \text{ für } j < i - m_e \text{ oder } j > i + m_r, \quad 1 \leq i \leq n$$

Die LR-Zerlegung einer Tridiagonalmatrix sei ohne Pivotisierung durchführbar. Dann gilt:

- L ist Tridiagonalmatrix
- R ist Tridiagonalmatrix

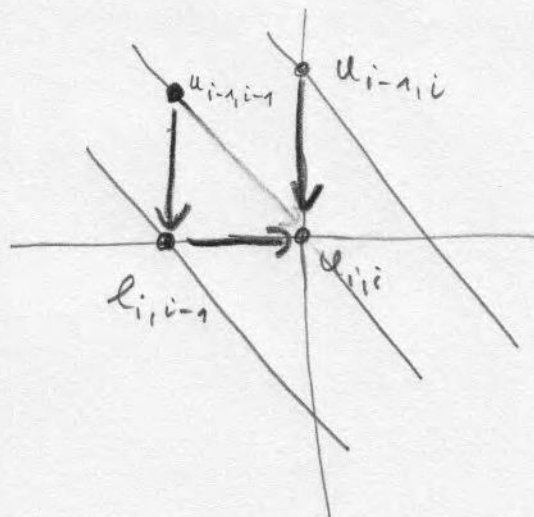


und damit:

Zeile i Spalte i : $l_{i,i-1} \cdot u_{i-1,i} + u_{i,i} = a_{i,i} \Rightarrow u_{i,i} = a_{i,i} - l_{i,i-1} u_{i-1,i}$

$$l_{i,i-1} \cdot u_{i-1,i-1} = a_{i,i-1} \Rightarrow l_{i,i-1} = \frac{a_{i,i-1}}{u_{i-1,i-1}}$$

$$1 \cdot u_{i-1,i} = a_{i-1,i} \Rightarrow u_{i-1,i} = a_{i-1,i}$$



Nichtreguläre Systeme

12
17.11.09

Es sei nun $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$ sowie $\text{Rang}(A)$ beliebig.

$Ax = b$ hat dann genau eine, unendlich viele oder gar keine Lösung.

Einige Grundbegriffe aus der linearen Algebra:

$$\text{Bild}(A) = \{y \in \mathbb{R}^m : y = Ax \text{ für ein } x \in \mathbb{R}^n\}$$

$$\text{Kern}(A) = \{x \in \mathbb{R}^n : Ax = 0\}$$

$$\text{Rang}(A) = \dim(\text{Bild}(A)) = \text{Rang}(A^T) = \dim(\text{Bild}(A^T))$$

(# l.u. Spalten = # l.u. Zeilen).

Orthogonales Komplement:

$$\text{Bild}(A)^\perp = \{y \in \mathbb{R}^m : (y, y')_2 = 0 \forall y' \in \text{Bild}(A)\}$$

$$(y, y')_2 = 0 \forall y' \in \text{Bild}(A)$$

$$\Leftrightarrow (y, Ax)_2 = 0 \forall x \in \mathbb{R}^n$$

$$\Leftrightarrow (A^T y)^T x = 0 \forall x \in \mathbb{R}^n$$

$$\Leftrightarrow y \in \text{Kern}(A^T)$$

und damit $\text{Bild}(A)^\perp = \text{Kern}(A^T)$.

Dies zeigt dann

$$\underbrace{\dim(\text{Bild}(A))}_{= \text{Rang}(A)} + \dim(\text{Kern}(A^T)) = m$$

$$\underbrace{\dim(\text{Bild}(A^T))}_{= \text{Rang}(A)} + \dim(\text{Kern}(A)) = n$$

Den Lösungsbegriff für lineare Gleichungssysteme kann man auf die folgende Weise erweitern.

Satz 5.17 (Least Squares Lösung)

a) Es existiert $\bar{x} \in \mathbb{R}^n$ so dass

$$\|A\bar{x} - b\|_2 = \min_{x \in \mathbb{R}^n} \|Ax - b\|_2.$$

b) Diese Bedingung ist äquivalent dazu dass \bar{x} Lösung von

$$A^T A \bar{x} = A^T b,$$

der sog. "Normalgleichung".

c) Falls $\text{Rang}(A) = n$ (damit ist zwingend $m \geq n$) ist \bar{x} eindeutig bestimmt, sonst hat jede weitere Lösung die Form $\bar{x} + y$ mit $y \in \text{Kern}(A)$.

Beweis.

(b) \Rightarrow (a) \bar{x} sei Lösung der Normalgleichung. Für ein beliebiges $x \in \mathbb{R}^n$ gilt:

$$\begin{aligned} \|Ax - b\|_2^2 &= \|A(x - \bar{x} + \bar{x}) - b\|_2^2 = \left\langle \underbrace{A\bar{x} - b}_{\in \text{Bild}(A)} + \underbrace{A(x - \bar{x})}_{\in \text{Kern}(A)}, \underbrace{A\bar{x} - b}_{\in \text{Bild}(A)} + \underbrace{A(x - \bar{x})}_{\in \text{Kern}(A)} \right\rangle_2 \\ &= \langle A\bar{x} - b, A\bar{x} - b \rangle_2 + 2 \langle \underbrace{A\bar{x} - b}_{\in \text{Bild}(A)}, \underbrace{A(x - \bar{x})}_{\in \text{Kern}(A)} \rangle_2 + \langle A(x - \bar{x}), A(x - \bar{x}) \rangle_2 \\ &= \|A\bar{x} - b\|_2^2 + \underbrace{\|A(x - \bar{x})\|_2^2}_{\geq 0} \\ &\geq \|A\bar{x} - b\|_2^2 \end{aligned}$$

da $A^T(A\bar{x} - b) = 0$

Damit erfüllt \bar{x} die Minimalitätsbedingung.

(a) \Rightarrow (b) Setze $F: \mathbb{R}^n \rightarrow \mathbb{R}$, $F(x) = \|Ax - b\|_2^2$. Notwendige Bedingung für ein Minimum ist $\nabla F(\bar{x}) = 0 \Leftrightarrow \frac{\partial F}{\partial x_k}(\bar{x}) = 0 \quad \forall k = 1, \dots, n$.

$$\begin{aligned} \frac{\partial F(\bar{x})}{\partial x_k} &= \frac{\partial}{\partial x_k} \langle Ax - b, Ax - b \rangle_2 \Big|_{x=\bar{x}} = \frac{\partial}{\partial x_k} \left(\sum_{i=1}^m \left[\sum_{j=1}^n a_{ij} x_j - b_i \right]^2 \right) \Big|_{x=\bar{x}} \\ &= \left(\sum_{i=1}^m 2 \left[\sum_{j=1}^n a_{ij} x_j - b_i \right] a_{ik} \right) \Big|_{x=\bar{x}} = 2 \sum_{i=1}^m (A^T)_{ki} \left[\sum_{j=1}^n a_{ij} \bar{x}_j - b_i \right] \\ &= 2 (A^T (A\bar{x} - b))_k \end{aligned}$$

Nachdifferenzieren.

Und damit $\nabla F(\bar{x}) = 2 A^T (A\bar{x} - b) \stackrel{!}{=} 0 \Leftrightarrow A^T A \bar{x} = A^T b$,
die Normalgleichung.

(c) Lösbarkeit der Normalgleichung.

$$\mathbb{R}^m = \text{Bild}(A) \oplus \text{Bild}(A)^\perp,$$

d.h. jedes $b \in \mathbb{R}^m$ besitzt eine eindeutige Zerlegung $b = s + r$

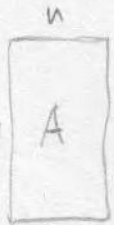
$$b = s + r \text{ mit } s \in \text{Bild}(A), r \in \text{Bild}(A)^\perp = \text{Kern}(A^T)$$

Da $s \in \text{Bild}(A)$ gibt es $\bar{x} \in \mathbb{R}^n$ mit $A\bar{x} = s$. Für dieses gilt dann

$$A^T A \bar{x} = A^T s = A^T s + \underbrace{A^T r}_{=0 \text{ da } r \in \text{Kern}(A^T)} = A^T b,$$

also löst dieses \bar{x} auch die Normalgleichung.

Sei $\text{Rang}(A) = n$, und damit $m \geq n$ (wg. $\text{Rang}(A) \leq \min(m, n)$).



Betrachte $A^T A x = 0$ also $\text{Kern}(A^T A)$.

$$A^T A x = 0 \Leftrightarrow A^T y = 0 \wedge y = Ax$$

Nun ist aber $\text{Kern}(A^T) \perp \text{Bild}(A)$ d.h. $\text{Kern}(A^T) \cap \text{Bild}(A) = \{0\}$.

Somit ist $A^T A$ regulär.

(Alternative: $A^T A$ ist symmetrisch und positiv definit wg. max. Rang).

Sei $\text{Rang}(A) < n$. Sei x_1 eine weitere Lösung der Normalgleichung (also $A^T A x_1 = A^T b$). Dann gilt

$$b = \underbrace{Ax_1}_{\in \text{Bild}(A)} + \underbrace{(b - Ax_1)}_{\in \text{Kern}(A^T), \text{ da } A^T(b - Ax_1) = A^T b - A^T A x_1 \stackrel{\downarrow}{=} 0} \quad x_1 \text{ Lsg der NG.}$$

Da die Zerlegung $\mathbb{R}^m = \text{Bild}(A) \oplus \text{Bild}(A)^\perp$ eindeutig ist muss

$$Ax_1 = A\bar{x} \text{ sein}$$

und damit

$$A(\bar{x} - x_1) = 0, \text{ also } \bar{x} - x_1 \in \text{Kern}(A). \quad \square$$

Bemerkung: $A^T A$ ist symmetrisch und positiv semidefinit.
Lösung prinzipiell mit Cholesky-Zerlegung möglich.

Anwendung: Gaußsche Ausgleichsrechnung.

15
17.11.09

Gegeben:

(i) n Funktionen $u_1, \dots, u_n: \mathbb{R} \rightarrow \mathbb{R}$, sowie

(ii) m Datenpunkte $(x_i, y_i) \in \mathbb{R}^2$, $1 \leq i \leq m \geq n$

Gesucht: n Koeffizienten c_1, \dots, c_n so dass

$$u(x) = \sum_{j=1}^n c_j u_j(x)$$

und

$$(*) \quad \sum_{i=1}^m (u(x_i) - y_i)^2 \rightarrow \text{minimal.}$$

Bei = ungl.
Eindeutig.

Dies lässt sich folgendermaßen formulieren:

$$c = (c_1, \dots, c_n)^T$$

$$y = (y_1, \dots, y_m)^T$$

$$a_{ij} = u_j(x_i)$$

Finde $c \in \mathbb{R}^n$ so dass $\|Ac - y\|_2^2$ minimal.

$$\text{Denn } \sum_{i=1}^m (u(x_i) - y_i)^2 = \sum_{i=1}^m \left(\underbrace{\sum_{j=1}^n c_j u_j(x_i) - y_i}_{(Ac - y)_i} \right)^2$$

$$= \|Ac - y\|_2^2.$$

Somit sind nach Satz 5.17 die gesuchten c_i Lösung der Normalgleichung

$$A^T A c = A^T y.$$

Lösung: z.B. mit Cholesky-Zerl.

$$\textcircled{Ü}: \text{ Zeige } \text{cond}_2(A^T A) = \text{cond}_2(A)^2.$$

6 Interpolation & Approximation

1.12.05
1

6.1 Einführung

Worum es geht: Darstellung und Auswertung von Funktionen im Rechner.

Anwendungen:

- Rekonstruktion eines funktionalen Zusammenhangs aus "gemessenen" Funktionswerten; Auswertung an Zwischenstellen.
- Teuer auszuwertende Funktionen effizienter auswerten
- Darstellung von Fonts (2D), Körpern (3D) im Rechner.
Voraussetzung für Simulation; Szenen in der Computergrafik.
- Lösungen von Differential- und Integralgleichungen
- Datenkompression
⇒ Beispiel auf den Folien.

Wir beschränken uns hier weitgehend auf Funktionen in einer Variablen, also etwa

$$f \in C[a, b].$$

$C[a, b]$ ist ein unendlichdimensionaler (Vektor-)raum von Funktionen. Im Rechner betrachten wir Funktionenklassen die durch eine endliche Zahl von Parametern bestimmt sind (müssen keine Teilräume sein!). z. B.

a) $p(x) = a_0 + a_1x + \dots + a_nx^n$ (Polynome)

b) $r(x) = \frac{a_0 + a_1x + \dots + a_nx^n}{b_0 + b_1x + \dots + b_mx^m}$ (rationale Funktionen)

$$c) \quad t(x) = \frac{1}{2} a_0 + \sum_{k=1}^n (a_k \cos(kx) + b_k \sin(kx)) \quad \boxed{1.12.09}$$

(trigonometrische Polynome)

$$d) \quad e(x) = \sum_{k=1}^n a_k \exp(b_k x) \quad (\text{Exponentialsumme})$$

Grundaufgabe der Approximation:

Gegeben eine Klasse von Funktionen P (siehe oben)
sowie eine Funktion f (z.B. in $C[a, b]$).

Finde $g \in P$ so dass der Fehler $f-g$ in
geeigneter Weise minimiert wird.

Beispiele:

$$a) \quad \left(\int_a^b (f-g)^2 dx \right)^{1/2} \rightarrow \min$$

$$b) \quad \max_{a \leq x \leq b} |f(x) - g(x)| \rightarrow \min$$

$$c) \quad \max_{i=0, \dots, n} |f(x_i) - g(x_i)| \rightarrow \min \quad \text{für } a \leq x_i \leq b, \quad i=0, \dots, n.$$

Man spricht von Interpolation falls g durch

$$g(x_i) = y_i := f(x_i) \quad i=0, \dots, n$$

festgelegt wird.

Dies ist ein Spezialfall der Approximationsaufgabe.

2.2 Polynominterpolation

1.12.09

P_n sei die Menge der Polynome über \mathbb{R} vom Grad kleiner gleich $n \in \mathbb{N}$:

$$P_n := \left\{ p(x) = \sum_{i=0}^n a_i x^i \mid a_i \in \mathbb{R} \right\}$$

P_n ist ein $n+1$ -dimensionaler Vektorraum.

Die Monome $1, x, x^2, \dots, x^n$ bilden eine Basis.

Zu gegebenen $n+1$ ^{paarweise verschiedenen} Stützstellen ist die Interpolationsaufgabe

$$p \in P_n : p(x_i) = y_i := f(x_i) \quad i=0, \dots, n$$

äquivalent zu dem linearen Gleichungssystem

$$\begin{bmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix}$$

- Falls $x_i = x_j$ $V[x_0, \dots, x_n]$

- $V[x_0, \dots, x_n]$ heißt Vandermonde-Matrix

- $x_i = x_j \Rightarrow$ Zeile $i =$ Zeile $j \Rightarrow V$ nicht regulär

- aber ist V für paarweise verschiedene x_i regulär? Ja! Siehe unten.
(Determinante lässt sich auch direkt ausrechnen).

- V ist schlecht konditioniert für große n $\text{cond}(V) \approx 2^n$!

Wie macht man es besser und einfacher?

Lagrange-Interpolation

4
1.12.09

Definition 6.1

Zu den paarweise verschiedenen Stützstellen $x_i, i=0, \dots, n$,
definiere die 10g. Lagrange polynome:

$$L_i^{(n)}(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x-x_j}{x_i-x_j}, \quad i=0, \dots, n$$

FKB

Diese Polynome haben die Eigenschaften:

a) $L_i^{(n)}$ hat Grad n .

klarung: $\prod_{\substack{j=0 \\ j \neq i}}^n (x-x_j)$ hat n Faktoren.

b) Es gilt

$$L_i^{(n)}(x_k) = \delta_{ik} = \begin{cases} 1 & i=k \\ 0 & i \neq k \end{cases}$$

$i \neq k$: Im Produkt $\prod_{\substack{j=0 \\ j \neq i}}^n (x-x_j)$ kommt $j=k \neq i$ vor und

damit $(x_k - x_k) = 0$ ein Faktor.

$i=k$: $L_i^{(n)}(x_i) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x_i - x_j}{x_i - x_j} = 1.$

c) Die $L_i^{(n)}$ bilden eine Basis von P_n

$L_k^{(n)}$ ist das einzige Polynom unter den $L_i^{(n)}$ mit $L_k(x_k) = 1$, daher ist es unabhängig von den $L_i^{(n)}, i \neq k$.

Da P_n $n+1$ dimensional bilden die $L_i^{(n)}, i=0, \dots, n$ eine Basis.

Lösung der Interpolationsaufgabe:

Für gegebene $(x_i, y_i), i=0, \dots, n$, erfüllt

$$p(x) = \sum_{i=0}^n y_i L_i^{(n)}(x)$$

die Interpolationsaufgabe: $p(x_k) = \sum_{i=0}^n y_i L_i^{(n)}(x_k) = \sum_{i=0}^n y_i \delta_{ik} = y_k.$

Satz 6.2 (Eindeutigkeit der Polynominterpolation)

5
1.12.09

Zu gegebenen paarweise verschiedenen Stützstellen x_0, \dots, x_n gibt es genau ein Polynom vom Grad n mit

$$p(x_i) = y_i \quad i = 0, \dots, n, \quad y_i \in \mathbb{R}.$$

Beweis:

kurz: $p = \sum_{i=0}^n y_i L_i^{(n)}$ interpoliert die gegebenen Werte. Da die $L_i^{(n)}$ eine Basis von P_n bilden ist diese Darstellung eindeutig.

lang: $p = \sum_{i=0}^n y_i L_i^{(n)}$ zeigt konstruktiv die Existenz des Interpolationspolynoms.

Eindeutigkeit: Ang. es gibt zwei Polynome p_1, p_2 , $p_1 \neq p_2$, aber $p_1(x_i) = p_2(x_i) = y_i \quad i = 0, \dots, n$.

$p := p_1 - p_2 \in P_n$ (Differenzpolynom) erfüllt $p(x_i) = 0, i = 0, \dots, n$, hat also $n+1$ Nullstellen, damit muss $p \equiv 0$ gelten. Dies folgt aus

Satz von Rolle: $u(x)$ sei auf $[a, b]$ stetig und in (a, b) differenzierbar sowie $u(a) = u(b) = 0$. Dann existiert mindestens ein $x \in (a, b)$ mit $u'(x) = 0$. ($u(x) = 0$ ist auch möglich).



Angenommen es sei p nicht das Nullpolynom, also $p = \alpha_m x^m$ mit $0 \leq m \leq n, \alpha_m \neq 0$. Dann gilt

$$p = p^{(0)} = \alpha_m x^m + \dots \quad \text{hat } n+1 \text{ Nullstellen}$$

$$p' = p^{(1)} = \alpha_m m x^{m-1} + \dots \quad \text{hat } n \text{ Nullstellen (Rolle!)}$$

$$p'' = p^{(2)} = \alpha_m m(m-1) x^{m-2} + \dots \quad \text{hat } n-1 \text{ Nullstellen}$$

$$\vdots$$
$$p^{(m)} = \alpha_m m! \quad \text{hat } n+1-m \geq 1 \text{ Nullstellen (da } m \leq n)$$

Dies geht nur für $\alpha_m = 0$ im Widerspruch zur Annahme. \square

Newton-Darstellung

6
1.12.09

Nachteil der Lagrange-Polynome: Hinzufügen einer Stützstelle

• Hinzufügen einer Stützstelle ändert alle bisherigen Basispolynome

- Eignet sich nicht zur „inkrementellen“ Erstellung des Interpolationspolynoms.

Besser ist hier die Newton-Darstellung mit den Basis-Polynomen:

$$N_0(x) = 1; \quad i = \underline{1}, \dots, n: \quad N_i(x) = \prod_{j=0}^{i-1} (x-x_j)$$

a) $N_i(x)$ ist ein Polynom vom Grad i

b) N_0, \dots, N_n bilden eine Basis von P_n .

c) $N_i(x_k) = 0$ für alle $i > k$, denn $(x-x_k)$ kommt in N_i für $i > k$ vor. $i-1 \geq k \Leftrightarrow i > k$

Gestaffelte Berechnung: Bestimme $p(x) = \sum_{i=0}^n a_i N_i(x)$ über

$$p(x_k) = \sum_{i=0}^n a_i N_i(x_k) = \sum_{i=0}^k a_i N_i(x_k) \stackrel{!}{=} y_i \quad i=0, \dots, n$$

also $a_0 = y_0$;

$$a_k = \left(y_k - \sum_{i=0}^{k-1} a_i N_i(x_k) \right) / N_k(x_k);$$

Alternativ lassen sich die Koeffizienten auch berechnen über

Satz 6.3 (Dividierte Differenzen)

Man definiert rekursiv die sog. „Dividierten Differenzen“

$$\forall i = 0, \dots, n \quad y[x_i] := y_i \quad (\text{die } n+1 \text{ Werte an Stützstellen})$$

$$\forall k = 1, \dots, n-i$$

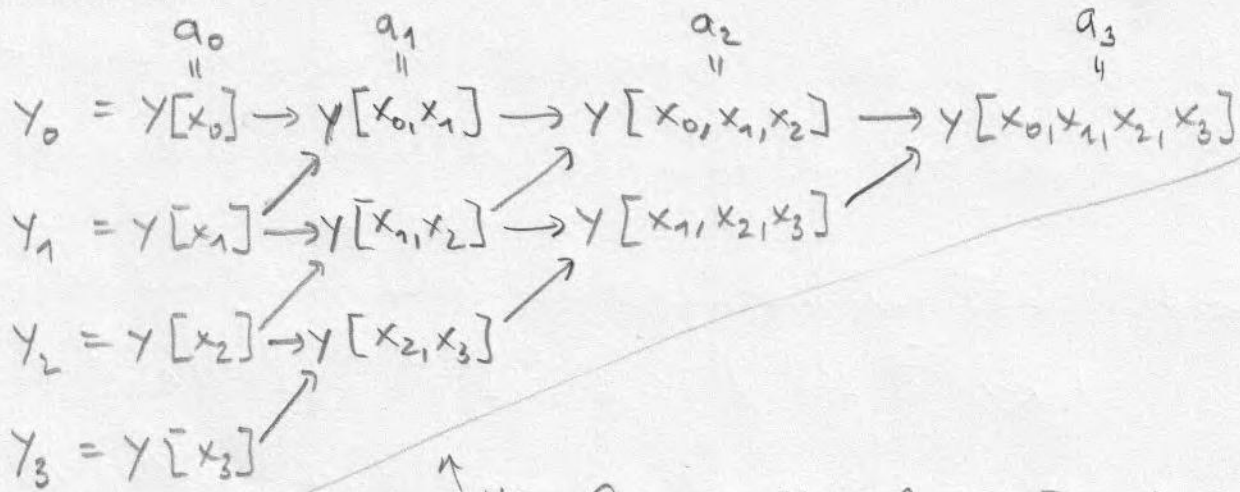
$$y[\underbrace{x_{i-1}, \dots, x_{i+k}}_{k+1 \text{ Stück}}] := \frac{y[\underbrace{x_{i+1}, \dots, x_{i+k}}_k] - y[\underbrace{x_{i-1}, \dots, x_{i+k-1}}_k]}{x_{i+k} - x_{i-1}}$$

Dann gilt

$$p(x) = \sum_{i=0}^n \underbrace{y[x_0, \dots, x_i]}_{=a_i} N_i(x).$$

Beweis: [Ramanujan, Satz 2.2]. ◻

Die dividierten Differenzen ordnet man in folgendem Tableau an:



↑ Hinzufügen von Y_4 erfordert Berechnung aller Koeffizienten in der Diagonale.

Diese Berechnung der Koeffizienten der Newton-Darstellung ist numerisch stabiler.

Effiziente und numerisch stabile Auswertung des Interpolationspolynoms an einzelnen Stellen (also $p(x)$) gelingt mit der Methode von Neville. Siehe [Ramacher, S. 27].

Interpolationsfehler

2.12.09

$y_i = f(x_i)$, $i=0, \dots, n$ sei die Auswertung einer Funktion f an $n+1$ paarweise verschiedenen Stützstellen

$P(x)$ sei das Polynom vom Grad n welches (x_i, y_i) , $i=0, \dots, n$ interpoliert.

Die Differenz erfüllt

$$e(x) = f(x) - p(x) = \begin{cases} 0 & x = x_i \\ ? & x \neq x_i \quad i=0, \dots, n. \end{cases}$$

Frage: Wie groß kann die Differenz werden?

Satz 6.4 Sei $f(x)$ $n+1$ mal stetig differenzierbar auf $[a, b]$ und

es sei $a \leq x_0 < x_1 < \dots < x_n \leq b$. (Löst auch Extrapolation zu) Dann gibt es zu jedem $x \in [a, b]$ ein $\xi_x \in (x_0, \dots, x_n, x)$ (= kleinstes Intervall welches alle Punkte enthält) so dass:

$$f(x) - p(x) = \frac{f^{(n+1)}(\xi_x)}{(n+1)!} \prod_{j=0}^n (x - x_j).$$

Beweis. Für $x \in \{x_0, \dots, x_n\}$ liefert $\prod_{j=0}^n (x - x_j) = 0$ somit ist ξ_x beliebig wählbar.

Sei also $x \in [a, b] \setminus \{x_0, \dots, x_n\}$. Zu diesem x definiere die Funktion

$$F_x(t) = f(t) - p(t) - \underbrace{\frac{f(x) - p(x)}{l(x)}}_{\substack{\text{eine Zahl} \\ \text{abhängig von } x}} l(t) \quad \text{mit } l(t) = \prod_{j=0}^n (t - x_j)$$

neue freie Variable

$F_x(t)$ hat die $n+2$ Nullstellen $\{x_0, \dots, x_n, x\}$, denn

$$i=0, \dots, n: F_x(x_i) = \underbrace{f(x_i) - p(x_i)}_{=0} - \frac{f(x) - p(x)}{l(x)} \underbrace{l(x_i)}_{=0}$$

$$x: F_x(x) = f(x) - p(x) - \frac{f(x) - p(x)}{l(x)} l(x) = 0$$

Satz von Rolle:
(siehe oben)

$F_x(t)$ $n+2$ Nullstellen (mindestens)
nach t ableiten! $F_x^{(n)}(t)$ $n+1$ Nullstellen (-"-)
 $F_x^{(n+1)}(t)$ 1 Nullstelle (-"-)

Zusätzlich gilt: Diese Nullstellen sind in (x_0, \dots, x_n, x)

Diese Nullstelle von $F_x^{(n+1)}$ sei $\xi_x \in (x_0, \dots, x_n, x)$ und es gilt

9
2.12.09

$$F_x^{(n+1)}(\xi_x) = f^{(n+1)}(\xi_x) - \underbrace{p^{(n+1)}(\xi_x)}_{=0 \text{ da } p \text{ Grad } n} - \frac{f(x) - p(x)}{l(x)} \underbrace{l^{(n+1)}(\xi_x)}_{l(t) = t^{n+1} + \dots}$$

$n+1$ mal differenzieren:
 $l^{(n+1)}(t) = (n+1)!$

$$= f^{(n+1)}(\xi_x) - \frac{f(x) - p(x)}{l(x)} (n+1)! \stackrel{!}{=} 0.$$

Auflösen nach $f(x) - p(x)$ liefert die Behauptung. □

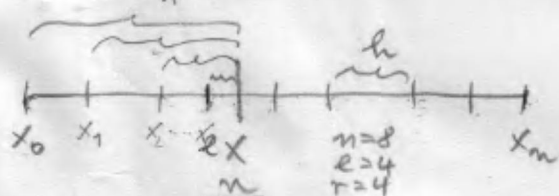
Diskussion des Interpolationsfehlers

Betrag bilden:

$$|f(x) - p(x)| \leq \frac{|f^{(n+1)}(\xi_x)|}{(n+1)!} \prod_{j=0}^n |x - x_j| \quad \text{für ein } \xi_x.$$

- Die Stützstellen seien äquidistant gewählt: $|x_i - x_{i+1}| = h$ und

$x_0 \leq x \leq x_n$:



l : # Punkte links von x

r : # Punkte rechts von x

$$\Rightarrow l + r = n + 1 \quad (\text{Gesamtzahl der Punkte})$$

und damit $\prod_{j=0}^n |x - x_j| \leq \underbrace{1h \cdot 2h \cdot \dots \cdot lh}_{\text{links von } x} \cdot \underbrace{1h \cdot 2h \cdot \dots \cdot rh}_{\text{rechts von } x}$

$$= l! r! h^{n+1}$$

also $|f(x) - p(x)| \leq |f^{(n+1)}(\xi_x)| \frac{l! r!}{(n+1)!} h^{n+1}$

≤ 1 da $l+r = n+1$

Für $|f^{(n+1)}|$ beschränkt und $n \rightarrow \infty$ geht $|f(x) - p(x)| \rightarrow 0$.

Allerdings sind die höheren Ableitungen auch einfache Funktionen für $n \rightarrow \infty$ oft nicht beschränkt sondern wachsen sehr schnell.

- Runacher: (Runge's Gegenbeispiel)

$$f(x) = \frac{1}{1+x^2} \quad |f^{(n)}(x)| \approx 2^n n! O\left(\frac{1}{|x|^{n+2}}\right)$$

⇒ siehe Folien, insbesondere am Rand des Intervalls.

- Beispiel 4.6. aus dem Numstoch-Skript als Ü-Aufgabe!

Bemerkung 6.5

Approximationssatz von Weierstraß: Jede Funktion $f \in C[a,b]$ kann beliebig gut gleichmäßig auf $[a,b]$ durch Polynome approximiert werden.

Obige Beobachtung ist kein Widerspruch, denn

- Approximation muss nicht durch Interpolation erfolgen (Der Beweis nutzt sog. Bernstein-Polynome).
- Mit nichtäquidistanten Stützstellen wird es auch schon sehr viel besser.

Bemerkung 6.6

Allgemein gilt für „Methoden hoher (Polynom-) Ordnung“, dass entsprechende Differenzierbarkeit von f vorliegen muss.

Konditionierung

$P(x; y)$ bezeichne das Interpolationspolynom zu den Ordinatenwerten $y = (y_0, \dots, y_n)^T$ und fixierten Abszissenwerten $(x_0, \dots, x_n)^T$.

$$\frac{P(x; y + \Delta y) - P(x; y)}{P(x; y)} = \left(\sum_{j=0}^n (y_j + \Delta y_j) L_j^{(n)}(x) - \sum_{j=0}^n y_j L_j^{(n)}(x) \right) / P(x; y)$$

$$= \sum_{j=0}^n \underbrace{\frac{L_j^{(n)}(x) y_j}{P(x; y)}}_{\text{Verstärk. faktor}} \cdot \underbrace{\frac{\Delta y_j}{y_j}}_{\text{relativer Eingabefehler}}$$

Für große n kann $L_j^{(n)}(x)$ sehr groß wachsen, dann ist die Interpolationsaufgabe schlecht konditioniert!

Numerische Differentiation

Problem:

- Berechne die Ableitung (der Ordnung k) einer tabellarisch gegebenen Funktion oder einer einer im Rechner als Prozedur gegebenen Funktion.

Idee: Erstelle Interpolationspolynom zu bestimmten Stützstellen, leite dieses ab und werte es aus.

Zunächst = Ableitungsordnung = Polynomgrad.

Lagrange-Polynome sind

$$L_i^{(n)}(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(x-x_j)}{(x_i-x_j)} = \underbrace{\left(\prod_{\substack{j=0 \\ j \neq i}}^n \frac{1}{(x_i-x_j)} \right)}_{=: \lambda_i \in \mathbb{R}} x^n + \alpha_{n-1} x^{n-1} + \dots + \alpha_0$$

 n -maliges differenzieren liefert:

$$\frac{d^n}{dx^n} L_i^{(n)}(x) = \lambda_i n! \quad (\text{konstant, unabhängig von } x)$$

Damit gilt für die n -te Ableitung eines Interpolationspolynoms vom Grad n :

$$\frac{d^n}{dx^n} \left(\sum_{j=0}^n y_j L_j^{(n)} \right)(x) = n! \sum_{j=0}^n y_j \lambda_j. \quad (\text{unabh. von } x)$$

Eine Aussage über den Fehler liefert:

Satz 6.7 Sei $f \in C^n[a, b]$ und $a = x_0 < \dots < x_n = b$. Dann gibt es ein $\xi \in (a, b)$ sodass

$$f^{(n)}(\xi) = n! \sum_{i=0}^n y_i \lambda_i.$$

Beweisskizze: n -malige Anwendung des Satzes von Rolle auf $g(x) = f(x) - p(x)$.

$$\begin{array}{l} g(x) = f(x) - p(x) \quad n+1 \text{ Nullst.} \\ g'(x) \quad n \\ g^{(n)}(x) \quad 1 \text{ Nullst.} \rightarrow f \end{array}$$

Um einfachere Formeln zu bekommen verwenden wir nun äquidistante Stützstellen, d.h. $x_i = x_0 + ih, 0 \leq i \leq n$.

Damit ist

$$\lambda_i = \frac{1}{\underbrace{(x_i - x_0) \cdots (x_i - x_{i-1})}_{i \text{ Stück positiv}} \underbrace{(x_i - x_{i+1}) \cdots (x_i - x_n)}_{n-i \text{ Stück negativ}}}$$

$x_0 + ih - x_0 - nh = -(n-i)h$

$$= \frac{1}{h^n (-1)^{n-i} i!(n-i)!} = \frac{(-1)^{n-i}}{h^n n!} \binom{n}{i}$$

← Binomialkoeffizient

und damit

$$f^{(m)}(x) \approx \frac{d^m}{dx^m} \left(\sum_{j=0}^n y_j L_j^{(m)}(x) \right) = \frac{1}{h^m} \sum_{i=0}^n (-1)^{n-i} \binom{n}{i} y_i$$

$\frac{n!}{i!(n-i)!}$ $n!$ kürzt
nich heraus!

Speziell:

$$f^{(1)}(x) \approx \frac{y_1 - y_0}{h}, \quad f^{(2)}(x) \approx \frac{y_2 - 2y_1 + y_0}{h^2}, \quad f^{(3)}(x) \approx \frac{y_3 - 3y_2 + 3y_1 - y_0}{h^3}$$

Bisher: m -te Ableitung aus Polynom vom Grad n (d.h. $n+1$ Werten)

Es geht auch: m -te Ableitung aus Polynom vom Grad $n > m$.

Wert hängt dann aber von der Auswertestelle ab.

Beispiel: $m=1, n=2$. Äquidistante Stützstellen $x_0, x_1 = x_0 + h, x_2 = x_0 + 2h$.

$$p'(x_1) = \frac{y_2 - y_0}{2h} \quad \text{„zentraler Differenzenquotient“}$$

Taylorreihenentwicklung (Übung!) zeigt

$$f'(x) = \frac{f(x+h) - f(x-h)}{2h} + O(h^2) \quad \text{für } f \in C^3,$$

$$f''(x) = \frac{f(x+h) - 2f(x) + f(x-h)}{h^2} + O(h^2) \quad \text{für } f \in C^4.$$

Beispiel 6.8

→ Beispiel 4.9 aus Numstoch zur

→ Übung!

Auslöschung bei numerischer Differentiation

und zur Motivation der Extrapolation!

Extrapolation zum Limes

15
3.12.09

Eine Größe $a(h)$ sei im Rechner für $h > 0$ berechenbar, nicht jedoch für $h = 0$. Man möchte

$$a(0) = \lim_{h \rightarrow 0} a(h)$$

mit guter Genauigkeit berechnen.

Beispiel 6.9

a) $a(0) = \lim_{x \rightarrow +0} \frac{\cos(x) - 1}{\sin(x)} \quad (= 0)$

b) Numerische Differentiation

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (\text{für kleine } h \text{ tritt Auslöschung ein})$$

c) Numerische Integration

$$\int_a^b f(x) dx = \lim_{N \rightarrow \infty} \sum_{i=1}^N \frac{1}{N} f\left(\frac{i-1}{2N} + \frac{i}{2N}\right) = \lim_{N \rightarrow \infty} \sum_{i=1}^N h f\left(\frac{2i-1}{2} h\right)$$

wobei $h := \frac{1}{N}$.

$h \rightarrow 0$ nicht möglich wg. Aufwand $\rightarrow \infty$

d) Numerische Lösung des Anfangswerts

$$y'(t) = f(t, y(t)); \quad y(0) = y_0$$

$$y(T) \approx y_N \quad \text{und} \quad y_n = y_{n-1} + h f(t, y_{n-1}); \quad h = 1/N;$$

$h \rightarrow 0$ bedeutet $N \rightarrow \infty$ und damit Aufwand $\rightarrow \infty$. ▣

Idee der Extrapolation:

Zu $h_0 > h_1 > \dots > h_n > 0$ bestimme Interpolationspolynom

$$p(h_i) = a(h_i) \quad i = 0, \dots, n$$

und berechne

$$a(0) \approx p(0) \quad (\text{Extrapolation, da } 0 \notin [h_n, \dots, h_0].)$$

graphisch:

Beispiel 6.10.

[Rauwacher]. Für $a(h) = \frac{\cos(h)-1}{\sin(h)}$ erhält man für

$h_0 = 1/8$	$a(h_0) = -6.258\ 151 \cdot 10^{-2}$	} halbiert sich, da h halbiert wird.
$h_1 = 1/16$	$a(h_1) = -3.126\ 018 \cdot 10^{-2}$	
$h_2 = 1/32$	$a(h_2) = -1.56\ 2627 \cdot 10^{-2}$	

und bei Extrapolation mit einem Polynom vom Grad 2:

$$a(0) \approx p_2(0) = -1.02 \cdot 10^{-5}$$

also sehr viel besser!

Übung: Extrapolation bei der numerischen Differentiation!

Warum funktioniert das so gut?

Die Funktion $a(x)$ sei $n+1$ mal stetig differenzierbar in einer genügend großen Umgebung von 0. Dann gibt es zu jedem $h > 0$ (in dieser Umgebung) ein $\xi_h \in [0, h]$ sodass: (Taylorreihe mit Restglied u. Lagrange)

$$\begin{aligned}
 a(h) &= a(0+h) = a(0) + h a'(0) + \dots + \frac{h^n}{n!} a^{(n)}(0) + \frac{h^{n+1}}{(n+1)!} a^{(n+1)}(\xi_h) \\
 &\stackrel{\text{gleich darstellen!}}{=} \underbrace{a(0) + a'(0)h + \dots + \frac{a^{(n)}(0)}{n!} h^n}_{\text{Polynom in } h \text{ vom Grad } n} + \underbrace{\frac{a^{(n+1)}(\xi_h)}{(n+1)!} h^{n+1}}_{\text{abhängig von } h!} \quad (6.1)
 \end{aligned}$$

Für $h > 0$ numerisch berechenbar.

$a^{(k)}(0)$ hängt nicht von h ab!

Idee: Für verschiedene h_i bilde Linearkombination:

$$\begin{aligned}
 \sum_{i=0}^n c_i a(h_i) &= \sum_{i=0}^n c_i \left(\sum_{j=0}^n a_j h_i^j \right) + \sum_{i=0}^n c_i \frac{a^{(n+1)}(\xi_{h_i})}{(n+1)!} h_i^{n+1} \\
 &= \sum_{j=0}^n a_j \left(\sum_{i=0}^n c_i h_i^j \right) + \text{Fehler} = a_0 + \text{Fehler} \\
 &\quad \text{unbekannt} \quad \text{Bestimmungsgleichung für die } c_i \quad \rightarrow \text{klein wg } h_i^{n+1}
 \end{aligned}$$

$$\sum_{i=0}^n c_i h_i^j = \begin{cases} 1 & j=0 \\ 0 & \text{sonst} \end{cases}$$

Als Gleichungssystem lautet die Bestimmungsgleichungen

15
3.12.09

$$V^T c = e^{(0)}$$

mit $e^{(0)} = (1, 0, \dots, 0)^T$

und $V = V[h_0, \dots, h_n]$ der Vandermondematrix.

Die Auswertung mit diesen Koeffizienten ist dann:

$$A := \sum_{i=0}^n c_i \underbrace{a(h_i)}_{y_i} = \underbrace{(V^{-T} e^{(0)})^T}_c y = (e^{(0)})^T V^{-1} y.$$

$y = (y_0, \dots, y_n)$

Für den Fehler gilt dann

$$|A - a(0)| = \left| \sum_{i=0}^n c_i \frac{a^{(n+1)}(\xi_i)}{(n+1)!} h_i^{(n+1)} \right|$$

fest. ge. für c!

$$V = \begin{bmatrix} 1 & h_0 & h_0^2 & \dots & h_0^n \\ 1 & h_1 & h_1^2 & \dots & h_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & h_n & h_n^2 & \dots & h_n^n \end{bmatrix}$$

$$= \begin{bmatrix} 1 & hr & (hr)^2 & \dots & (hr)^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & hr^n & (hr^n)^2 & \dots & (hr^n)^n \end{bmatrix}$$

$$= \begin{bmatrix} 1 & r & r^2 & \dots & r^n \\ 1 & r^2 & r^4 & \dots & r^{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & r^n & r^{2n} & \dots & r^{n^2} \end{bmatrix} \begin{bmatrix} h & 0 & \dots & 0 \\ 0 & h & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & h \end{bmatrix}$$

$$h_i = h r^i \rightarrow \leq \sum_{i=0}^n \underbrace{\|V^{-T} e^{(0)}\|_\infty}_{\max c_i} \frac{|a^{(n+1)}(\xi_i)|}{(n+1)!} h^{n+1} r^{i(n+1)}$$

z.B. $r = 1/2$

$$\leq \|V^{-T}\|_\infty \underbrace{\|e^{(0)}\|_\infty}_{=1} |a^{(n+1)}(\xi_{\max})| \frac{h^{n+1}}{(n+1)!} \underbrace{\sum_{i=0}^n r^{i(n+1)}}_{\text{geometrische Reihe}}$$

$$\leq \|V^{-T}\|_\infty \underbrace{|a^{(n+1)}(\xi_{\max})|}_{\text{beschränkt}} \frac{h^{n+1}}{(n+1)!} \underbrace{(1+r^{n+1})}_{= r^{n+1} + 1}$$

$\leq \frac{5}{4}$ für $r \leq \frac{1}{2}, n \geq 1$

Nun betrachte das Interpolationspolynom

$$p(h_i) = \sum_{j=0}^n b_j h_i^j \stackrel{!}{=} a(h_i) = y_i$$

und somit $Vb = y$ für dessen Koeffizienten.

Auswerten an der Stelle Null bedeutet

$$p(0) = b_0 = (e^{(0)})^T b = (e^{(0)})^T V^{-1} y = A \text{ von oben!}$$

Oben beschriebene Elimination der Fehlerformel entspricht also genau der Extrapolationsmethode!

Entscheidend ist also die Fehlerdarstellung (6.1).

16
3.12.09

Und es geht noch besser! Wir betrachten die Näherung für $f''(x)$:
(n gerade:)

$$f(x+h) = f(x) + h f'(x) + \frac{h^2}{2} f''(x) + \frac{h^3}{3!} f^{(3)}(x) + \dots + \frac{h^{2n+2}}{(2n+2)!} f^{(2n+2)}(x) + \frac{h^{2n+4}}{(2n+4)!} f^{(2n+4)}(\xi_+)$$

$$f(x-h) = f(x) - h f'(x) + \frac{h^2}{2} f''(x) - \frac{h^3}{3!} f^{(3)}(x) + \dots + \frac{h^{2n+2}}{(2n+2)!} f^{(2n+2)}(x) - \frac{h^{2n+4}}{(2n+4)!} f^{(2n+4)}(\xi_-)$$

$$f(x+h) + f(x-h) = 2f(x) + h^2 f''(x) + \dots + \frac{h^{2n+2}}{2(2n+2)!} f^{(2n+2)}(x) + \frac{h^{2n+4}}{(2n+4)!} [f^{(2n+4)}(\xi_+) + f^{(2n+4)}(\xi_-)]$$

$$a(h) := \frac{f(x+h) - 2f(x) + f(x-h)}{h^2} =$$

$$= f''(x) + \frac{h^2}{2(4!)} f^{(4)}(x) + \dots + \frac{h^{2n}}{2(2n+2)!} f^{(2n+2)}(x) + \frac{h^{2n+2}}{(2n+4)!} [f^{(2n+4)}(\xi_+) + f^{(2n+4)}(\xi_-)]$$

$$= P_x(h^2) + O(h^{2n+2})$$

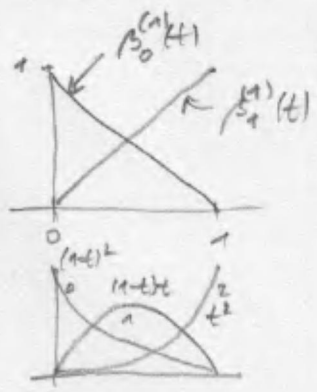
Da alle ungeraden Terme „von selber“ wegfallen kann man mit der gleichen Anzahl von Auswertungen doppelt so viele Fehlerterme eliminieren!

Natürlich muss f dazu entsprechend oft differenzierbar sein.

6.4 Bernstein-Polynome und Kurvendarstellung

7.12.09

Wir gehen nun über von der Interpolation zur Approximation.
 Speziell für Kurven, d.h. Funktionen $u(t): [a, b] \rightarrow \mathbb{R}^d$,
 $d=2,3$, haben sich Bernstein-Polynome bewährt.



Definition 6.11 (Bernstein-Polynome)

Die Polynome

$$\beta_i^{(n)}(t) = \binom{n}{i} (1-t)^{n-i} t^i \quad i=0, \dots, n,$$

Vom Grad n heißen Bernstein-Polynome auf $[0, 1]$.

Mittels der Transformation $\varphi: [a, b] \rightarrow [0, 1]$, $\varphi(u) = \frac{u-a}{b-a}$, definiert man die Bernstein-Polynome auf einem allgemeinen Intervall $[a, b]$:

$$\begin{aligned} \beta_{i, [a, b]}^{(n)} &= \beta_i^{(n)}(\varphi(u)) = \binom{n}{i} \left(1 - \frac{u-a}{b-a}\right)^{n-i} \left(\frac{u-a}{b-a}\right)^i \\ &= \binom{n}{i} \frac{1}{(b-a)^n} (b-u)^{n-i} (u-a)^i. \end{aligned}$$

Satz 6.12 (Eigenschaften der Bernstein-Polynome)

(a) $\sum_{i=0}^n \beta_i^{(n)}(t) = 1$

(b) Beweis: binomischer Lehrsatz: $1 = (1-t+t)^n = \sum_{i=0}^n \binom{n}{i} (1-t)^{n-i} t^i = \sum_{i=0}^n \beta_i^{(n)}(t)$

(c) $t=0$ ist i -fache Nullstelle von $\beta_i^{(n)}$

Beweis: Es sei $0 \leq j \leq i$, dann ist wg Produktregel $\frac{d^j}{dt^j} \beta_i^{(n)}(t) = \sum_{k=0}^j g^{(k)}(t) t^{i-k}$
 und damit $\frac{d^j}{dt^j} \beta_i^{(n)}(0) = 0$ falls $i-j > 0 \Leftrightarrow 0 \leq j < i$

- d.h. $i=0$: keine Nullstelle
- $i=1$: $0 \leq j < 1$ eine Nullstelle
- $i=2$: $0 \leq j < 2$ Nullstelle, Nullstelle der Ableitung

(c) $t=1$ ist $n-i$ -fache Nullstelle von $\beta_i^{(n)}$

$\Rightarrow t=0, t=1$ sind die einzigen Nullstellen.
Beweis: analog zu (b).

(d) Symmetrie: $\beta_i^{(n)}(t) = \beta_{n-i}^{(n)}(1-t)$

Bew: Einsetzen

(e) Positivität - $0 \leq \beta_i^{(n)} \leq 1$ für $t \in [0, 1]$
- $\beta_i^{(n)}(t) > 0$ für $t \in (0, 1)$

Beweis: $\forall t \in [0, 1] \Rightarrow t \geq 0, 1-t \geq 0$ also $\beta_i^{(n)}(t) \geq 0$
 $t \in (0, 1) \Rightarrow t > 0, 1-t > 0$ also $\beta_i^{(n)}(t) > 0$

und $\sum_{i=0}^n \beta_i^{(n)}(t) = 1 \Leftrightarrow \beta_i^{(n)}(t) = 1 - \sum_{\substack{j=0 \\ j \neq i}}^n \beta_j^{(n)}(t) \leq 1$.

(f) $\beta_i^{(n)}$ hat in $[0, 1]$ genau ein Maximum in i/n

$$\frac{d}{dt} \beta_i^{(n)}(t) = \binom{n}{i} \left[-(n-i)(1-t)^{n-i-1} t^i + i(1-t)^{n-i} t^{i-1} \right]$$

$$= \binom{n}{i} (1-t)^{n-i-1} t^{i-1} \left(\underbrace{(1-t)i - (n-i)t}_{i-it-nt+it} \right)$$

$$= \binom{n}{i} \underbrace{(1-t)^{n-i-1}}_{n-i-1\text{-fache Nullst. bei } t=1} \underbrace{t^{i-1}}_{i-1\text{-fache Nullstelle bei } t=0} \underbrace{(i-nt)}_{\text{Nullstelle bei } t=i/n}$$

$n-i-1 + i-1 + 1 = n-1$ Nullstellen.
da $\frac{d}{dt} \beta_i^{(n)}$ Polynom von Grad $n-1$ sind dies alle!

Aus (b) und (c) $\xrightarrow{\text{und (e)}}$ folgt dass bei i/n ein Maximum vorliegt.

(g) Die $\{\beta_i^{(n)}\}_{i=0}^n$ sind linear unabhängig und bilden eine Basis von P_n .

Bew: zu zeigen $\sum_{i=0}^n b_i \beta_i^{(n)}(t) = 0 \forall t \in \mathbb{R} \Rightarrow b_i = 0$.

da $\sum b_i \beta_i$ Nullfunkt.

Betrachte Ableitungen: $\frac{d^j}{dt^j} \sum_{i=0}^n b_i \beta_i^{(n)} = \sum_{i=0}^n b_i \frac{d^j}{dt^j} \beta_i^{(n)}(t) = 0 \forall t \in \mathbb{R}$

Setze $j=0, t=0$: Es ist nur $\beta_0^{(n)}(0) \neq 0$, alle anderen haben dort Nullst. $\Rightarrow b_0 = 0$

$j=1, t=0$: Es ist nur $\frac{d}{dt} \beta_1^{(n)}(0) \neq 0 \Rightarrow b_1 = 0$

usw.

(h) Die Bernstein-Polynome erlauben folgende rekursive Darstellung über den Grad n :

$$\beta_i^{(n)}(t) = \begin{cases} (1-t) \beta_0^{(n-1)}(t) & i=0 \\ t \beta_{i-1}^{(n-1)}(t) + (1-t) \beta_i^{(n-1)}(t) & 0 < i < n \\ t \beta_{n-1}^{(n-1)}(t) & i=n \end{cases}$$

Bew. $i=0, i=n$ sieht man durch einsetzen.

$0 < i < n$

$$t \binom{n-1}{i-1} (1-t)^{n-1-(i-1)} t^{i-1} + (1-t) \binom{n-1}{i} (1-t)^{n-1-i} t^i$$

$$= \binom{n-1}{i-1} (1-t)^{n-i} t^i + \binom{n-1}{i} (1-t)^{n-i} t^i = \left[\binom{n-1}{i-1} + \binom{n-1}{i} \right] (1-t)^{n-i} t^i$$

$= \binom{n}{i}$ Rekursionsformel für Binomialk.

(i) Für die erste Ableitung gilt die Rekursionsformel:

$$\frac{d}{dt} \beta_i^{(n)}(t) = \begin{cases} -n \beta_0^{(n-1)}(t) & i=0 \\ n [\beta_{i-1}^{(n-1)}(t) - \beta_i^{(n-1)}(t)] & 0 < i < n \\ n \beta_{n-1}^{(n-1)}(t) & i=n \end{cases}$$

Kein Beweis?

Achtung: Rechts stehen keine Ableitungen sondern Bernstein-Polynome vom Grad $n-1$!

Kurvendarstellung mittels Bernstein-Polynomen beschreibt:

Definition 6.13 (Bezier-Kurven)

Für gegebene Punkte $b_0, \dots, b_n \in \mathbb{R}^d$ heißt das vektorwertige Polynom

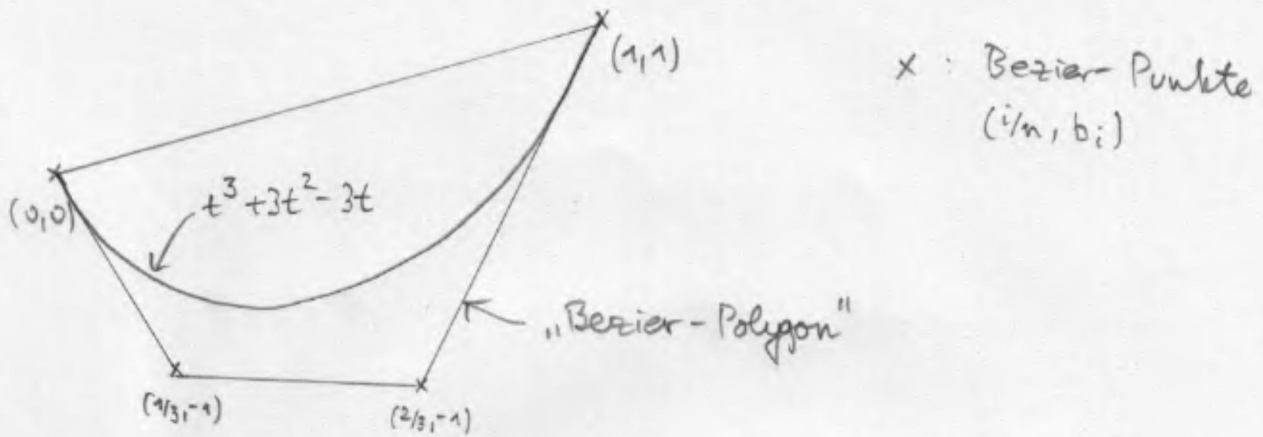
$$B(t) = \sum_{i=0}^n b_i \beta_i^{(n)}(t)$$

Bezier-Kurve.

Beispiel 6.14

Betrachte $b_0 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$, $b_1 = \begin{pmatrix} 1/3 \\ -1 \end{pmatrix}$, $b_2 = \begin{pmatrix} 2/3 \\ -1 \end{pmatrix}$, $b_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

Die zugehörige Bezier-Kurve: $B(t) = \sum_{i=0}^3 b_i \beta_i^{(3)}(t) = \begin{pmatrix} t \\ t^3 + 3t^2 - 3t \end{pmatrix}$



Die Verbindung der Punkte $(i/n, b_i)$ nennt man „Bezier-Polygon“.

Es gelten folgende Eigenschaften:

- Das Bezier-Polynom liegt in der konvexen Hülle der Bezier-Punkte.

$$B(t) = \sum_{i=0}^n b_i \beta_i^{(n)}(t) \quad \text{ist wegen} \quad 0 \leq \beta_i^{(n)}(t) \leq 1 \quad \text{und} \quad \sum_{i=0}^n \beta_i^{(n)}(t) = 1$$

eine Konvexkombination.



- Es ist $B(0) = b_0$ und $B(1) = b_n$. (Liegt an den Nullstellen 6.12 b,c)

- Die Ableitung (Tangente an die Kurve) hat den Endpunkten die Richtung $(b_1 - b_0)$ bzw. $(b_n - b_{n-1})$.

Nutze rekursive Darstellung aus Satz 6.12 (i):

$$\begin{aligned} \frac{d}{dt} B(t) &= \sum_{i=0}^n b_i \frac{d}{dt} \beta_i^{(n)}(t) = \\ &= b_0 \left[-n \underbrace{\beta_0^{(n-1)}(t)}_{=1} \right] + b_1 \left[n \left(\underbrace{\beta_0^{(n-1)}(t)}_{=1} - \underbrace{\beta_1^{(n-1)}(t)}_{=1} \right) \right] + \dots \\ &\quad \dots + b_{n-1} \left[n \left(\underbrace{\beta_{n-2}^{(n-1)}(t)}_{=1} - \underbrace{\beta_{n-1}^{(n-1)}(t)}_{=1} \right) \right] + b_n \left[n \underbrace{\beta_{n-1}^{(n-1)}(t)}_{=1} \right] \end{aligned}$$

$t=0$: $\beta_0^{(n-1)}(0) = 1$: $-n b_0 + b_1 \cdot n = n(b_1 - b_0)$

$t=1$: $\beta_{n-1}^{(n-1)}(1) = 1$: $b_{n-1} \cdot (-n) + b_n n = n(b_n - b_{n-1})$.

$= \begin{cases} b=0 & \text{Fall-} \\ b=1 & \text{untersch.} \end{cases}$

Bernstein-Polynome erlauben eine gute Kontrolle der Werte des Polynoms durch dessen Koeffizienten.

Vergleichen wir dies mit anderen Polynombasen:

- Monome, Newton: Koeffizienten erlauben keine einfache Kontrolle der Werte des Polynoms
- Lagrange: Erlaubt exakte Kontrolle an den Stützpunkten, dazwischen aber (sehr) große Abweichungen.
- Bernstein: $\min\{b_0, \dots, b_n\} \leq \sum_{i=0}^n b_i \beta_i^{(n)}(t) \leq \max\{b_0, \dots, b_n\}$

Eine effiziente und numerisch stabile Auswertung einer Bézierkurve erlaubt der

Algorithmus von de Casteljau

Mit Hilfe der Rekursionsformel ^{aus Satz} 6.12 (b) erhält man

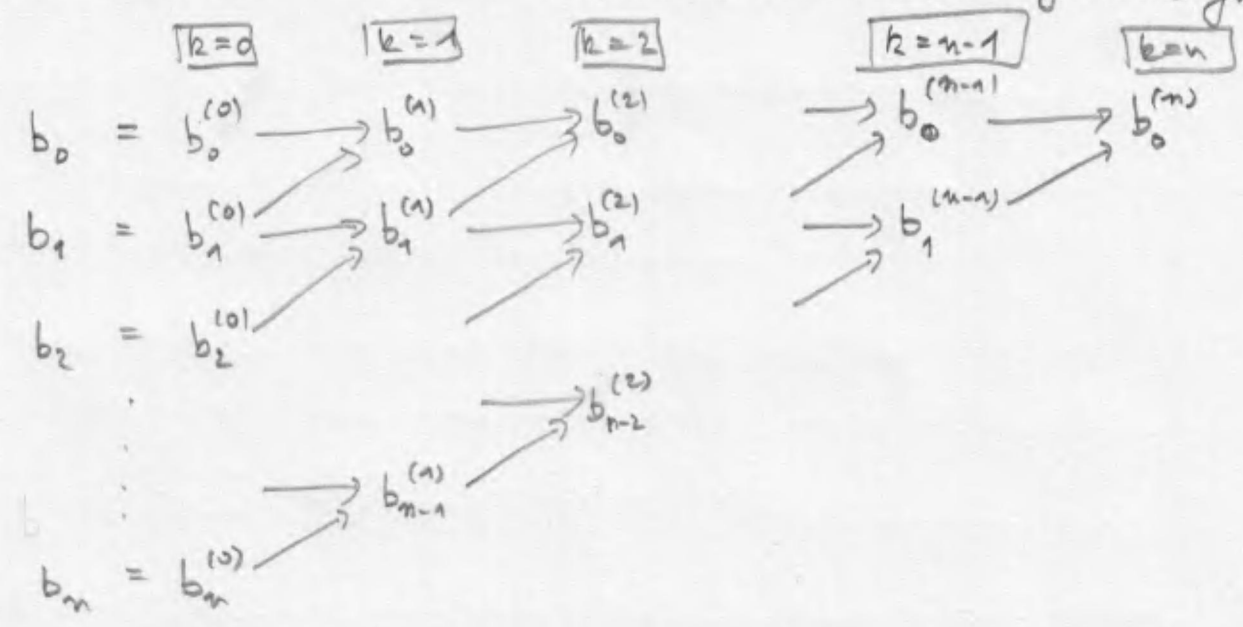
$$\begin{aligned}
 B(t) &= \sum_{i=0}^n b_i^{(0)} \beta_i^{(n)}(t) \quad \text{Eingangskoeffizienten} \\
 &= b_0^{(0)} (1-t) \beta_0^{(n-1)}(t) + b_1^{(0)} [t \beta_0^{(n-1)}(t) + (1-t) \beta_1^{(n-1)}(t)] \\
 &\quad + b_2^{(0)} [t \beta_1^{(n-1)}(t) + (1-t) \beta_2^{(n-1)}(t)] + \dots \\
 &\quad + b_{n-1}^{(0)} [t \beta_{n-2}^{(n-1)}(t) + (1-t) \beta_{n-1}^{(n-1)}(t)] + b_n^{(0)} t \beta_{n-1}^{(n-1)}(t) \\
 &= \sum_{i=0}^{n-1} [b_i^{(0)} (1-t) + b_{i+1}^{(0)} t] \beta_i^{(n-1)}(t) \\
 &\quad \quad \quad =: b_i^{(1)}
 \end{aligned}$$

d.h. wir haben nun n neue Koeffizienten $b_i^{(1)}$ für ein Polynom von Grad n
Allgemein:

Gegeben $b_i^{(0)} = b_i$ $k=0, 0 \leq i \leq n$

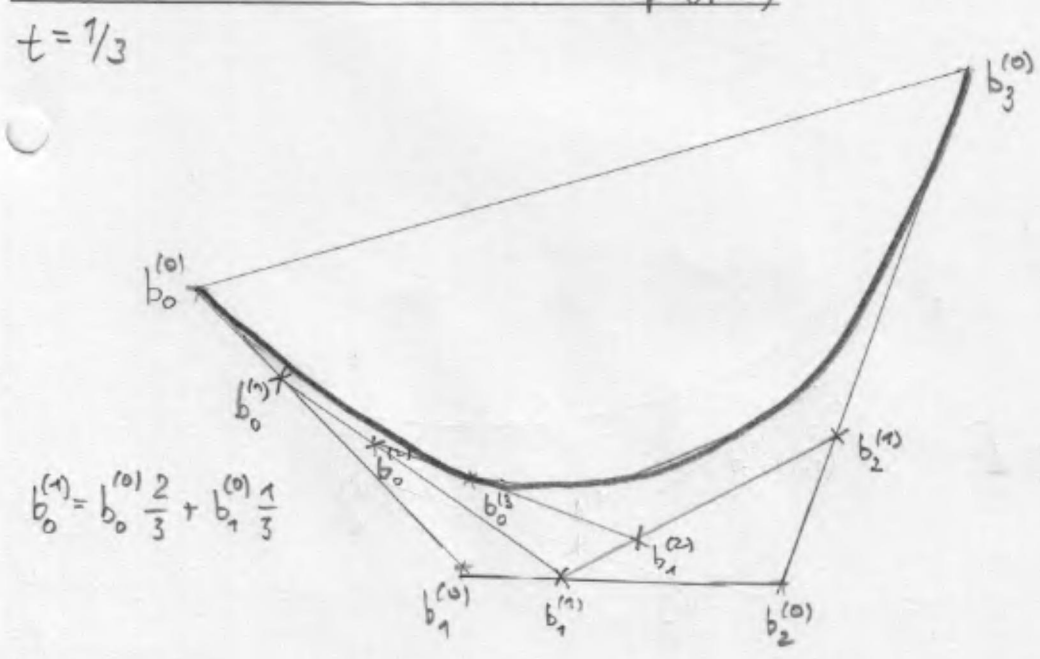
Setze $b_i^{(k)} = b_i^{(k-1)} (1-t) + b_{i+1}^{(k-1)} t$ $0 < k \leq n, 0 \leq i \leq n-k$

Die Auswertung von $B(t)$ erfolgt in einer ähnlichen Struktur wie beim Neville-Verfahren zur Auswertung von Polynomen:



Und damit gilt dann $B(t) = b_0^{(n)} \underbrace{(\beta_0^{(0)}(t))}_{=1} = b_0^{(n)}$
für alle t

Geometrisch interpretiert man das folgendermaßen
(Bernstein-Kurve aus Beispiel 6.14)
Beispiel 6.15 (Fortb. von Bsp. 6.14)



Wir kehren wieder zurück zur Interpolation.

Bis jetzt

- # Stützstellen = Polynomgrad + 1
- Großer Polynomgrad = viele Stützstellen \Rightarrow ^{evtl.} starke Abweichung zwischen den Stützstellen.

Idee: Stückweise Polynome niedrigen Grades.

Definition 6.16

Sei $X = (x_0, x_1, \dots, x_n)$ mit $a = x_0 < x_1 < \dots < x_n = b$ eine Zerlegung des Intervalles $[a, b]$ und sei $m \in \mathbb{N}$. Die Menge

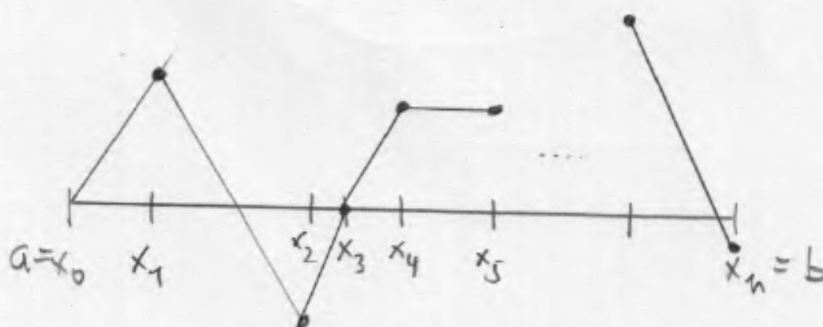
$$S^m(X) = \left\{ s \in C^{m-1}([a, b]) : s|_{[x_i, x_{i+1}]} \in P_m, 0 \leq i < n \right\}$$

heißt Spline-Raum vom Grad m über der Zerlegung X .

Beispiel 6.17

$S^1(X)$ bedeutet:

- $s \in S^1(X)$ ist Polynom vom Grad 1 auf jedem Teilintervall $[x_i, x_{i+1}]$
- $S^1(X) \subset C^0([a, b])$, als $s \in S^1(X)$ stetig



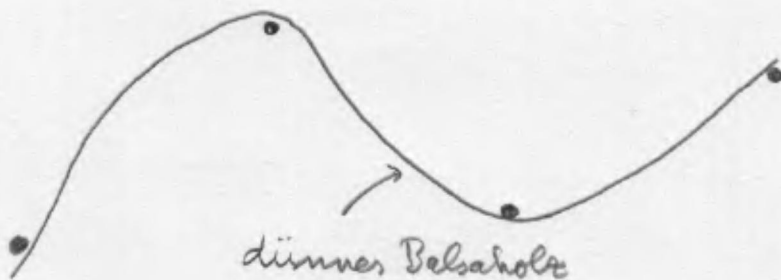
Stetigkeit $\Rightarrow s \in S^1(X)$ ist eindeutig durch Werte in den Stützstellen beschreibbar
 $\dim S^1(X) = |X|$

Kubische Splines

In der Praxis ist $S^3(X)$ sehr beliebt.

$S^3(X)$ heißt Raum der kubischen Splines.

Geschichte: „Straklatte“ zur Konstruktion glatter Kurven im Schiffs- und Flugzeugbau.



dünnes Balsaholz

biegt sich unter Energieminimierung

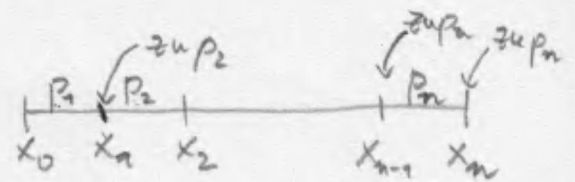
$$\int_a^b \frac{|y''(t)|^2}{1+|y'(t)|^2} dt \approx \int_a^b |y''(t)|^2 dt \rightarrow \min.$$

↑ Krümmung. ↑ $|y'(t)| \ll 1$

Konstruktion von $S \in S^3(X)$. Die Funktion setzt sich stückweise aus n Polynomen zusammen:

$$S(x) = \begin{cases} P_i(x) & x \in [x_{i-1}, x_i) \quad i \in \{1, \dots, n\} \\ P_n(x_n) & x = x_n \end{cases}$$

Für die Polynome P_i gelten folg. Bed:



(a) Interpolationsbedingung (Stetigkeit)

$$i=1, \dots, n: \left. \begin{aligned} P_i(x_{i-1}) &= y_{i-1} \\ P_i(x_i) &= y_i \end{aligned} \right\} 2n \text{ Bedingungen}$$

(b) Stetigkeit der ersten und zweiten Ableitung an inneren Punkten:

$$i=1, \dots, n-1 \left. \begin{aligned} P_i'(x_i) &= P_{i+1}'(x_i) \\ P_i''(x_i) &= P_{i+1}''(x_i) \end{aligned} \right\} 2(n-1) = 2n-2 \text{ Bedingungen}$$

⇒ zusammen $4n-2$ Bedingungen.

Pro Polynom p_i (vom Grad 3) hat man 4, also insgesamt $4n$ Freiheitsgrade.

25
8.12.09

Die fehlenden 2 Bedingungen erhält man durch Randbedingungen an den Stellen x_0 und x_n . Dabei gibt es verschiedene Varianten:

(c) Randbedingungen. Eine der folgenden Varianten:

i) Natürliche Randbedingungen:

$$p_1''(x_0) = 0$$

$$p_n''(x_n) = 0$$

ii) Hermite-Randbedingungen

$$p_1'(x_0) = f'(x_0)$$

$$p_n'(x_n) = f'(x_n)$$

iii) Periodische Randbedingungen

$$p_1'(x_0) = p_n'(x_n)$$

$$p_1''(x_0) = p_n''(x_n)$$

Wir behandeln im folgenden nur die natürliche RB (c)(i).

Satz 6.18 (Berechnung kubischer Splines)

26
8.12.09

Wir schreiben die Teilpolynome des Splines in der Form

$$p_i(x) = a_0^{(i)} + a_1^{(i)}(x-x_i) + a_2^{(i)}(x-x_i)^2 + a_3^{(i)}(x-x_i)^3 \quad i=1, \dots, n.$$

Die $a_2^{(i)}$ sind dann die Lösung des linearen Gleichungssystems der Dimension $n-1$:

$$h_i a_2^{(i-1)} + 2(h_i + h_{i+1}) a_2^{(i)} + h_{i+1} a_2^{(i+1)} = 3 \left(\frac{Y_{i+1} - Y_i}{h_{i+1}} - \frac{Y_i - Y_{i-1}}{h_i} \right) \quad i=1, \dots, n-1, \quad (6.1)$$

wobei $a_2^{(0)} = a_2^{(n)} = 0$ (natürliche Randbedingung!) und $h_i = x_i - x_{i-1}$.

Die restlichen Koeffizienten ergeben sich zu:

$$a_0^{(i)} = Y_i \quad (6.2)$$

$$a_1^{(i)} = \frac{Y_i - Y_{i-1}}{h_i} + \frac{h_i}{3} (2a_2^{(i)} + a_2^{(i-1)}) \quad (6.3)$$

$$a_3^{(i)} = \frac{a_2^{(i)} - a_2^{(i-1)}}{3h_i} \quad (6.4)$$

Beweis:

(i) Berechne Ableitungen der Teilpolynome

$$p_i'(x) = a_1^{(i)} + 2a_2^{(i)}(x-x_i) + 3a_3^{(i)}(x-x_i)^2 \quad (6.5)$$

$$p_i''(x) = 2a_2^{(i)} + 6a_3^{(i)}(x-x_i) \quad (6.6)$$

(ii) Interpolationsbed. nutzen. Einsetzen von x_i :

$$Y_i = p_i(x_i) = a_0^{(i)} \Rightarrow \boxed{a_0^{(i)} = Y_i} \quad i=1, \dots, n. \quad (6.7)$$

Das ist (6.2).

Einsetzen von x_{i-1} :

$$Y_{i-1} = p_i(x_{i-1}) = \underbrace{a_0^{(i)}}_{Y_i} + a_1^{(i)}(h_i) + a_2^{(i)}h_i^2 + a_3^{(i)}(-h_i^3) \quad i=1, \dots, n$$

$$\Leftrightarrow Y_{i-1} - Y_i = -h_i a_1^{(i)} + h_i^2 a_2^{(i)} - h_i^3 a_3^{(i)} \quad (6.8)$$

(iii) Randbedingungen einsetzen. Wir behandeln nur natürliche. 27
8.12.09

$$0 = p_1''(x_0) = 2a_2^{(1)} - 6a_3^{(1)}h_1 \quad (6.9)$$

$$0 = p_n''(x_n) = 2a_2^{(n)} \Rightarrow \boxed{a_2^{(n)} = 0} \quad (6.10)$$

(iv) Stetigkeit der ersten Ableitung

$$p_i'(x_i) = p_{i+1}'(x_i) \quad i = 1, \dots, n-1$$

$$\Leftrightarrow a_1^{(i)} = a_1^{(i+1)} - 2a_2^{(i+1)}h_{i+1} + 3a_3^{(i+1)}h_{i+1}^2 \quad (6.11)$$

(v) Stetigkeit der zweiten Ableitung

$$p_i''(x_i) = p_{i+1}''(x_i) \quad i = 1, \dots, n-1$$

$$\Leftrightarrow 2a_2^{(i)} = 2a_2^{(i+1)} - 6a_3^{(i+1)}h_{i+1} \quad (6.12)$$

(vi) Drücke $a_3^{(i)}$ durch $a_2^{(i)}$ aus; dh. löse (6.12) nach $a_3^{(i+1)}$ auf.

$$a_3^{(i+1)} = \frac{a_2^{(i+1)} - a_2^{(i)}}{3h_{i+1}} \quad i = 1, \dots, n-1$$

umnummerieren

$$\Leftrightarrow \boxed{a_3^{(i)} = \frac{a_2^{(i)} - a_2^{(i-1)}}{3h_i}} \quad i = 2, \dots, n \text{ aus (v)} \quad (6.13)$$

$i = 1$ aus (6.9) wenn man

$$\boxed{a_2^{(0)} := 0} \text{ setzt.}$$

also $i = 1, \dots, n$

Das ist (6.4)

(vii) $a_1^{(i)}$ durch $a_2^{(i)}$ ausdrücken. Dazu löse (6.8) nach a_1 auf:

$$a_1^{(i)} = \frac{y_i - y_{i-1}}{h_i} + h_i a_2^{(i)} - h_i^2 a_3^{(i)} \quad i = 1, \dots, n$$

$$\stackrel{(6.13)}{\text{aus.}} = \frac{y_i - y_{i-1}}{h_i} + h_i a_2^{(i)} - h_i^2 \left(\frac{a_2^{(i)} - a_2^{(i-1)}}{3h_i} \right)$$

$$\boxed{a_1^{(i)} = \frac{y_i - y_{i-1}}{h_i} + \frac{h_i}{3} (2a_2^{(i)} + a_2^{(i-1)})} \quad i = 1, \dots, n$$

Das ist (6.3)

(viii) Nun setze $a_2^{(i)}$ und $a_2^{(i)}$ in die verbleibende Gleichung (6.11) ein

28
8.12.09

$$\frac{y_i - y_{i-1}}{h_i} + \frac{h_i}{3} \left(\underbrace{2a_2^{(i)}}_{\sqrt{}} + \underbrace{a_2^{(i-1)}}_{\sqrt{}} \right) = \frac{y_{i+1} - y_i}{h_{i+1}} + \frac{h_{i+1}}{3} \left(\underbrace{2a_2^{(i+1)}}_{\sqrt{}} + \underbrace{a_2^{(i)}}_{\sqrt{}} \right) - 2a_2^{(i+1)} h_{i+1} + 3h_{i+1}^2 \left(\frac{a_2^{(i+1)}}{3h_{i+1}} - a_2^{(i)} \right) \quad i=1, \dots, n-1$$

(\Rightarrow)

$$a_2^{(i-1)} \left(+ \frac{h_i}{3} \right) + a_2^{(i)} \left(\frac{2h_i}{3} - \frac{h_{i+1}}{3} + h_{i+1} \right) + a_2^{(i+1)} \left(- \frac{2h_{i+1}}{3} + 2h_{i+1} - h_{i+1} \right) = \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{y_i - y_{i-1}}{h_i}$$

mal 3

$$\Leftrightarrow h_i a_2^{(i-1)} + 2(h_i + h_{i+1}) a_2^{(i)} + h_{i+1} a_2^{(i+1)} = 3 \left(\frac{y_{i+1} - y_i}{h_{i+1}} - \frac{y_i - y_{i-1}}{h_i} \right)$$

Und das ist (6.1).

Beachte, dass in (iii) und (vi) $a_2^{(0)} = a_2^{(n)} = 0$ gesetzt wurde. \blacksquare

Zur Lösung des Tridiagonalsystem:

- GEM/LR-Zerlegung hat in diesem Fall $O(n)$ Aufwand!
→ sehr schnell.
- Das Gleichungssystem ist symmetrisch und strikt diagonal dominant

$$\sum_{\substack{j=1 \\ j \neq i}}^{n-1} |a_{ij}| < |a_{ii}|$$

\Rightarrow Regularität, stabile LR-Zerlegung ohne Pivotisierung

Beispiel 6.19 Aus Numstoch Beispiel 6.2.

Satz 6.20 (Fehlerabschätzung)

29
8.12.09

Sei $f \in C^4([a, b])$. Erfüllt der kubische Spline

$$s''(a) = f''(a) \quad \text{und} \quad s''(b) = f''(b) \quad (\text{hatten wir oben nicht})$$

so gilt

$$\max_{a \leq x \leq b} |f(x) - s(x)| \leq \frac{1}{2} h^4 \max_{a \leq x \leq b} |f^{(4)}(x)|$$

$$\text{für } h := \max_{1 \leq i \leq n} |x_i - x_{i-1}|.$$

Beweis: Warner / Schaback II, siehe [Rannacher].
(wesentlich) □

Selbst unter noch schwächeren Voraussetzungen konvergiert die Spline-Interpolierende gleichmäßig gegen f .

Beispiel 6.21 Zur Konvergenzordnung Beispiel 6.5 aus Num Stoch.

6.6 Trigonometrische Interpolation

11.01.10

Problem: Interpolation „periodischer“ Funktionen, d.h. es gebe ein $\omega \in \mathbb{R}$, $\omega > 0$ so dass

$$f(x+\omega) = f(x) \quad \forall x \in \mathbb{R}.$$

Es bietet sich die Interpolation mit „trigonometrischen Summen“

$$t_n(x) = \frac{a_0}{2} + \sum_{k=1}^m \left\{ a_k \cos\left(\frac{kx2\pi}{\omega}\right) + b_k \sin\left(\frac{kx2\pi}{\omega}\right) \right\} \quad (6.14)$$

an, denn jeder Summand ist ω -periodisch: $\cos(k(x+\omega)2\pi/\omega) = \cos(kx2\pi/\omega + k2\pi) = \cos(kx2\pi/\omega)$

(6.14) hat $2m+1$ Parameter, deshalb setze

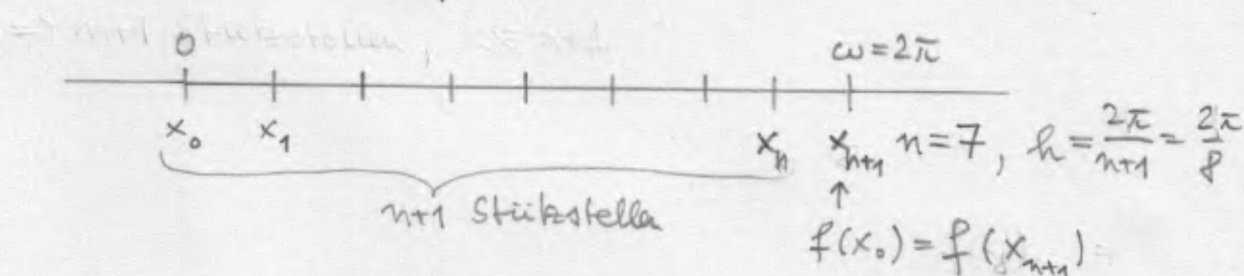
$$n := 2m$$

(bisher: Polynom vom Grad n hatte $n+1$ freie Parameter).

Ab jetzt: Nehme an, dass $\omega = 2\pi$ (spricht Transformation des Arguments).

Ausserdem: Verwende wieder äquidistante Stützstellen

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



Es zeigt sich: Die Interpolationsaufgabe $t_n(x_k) = f(x_k) \quad k=0, \dots, n$ ist zunächst einfacher im Körper \mathbb{C} zu lösen!

D.h. betrachte das komplexe trigonometrische Polynom

$$t_n^*(x) = \sum_{k=0}^n c_k e^{ikx}$$

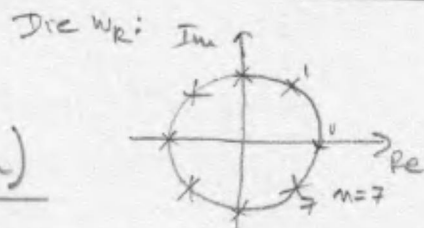
mit $i = \sqrt{-1}$ imaginäre Einheit

$c_k \in \mathbb{C}$ komplexe Koeffizienten

Eulersche Identität $e^{i\varphi} = \cos \varphi + i \sin \varphi$ für $\varphi \in \mathbb{R}$.

Wir untersuchen zunächst die Eigenschaften der komplexen Exponentialfunktion:

11.01.10



Hilfssatz 6.22 (Komplexe Einheitswurzeln)

Setze $w_k := e^{ix_k} = e^{i \frac{2\pi k}{n+1}}$ für alle $k \in \mathbb{Z}$ und gegebenes $n \in \mathbb{N}$.
 $w_k \in \mathbb{C}$ heißt „ k -te Einheitswurzel“ und hat folgende Eigenschaften

a) $w_k^{n+1} - 1 = 0$ für alle $k \in \mathbb{Z}$.

M. a. W. die w_k sind Lösungen der Gleichung $w^{n+1} - 1 = 0$ in \mathbb{C} .

Bew: $w_k^{n+1} = \left(e^{i \frac{2\pi k}{n+1}} \right)^{n+1} = e^{i 2\pi k} = \underbrace{\cos 2\pi k}_{=1} + i \underbrace{\sin 2\pi k}_{=0} = 1.$

b) $w_R^j = w_j^k \quad \forall j, k \in \mathbb{Z}$

Bew: $w_R^j = \left(e^{i \frac{2\pi k}{n+1}} \right)^j = e^{i \frac{2\pi k j}{n+1}} = \left(e^{i \frac{2\pi j}{n+1}} \right)^k = w_j^k$

c) $w_R^{-j} = w_j^{-k} \quad \forall j, k \in \mathbb{Z}$

Beweis genau wie b, ist aber nicht identisch zu b).

d) $w_R^j = w_R^{-j \bmod (n+1)} = w_{k \bmod (n+1)}^j = w_{k \bmod (n+1)}^j \quad \forall j, k \in \mathbb{Z}$.

zeige nur erste Identität, Rest analog:

Sei $j = r(n+1) + s$ mit $0 \leq s < n+1$.

$w_R^j = e^{i \frac{2\pi k j}{n+1}} = e^{i \frac{2\pi [k r(n+1) + k s]}{n+1}} = \underbrace{e^{2\pi k r}}_{=1} \cdot \underbrace{e^{i \frac{2\pi k s}{n+1}}}_{(w_R^k)^s} = w_R^{j \bmod (n+1)}$

e) $\sum_{j=0}^n w_R^j = \begin{cases} n+1 & k \bmod (n+1) = 0 \\ 0 & \text{sonst} \end{cases}$

Sei $k=0$: $w_0^j = e^0 = 1 \quad \forall j$, also $\sum_{j=0}^n 1 = n+1$.

$k \neq 0$: w_k ist nach a) Lösung von „Teleskopsumme“

$w^{n+1} - 1 = (w-1)(w^n + w^{n-1} + \dots + 1) = 0$

Für $k \neq 0$ ist $w_k \neq 1$ also $w_k - 1 \neq 0$ somit muss der zweite Faktor

also $\sum_{j=0}^n w_R^j = 0$ sein.

Satz 6.23 (Komplexwertige trigonometrische Interpolation)

Zu gegebenen Zahlen $y_0, \dots, y_n \in \mathbb{C}$ gibt es genau eine Funktion der Gestalt

$$t_n^*(x) = \sum_{k=0}^n c_k e^{ikx},$$

die den Interpolationsbedingungen

$$t_n^*(x_j) = y_j \quad j=0, \dots, n, \quad x_j = \frac{2\pi j}{n+1}$$

genügt. Die komplexen Koeffizienten sind bestimmt durch

$$c_k = \frac{1}{n+1} \sum_{j=0}^n y_j \underbrace{e^{-ijx_k}}_{= e^{-i \frac{2\pi jk}{n+1}} = w_k^{-j}} \quad \forall k=0, \dots, n. \quad (6.15)$$

Beweis: Mit der Abkürzung $w = e^{ix}$ gibt

$$t_n^*(x) = \sum_{k=0}^n c_k e^{ikx} = \sum_{k=0}^n c_k \underbrace{(e^{ix})^k}_{= w^k} = \sum_{k=0}^n c_k w^k = P_n(w).$$

Jedem t_n^* entspricht also ein genau komplexes Polynom P_n vom Grad n . Transformation $w = e^{ix}$ eindeutig auf $0 \dots 2\pi$

Übertragung der Interpolationsbedingungen für t_n^* auf P_n ergibt:

$$t_n^*(x_j) = P_n(\underbrace{e^{ix_j}}_{= w_j}) = y_j \quad \forall j=0, \dots, n.$$

Da die Polynominterpolation zu paarw. vord. Stützstellen (auch im Komplexen) eindeutig ist, gibt es genau ein solches P_n , also auch t_n^* .

Bleibt die Berechnung der Koeffizienten c_k . Diese ergeben sich durch das lineare Gleichungssystem

$$P_n(e^{ix_j}) = \sum_{l=0}^n c_l (e^{ix_j})^l = \sum_{l=0}^n c_l e^{i \frac{2\pi lj}{n+1}} = \sum_{l=0}^n c_l w_j^l = y_j \quad j=0, \dots, n.$$

(n+1) x (n+1) LGS für die Koeffizienten

Das Gleichungssystem kann man explizit auflösen

12.01.10

Für ein $k \in 0, \dots, n$:

$$\sum_{j=0}^n w_k^{-j} \left(\sum_{l=0}^n c_l w_j^l \right) = \sum_{l=0}^n c_l \left(\sum_{j=0}^n w_j^{l-k} \right) = \sum_{j=0}^n w_k^{-j} y_j$$

\uparrow
 $= y_j$ 1) \sum vertauscht 2) $w_k^{-j} = w_j^{-k}$

$$\sum_{j=0}^n w_j^{l-k} = \begin{cases} n+1 & l-k=0 \\ 0 & l-k \neq 0 \end{cases}$$

und damit

$$c_k^{(n+1)} = \sum_{j=0}^n y_j w_k^{-j} \Leftrightarrow c_k = \frac{1}{n+1} \sum_{j=0}^n y_j e^{-ijx_k}$$

Damit ist auch die reelle Interpolationsaufgabe gelöst, indem man $y_j \in \mathbb{R}$ annimmt.

Zu zeigen ist nun, wie die Koeffizienten a_k, b_k berechnet werden.

Satz 6.24 (Diskrete Fourier-Analyse)

Für $n \in \mathbb{N}_0$ gibt es zu gegebenen reellen Zahlen y_0, \dots, y_n genau ein trigonometrisches Polynom der Form

$$t_n(x) = \frac{a_0}{2} + \sum_{k=1}^m \{ a_k \cos(kx) + b_k \sin(kx) \} + \frac{\theta}{2} a_{m+1} \cos((m+1)x)$$

mit $t_n(x_j) = y_j$, $j=0, \dots, n$, $x_j = \frac{2\pi j}{n+1}$, sowie

$$\theta = 0, m = \frac{n}{2} \quad n \text{ gerade} \rightarrow a_0, \dots, a_m, b_1, \dots, b_m$$

$$\theta = 1, m = \frac{n-1}{2} \quad n \text{ ungerade} \rightarrow a_0, \dots, a_{m+1}, b_1, \dots, b_m$$

(n gerade: $2 \cdot (n/2) + 1 = 2m + 1$, n ungerade: $2 \cdot (n-1)/2 + 2 = n+1$).

Die Koeffizienten werden bestimmt durch

$$a_k = \frac{2}{n+1} \sum_{j=0}^n y_j \cos(jx_k), \quad b_k = \frac{2}{n+1} \sum_{j=0}^n y_j \sin(jx_k)$$

Beweisskizze: (ausführlich: siehe Ra-Skript)

5
12.01.10

Man bestimmt die Koeffizienten c_k des komplexen trigonometrischen Polynoms zu reellen Daten $y_j \in \mathbb{R}$ und setzt dann

$$a_0 = 2c_0$$

$$a_k = c_k + c_{n+1-k} \quad k=1, \dots, m$$

$$b_k = i(c_k - c_{n+1-k}) \quad k=1, \dots, m$$

$$a_{n+1} = 2c_{n+1} \quad n=2m+1 \text{ (} n \text{ ungerade).}$$

Dann rechnet man nach, dass $t_n^*(x_j) = t_n(x_j) = y_j \in \mathbb{R}$.

Wir wollen noch nachrechnen, dass die a_k reell sind:

$$\begin{aligned} a_k &= c_k + c_{n+1-k} \\ (6.15) \rightarrow &= \frac{1}{n+1} \sum_{j=0}^n y_j (e^{-ijx_k} + e^{-ijx_{n+1-k}}) \\ &= \frac{1}{n+1} \sum_{j=0}^n y_j (e^{-ijx_k} + e^{ijx_k}) \\ &= \frac{1}{n+1} \sum_{j=0}^n y_j (\underbrace{\cos(-jx_k)}_{=\cos(jx_k)} + i \underbrace{\sin(-jx_k)}_{=-\sin(jx_k)} + \cos(jx_k) + i \sin(jx_k)) \\ &= \frac{1}{n+1} \sum_{j=0}^n y_j \cdot 2 \cos(jx_k). \end{aligned}$$

$e^{-i \frac{2\pi j(n+1-k)}{n+1}} = e^{i \frac{2\pi jk}{n+1}} = e^{ijx_k}$

Schnelle Fouriertransformation

6
12.01.10

Einer der berühmtesten Algorithmen der angewandten Mathematik / Informatik!

Entwickelt von James Cooley und John Tukey 1965.

Zurück zur komplexen trigonometrischen Interpolation:

$$c_k = \frac{1}{N} \sum_{j=0}^{N-1} y_j e^{-i \frac{2\pi j k}{N}} \quad k=0, \dots, N-1 \quad \text{(Hintransformation)} \quad (6.15a)$$
$$y_j = \sum_{k=0}^{N-1} c_k e^{i \frac{2\pi k j}{N}} \quad j=0, \dots, N-1 \quad \text{(Rücktransformation)} \quad (6.15b)$$

- Hier wurde $N = n+1$ gesetzt.
- Die Berechnung der y_j aus den c_k ist schlicht die Definition des komplexen trig. Polynoms
- 6.15 bezeichnet man als diskrete Fouriertransformation.

Aufwand: Setze $c := (c_0, \dots, c_{N-1})^T$, $y := (y_0, \dots, y_{N-1})$

$$c := (c_0, \dots, c_{N-1})^T, \quad y := (y_0, \dots, y_{N-1})$$

Dann ist (6.15a) äquivalent zu

$$c = \frac{1}{N} W y$$

$$\text{mit } (W)_{kj} = e^{-i \frac{2\pi j k}{N}} = W_k^{-j} \quad (\text{Einheitswurzel})$$

Die Rücktransformation (6.15b) lautet

$$y = U c \quad \text{mit } (U)_{j,k} = W_j^k$$

Wegen $U \frac{1}{N} W = I$ ist $W^{-1} = \frac{1}{N} U$, man hat also explizite Fact. der Inversa.

Da W, U voll besetzt sind beträgt der Aufwand für Hin- und Rücktransformation jeweils $O(N^2)$.

Betrachte (6.15a) ohne den Vorfaktor $\frac{1}{N}$, nenne das \tilde{c}_k .
Es sei N gerade, dann gilt:

$$\tilde{c}_k = \sum_{j=0}^{N-1} y_j e^{-i \frac{2\pi jk}{N}}$$

Aufspalten
der Summe in
gerade und unger.
Indizes

$$\underbrace{\sum_{j=0}^{N/2-1} y_{2j} e^{-i \frac{2\pi 2jk}{N}}}_{\text{gerader Teil}} + \underbrace{\sum_{j=0}^{N/2-1} y_{2j+1} e^{-i \frac{2\pi (2j+1)k}{N}}}_{\text{ungerader Teil}}$$

bringe
N/2 in
Nenn

$$\underbrace{\sum_{j=0}^{N/2-1} y_{2j} e^{-i \frac{2\pi jk}{N/2}}}_{=: \tilde{c}_k^g} + e^{-i \frac{2\pi k}{N}} \underbrace{\sum_{j=0}^{N/2-1} y_{2j+1} e^{-i \frac{2\pi jk}{N/2}}}_{=: \tilde{c}_k^u}$$

Wegen der N/2-Periodizität der Einheitswurzeln $e^{-i \frac{2\pi k}{N/2}}$ gilt

$$\tilde{c}_{k+N/2}^g = \tilde{c}_k^g \quad \text{und} \quad \tilde{c}_{k+N/2}^u = \tilde{c}_k^u \quad \forall k=0, \dots, N/2-1.$$

Damit gilt:

- $\tilde{c}_k^g, \tilde{c}_k^u, k=0, \dots, N/2-1$ berechnen sich jeweils durch eine DFT der Länge $N/2$
- Daraus berechnet man dann die ursprünglich gesuchten Koeffizienten

$$\tilde{c}_k = \tilde{c}_k^g + e^{-i \frac{2\pi k}{N}} \tilde{c}_k^u \quad 0 \leq k < N/2 \quad (6.16a)$$

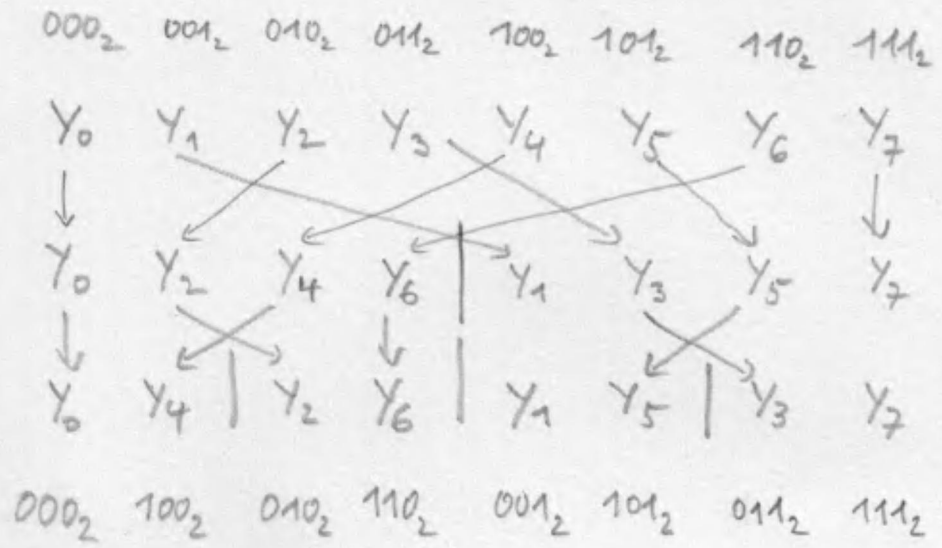
$$\tilde{c}_k = \tilde{c}_{k-N/2} + e^{-i \frac{2\pi k}{N}} \tilde{c}_{k-N/2}^u \quad N/2 \leq k < N \quad (6.16b)$$

- Falls $N/2$ wieder gerade kann man das Prinzip rekursiv fortsetzen.
- Ist $N=2^d$ eine Zweierpotenz erhält man schließlich eine DFT der Länge 2, die einfach durchführbar ist.

Beispiel $N=8$.

"Abstiegsphase" der Rekursion: Umsortieren der Eingabedaten.

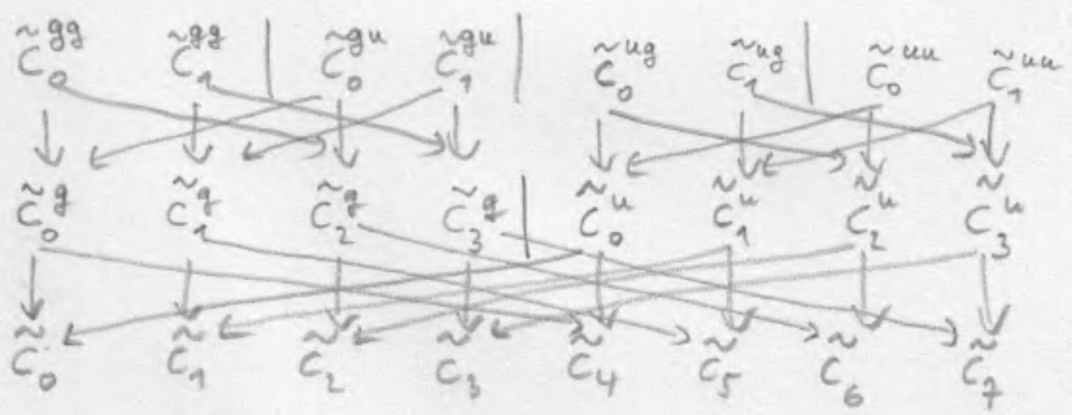
Am besten schreibt man die Indizes zur Basis 2



Permutation $(b_{d-1} \dots b_0)_2 \rightarrow (b_0 \dots b_{d-1})_2$ heißt "bit reversal"

"Aufstiegsphase": Rekombination der Koeffizienten nach (6.16)

↑ von unten nach oben



Dieses Verknüpfungsmuster nennt man "perfect shuffle".

Aufwand der FFT: $N = 2^d$ Zweipötenz
 $\Rightarrow d = \log_2 N$.

12.01.10

$$A(N) = 2A\left(\frac{N}{2}\right) + cN$$

Gleitkom-
operationen
für DFT der
Länge N

2 Transf. der
Länge $N/2$ Ber.

Aufwand für (6.16)

$$= 2 \left[2A\left(\frac{N}{4}\right) + c\frac{N}{2} \right] + cN$$

$$= 4A\left(\frac{N}{4}\right) + cN + cN$$

$$= 4 \left(2A\left(\frac{N}{8}\right) + c\frac{N}{4} \right) + cN + cN$$

$$= 8A\left(\frac{N}{8}\right) + cN + cN + cN$$

d mal

$$= 2^d A(1) + \underbrace{cN + \dots + cN}_{d-1 \text{ mal}} = cdN$$

$$= O(N \log N)$$

deutlich schneller für große N .

Praktisches zur DFT:

- Spektralanalyse
- Beispiele auf dem Computer.

Wir betrachten Approximation von Funktionen in
Prähilberträumen:

Definition 6.25 Ein Vektorraum von Funktionen über \mathbb{R} oder \mathbb{C} mit Skalarprodukt heißt Prähilbertraum.

Beispiele:

a) Raum der stetigen Funktionen $C(a,b)$ mit dem Skalarprodukt

$$(f, g) = \int_a^b f(x) \overline{g(x)} dx.$$

b) Raum der quadratintegrierbaren Funktionen

$$L^2(a,b) = \left\{ f : \int_a^b |f(t)|^2 dx < \infty \right\}$$

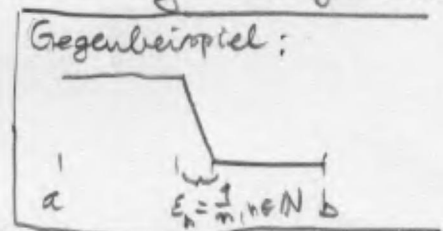
Dabei ist $\int_a^b dx$ das „Lebesgue-Integral“ (\rightarrow siehe Funktional analysis),
jede Cauchy-Folge konvergiert.

Im Gegensatz zu $C(a,b)$ ist $L^2(a,b)$ vollständig bezüglich der Norm

$$\|f\| = \sqrt{(f, f)}$$

und damit ein Hilbertraum.

$L^2(a,b)$ ist das Analogon des \mathbb{R}^n für Funktionenräume.



c) Gegeben $\Psi = \{\psi_1, \dots, \psi_N\}$, $\psi_i \in C(a,b)$ (z.B.)

$S = \text{span } \Psi = \left\{ f : f = \sum_{i=1}^N c_i \psi_i \right\}$ ist endlichdimensionaler Prähilbertraum.

Im folgenden sei H Prähilbertraum und $S \subseteq H$ ein endlichdimensionaler Teilraum.

Wir betrachten folgende Aufgabe:

11
12.01.10

Zu $f \in H$ finde $g \in S$ so dass

$$\|f - g\| \rightarrow \min \quad (6.17)$$

wobei $\|f\| = (f, f)^{1/2}$ die durch das Skalarprodukt induzierte Norm ist.

Satz 6.26 (Allgemeine Gauß-Approximation)

Die Aufgabe (6.17) hat genau eine Lösung $g \in S$.

Dies ist charakterisiert durch

$$(g, \varphi) = (f, \varphi) \quad \forall \varphi \in S. \quad (6.18)$$

Beweis

a) Sei $\|f - g\|$ minimal für $g \in S$.

$$F_\varphi(t) := \|f - (g + t\varphi)\|^2, \quad t \in \mathbb{R}$$

ist für beliebiges $\varphi \in S$ eine stetig differenzierbare Funktion mit Minimum bei $t=0$.

Also

$$\begin{aligned} \frac{d}{dt} F_\varphi(t) \Big|_{t=0} &= \frac{d}{dt} (f - g - t\varphi, f - g - t\varphi) \Big|_{t=0} \\ &= \frac{d}{dt} \left[(f - g, f - g) - 2t(f - g, \varphi) + t^2(\varphi, \varphi) \right] \Big|_{t=0} \\ &= \left[-2(f - g, \varphi) + 2t(\varphi, \varphi) \right] \Big|_{t=0} \\ &= -2(f - g, \varphi) \stackrel{!}{=} 0 \quad (\text{für alle } \varphi \in S) \end{aligned}$$

Dies ist (6.18).

Geometrisch: „Fehler“ $f - g$ steht senkrecht auf allen $\varphi \in S$.

b) Sei nun $(f - g, \varphi) = 0 \quad \forall \varphi \in S$. Für beliebiges $g' \in S$ gilt dann

$$\|f - g'\|^2 = \|f - g - \varphi\|^2 = (f - g - \varphi, f - g - \varphi)$$

$$g' = g + \underbrace{g' - g}_{=: \varphi} = g + \varphi \quad = (f - g, f - g) - 2t \underbrace{(f - g, \varphi)}_{= 0 \text{ n. Vor.}} + (\varphi, \varphi)$$

$$= \|f - g\|^2 + \|\varphi\|^2 \geq \|f - g\|^2 \quad \text{also } g \text{ Minimum.}$$

Damit ist die Äquivalenz gezeigt:

12
12.01.10

$$\|f - g\| \rightarrow \min \Leftrightarrow (g, \varphi) = (f, \varphi) \quad \forall \varphi \in S.$$

c) Eindeutigkeit des Minimums.

Angenommen es gibt zwei Minima g_1, g_2 mit $g_1 \neq g_2$. Dann ist

$$\begin{aligned} \|f - g_1\|^2 &= \|f - g_2 + \underbrace{g_1 - g_2}_{=:\varphi}\|^2 \\ &= \|f - g_2\|^2 + \underbrace{\|g_1 - g_2\|^2}_{> 0 \text{ da } g_1 \neq g_2} > \|f - g_2\|^2 \quad \downarrow \text{ zu } g_1 \text{ Minimum.} \end{aligned}$$

d) Existenz des Minimums. Dies zeigen wir konstruktiv und erhalten damit auch einen Algorithmus zur Berechnung von g .

Da S endlich-dimensional gibt es eine Basis $\Psi = \{\psi_1, \dots, \psi_N\}$, $N = \dim S$. Das gesuchte Element g hat die Darstellung

$$g = \sum_{j=1}^N \alpha_j \psi_j.$$

Einsetzen der Basisdarstellung in (6.18) liefert

$$(g, \varphi) = (f, \varphi) \quad \forall \varphi \in S$$

$$\Leftrightarrow \left(\sum_{j=1}^N \alpha_j \psi_j, \psi_i \right) = (f, \psi_i) \quad i = 1, \dots, N$$

$$\Leftrightarrow \sum_{j=1}^N \alpha_j (\psi_j, \psi_i) = (f, \psi_i) \quad i = 1, \dots, N$$

$$\Leftrightarrow A \alpha = b$$

mit $(A)_{ij} = (\psi_j, \psi_i)$ Massenmatrix oder Gramsche Matrix

$$(b)_i = (f, \psi_i)$$

A ist symmetrisch und positiv definit, denn

$$i) \quad (A)_{ij} = (\psi_j, \psi_i) \underset{\substack{\uparrow \\ \text{Sym. des SP.}}}{=} (\psi_i, \psi_j) = (A)_{ji}$$

$$ii) \quad \alpha^T A \alpha = \sum_{i=1}^N \alpha_i \left(\sum_{j=1}^N (A)_{ij} \alpha_j \right) = \sum_{i=1}^N \alpha_i \left(\sum_{j=1}^N (\psi_j, \psi_i) \alpha_j \right)$$

$$= \sum_{i=1}^N \alpha_i \left(\sum_{j=1}^N \alpha_j \psi_j, \psi_i \right) = \left(\underbrace{\sum_{j=1}^N \alpha_j \psi_j}_{=: g}, \underbrace{\sum_{i=1}^N \alpha_i \psi_i}_{=: g} \right)$$

$$= (g, g) > 0 \quad \text{da } g \neq 0 \text{ wenn } \alpha \neq 0.$$

(Wir nehmen hier $K = \mathbb{R}$ an, sonst hermitesch).

Da A s.p.d. hat das LGS $Ax = b$ genau eine Lösung für jedes b.

Nach a)-c) ist $\sum \alpha_j \psi_j$ die eindeutig bestimmte Bestapproximation.

Approximation mit Orthonormalbasen

Besonders einfach wird die Lösung der Approximationsaufgabe wenn Ψ eine Orthonormalbasis ist, d.h. $(\psi_i, \psi_j) = \delta_{ij}$.

Dann gilt

$$\sum_{j=1}^N \alpha_j (\psi_j, \psi_i) = \alpha_i \stackrel{!}{=} (f, \psi_i) \quad i = 1, \dots, N$$

und somit

$$g = \sum_{j=1}^N \alpha_j \psi_j = \sum_{j=1}^N (f, \psi_j) \psi_j$$

Beispiel 6.27 (Fourierreihe)

Für $N = 2m + 1$, $m \in \mathbb{N}$ ist

$$\Psi_F = \left\{ \frac{1}{\sqrt{2\pi}}, \frac{1}{\sqrt{\pi}} \cos x, \dots, \frac{1}{\sqrt{\pi}} \cos(mx), \frac{1}{\sqrt{\pi}} \sin x, \dots, \frac{1}{\sqrt{\pi}} \sin(mx) \right\}$$

eine Orthonormalbasis von Funktionen auf dem Intervall $[-\pi, \pi]$. (nachrechnen!).

Damit gilt dann

14
12.01.10

$$g(x) = \frac{a_0}{2} + \sum_{k=1}^m \{ a_k \cos(kx) + b_k \sin(kx) \}$$

und
$$a_k = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos(kx) dx, \quad k=0, \dots, m,$$

$$b_k = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin(kx) dx, \quad k=1, \dots, m.$$

Für unendlich viele Glieder ($m=\infty$) nennt man die Reihe Fourier-Reihe. Diese konvergiert gegen ein Element aus $L^2(-\pi, \pi)$.

Ü: DFT lässt sich als näherungsweise Auswertung der Integrale mittels Trapezregel verstehen.

Fehlerkontrolle

Bisher: $S \subset H$ fest gewählt; Berechnung des „optimalen“ g .
Der Fehler $\|f-g\|$ wird akzeptiert.

Bsp.: Fourierreihe für $m \rightarrow \infty$ sollte $\|f-g\|$ immer kleiner werden.

Verfeinerung der Approximationsaufgabe:

Specialfall „Kompression“:
H ist selbst schon endlichdimensional!

Finde $S \subset H$ mit $\dim S$ möglichst klein, so dass

$$\|f-g\| \leq \text{ToL} \quad (= \text{vorgegebene Zahl}).$$

Mittels Orthonormalbasen lässt sich der Fehler recht einfach messen.

Sei $S_N \subset H$, $N \in \mathbb{N}$, eine Folge von Approximationsräumen mit $\dim S_N = N$. Oft gilt ^{sogar} $S_N \subset S_{N+1}$.

Sei nun g_N die Bestapproximation von f in S_N .

Weiter sei ψ_N eine Orthonormalbasis von S_N (oft gilt $\psi_N \subset \psi_{N+1}$).

Dann gilt für den Fehler:

15
12.01.10

$$\begin{aligned} 0 \leq \|f - g_N\|^2 &= (f - g_N, f - g_N) \\ &= (f, f) - 2(f, g_N) + (g_N, g_N) \\ \text{Basis} & \\ \text{einsetzen} & \quad \rightarrow = (f, f) - 2\left(f, \sum_{i=1}^N \underbrace{(f, \psi_i)}_{x_i} \psi_i\right) + \left(\sum_{i=1}^N (f, \psi_i) \psi_i, \sum_{j=1}^N (f, \psi_j) \psi_j\right) \\ &= (f, f) - 2 \sum_{i=1}^N (f, \psi_i) (f, \psi_i) + \sum_{i=1}^N \sum_{j=1}^N \underbrace{(f, \psi_i)}_{x_i} \underbrace{(f, \psi_j)}_{x_j} \underbrace{(\psi_i, \psi_j)}_{=\delta_{ij}} \\ &= (f, f) - 2 \sum_{i=1}^N (f, \psi_i)^2 + \sum_{i=1}^N (f, \psi_i)^2 \\ &= (f, f) - \sum_{i=1}^N (f, \psi_i)^2 \end{aligned}$$

Wir wollen

$$\|f - g_N\|^2 = (f, f) - \sum_{i=1}^N (f, \psi_i)^2 \leq \underbrace{\text{TOL}}_{\text{relative Toleranz}} \cdot (f, f)$$

$$\Leftrightarrow \sum_{i=1}^N (f, \psi_i)^2 \geq (1 - \text{TOL}) (f, f)$$

relative Toleranz

$$\text{TOL} \in [0, 1)$$

$$\text{TOL} = 1 \Rightarrow N = 1 \quad (\text{b. 18})$$

$$\text{TOL} = 0 \Rightarrow N = \infty \text{ und } g = f.$$

(Kompression: $H = S_M$, dann ist $N = M$ und $\|f - g_N\| = 0$ möglich)

Bemerkung:

a) (f, f) sei berechenbar (zumindest mit ausreichender Genauigkeit, bei Kompression hat man das).

b) Falls $\Psi_N \subset \Psi_{N+1}$ sind einfach so viele Basisfunktionen hinzuzufügen bis (b. 18) erreicht ist.

c) Für die Fourierreihe gilt $\Psi_N \subset \Psi_{N+1}$

d) Aus der obigen Fehlerdarstellung folgt bei $\Psi_N \subset \Psi_{N+1}$

unmittelbar

$$\|f - g_{N+1}\|^2 = (f, f) - \sum_{i=1}^{N+1} (f, \psi_i)^2 = (f, f) - \underbrace{\sum_{i=1}^N (f, \psi_i)^2}_{\|f - g_N\|^2} - \underbrace{(f, \psi_{N+1})^2}_{\geq 0}$$

$$\leq \|f - g_N\|^2.$$

Der Fehler nimmt also höchstens ab!

Es bezeichne

$$\text{tr}(f) = \{x : f(x) \neq 0\}$$

den Träger einer Funktion.

In Oben gewählte ONB aus Sinus u. Cosinusfunktion haben alle ψ_i im wesentl. globalen Träger (d.h. $\overline{\text{tr}(f)} = [a, b]$)

Oft variiert die zu approximierende Funktion aber nur lokal sehr stark. In diesem Fall möchte man ein Funktionensystem mit folgenden Eigenschaften:

- a) $(\psi_i, \psi_j) = \delta_{ij}$ Orthonormalität
- b) $\psi_N \subset \psi_{N+1}$ Geschachteltheit \rightarrow „inkrementelles Hinzufügen“
- c) $\text{diam}(\text{tr}(\psi_i)) \rightarrow 0$ für $i \rightarrow \infty$

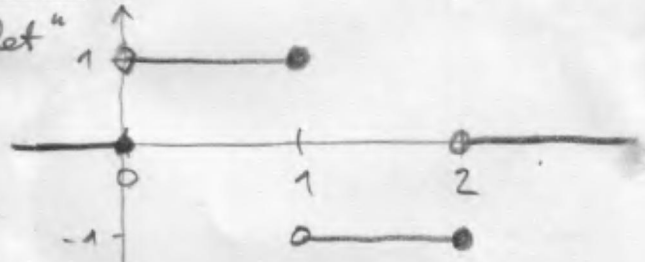
Diese Eigenschaften besitzen sog. Waveletfunktionen.

Wir untersuchen das sog. Haar-Wavelet.

Definition 6.28 (Haar-Wavelet)

Wir definieren das „mother wavelet“

$$\psi(x) = \begin{cases} 1 & 0 < x \leq 1 \\ -1 & 1 \leq x \leq 2 \\ 0 & \text{sonst} \end{cases}$$



und die Abschneidefunktion auf $(0, 1]$:

$$\chi(x) = \begin{cases} 1 & 0 < x \leq 1 \\ 0 & \text{sonst} \end{cases}$$

Für $l \in \mathbb{N}_0$ (Stufe) und $0 \leq i < 2^{l-1}$ (Index) ist das Haar-Wavelet gegeben durch

$$\psi_i^l(x) = \max(\sqrt{2^{l-1}}, 1) \cdot \psi(2^l x - 2i) \cdot \chi(x).$$

Die Haar-Waveletbasis der Stufe l ist

$$\underline{\psi}^l = \{\psi_0^0\} \cup \bigcup_{j=1}^l \bigcup_{i=0}^{2^j-1} \{\psi_i^j\}.$$

Wir fassen einige Eigenschaften der Haar-Wavelets zusammen:

17
16.01.10

a) Für $l > 0$ gilt

$$\psi_i^l(x) = \begin{cases} 0 & x \leq \frac{2i}{2^l} \vee x > \frac{2i+2}{2^l} \\ \sqrt{2^{l-1}} & \frac{2i}{2^l} < x \leq \frac{2i+1}{2^l} \\ -\sqrt{2^{l-1}} & \frac{2i+1}{2^l} < x \leq \frac{2i+2}{2^l} \end{cases}$$

Betrachte Argument des mother wavelet

$$2^l x - 2i \leq 0 \Leftrightarrow x \leq \frac{2i}{2^l} \rightarrow \psi_i^l(x) = 0 \text{ für } x \leq \frac{2i}{2^l}$$

$$2^l x - 2i \leq 1 \Leftrightarrow x \leq \frac{2i+1}{2^l} \rightarrow \psi_i^l(x) = 1 \text{ für } \frac{2i}{2^l} < x \leq \frac{2i+1}{2^l}$$

$$2^l x - 2i \leq 2 \Leftrightarrow x \leq \frac{2i+2}{2^l} \rightarrow \psi_i^l(x) = -1 \text{ für } \frac{2i+1}{2^l} < x \leq \frac{2i+2}{2^l}$$

und $\psi_i^l(x) = 0$ für $x > \frac{2i+2}{2^l}$

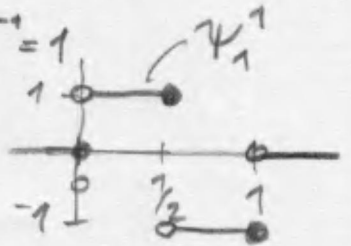
Der letzte Index ist $2^{l-1} - 1 \in \mathbb{N}_0$ (für $l > 0$!) und damit

$$\frac{2i+2}{2^l} \leq \frac{2(2^{l-1}-1)+2}{2^l} = \frac{2^l - 2 + 2}{2^l} = 1.$$

Für $l > 0$ ist also $\psi(2^l x - 2i) = 0$ für $x \leq 0$ und $x > 1$ und somit $\chi(x)$ überflüssig

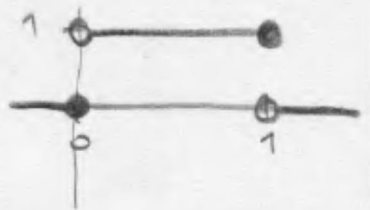
b) Speziell für $l=1$ (in a) enthalten) gilt $0 \leq i < 2^{1-1} = 1$

$$\psi_0^1(x) = \underbrace{\max(\sqrt{2^{1-1}}, 1)}_{=1} \cdot \underbrace{\psi(2x)}_{=0 \text{ für } x \leq 0, x > 1} \cdot \chi(x) = \psi(2x)$$



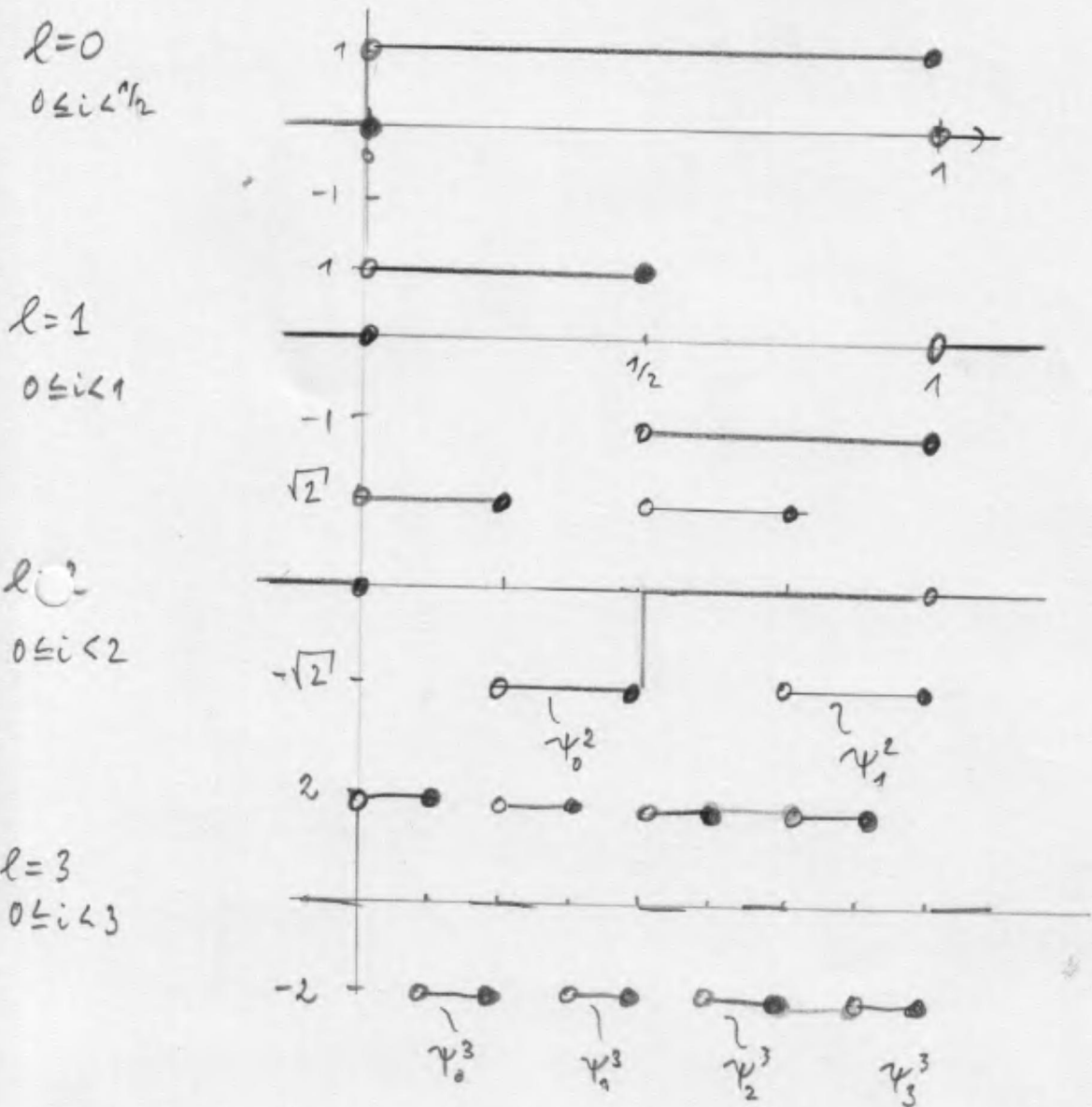
c) Für $l=0$ gilt $0 \leq i < 2^{-1} = \frac{1}{2} \rightarrow$ Ein Index $i=0$

$$\psi_0^0(x) = \underbrace{\max(\sqrt{1/2}, 1)}_{=1} \cdot \underbrace{\psi(x)}_{=1 \text{ für } 0 \leq x \leq 1} \cdot \chi(x) = \chi(x)$$



d) Graphische Veranschaulichung der Haar-Wavelets

16.01.10



$$\Rightarrow \text{tr}(\psi_i^l) = \text{tr}(\psi_{2i}^{l+1}) \cup \text{tr}(\psi_{2i+1}^{l+1})$$

Die Wavelets bilden eine Baumstruktur bezüglich der Inklusion der Träger.

ψ_j^{l+1} ist Kind von ψ_i^l genau dann wenn $\text{tr}(\psi_j^{l+1}) \subseteq \text{tr}(\psi_i^l)$



vollständiger Binärer Baum.
 \Rightarrow effiziente Auswertung!

e) Orthogonalitätseigenschaften

Es gilt $(\psi_i^l, \psi_j^k) = \begin{cases} 1 & i=j \wedge l=k \\ 0 & \text{sonst.} \end{cases}$

1) $i=j \wedge l=k$

Für $l=0$ stimmt es offensichtlich $\frac{2i+2}{2^e}$

Für $l>0$ gilt $(\psi_i^l, \psi_i^l) = \int_{\frac{2i}{2^e}}^{\frac{2i+2}{2^e}} (\sqrt{2^{e-1}})^2 dx = 2^{e-1} \cdot \frac{1}{2^{e-1}} = 1$

2) Sei $l=k$ aber $i \neq j$

dann gilt $\text{tr}(\psi_i^l) \cap \text{tr}(\psi_j^l) = \emptyset$ also $(\psi_i^l, \psi_j^l) = 0$

3) o.B.d.A. sei $l > k$ und somit insb. $l > 0$

Dann ist ψ_j^k konstant auf $\text{tr}(\psi_i^l)$ und somit

$(\psi_i^l, \psi_j^k) = c \int_{\frac{2i}{2^e}}^{\frac{2i+2}{2^e}} \psi_i^l(x) dx = 0$ da ψ_i^l

↑
Wert von ψ_j^k $\frac{2i}{2^e}$

f) Offensichtlich erfüllt die Konstruktion

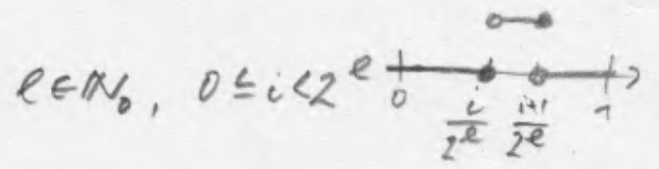
$\Psi^l \subset \Psi^{l+1}$

da immer nur Basisfunktionen hinzugenommen werden.

Damit sind alle geforderten Eigenschaften erfüllt.

g) Um die durch Ψ^l darstellbaren Funktionen noch anschaulicher zu machen definieren wir die Funktionen

$\varphi_i^l(x) = \begin{cases} 1 & \frac{i}{2^e} < x \leq \frac{i+1}{2^e} \\ 0 & \text{sonst} \end{cases}$

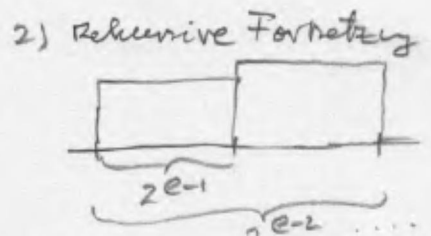
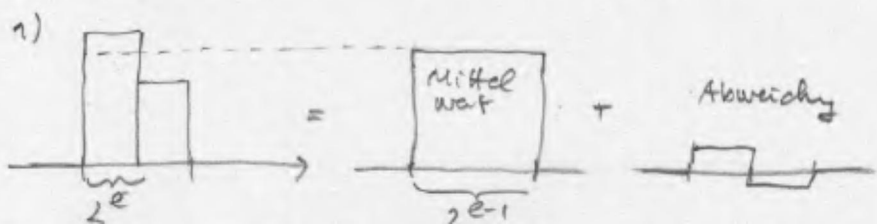


$\Phi^l = \bigcup_{i=0}^{2^e-1} \{\varphi_i^l\}$

$S^l = \text{span } \Phi^l$ ist der Raum der stückweise konstanten Funktionen auf dem Intervall $(0,1]$ bezüglich der Unterteilung $(\frac{i}{2^e}, \frac{i+1}{2^e}]$.

Es gilt $\text{span } \Psi^l = \text{span } \Phi^l = S^l$.

Beweisidee:



Die Beweisidee zeigt auch wie man die Darstellung einer Funktion f bezüglich der Basis Φ^l in eine bezüglich der Basis Ψ^l umrechnet.

Umgekehrt bedeutet $f(x) = \sum_{i=0}^l c_i \psi_i^l$ die Umrechnung in die Basis Φ^l (diesen Punkt kann man noch ausführlicher machen)

h) Effiziente Auswertung.

Hat man $f(x) = \sum_{j=0}^l \sum_{i=0}^{2^{j-1}} c_i^j \psi_i^j$ zu berechnen,

so spielen für ein $x \in (0, 1]$ wegen der lokalen Träger nur wenige x eine Rolle.

Algorithmus:

```

f = c_0^0 \cdot \psi_0^0(x);
j = 1, i = 0;
while (true) {
    f = f + c_i^j \cdot \psi_i^j(x);
    if (j == l) break;
    i = { 2i   x \in tr(\psi_{2i}^{j+1})
        { 2i+1 x \in tr(\psi_{2i+1}^{j+1})
    j = j + 1;
}
    
```

Datenkompression mit Wavelets

Anwendung von (6.18). Geg: Darst von $f \in S^l$ in Basis Ψ^l .

Gesucht: Teilraum $\tilde{S} \subset S^l$ so dass $\|f - \tilde{f}\| \leq \text{TOL} \cdot (f, f)$ und \tilde{S} möglichst klein.

Algorithmus:

```

I = { (0, 0) };
I \subseteq { (i, l) | l \in \mathbb{N}_0, 0 \leq i < 2^{l-1} }
S = (f, \psi_0^0)^2
while (e < (1 - \text{TOL}) \cdot (f, f)) {
    Wähle Kind (k, l+1) eines (i, l) \in I mit (f, \psi_k^{l+1}) maximal.
    I = I \cup { (k, l+1) };
    S = S + (f, \psi_k^{l+1})^2;
}
    
```

Charakterisiert eine beliebige Menge von Wavelet Basen.

Ausblick

21
19.01.10

- Eine leichte Änderung von (6.18) :

Finde $u \in S$ so dass

$$\int_{\Omega} \nabla u \cdot \nabla \varphi \, dx = \int_a^b f \varphi \, dx \quad \forall \varphi \in S$$

führt auf die (numerische) Lösung der partiellen Differentialgleichung

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = f \quad \text{in } \Omega,$$

den Physiker als „Laplace-Gleichung“ bekannt.

Mehr davon in Numerik 2

- Die Welt ist nicht eindimensional.

Einen kurzen Ausflug ins Mehrdimensionale bei der Integration.

7 Numerische Integration

18.01.10

auch: „numerische Quadratur“.

Wir behandeln die numerische Berechnung bestimmter Integrale in einer Raumdimension:

$$I_{(a,b)} = \int_a^b f(x) dx.$$

Alle von uns behandelten Verfahren führen auf folgende Form

$$I(f) = \sum_{i=0}^n w_i f(x_i) + \text{Fehler.}$$

Hierbei sind

$w_i \in \mathbb{R}$ die Gewichte und

$x_i \in \mathbb{R}$ die Stützstellen.

7.1 Newton-Cotes Formeln

sind sog. interpolatorische Quadraturformeln.

Idee: Stelle Interpolationspolynom p zu gewissen Stützstellen auf und berechne das Integral über p exakt.

Formel: Stützstellen $(x_i, f(x_i))$ $i=0, \dots, n$

Interpolationspolynom in Lagrange-Darstellung

$$p_n(x) = \sum_{i=0}^n f(x_i) L_i^{(n)}(x) \quad L_i^{(n)}(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(x-x_j)}{(x_i-x_j)}$$

$$\text{also } I(f) \approx I^{(n)}(f) = \int_a^b p_n(x) dx = \sum_{i=0}^n f(x_i) \int_a^b L_i^{(n)}(x) dx \quad (7.1)$$

Bevor wir explizite Formeln für die w_i angeben noch eine

Definition 7.1 (Ordnung einer Quadraturformel)

18.01.10


Eine Quadraturformel $I^{(m)}(f)$ hat mindestens die Ordnung m , wenn sie Polynome vom Grad $m-1$ exakt integriert. \square

- Hier: $n+1$ Stützstellen \Rightarrow Polynom vom Grad n exakt \Rightarrow Ordnung mind. n .
- Später: Bei geschickter Wahl der Stützstellen max. Ordnung $2n+1$ bei $n+1$ Stützst.
- Für $f \equiv 1$ gilt $p \equiv 1$ und damit $\int_a^b 1 dx = b-a = \sum_{i=0}^n w_i$.

Newton-Cotes-Formeln nutzen äquidistante Stützstellen:

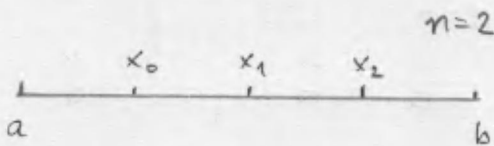
Variante a) Abgeschlossene Formeln: $a, b \in$ Stützstellen

$n=4$


$$x_i = a + iH \quad i=0, \dots, n, \quad H = \frac{b-a}{n}$$

Variante b) Offene Formeln: $a, b \notin$ Stützstellen

$n=2$


$$x_i = a + (i+1)H \quad i=0, \dots, n, \quad H = \frac{b-a}{n+2}$$

Berechnung der Gewichte (abgeschlossene Form):

$$I^{(n)}(f) = \sum_{i=0}^n f(x_i) \int_a^b L_i^{(n)}(x) dx = (b-a) \sum_{i=0}^n \underbrace{\left(\frac{1}{b-a} \int_a^b L_i^{(n)}(x) dx \right)}_{=: w_i \text{ unabhängig von } a, b!} f(x_i)$$

Mittels Substitution $x = g(s) = a + sH \Rightarrow s = g^{-1}(x) = \frac{x-a}{H}$, $g'(s) = H$ ergibt sich:

$$\begin{aligned} w_i &= \frac{1}{b-a} \int_a^b L_i^{(n)}(x) dx = \frac{1}{b-a} \int_{g^{-1}(a)}^{g^{-1}(b)} L_i^{(n)}(a+sH) \underbrace{g'(s)}_{=H} ds \\ &= \frac{1}{b-a} \frac{b-a}{H} \int_0^1 \prod_{\substack{j=0 \\ j \neq i}}^n \frac{[a+sH - (a+jH)]}{[a+iH - (a+jH)]} ds \\ &= \frac{1}{n} \int_0^1 \prod_{\substack{j=0 \\ j \neq i}}^n \frac{s-j}{i-j} ds. \end{aligned}$$

Berechnet man diese Gewichte so erhält man:

19.01.10

a) Abgeschlossene Formeln $n=1, 2, 3$, $H = \frac{b-a}{n}$

$$I^{(1)}(f) = \frac{b-a}{2} \{ f(a) + f(b) \} \quad \text{Trapez, Sehnens-Trapezregel}$$

$$I^{(2)}(f) = \frac{b-a}{6} \{ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \} \quad \text{Simpsonregel / Keplersche Fassregel}$$

$$I^{(3)}(f) = \frac{b-a}{8} \{ f(a) + 3f(a+H) + 3f(b-H) + f(b) \} \quad \frac{3}{8}\text{-Regel}$$

b) Offene Formeln $n=0, 1, 2$, $H = \frac{b-a}{n+2}$

$$I^{(0)}(f) = (b-a) f\left(\frac{a+b}{2}\right) \quad \text{Mittelpunkt, Tangenten-Trapez, Rechteckregel}$$

$$I^{(1)}(f) = \frac{b-a}{2} \{ f(a+H) + f(b-H) \}$$

$$I^{(2)}(f) = \frac{b-a}{3} \{ 2f(a+H) - f\left(\frac{a+b}{2}\right) + 2f(b-H) \}$$

Bemerkung 7.2

Ab $n=7$ für abgeschlossene und $n=2$ für offene Formeln treten negative Gewichte w_i auf. Dies ist ungünstig weil

- Für $f(x) \geq 0$ $\forall x$ garantieren positive w_i dass $I^{(n)}(f) \geq 0$, sonst nicht.

- Erhöhte Gefahr der Auslöschung

- Konditionierung wird schlechter:

Sei $\tilde{f}(x_i) = f(x_i) + \Delta y_i$ mit $|\Delta y_i| \leq \varepsilon$ so gilt:

$$I^{(n)}(\tilde{f}) = \sum_{i=0}^n w_i (f(x_i) + \Delta y_i) = I^{(n)}(f) + \sum_{i=0}^n w_i \Delta y_i$$

also

$$|I^{(n)}(\tilde{f}) - I^{(n)}(f)| = \left| \sum_{i=0}^n w_i \Delta y_i \right| \leq \varepsilon \sum_{i=0}^n |w_i|$$

Sind alle w_i positiv so gilt $\sum_{i=0}^n |w_i| = \sum_{i=0}^n w_i = b-a$.

Somit kann Kondition schlechter werden. □

Satz 7.3 (Restglieder)

4
19.01.10

Den begangenen Fehler kann man folgendermaßen abschätzen:

(i) Trapezregel

$$I(f) - \frac{b-a}{2} \{f(a) + f(b)\} = -\frac{(b-a)^3}{12} f''(\xi), \quad f \in C^2[a, b].$$

(ii) Simpson-Regel

$$I(f) - \frac{b-a}{6} \left\{ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right\} = -\frac{(b-a)^5}{2880} f^{(4)}(\xi), \quad f \in C^4[a, b].$$

(iii) Mittelpunkregel:

$$I(f) - (b-a) f\left(\frac{a+b}{2}\right) = \frac{(b-a)^3}{24} f''(\xi), \quad f \in C^2[a, b].$$

Wobei $\xi \in [a, b]$.

Beweis:

I-Fall:

$$\int_a^b f(x) - p_n(x) dx \stackrel{!}{=} \int_a^b \frac{f^{(n+1)}(\eta(x))}{(n+1)!} \prod_{j=0}^n (x-x_j)$$

Speziell für die Trapezregel gilt ($n=1$)

$$I(f) - I^1(f) = \frac{1}{2} \int_a^b f'''(\eta(x)) \underbrace{(x-a)(x-b)}_{=: g(x) \leq 0} dx$$

Mittelwert-
satz der
Integralrech.

$$\rightarrow = \frac{1}{2} f'''(\xi) \int_a^b (x-a)(x-b) dx = -\frac{(b-a)^3}{12} f'''(\xi), \quad \xi \in [a, b].$$

Er fordert $g(x) \geq 0$
oder $g(x) \leq 0$

Die beiden anderen Fälle sind etwas schwieriger wg. Vorzeichenwechsel von g , siehe Kammacher. E

- Mittelpunkregel hat den halben Fehler der Trapezregel bei nur einer Auswertung von f (tuner).

- Restglieder haben immer die typische Form $c(b-a)^{m+1} f^{(m)}(\xi)$

Kopplung!

7.2 Summierte Quadraturformeln

19.01.10

Erhöhen des Polynomgrades ist wenig sinnvoll, da

- negative Gewichte auftreten
- Lagrange-Interpol. nicht punktweise konvergiert
- Entsprechende Differenzierbarkeit von f gegeben sein muss

Idee: - Unterteile $[a, b]$ in N Teilintervalle

$$[x_i, x_{i+1}] \quad x_i = a + ih, \quad i = 0, \dots, N-1, \quad h = \frac{b-a}{N}$$

- Wende eine der obigen Formeln in jedem Teilintervall an
- Das Ergebnis nennt man „summierte Quadraturformel“

Satz 7.4 (Restglied für summierte Quadraturen)

Für die je Teilintervall verwendete Quadraturformel gelte

$$\frac{I}{[x_i, x_{i+1}]}(f) - I_{[x_i, x_{i+1}]}^{(n)}(f) = \alpha_n h^{m+2} f^{(m+1)}(\xi_i), \quad \xi_i \in [x_i, x_{i+1}].$$

α_n unabh. von $[x_i, x_{i+1}]$

Dann gilt für die summierte Quadraturformel

$$I(f) - \underbrace{I_h^{(n)}}_{\text{summierte Formel}}(f) = \alpha_n (b-a) h^{m+1} f^{(m+1)}(\xi), \quad \xi \in [a, b].$$

Trapezregel: $n = m = 1$, also $O(h^2)$. Die Ordnung ist (mind.) 2 (dies motiviert die Def. der Ordnung).

Simpson: $n = 2, m = 3$, also $O(h^4)$. Die Ordnung ist mindestens 5 (Pol. v. Grad 2)

Beweis: Zwischenwertsatz aus der Analysis:

Sei $g(x)$ stetig auf $[a, b]$ dann \exists zu jedem $u \in [\min(g(a), g(b)), \max(g(a), g(b))]$ mind. ein $\eta \in [a, b]$ mit $g(\eta) = u$.

Aus dem Zwischenwertsatz folgt: Sei $\xi_i \in [a, b]$, $i = 0, \dots, N-1$ und g stetig auf $[a, b]$

Dann nimmt g jeden Wert zwischen $\min_i g(\xi_i)$ und $\max_i g(\xi_i)$ an.

Wegen $\min_i g(\xi_i) \leq \frac{1}{N} \sum_{i=0}^{N-1} g(\xi_i) \leq \max_i g(\xi_i)$ existiert $\xi \in [a, b]$ mit $\sum_{i=0}^{N-1} g(\xi_i) = Ng(\xi)$.

$$\text{Also: } I(f) - \underbrace{I_h^{(n)}}_{\text{Var.}}(f) = \sum_{i=0}^{N-1} \alpha_n h^{m+2} f^{(m+1)}(\xi_i) = \alpha_n h^{m+2} \sum_{i=0}^{N-1} f^{(m+1)}(\xi_i) = \alpha_n h^{m+2} N f^{(m+1)}(\xi)$$

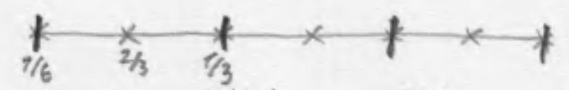
$$h = \frac{b-a}{N} \rightarrow = \alpha_n h^{m+2} \frac{b-a}{h} f^{(m+1)}(\xi) = \alpha_n (b-a) h^{m+1} f^{(m+1)}(\xi).$$

(i) Summierte Trapezregel

$$I_h^{(1)}(f) = \sum_{i=0}^{N-1} \frac{\overbrace{x_{i+1}-x_i}^{=h}}{2} \{f(x_i) + f(x_{i+1})\} = h \left\{ \frac{f(a)}{2} + \sum_{i=1}^{N-1} f(x_i) + \frac{f(b)}{2} \right\}$$

$$I(f) - I_h^{(1)}(f) = -\frac{b-a}{12} h^2 f''(\xi), \quad \xi \in [a, b].$$

(ii) Summierte Simpson Regel



$$I_h^{(2)}(f) = \sum_{i=0}^{N-1} \frac{\overbrace{x_{i+1}-x_i}^{=h}}{6} \left\{ f(x_i) + 4f\left(\frac{x_i+x_{i+1}}{2}\right) + f(x_{i+1}) \right\} = h \left\{ \frac{f(a)}{6} + \sum_{i=1}^{N-1} \frac{f(x_i)}{3} + \frac{2}{3} \sum_{i=0}^{N-1} f\left(\frac{x_i+x_{i+1}}{2}\right) + \frac{f(b)}{6} \right\}$$

$$I(f) - I_h^{(2)}(f) = -\frac{b-a}{2880} h^4 f^{(4)}(\xi), \quad \xi \in [a, b].$$

(iii) Summierte Mittelpunktregel

$$I_h^{(0)}(f) = \sum_{i=0}^{N-1} (x_{i+1}-x_i) f\left(\frac{x_i+x_{i+1}}{2}\right) = h \sum_{i=0}^{N-1} f\left(\frac{x_i+x_{i+1}}{2}\right)$$

$$I(f) - I_h^{(0)}(f) = \frac{b-a}{24} h^2 f''(\xi), \quad \xi \in [a, b].$$

Übung: Man zeige folgende Formeln:

$$\underbrace{I_h^{(2)}(f)}_{\text{Simpson } O(h^4)} = \frac{1}{3} \underbrace{I_h^{(1)}(f)}_{\text{Trapez je } O(h^2)} + \frac{2}{3} \underbrace{I_h^{(0)}(f)}_{\text{Mittelpunktregel}}$$

$$\underbrace{I_{h/2}^{(1)}}_{\text{halbe Gitterweite}} = \frac{1}{2} I_h^{(1)} + \frac{1}{2} I_h^{(0)}(f)$$

Diese Formeln können zur Fehlerkontrolle verwendet werden.

$$\sim |I(f) - I_h(f)| \leq \text{TOL.}$$

Beispiel 7.5 (Zur Quadratur)

Siehe Rechner.

7.3 Quadraturen höherer Ordnung

19.01.10

Wie gesehen ist mit Newton-Cotes maximal $n=7$ ($n=2$, offen) sinnvoll.

Wie erreicht man höhere Ordnung?

Romberg-Integration (= Extrapolation zum Limes mit Trapezregel)

Mit den summierten Formeln ist man an $\lim_{h \rightarrow 0} I_h^{(n)}(f)$ interessiert. Extrapolation zum Limes ist anwendbar.

Herzstück ist die Euler-Maclaurinsche Summenformel:

$$I(f) - \underbrace{I_h^{(n)}(f)}_{\text{Trapezsumme!}} = \sum_{k=1}^m h^{2k} \frac{B_{2k}}{(2k)!} (f^{(2k-1)}(b) - f^{(2k-1)}(a)) \\ + h^{2m+2} \frac{B_{2m+2}}{(2m+2)!} (b-a) f^{(2m+2)}(\xi), \quad \xi \in [a, b].$$

Die Trapezsumme hat also eine Fehlerentwicklung in gerade Potenzen. Damit ist die Extrapolation besonders effizient.

B_i : Bernoulli Zahlen.

Gauß-Integration

Frage: Lässt sich die Genauigkeit von Quadraturformeln durch nichtäquidistante Stützstellen verbessern?

Idee: Wähle w_i, x_i so dass Polynome von möglichst hohem Grad exakt integriert werden (Ordnungsmaximierung).

Satz 7.6 Die maximale Ordnung einer Quadraturformel mit $n+1$ Stützstellen ist $2n+2$ (d.h. Polynome vom Grad $2n+1$ werden exakt integriert).

Beweis: Ang. man könnte Polynome vom Grad $2n+2$ bei $n+1$ Stützst. exakt integrieren (Ordng $2n+3$). Betrachte

$$q(x) = \prod_{i=0}^n (x-x_i)^2$$

- $q(x)$ hat Grad $2n+2$

- $q(x) \geq 0 \forall x$ und $q(x) \neq 0$ also $\int_{-1}^1 q(x) dx > 0$

- Andererseits: $q(x_i) = 0$ und damit $\sum_{i=0}^n q(x_i) w_i = 0$, d.h. q wird nicht exakt integriert ∇ \blacksquare .

Die maximal mögliche Ordnung wird mit der Gauß-Quadratur erreicht:

19.10.10

Satz 7.7 (Gauß-Quadratur)

Es gibt genau eine interpolatorische Quadraturformel zu $n+1$ paarweise verschiedenen Stützstellen in $[-1, 1]$ mit der Ordnung $2n+2$.

Ihre Stützstellen sind die Nullstellen $\lambda_0, \dots, \lambda_n \in (-1, 1)$ des $(n+1)$ -ten Legendrepolynoms L_{n+1}

$$L_0(x) = 1, \quad L_1(x) = x, \quad L_{n+1}(x) = \frac{2n+1}{n+1} L_n(x) - \frac{n}{n+1} L_{n-1}(x).$$

Die Gewichte erhält man mittels

$$w_i = \int_{-1}^1 \prod_{\substack{j=0 \\ j \neq i}}^n \left(\frac{x - \lambda_j}{\lambda_i - \lambda_j} \right)^2 dx > 0 \quad i=0, \dots, n.$$

Bem: Die Legendrepolynome bilden Orthogonalsystem $\int_{-1}^1 L_n(x) L_m(x) dx = 0$ $n \neq m$.

Beispiele: $h = \frac{b-a}{2}$, $c = \frac{b+a}{2}$

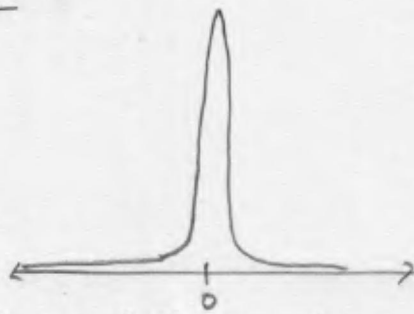
$n=1$: $I^{(1)}(f) = \frac{b-a}{2} \left\{ f\left(c - \sqrt{\frac{1}{3}}h\right) + f\left(c + \sqrt{\frac{1}{3}}h\right) \right\}$ Ordnung 4.

$n=2$: $I^{(2)}(f) = \frac{b-a}{18} \left\{ 5f\left(c - \sqrt{\frac{3}{5}}h\right) + 8f(c) + 5f\left(c + \sqrt{\frac{3}{5}}h\right) \right\}$ Ord. 6.

Entsprechend lassen sich summierte Formeln definieren.

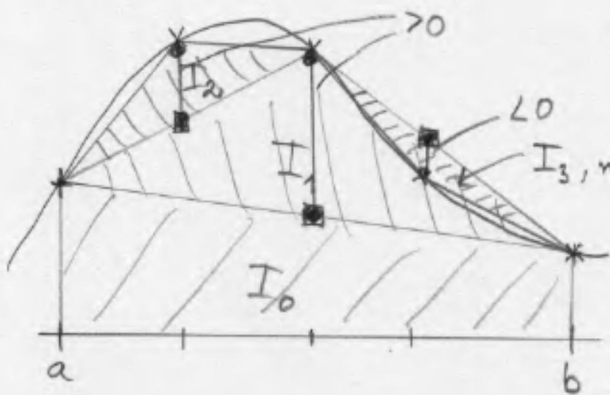
Adaptive Quadratur

Betr. $f(x) = \frac{1}{10^{-5} + x^2}$



Summierte Quadratur mit fester Schrittweite ineffizient.

Prinzip von Archimedes:



1) $I(f) = I_0 + I_1 + I_2 + \dots$

2) Breche rekursive Unterteilung ab falls $|I_j|$ klein genug.

Mehrdimensionale Quadratur

Die Welt ist nicht eindimensional.

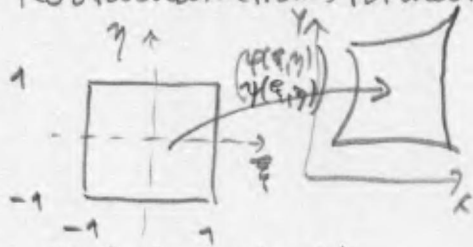
Für Rechtecke ($d=2$), Quader ($d=3$), ... lassen sich obige Formeln leicht erweitern:

$$\int_c^b \int_a^d f(x,y) dx dy \approx \int_c^b \sum_{i=0}^n f(x_i, y) w_i dy = \sum_{i=0}^n \int_c^b f(x_i, y) dy w_i$$

$$\approx \sum_{i=0}^n \sum_{j=0}^n f(x_i, y_j) \underbrace{w_i w_j}_{=w_{ij}}$$

Allerdings sind nicht alle Gebiete Rechtecke (anders als in 1D!)

Koordinatentransformation:



„Einheitsrechteck“

Abb $\begin{pmatrix} \varphi(\xi, \eta) \\ \psi(\xi, \eta) \end{pmatrix} : [-1, 1] \times [-1, 1] \rightarrow \Omega \subset \mathbb{R}^2$

$$\int_{\Omega} f(x,y) dx dy = \int_{-1}^1 \int_{-1}^1 f(\varphi(\xi, \eta), \psi(\xi, \eta)) \left| \frac{\partial(\varphi, \psi)}{\partial(\xi, \eta)} \right| d\xi d\eta$$

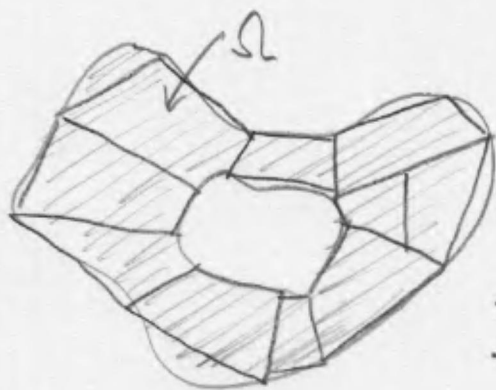
$$\left| \frac{\partial(\varphi, \psi)}{\partial(\xi, \eta)} \right| = \det \begin{bmatrix} \frac{\partial \varphi}{\partial \xi}(\xi, \eta) & \frac{\partial \psi}{\partial \xi}(\xi, \eta) \\ \frac{\partial \varphi}{\partial \eta}(\xi, \eta) & \frac{\partial \psi}{\partial \eta}(\xi, \eta) \end{bmatrix}$$

Transponierte Jacobimatrix.

Summierte Formeln in mehreren Raumdimensionen.

10
19.01.10

Bei komplexeren, z. B. Gebieten mit Löchern reicht das nicht



- Zerlegung in Teilgebiete die sich auf Rechtecke transformieren lassen.

- "Gittergenerierung" nicht trivial und schwierig automatisch zu machen
- Erfordert Beschreibung des Gebietes Ω .
- Zusätzlicher Geometriefehler durch nicht exakte Approximation der Geometrie

- Man kann auch Simplexes (= Dreieck, Tetraeder) zur Unterteilung verwenden

Fluch der Dimension

Ist d sehr groß so sind die hier behandelten Methoden nicht brauchbar.

Betrachte $\Omega = [0, 1]^d$. Zerlegt man $[0, 1]$ in zwei Teilintervalle je Richtung so hat man den d -dimensionalen Würfel in 2^d Teilwürfel zerlegt.

\Rightarrow Der Aufwand steigt exponentiell in d an. Dies bezeichnet man als "Fluch der Dimension".

Eine Möglichkeit ist dann die Monte-Carlo Integration

$$I(f) \approx \frac{C}{N} \sum_{i=1}^N f(\xi_i) \quad \text{mit Zufallszahl } \xi_i \in \Omega.$$

8 Iterative Lösung von Gleichungssystemen

25.01.10

In diesem Abschnitt betrachten wir die Lösung von algebraischen Gleichungen

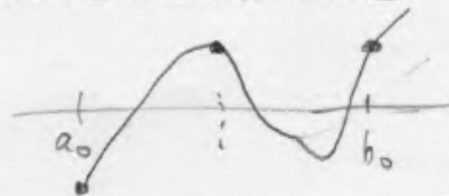
$$f(x) = 0 \quad \text{mit} \quad f: \mathbb{R}^n \rightarrow \mathbb{R}^n.$$

Dabei beschränken wir uns zunächst auf den Fall $n=1$ (Skalar).

Intervallshachtelung

Idee: Angenommen man kennt ein Teilintervall $[a_0, b_0]$ so dass $f(a_0)f(b_0) < 0$ (unterschiedliche Vorzeichen) und f sei stetig.

Dann hat f nach dem Zwischenwertsatz mind. eine Nullstelle in $[a_0, b_0]$.



Algorithmus:

Geg. $I_0 = [a_0, b_0]$ mit $f(a_0)f(b_0) < 0$ und $\epsilon > 0$

$t=0$;

while ($b_t - a_t > \epsilon$) {

$$x_t = (a_t + b_t) / 2;$$

if ($f(x_t) == 0$) {

$$a_t = x_t - \epsilon; \quad b_t = x_t;$$

} else if ($f(a_t)f(x_t) < 0$) {

$$a_{t+1} = a_t; \quad b_{t+1} = x_t;$$

// Nullstelle in $[a_t, x_t]$

} else {

$$a_{t+1} = x_t; \quad b_{t+1} = b_t;$$

// es ist $f(x_t)f(b_t) < 0$ da
VZ von $f(a_t) = \text{VZ von } f(x_t)$

} $t = t + 1$;

}

Analyse:

25.01.10

$$\text{Es gilt } a_t \leq a_{t+1} < b_{t+1} \leq b_t$$

$$\text{und } |b_{t+1} - a_{t+1}| \leq \frac{1}{2} |b_t - a_t| = \left(\frac{1}{2}\right)^{t+1} |b_0 - a_0|$$

(solange nicht $f(x_t) \equiv 0$).

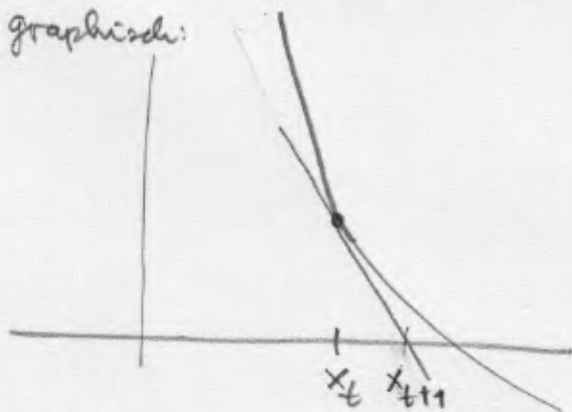
Bemerkung:

- Konvergenz ist linear mit Rate $\frac{1}{2}$.
- Sehr gut geeignet für monotone Funktionen
- nur für reelle Funktionen im \mathbb{R}^1 geeignet.

§. 1 Newton Verfahren

Die Funktion sei (mindestens) einmal stetig differenzierbar.

Idee: graphisch:



geg. x_t . Da $f \in C^1$ gibt es „Tangente“

$$T_t(x) = f'(x_t)(x - x_t) + f(x_t)$$

Nullstelle der Tangente:

$$T_t(x) = 0 \Leftrightarrow x = x_t - \frac{f(x_t)}{f'(x_t)}$$

Dies führt zur Iterationsvorschrift

$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)}$$

Offensichtlich ist $|f'(x_t)| > 0$ erforderlich, d. h. wir setzen voraus, dass die Nullstelle einfach ist.

Das Newton-Verfahren lässt sich auf Systeme

$f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ erweitern:

Es existiere die Taylorentwicklung von f :

$$f_i(x_t + \Delta x) = f_i(x_t) + \sum_{j=1}^n \frac{\partial f_i}{\partial x_j}(x_t) \Delta x_j + R_i(x_t, \Delta x) \quad i=1, \dots, n$$

in vektorieller Schreibweise

$$f(x_t + \Delta x) = f(x_t) + J(x_t) \Delta x + R(x_t, \Delta x)$$

$$(J(x_t))_{ij} = \frac{\partial f_i}{\partial x_j}(x_t) \quad \text{„Jacobimatrix“}$$

Ignorieren des Restgliedes entspricht „Linearisierung von f “.

$$f(x_t) + J(x_t) \Delta x \stackrel{!}{=} 0$$

$$\Leftrightarrow \Delta x = -J(x_t)^{-1} f(x_t)$$

führt zur Iteration

$$x_{t+1} = x_t - J(x_t)^{-1} f(x_t)$$

Jeder Schritt erfordert Lösung eines LGS mit der Jacobimatrix!

Nun untersuchen wir die Konvergenz des Newton-Verfahrens.
Allerdings nur im \mathbb{R}^1 .

Satz 8.1 (Newton-Verfahren)

4
25.01.10

Die Fkt $f \in C^2[a,b]$ habe in (a,b) (Inneren!) eine Nullstelle z und es sei

$$m := \min_{a \leq x \leq b} |f'(x)| > 0, \quad M := \max_{a \leq x \leq b} |f''(x)|.$$

Es sei $\varrho > 0$ so gewählt, dass

$$\varrho := \frac{M}{2m} \varrho < 1, \quad K_\varrho(z) := \{x \in \mathbb{R} : |x-z| \leq \varrho\} \subset [a,b].$$

Dann sind für jeden Startwert $x_0 \in K_\varrho(z)$ die Newton-Iterierten $x_t \in K_\varrho(z)$ definiert und konvergieren gegen die Nullstelle z .

Dabei gilt die a-priori Fehlerabschätzung

$$|x_t - z| \leq \frac{2m}{M} \varrho^{(2^t)}, \quad t \in \mathbb{N}$$

a-priori:
nur Abh. von den
Voraussetzungen

Und die a-posteriori Fehlerabschätzung

$$|x_t - z| \leq \frac{1}{m} |f(x_t)| \leq \frac{M}{2m} |x_t - x_{t-1}|^2, \quad t \in \mathbb{N}.$$

a-posteriori:
auch Abh. von bereits
berechnete Iterierten.

Beweis: Vorbereitungen

i) Mittelwert der Differentialrechnung liefert für alle $x, y \in [a,b], x \neq y$:

$$\left| \frac{f(x) - f(y)}{x - y} \right| = |f'(\xi)| \geq m \quad \Leftrightarrow \quad |x - y| \leq \frac{1}{m} |f(x) - f(y)|$$

\uparrow
Vor.

- f ist Lipschitz-stetig

- die Nullstelle z ist eindeutig, da sonst $0 < |z_1 - z_2| \leq \frac{1}{m} |f(z_1) - f(z_2)| = 0$ \downarrow

ii) Da $f \in C^2[a,b]$ gilt folg. Taylordarstellung:

$$f(y) = f(x) + (y-x)f'(x) + \underbrace{\int_x^y (x-t)f''(t) dt}_{=: R(y;x)} \quad \text{Restglied.}$$

Transformation des Integrals mit

$$\varphi(s) = x + s(y-x)$$

$$\varphi: [0, 1] \rightarrow [x, y]$$

liefert für das Restglied

$$\begin{aligned}
R(y; x) &= \int_x^y (x-t) f''(t) dt = \int_0^1 (x-\varphi(s)) f''(\varphi(s)) \varphi'(s) ds \\
&= \int_0^1 (\underbrace{x-x}_{=0} - s(y-x)) f''(x+s(y-x)) (y-x) ds \\
&= -(y-x)^2 \int_0^1 s f''(x+s(y-x)) ds.
\end{aligned}$$

Und damit

$$|R(y; x)| \leq (y-x)^2 \int_0^1 s \underbrace{|f''(x+s(y-x))|}_{\leq M \text{ n. Vor}} ds \leq \frac{M}{2} |y-x|^2.$$

iii) Nun setze $g(x) := x - \frac{f(x)}{f'(x)}$ (d.h. $x_{t+1} = g(x_t)$)

Dann gilt:

$$g(x) - z = x - \frac{f(x)}{f'(x)} - z = -\frac{1}{f'(x)} \underbrace{\left\{ f(x) + (z-x)f'(x) \right\}}_{= -R(z; x)}$$

Für $x \in K_g(z)$ gilt dann

$$\text{wg: } \underbrace{f(z)}_{=0} = \underbrace{f(x) + (z-x)f'(x)}_{\text{wg} < 1} + R(z; x)$$

$$|g(x) - z| = \left| \frac{1}{f'(x)} R(z; x) \right| \leq \frac{1}{m} \frac{M}{2} |z-x|^2 = \frac{M}{2m} \underbrace{|x-z|}_{\leq \rho} \underbrace{|x-z|}_{\leq \rho} < \rho$$

\uparrow min. macht ρ größer
 $\underbrace{< 1}_{\text{n. Wahl von } \rho} \quad \text{da } x \in K_g(z)$

Somit folgt aus $x \in K_g(z)$, dass auch $g(x) \in K_g(z)$.
 g bildet die Menge $K_g(z)$ auf sich selbst ab.

Die Newton-Iterationen sind $x_{t+1} = g(x_t)$.

(*) von oben

25.01.10

Setze $s_t := \frac{M}{2m} |x_t - z|$, Dann gilt mit der Absch. von oben

$$s_t = \frac{M}{2m} |x_t - z| = \frac{M}{2m} |g(x_{t-1}) - z| \stackrel{(*)}{\leq} \frac{M}{2m} \frac{M}{2m} |x_{t-1} - z|^2 = s_{t-1}^2$$

Somit gilt nach t Schritten:

$$s_t \leq s_{t-1}^2 \leq s_{t-2}^4 \leq \dots \leq s_{\underbrace{t-t}_{=0}}^{(2^t)} = s_0^{(2^t)}$$

und damit wegen $|x_t - z| = \frac{2m}{M} s_t$ und $s_0 = \frac{M}{2m} \underbrace{|x_0 - z|}_{\leq \rho} \leq \rho < 1$

$$|x_t - z| = \frac{2m}{M} s_t \leq \frac{2m}{M} s_0^{(2^t)} \leq \frac{2m}{M} \rho^{(2^t)}$$

was zu zeigen war.

A-posteriori Abschätzung folgt aus Taylor-Formel für x_t, x_{t-1}

$$f(x_t) = \underbrace{f(x_{t-1}) + (x_t - x_{t-1}) f'(x_t)}_{=0 \text{ nach Konstruktion!}} + R(x_t; x_{t-1})$$

und

$$|x_t - z| \leq \frac{1}{m} |f(x_t) - \underbrace{f(z)}_{=0}| = \frac{1}{m} |R(x_t; x_{t-1})| \leq \frac{M}{2m} |x_t - x_{t-1}|^2$$

Lipschitz
vom Anfang

□

Beispiel 8.2 (Wurzelberechnung mit Newton-Verfahren)

25.01.10

$$a > 0, n \geq 1, \text{ löse } x^n = a \Leftrightarrow f(x) = x^n - a = 0, f'(x) = n x^{n-1}$$

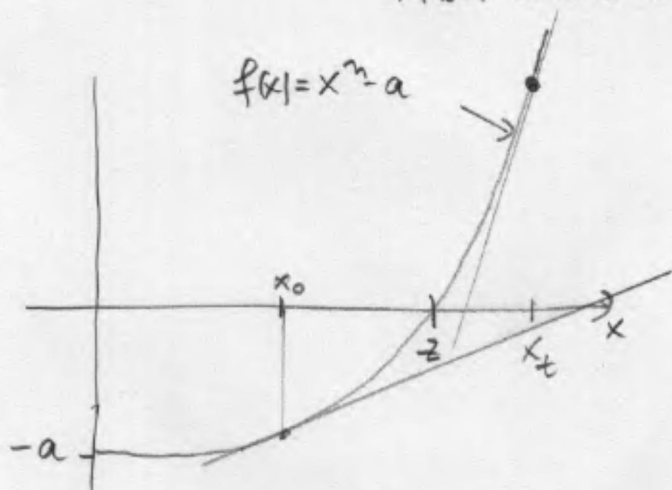
Also

$$\begin{aligned} x_{t+1} &= x_t - \frac{f(x_t)}{f'(x_t)} = x_t - \frac{x_t^n - a}{n x_t^{n-1}} = \frac{n x_t^n - x_t^n + a}{n x_t^{n-1}} = \frac{(n-1)x_t^n + a}{n x_t^{n-1}} \\ &= \frac{1}{n} \left\{ (n-1)x_t + \frac{a}{x_t^{n-1}} \right\} \end{aligned}$$

Satz 8.1 behauptet: Iteration konvergiert, falls x_0 nahe genug an z .

Hier gilt jedoch: Iteration konvergiert global, d.h. für alle $x_0 > 0$.

Aber nicht unbedingt quadratisch von Beginn an



1) für $x_t \geq z$ gilt

$$|x_{t+1} - z| < |x_t - z|$$

$$\text{da } f(x_t) > 0 \text{ und } f'(x_t) > \frac{f(x)}{x_t - z}$$

2) $0 < x_0 < z$

damit ist $x_1 > z$

$$\text{da } f(x_0) < 0 \text{ und } f'(x_0) < \frac{-f(x_0)}{z - x_0}$$

Man zeigt: für $n=2$ ist für $|x_0 - \sqrt{a}| \leq 2\sqrt{a}$ die Konvergenz quadratisch.

Bemerkungen zum Newton-Verfahren.

- Das Newton-Verfahren konvergiert nur lokal, d.h. wenn $|x_0 - z| \leq \delta$, \rightarrow "Einzugsbereich". wobei

-- δ i.d.R. unbekannt

-- δ möglicherweise sehr klein ist. oben: $\frac{M}{2m} \delta < 1 \Rightarrow \delta < \frac{2m}{M} \leftarrow \min f'$

- Newton-Verfahren konvergiert quadratisch.

$$|x_t - z| \leq c |x_{t-1} - z|^2. \text{ zum Vgl Intervallsch: } |x_t - z| \leq \frac{1}{2} |x_{t-1} - z|^2$$

- gedämpftes Newton-Verfahren: Verbesserung der Konvergenz außerhalb des Einzugsbereichs:

$$x_{t+1} = x_t - \lambda_t \frac{f(x_t)}{f'(x_t)}, \lambda_t \in (0, 1]$$

Wahl von λ_t
"Dämpfungsstrategie"

- Mehrfache Nullstellen

Sei z ^{zunächst} zweifache Nullstelle, d.h. $f(z) = f'(z) = 0$ und $f''(z) \neq 0$.

Wegen

$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)} = x_t - \frac{f(x_t) - f(z)}{f'(x_t) - f'(z)} \stackrel{\text{Erweitern}}{=} x_t - \frac{\frac{f(x_t) - f(z)}{x_t - z}}{\frac{f'(x_t) - f'(z)}{x_t - z}} = x_t - \frac{f'(z)}{f''(\xi_t)}$$

und $f''(z) \neq 0$ bleibt die Iteration für $x_t \rightarrow z$ (und damit $\xi_t \rightarrow z$) wohldefiniert.

Man zeigt: Für p -fache Nullstelle zeigt

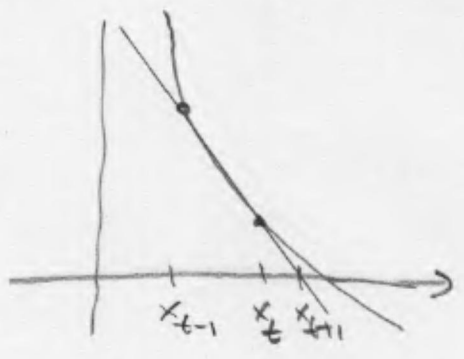
$$x_{t+1} = x_t - p \frac{f(x_t)}{f'(x_t)}$$

quadratische Konvergenz.

- Sekanten-Methode

Berechnung der Ableitung unter Umständen teuer.

Idee: Ersetze Tangente durch eine Sekante



$$s(x) = f(x_t) + (x - x_t) \frac{f(x_t) - f(x_{t-1})}{x_t - x_{t-1}}$$

Ansatz
 $\Rightarrow s(x) \stackrel{!}{=} 0$

führt auf Iteration

$$x_{t+1} = x_t - f(x_t) \frac{x_t - x_{t-1}}{f(x_t) - f(x_{t-1})}$$

Konvergenz: lokal mit

$$|x_t - z| \leq \frac{2m}{M} q^{t_0} , t \in \mathbb{N}$$

$\rho_0 = \rho_1 = 1$
 $\rho_{t+1} = \rho_t + \rho_{t-1}$ "Fibonacci Zahlen"

Nur eine f -Auswertung pro Iteration notwendig

$\rho_t \sim 0.723 \cdot (1.618)^t$: Konvergenzordnung 1,618 also zw. 1 u. 2

Problem: Auslöschung in $\frac{x_t - x_{t-1}}{f(x_t) - f(x_{t-1})}$

8.2 Sukzessive Approximation

25.01.10

Mit $g(x) = x - \frac{f(x)}{f'(x)}$ hat das Newton-Verfahren die Form

$$x_{t+1} = g(x_t).$$

Da die Nullstelle z von wg $f(z) = 0$ einen Fixpunkt der Iteration $x_{t+1} = g(x_t)$ ist nennt man das auch Fixpunktiteration.

Hier untersuchen wir nun allgemeine Iterationen dieser Art.

Z. B. könnte die Berechnung von $f'(x)$ sehr teuer sein und man wertet f' nur einmal „in der Nähe“ von z aus:

$$x_{t+1} = x - \frac{f(x)}{f'(z)}$$

Frage: Wam konvergiert so eine Iteration. Insbesondere wollen wir auch $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ zulassen.

Antwort gibt der sog. „Banachsche Fixpunktsatz“.

Satz 8.3 (Sukzessive Approximation)

Sei $G \subseteq \mathbb{R}^n$ eine nichtleere, abgeschlossene Punktmenge und $g: G \rightarrow G$ Lipschitz-stetig mit Konstante $q < 1$, d. h.

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

Hierbei ist $\|\cdot\|$ eine Vektornorm im \mathbb{R}^n und g nennt man eine „Kontraktion“. Dann existiert genau ein Fixpunkt $z \in G$ von g und für jeden Startpunkt $x^{(0)} \in G$ konvergiert die Folge der Iterierten $x^{(t+1)} = g(x^{(t)})$ gegen z .

Es gelten die a posteriori und a priori Fehlerabschätzungen

$$\|x^{(t)} - z\| \leq \frac{q}{1-q} \|x^{(t)} - x^{(t-1)}\| \leq \frac{q^t}{1-q} \|x^{(1)} - x^{(0)}\|.$$

(Wir schreiben den Iterationsindex oben in Klammer damit bei Vektoren unten Platz für den Komponentenindex bleibt)

Beweis:

10
25.01.10

i) Da $g: G \rightarrow G$ ist $x^{(t)} = g(x^{(t-1)}) = g(g(x^{(t-2)})) = \dots = \underbrace{g(\dots g(x^{(0)}))}_{t \text{ mal}}$
wohldefiniert.

ii) Weiter ist

$$\begin{aligned} \|x^{(t+1)} - x^{(t)}\| &= \|g(x^{(t)}) - g(x^{(t-1)})\| \leq q \|x^{(t)} - x^{(t-1)}\| \\ &\vdots \\ &\leq q^t \|x^{(1)} - x^{(0)}\| \end{aligned}$$

iii) Zeige nun, dass die $x^{(t)}$ eine Cauchy-Folge bilden.
Sei $m \geq 1$ und $\epsilon > 0$ gegeben. Es ist

$$\begin{aligned} \|x^{(t+m)} - x^{(t)}\| &\leq \|x^{(t+m)} - x^{(t+m-1)} + x^{(t+m-1)} - x^{(t+m-2)} + \dots + x^{(t+1)} - x^{(t)}\| \\ \text{Dreiecks-} &\rightarrow \leq \|x^{(t+m)} - x^{(t+m-1)}\| + \|x^{(t+m-1)} - x^{(t+m-2)}\| + \dots + \|x^{(t+1)} - x^{(t)}\| \\ \text{ungl.} & \quad \text{(ii)} \rightarrow \leq q^{t+m-1} \|x^{(1)} - x^{(0)}\| + q^{t+m-2} \|x^{(1)} - x^{(0)}\| + \dots + q^t \|x^{(1)} - x^{(0)}\| \\ \text{Ausklammern} & \rightarrow = (q^{t+m-1} + q^{t+m-2} + \dots + q^t) \|x^{(1)} - x^{(0)}\| \\ \text{geom. Reihe} & \rightarrow \leq q^t \frac{1-q^m}{1-q} \|x^{(1)} - x^{(0)}\| \leq \epsilon \text{ für } t \geq t(\epsilon) \end{aligned}$$

hinreichend groß.

\mathbb{R}^n ist vollständig, jede Cauchy-Folge konvergiert.
Also existiert $z = \lim_{t \rightarrow \infty} x^{(t)}$ und $z \in G$, da G abgeschlossen.

Schließlich ist $z = g(z)$. (Das zeigt man so:

$$\begin{aligned} \|z - g(z)\| &= \|z - x^{(t)} + x^{(t)} - g(z)\| \\ &\stackrel{t \text{ bel.}}{\leq} \|z - x^{(t)}\| + q \|x^{(t-1)} - z\| \rightarrow 0 \\ &\quad \rightarrow 0 \text{ für } t \rightarrow \infty \quad \rightarrow 0 \text{ für } t \rightarrow \infty \end{aligned}$$

iv) Fehlerabschätzung

$$\begin{aligned} \|x^{(t+m)} - x^{(t)}\| &\leq \|x^{(t+m)} - x^{(t+m-1)}\| + \dots + \|x^{(t+1)} - x^{(t)}\| \quad (\text{wie oben}) \\ &\stackrel{\downarrow m \text{ mal}}{\leq} q^m \|x^{(t)} - x^{(t-1)}\| + \dots + q \|x^{(t)} - x^{(t-1)}\| \\ &= (q^m + q^{m-1} + \dots + q) \|x^{(t)} - x^{(t-1)}\| \leq \frac{q}{1-q} \|x^{(t)} - x^{(t-1)}\| \\ &\quad \text{absh. durch geom. Reihe} \end{aligned}$$

Für $m \rightarrow \infty$ gilt $x^{(t+m)} \rightarrow z$, rechte Seite ist unabh. von m , also

$$\|z - x^{(t)}\| \leq \frac{q}{1-q} \|x^{(t)} - x^{(t-1)}\| \leq \frac{q}{1-q} q^{t-1} \|x^{(1)} - x^{(0)}\| = \frac{q^t}{1-q} \|x^{(1)} - x^{(0)}\|$$

Kann man benutzen um Abstand zur exakten Lösung zu schätzen



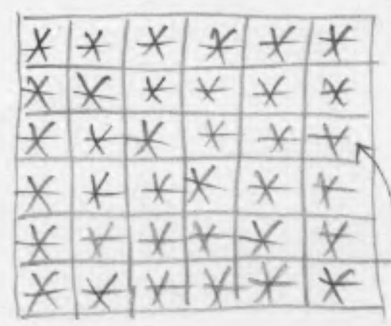
8.3 Iterationsverfahren zur Lösung linearer Gleichungssysteme

Wir kehren zurück zur Lösung von linearen Gleichungssystemen

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad b \in \mathbb{R}^n. \quad A \text{ sei regulär.}$$

Definition 8.4 Eine Menge von Matrizen $\{A^{(n)} \mid n \in \mathbb{N}\}$ heißt dünn besetzt falls $|\{a_{ij}^{(n)} \mid a_{ij}^{(n)} \neq 0\}| = \text{nnz}(A^{(n)}) = O(n)$. \square

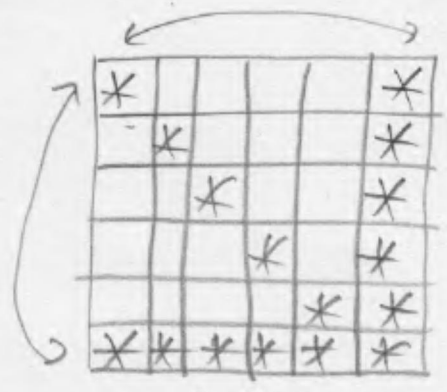
Gauß-Elimination ist für dünn besetzte Matrizen oft schlecht geeignet aufgrund von fill-in.



$$\text{nnz}(A_n) = 3n - 2$$

fill-in

fill-in
Minimierung
durch
Umordnen



no fill-in

Lösen von $Ax = b \iff$ „Nullstellensuche“ $f(x) = b - Ax = 0$.

Definiere Iteration

$$\begin{aligned} x^{(t+1)} &= g(x^{(t)}) = x^{(t)} + C^{-1} f(x^{(t)}) \\ &= x^{(t)} + C^{-1} (b - Ax^{(t)}) \\ &= \underbrace{(I - C^{-1}A)}_{=: B \text{ „Iterationsmatrix“}} x^{(t)} + C^{-1}b \end{aligned}$$

Für $x := A^{-1}b$ gilt

$$g(x) = (I - C^{-1}A) \underbrace{A^{-1}b}_x + C^{-1}b = A^{-1}b - C^{-1}b + C^{-1}b = A^{-1}b = x$$

Also x Fixpunkt von g .

Für die Lipschitzkonstante der Funktion g gilt

12
25.01.10

$$\begin{aligned}\|g(x) - g(y)\| &= \|Bx + C^{-1}b - By - C^{-1}b\| = \|B(x-y)\| \\ &\leq \|B\| \|x-y\|\end{aligned}$$

Falls $\|B\| < 1$ ($\|\cdot\|$ verträgliche Matrixnorm) ist g Kontraktion auf \mathbb{R}^n .

Beispiele für Iterationsverfahren.

Setze $A = L + D + U$ L strikte untere Dreiecksmatrix
 D Diagonalmatrix
 U obere Dreiecksmatrix

$C = D$, also $x^{(t+1)} = x^{(t)} + D^{-1}(b - Ax^{(t)})$ „Jacobi-Verfahren“

$C = L + D$, also $x^{(t+1)} = x^{(t)} + (L + D)^{-1}(b - Ax^{(t)})$ „Gauß-Seidel Verf.“

Iterationsverfahren konvergieren in der Regel nur für bestimmte Klassen von Matrizen. Hier ein Beispiel.

Definition 8.5 Eine Matrix heißt strikt diagonaldominant

falls

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

Beispiel: Splines, Richardson-Verfahren.

Satz 8.6 Das Jacobi-Verfahren konvergiert für strikt diagonaldominante Matrizen.

Beweis. $B = I - D^{-1}A$. Zeige $\|B\|_{\infty} < 1$, (Zeilensummenorm).

$$\|B\|_{\infty} = \|I - D^{-1}A\|_{\infty} = \|I - D^{-1}(L + D + U)\|_{\infty} = \| -D^{-1}(L + U) \|_{\infty}$$

$$\|B\|_{\infty} = \|D^{-1}(L + U)\|_{\infty} = \max_{i=1, \dots, n} \left(\sum_{j \neq i} \left| \frac{a_{ij}}{a_{ii}} \right| \right) = \max_{i=1, \dots, n} \frac{1}{|a_{ii}|} \underbrace{\sum_{j \neq i} |a_{ij}|}_{< |a_{ii}| \text{ n. Ver.}} < 1. \quad \square$$

Es gibt viele weitere solche Aussagen für sym. pos. def. Matrizen, schwach diagonaldom. Matrizen, M-Matrizen, ...

Aufwand für Iterationsverfahren

1) Aufwand für eine Iteration $x^{(t+1)} = x^{(t)} + C^{-1}(b - Ax^{(t)})$
 sei $\alpha(n)$. Typischerweise $\alpha(n) = O(n)$.

$$2) \|x^{(t)} - x\| \leq \|B\| \|x^{(t-1)} - x\|$$

$$\text{also } \|x^{(t)} - x\| \leq \|B\|^t \|x^{(0)} - x\|$$

$$\text{brauche } \|B\|^t \leq \varepsilon \Leftrightarrow t \underbrace{\log \|B\|}_{< 0} \leq \underbrace{\log \varepsilon}_{< 0} \Leftrightarrow t \geq \underbrace{\frac{\log \varepsilon}{\log \|B\|}}_{> 0}$$

$$\text{Gesamtaufwand: } t \cdot \alpha(n) = \frac{\log \varepsilon}{\log \|B\|} \alpha(n)$$

$\|B\|$ problemabhängig, je nach Verfahren auch von n abhängig.

Es gibt Verfahren, die relevante Probleme (z.B. Rdu-matrix)

im Gesamtaufwand $O(n)$ lösen können!