

## Gleitkomma-Arithmetik

Für  $e_{min}, e_{max} \in \mathbb{Z}$ ,  $e_{min} < e_{max}$  ist ein Gleitkommasystem wie folgt definiert:

$$\mathcal{F} = \mathcal{F}(\beta, t, e_{min}, e_{max}) \\ = \{\pm m\beta^{e-t} \mid m \in \mathbb{N}, \beta^{t-1} \leq m \leq \beta^t - 1 \vee m = 0, \\ e_{min} \leq e \leq e_{max}\}$$

$$x \in \mathcal{F} \setminus \{0\} \Rightarrow \beta^{e_{min}-1} \leq |x| \leq \beta^{e_{max}}(1 - \beta^{-1}).$$

## Normalisierte Darstellung

Für  $d_1 \neq 0$ ,  $0 < d_1 \leq \beta - 1$ :

$$x = \pm \beta^e \left( \frac{d_1}{\beta^1} + \frac{d_2}{\beta^2} + \dots + \frac{d_t}{\beta^t} \right) =: \pm 0.d_1 d_2 \dots d_t \cdot \beta^e$$

## Relative Maschinengenauigkeit

$fl(x) \in \mathcal{F}$  ist die  $x \in \mathbb{R}$  am nächsten liegende Gleitkommazahl.

Für relative Maschinengenauigkeit  $\epsilon := \frac{1}{2}\beta^{1-t}$ :

$$\frac{|fl(x) - x|}{|x|} < \epsilon, \quad \frac{|fl(x) - x|}{|fl(x)|} \leq \epsilon$$

## Arithmetische Grundoperationen

Für  $x, y \in \mathcal{F}$  sind Operationen  $o \in \{x, -, *, \div\}$  bzgl. eines Gleitkommasystems definiert:

$$\bar{o}(x, y) := fl(o(x, y))$$

Zu beachten ist hier die Ungültigkeit der Assoziativ- und Distributivgesetze.

## Kondition mathematischer Probleme

Ein mathematisches Problem  $(f, x)$  ist die Auswertung von  $f(x)$  an  $x \in E$  wobei  $f: E \subset X \rightarrow R \subset Y$ .  $X$  und  $Y$  sind normierte Räume,  $E$  Menge der Eingaben,  $R$  Menge der Resultate.

Ein Algorithmus für  $f$  ist Abbildung  $\tilde{f}: E \subset X \rightarrow Y$  s.d.  $\tilde{f}(x)$  in endlich vielen Schritten auswertbar ist und  $\tilde{f}(x) \approx f(x)$  gilt.

Die Konditionszahl eines mathematischen Problems  $(f, x)$  ist die kleinste Zahl  $\kappa_f(x) \geq 0$  mit:

$$\frac{\|f(x + \Delta x) - f(x)\|_Y}{\|f(x)\|_Y} \leq \kappa_f(x) \frac{\|\Delta x\|_X}{\|x\|_X} + o(\|\Delta x\|_X)$$

Für  $\|\Delta x\|_X \rightarrow 0$ .

$$\kappa_f(x) = \limsup_{\delta \rightarrow 0} \left\{ \frac{\|f(x + \Delta x) - f(x)\|_Y \|x\|_X}{\|f(x)\|_Y \|\Delta x\|_X} \right\}$$

Für  $\Delta x \in X$ ,  $x + \Delta x \in E$ ,  $\|\Delta x\|_X \leq \delta$ .

Ein Problem  $(f, x)$  ist gut konditioniert für *kleine* und schlecht konditioniert für *große* Konditionszahlen  $\kappa_f(x)$ .

Existiert  $\limsup$  nicht wird  $\kappa_f(x) = \infty$  gesetzt und das Problem als *schlecht gestellt* bezeichnet.

## Kondition stetig differenzierbarer Fkt.

Für  $f \in C^1(E, \mathbb{R}^m)$  in Umgebung  $E \subseteq \mathbb{R}^n$  von  $x$ :

$$\kappa_f(x) = \frac{\|f'(x)\|_\infty \cdot \|x\|_X}{\|f(x)\|_Y}$$

## Stabilität numerischer Algorithmen

Elementaroperationen eines Computers sind mit dem relativen Fehler  $\epsilon$  behaftet. Zusätzlich existiert ein Eingabefehler derselben Größenordnung. Es ist also in jedem Fall mit einem relativen Fehler  $\kappa_f(x)\epsilon$  zu rechnen.

## Vorwärtsanalyse

*Stabilitätsindikator der Vorwärtsanalyse* eines Algorithmus'  $\tilde{f}$  zur Lösung von  $(f, x)$  ist minimales  $\sigma = \sigma(x) \geq 0$  für  $\{x^\epsilon\}_{\epsilon > 0}$  mit  $\|x - x^\epsilon\|_X \leq \epsilon \|x\|_X$ :

$$\frac{\|\tilde{f}(x^\epsilon) - f(x^\epsilon)\|_Y}{\|f(x^\epsilon)\|_Y} \leq \sigma \kappa_f(x^\epsilon) \epsilon + o(\epsilon) \text{ für } \epsilon \rightarrow 0$$

Algorithmus  $\tilde{f}$  ist stabil im Sinne der Vorwärtsanalyse, wenn  $\sigma \leq \# \text{Elementaroperationen}$ .

## Rückwärtsanalyse

*Stabilitätsindikator der Rückwärtsanalyse* ist minimales  $\rho = \rho(x) \geq 0$  für  $\{x^\epsilon\}_{\epsilon > 0}$  s.d. für bel.  $\|x^\epsilon - x\|_X \leq \epsilon \|x\|_X$  Schar  $\{\hat{x}^\epsilon\}_{\epsilon > 0}$  ex. mit  $f(\hat{x}^\epsilon) = \tilde{f}(x^\epsilon)$ :

$$\frac{\|\hat{x}^\epsilon - x^\epsilon\|_X}{\|x^\epsilon\|_X} \leq \rho \epsilon + o(\epsilon) \text{ für } \epsilon \rightarrow 0$$

Vorwärtsstabilität folgt aus Rückwärtsstabilität.

## Vektor- und Matrixnormen

### Induzierte Matrixnorm / Operatornorm

Für Normen  $\|\cdot\|_o, \|\cdot\|_\star$  auf  $\mathbb{K}^n$  bzw.  $\mathbb{K}^m$  ist eine Matrixnorm  $\|\cdot\|: \mathbb{K}^{m \times n} \rightarrow [0, \infty)$  auf dem Vektorraum der  $m \times n$ -Matrizen definiert:

$$\|A\| := \max_{v \in \mathbb{K}^n \setminus \{0\}} \frac{\|Av\|_\star}{\|v\|_o} = \max_{\{v \in \mathbb{K}^n \mid \|v\|_o = 1\}} \|Av\|_\star$$

### Eigenschaften

Für  $A \in \mathbb{K}^{m \times n}$  gilt  $\forall v \in \mathbb{K}^n: \|Av\|_\star \leq \|A\| \cdot \|v\|_o$ .  
Submultiplikativität:  $\|AB\| \leq \|A\| \cdot \|B\|$

### Matrix-p-Normen

Induzierte Matrixnorm bei Wahl der  $p$ -Normen über  $\mathbb{K}^n$  bzw.  $\mathbb{K}^m$ :

$$\|A\|_p := \max_{v \in \mathbb{K}^n} \frac{\|Av\|_p}{\|v\|_p} = \max_{\|v\|_p = 1} \|Av\|_p \text{ für } 1 \leq p \leq \infty$$

## Spaltensummennorm

Für  $A = (a_1, \dots, a_n)$  mit  $a_j \in \mathbb{K}^m$ :

$$\|A\|_1 = \max_{1 \leq j \leq n} \|a_j\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{i,j}|$$

## Zeilensummennorm

$$\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{i,j}|$$

## Spektralnorm

Die Matrix-2-Norm wird so genannt, da  $\|A\|_2 = \sqrt{\lambda_{max}(A^H A)}$  für  $\lambda_{max}(A^H A)$  als Bezeichner des größten Eigenwerts von  $A^H A \in \mathbb{K}^{n \times n}$ .

$$\|A\|_2 = \|A^H\|_2, \|A^H A\|_2 = \|A\|_2^2 \\ \|QA\|_2 = \|A\|_2 \text{ für unitäre } Q.$$

## Induzierte Normen

Für  $A \in \mathbb{R}^{n \times n}$  mit  $A > 0$  ist  $\langle z, z \rangle_A := \langle Az, z \rangle_2$  ein Skalarprodukt auf  $\mathbb{R}^n$ . Dieses induziert die Energienorm  $\|z\|_A = \sqrt{\langle z, z \rangle_A}$ .

## Kondition einer Matrix

Für  $A \in \mathbb{K}^{n \times n} \in GL_n \mathbb{R}$ ,  $\|\cdot\|$  induzierte Matrixnorm:

$$\kappa(A) = \|A\| \cdot \|A^{-1}\|$$

$$1 = \|Id\| = \|AA^{-1}\| \leq \|A\| \cdot \|A^{-1}\| = \kappa(A)$$

## Besondere Matrizen

### Diagonaldominante Matrizen

$A \in \mathbb{R}^{n \times n}$  ist diagonaldominant, falls:

$$\forall i \in \{1, \dots, n\}: |a_{i,i}| > \sum_{j=1, j \neq i}^n |a_{i,j}|$$

Insbesondere sind solche  $A$  regulär. Gilt nur  $\geq$  so heißt die Matrix schwach diagonaldominant.

### Positiv definite Matrizen

$A \in \mathbb{R}^{n \times n}$  ist positiv definit d.h.  $A > 0$  falls  $A = A^T$  und  $\forall x \in \mathbb{R}^n \setminus \{0\}: x^T A x > 0$ .

### Hessenberg-Matrizen

Fast obere / untere Dreiecksmatrix wobei 1. untere / obere Nebendiagonale besetzt sein kann.

### Neumannsche-Reihe

$$(Id - M)^{-1} = \sum_{k=0}^{\infty} M^k$$

### Bezüglich $A > 0$ konjugierte Vektoren

Vektoren  $p, q \in \mathbb{R}^n$  sind konjugiert bzgl.  $A > 0$  d.h.  $A$ -orthogonal, falls  $Ap \perp q$ , also  $\langle Ap, q \rangle_2 = \langle p, q \rangle_A = 0$ .

## Direkte Verfahren zur LGS Lösung

### Cramersche Regel

Sei  $A = (a_{i,j})_{ij} \in GL_n(\mathbb{R})$ ,  $b \in \mathbb{R}^n$ ,  $A[j] = (a_1, \dots, a_{j-1}, b, a_{j+1}, \dots, a_n) \in \mathbb{R}^{n \times n}$ ,  $a_k$   $k$ -ter Spaltenvektor von  $A$ . Dann bildet  $x_j = \frac{\det(A[j])}{\det(A)}$  für  $j = 1, \dots, n$  die eindeutige Lösung  $x \in \mathbb{R}^n$  s.d.  $Ax = b$ . Aufgrund des hohen Aufwands von allg. mehr als  $(n+1)!$  arithmetischen Operationen ist die Cramersche Regel nur von theoretischer Bedeutung.

### Lösung gestaffelter Systeme

Obere Dreiecksmatrizen können mittels Rückwärtssubstitution, untere Dreiecksmatrizen mittels Vorwärtssubstitution in  $\mathcal{O}(n^2)$  gelöst werden.

### LR-Zerlegung

$A = LR$  wobei  $L$  untere Dreiecksmatrix mit 1-Diagonale und  $R$  obere Dreiecksmatrix.

### Berechnung LR-Zerlegung

Die LR-Zerlegung existiert insofern die Diagonaleinträge nicht verschwinden. Insbesondere gilt dies für diagonaldominante Matrizen.

- Spaltenweises Nullen der unteren Einträge mittels *Gauß*, Matrizen  $L_1, \dots, L_{n-1}$
- $L = L_1^{-1} \dots L_{n-1}^{-1}$ ,  $R = L_{n-1} \dots L_1 A$

### Lösung $Ax = b$ mittels LR-Zerlegung

- $A = LR$  berechnen
- $Lz = b$  Vorwärtssubstitution
- $Rx = z$  Rückwärtssubstitution

### Spaltenpivotsuche

Die normale LR-Zerlegung ist nur Vorwärts- und nicht Rückwärtsstabil. Dies kann durch Spaltenpivotsuche verbessert werden. Hier wird in jedem Schritt mittels einer Permutationsmatrix immer mit der größten verbleibenden Zeile eliminiert.

Für alle regulären Matrizen existiert eine Spaltenpivotsuchen LR-Zerlegung so, dass  $PA = LR$ .

### Cholesky-Zerlegung

Für symmetrische  $A > 0$  existiert untere Dreiecksmatrix  $L$  mit positiver Diagonale, so dass  $A = LL^T$ . Sym.  $A > 0$  können eindeutig als  $A = LDL^T$  geschrieben werden s.d.  $L$  untere Dreiecksmatrix mit 1-Diagonale,  $D$  positive Diagonalmatrix.  $D = D^{1/2} D^{1/2}$  und  $G = LD^{1/2}$  ergibt die äquivalente Zerlegung  $A = GG^T$ .



## Interpolation

Interpolation von  $f$  mit  $\varphi$  s.d.  $\varphi(t_i) = f(t_i)$  für  $i = 0, \dots, n$ . Approximation von  $f$  mit  $\varphi$  s.d.  $\|\varphi - f\|$  möglichst klein in geeigneter Norm.

### Klassische Polynom-Interpolation

Zu gegebenen Knoten  $t_0 < \dots < t_n$  und Stützwerten  $f_i = f(t_i)$  für  $i = 0, \dots, n$  wird Polynom  $\in \Pi_n$  gesucht s.d.  $P(t_i) = f_i$  für  $i = 0, \dots, n$ . Zu  $n + 1$  Stützwerten  $f_i$  und paarweise verschiedenen Knoten  $t_i$  existiert dabei genau ein solches Interpolationspolynom  $P = P(f|t_0, \dots, t_n) \in \Pi_n$ .

### Interpolationsfehler

$$\|f - P(f|t_0, \dots, t_n)\|_\infty \leq \sup_{\tau \in [a,b]} \frac{|f^{(n+1)}(\tau)|}{(n+1)!} \|\omega_{n+1}\|_\infty$$

$\omega_{n+1} \in \Pi_{n+1}$  ist das *Newton-Polynom* bzgl.  $t_0, \dots, t_n$  mit  $\omega_{n+1}(t) := \prod_{i=0}^n (t - t_i)$ .

### Vandermonde-Matrix

$$\begin{pmatrix} 1 & t_0 & t_0^2 & \dots & t_0^n \\ 1 & t_1 & t_1^2 & \dots & t_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_n & t_n^2 & \dots & t_n^n \end{pmatrix} \begin{pmatrix} a_0 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} f_0 \\ \vdots \\ f_n \end{pmatrix}$$

Die Lösung der Vandermonde-Matrix beschreibt  $P(f|t_0, \dots, t_n) \in \Pi_n$ , was jedoch zu aufwändig ist.

### Lagrange-Basis

Basis  $\{L_{n,0}, \dots, L_{n,n}\}$  von  $\Pi_n$  abhg.  $t_0 < \dots < t_n$  wg.

$$L_{n,k}(t_i) = \delta_{k,i} = \begin{cases} 1 & k = i \\ 0 & \text{sonst} \end{cases}$$

$$L_{n,k}(t) := \prod_{j=0, j \neq k}^n \frac{t - t_j}{t_k - t_j}$$

Es gilt also:  $P(f|t_0, \dots, t_n) = \sum_{k=0}^n f_k \cdot L_{n,k}$

Ein Lagrange Polynom zu Stützstelle  $t_k$  nimmt an dieser 1, an allen anderen Stützstellen 0 an.

### Lemma von Aitken

$$P = P(f|t_0, \dots, t_n)(t) =$$

$$\frac{(t_0 - t)P(f|t_1, \dots, t_n)(t) - (t_n - t)P(f|t_0, \dots, t_{n-1})(t)}{t_0 - t_n}$$

### Schema von Neville

Sei  $t \in \mathbb{R}$  fest,  $P_{i,k}(t) = P_i, k = P(f|t_{i-k}, \dots, t_i)(t)$  für  $i \geq k$ . Dann ist insb.  $P_{i,0} = f_i$  und  $P_{n,n} = P(f|t_0, \dots, t_n)(t)$  kann rekursiv mit dem *Schema von Neville* berechnet werden:

$$P_{i,k} = \frac{(t_{i-k} - t)P_{i,k-1} - (t_i - t)P_{i-1,k-1}}{t_{i-k} - t_i} = \frac{t - t_{i-k}}{t_i - t_{i-k}} P_{i,k-1} - \frac{t - t_i}{t_i - t_{i-k}} P_{i-1,k-1}$$

### Tschebyscheff-Knoten

Für  $i = 0, \dots, n$ :

$$t_i^{[a,b]} = \frac{b+a}{2} + \frac{b-a}{2} \cos\left(\frac{2i+1}{2n+2}\pi\right)$$

Diese Knotenfolge liegt dichter zu den Intervallgrenzen hin und ergibt eine bessere Interpolation als äquidistante Knoten.

### Satz von Faber

Zu jeder Folge von Knoten  $\{t_0^{(n)}, \dots, t_n^{(n)}\}_{n \in \mathbb{N}}$  in  $[a, b]$  gibt es ein  $f \in C([a, b])$  so, dass  $\{P(f|t_0^{(n)}, \dots, t_n^{(n)})\}_{n \in \mathbb{N}}$  für  $n \rightarrow \infty$  nicht glm. gegen  $f$  konvergiert.

## Splines

Nachteile der Polynom-Interpolation bei einer größeren Anzahl von Knoten:

1. Starke Oszillation des Polynoms
2. Konvergenz des Polynoms gegen die interpolierte Funktion ist nicht gewährleistet

Sei  $\Delta = \{t_0, \dots, t_{l+1}\}$  ein Gitter paarweise verschiedener Knoten  $a = t_0 < \dots < t_{l+1} = b$ .  $s \in \mathcal{C}^{k-2}(a, b)$  ist Spline der Ordnung  $k \in \mathbb{N}$  bzgl.  $\Delta$  wenn sie auf jedem Intervall  $[t_i, t_{i+1}]$  mit einem Polynom  $s_i \in \Pi_{k-1}$  übereinstimmt.

### Spline-Raum

$S_{k,\Delta}$  ist Raum aller Splines der Ordnung  $k$  bzgl.  $\Delta$ . Der Spline-Raum  $S_{k,\Delta}$  ist ein reeller Vektorraum mit  $\Pi_{k-1}[a, b] \subset S_{k,\Delta}$ .

Zusätzlich gilt auch  $(t - t_i)_+^{k-1} \in S_{k,\Delta}$

### Abgebrochene Potenzen

Abgebrochene Potenzen vom Grad  $k - 1$ :

$$(t - t_i)_+^{k-1} := \begin{cases} (t - t_i)^{k-1} & : t \geq t_i \\ 0 & : t < t_i \end{cases}$$

Für  $t_i \in \Delta$ ,  $i \neq l + 1$

### Basis des Spline-Raumes

$\mathcal{B} = \{1, t, \dots, t^{k-1}, (t - t_1)_+^{k-1}, \dots, (t - t_l)_+^{k-1}\}$  ist eine Basis von  $S_{k,\Delta}$  mit  $\dim(S_{k,\Delta}) = k + l$ .

## Spline-Interpolation

Interpolation einer Funktion  $f$  bzgl. eines Gitters  $\Delta = \{t_0, \dots, t_{l+1}\}$  durch Spline der Ordnung  $k$ .

Im linearen Fall mit  $k = 2$  stimmt die Anzahl der Knoten  $l + 2$  mit  $\dim(S_{2,\Delta}) = l + 2$  überein. Es gibt also genau einen Spline der  $(t_i, f(t_i))$  interpoliert.

$$I_2 f = \sum_{i=0}^{l+1} f(t_i) B_i$$

Für  $B_i \in S_{2,\Delta}$  mit  $B_i(t_k) = \delta_{i,k}$

### Kubische Splines

Kubische Splines der Ordnung 4 eignen sich für die Darstellung von Kurven, da das menschliche Auge diese als glatt empfindet.

Die Interpolationsbedingungen reichen zur eindeutigen Bestimmung eines interpolierenden Splines aus  $S_{4,\Delta}$  nicht aus.

Wegen  $\dim(S_{4,\Delta}) - (l + 2) = l + 4 - (l + 2) = 2$  bleiben zwei Freiheitsgrade unbestimmt.

Eine zusätzliche Bedingung ist, dass der interpolierende kubische Spline die minimale Krümmung aller interpolierenden  $\mathcal{C}^2$ -Funktionen besitzen soll.

### Krümmung einer Funktion

Krümmung von  $y : [a, b] \rightarrow \mathbb{R}, y \in \mathcal{C}^2$ :

$$\kappa(t) := \frac{y''(t)}{(1 + y'(t))^2}$$

$1/\kappa(t)$  ist der Radius des *Krümmungskreises*.

Das Krümmungsverhalten von  $y$  über ganz  $[a, b]$  ist durch ein Integral messbar:

$$\|y''\|_2 := \left( \int_a^b y''(t)^2 dt \right)^{1/2}$$

### Randbedingungen

Für  $s \in S_{4,\Delta}$ ,  $\Delta = \{t_0, \dots, t_{l+1}\}$  sind mögliche Randbedingungen:

- (a)  $s'(a) = f'(a)$  und  $s'(b) = f'(b)$ : *vollständige Spline-Interpolation*
- (b)  $s''(a) = s''(b) = 0$ : *natürliche Interpolation*
- (c)  $s'(a) = s'(b)$  und  $s''(a) = s''(b)$  falls  $f$   $b - a$  periodisch: *periodische Interpolation*

Ist eine dieser Randbedingungen erfüllt, so ist  $s$  eindeutig bestimmt. Ferner gilt für alle anderen interpolierenden  $y \in \mathcal{C}^2(a, b)$ :  $\|s''\|_2 < \|y''\|_2$

## Momente von Splines

Sei  $h_{j+1} := t_{j+1} - t_j$  Länge von  $[t_j, t_{j+1}]$ .  $M_j = s''(t_j)$  für  $j = 0, \dots, l + 1$  sind die Momente des Splines  $s \in S_{4,\Delta}$ . Aus den Momenten kann der Spline vollständig rekonstruiert werden. Da  $s_j := s|_{[t_j, t_{j+1}]}$  kubisch ist gilt für  $s_j'' = s''|_{[t_j, t_{j+1}]}$ :

$$s_j''(t) = M_j \frac{t_{j+1} - t}{h_{j+1}} + M_{j+1} \frac{t - t_j}{h_{j+1}}$$

$$s_j'(t) = -M_j \frac{(t_{j+1} - t)^2}{2h_{j+1}} + M_{j+1} \frac{(t - t_j)^2}{2h_{j+1}} + A_j$$

$$s_j(t) = M_j \frac{(t_{j+1} - t)^3}{6h_{j+1}} + M_{j+1} \frac{(t - t_j)^3}{6h_{j+1}} + A_j(t - t_j) + B_j$$

Die Integrationskonstanten  $A_j, B_j$  lassen sich aus den Interpolationsbedingungen berechnen.

## Ein paar Matlab Grundlagen

```
A = [ 1 0 ; 0 1 ] % = eye(2)
b = [ 3 4 ]' % = [ 3; 4 ]
b(1) = 4 % => b = [ 4 4 ]'
A(2,1) = 2 % => A = [ 1 0 ; 2 1 ]
c = 1:3 % = [ 1 2 3 ]
c = 1:2:6 % = [ 1 3 5 ]
A(2,:) % = [ 3 1 ]
A(:,1) % = [ 1 ; 3 ]
A.^2 % = [ 1 0 ; 9 1 ]
```

## Beispiel: LR-Zerlegung

```
function [A, L, R] = lr(A)
[w,h] = size(A);

if w ~= h
    error('A is not a square matrix')
end

for k = 1:w-1
    if abs(A(k,k)) < eps
        error('No LU-decomposition');
    end

    A(k+1:w,k) = A(k+1:w,k) / A(k,k);
    A(k+1:w,k+1:w) = A(k+1:w,k+1:w)
        - A(k+1:w,k)
            * A(k,k+1:w);
    end

    L = tril(A,-1) + eye(w);
    R = triu(A);
end
```