

Rechnerarithmetik

Zahlendarstellung

Maschinenzahlen Eine maschinendarstellbare Zahl zur Basis B ist ein Element der Menge:

$$M = \{x \in \mathbb{R} \mid x = \pm 0.m_1m_2m_3 \dots m_n \cdot B^{\pm e_1e_2 \dots e_l}\} \cup \{0\}$$

- Mit:
- $m_1 \neq 0$ (Normalisierungsbedingung)
 - $m_i, e_i \in \{0, 1, \dots, B - 1\}$ für $i \neq 0$
 - $B \in \mathbb{N}, B > 1$ (Basis)

Zahlenwert Der Wert $\hat{\omega}$ einer Maschinenzahl berechnet sich durch:

$$\hat{\omega} = \sum_{i=1}^n m_i B^{\hat{e}-i}, \quad \text{mit} \quad \hat{e} = \sum_{i=1}^l e_i B^{l-i}$$

Werteberechnung Berechnung einer vierstelligen Zahl zur Basis 4:

$$\underbrace{0.3211}_{n=4} \cdot \underbrace{4^{12}}_{l=2}$$

- Exponent: $\hat{e} = 1 \cdot 4^1 + 2 \cdot 4^0 = 6$
- Wert: $\hat{\omega} = 3 \cdot 4^5 + 2 \cdot 4^4 + 1 \cdot 4^3 + 1 \cdot 4^2 = 3664$

IEEE-754 Standard

Der IEEE-754 Standard definiert zwei wichtige Gleitpunktformate:

Single Precision (32 Bit)

- Vorzeichen (V): 1 Bit
- Exponent (E): 8 Bit (Bias 127)
- Mantisse (M): 23 Bit + 1 hidden bit

Double Precision (64 Bit)

- Vorzeichen (V): 1 Bit
- Exponent (E): 11 Bit (Bias 1023)
- Mantisse (M): 52 Bit + 1 hidden bit

Darstellungsbereich Für jedes Gleitpunktsystem existieren:

- Grösste darstellbare Zahl: $x_{\max} = (1 - B^{-n}) \cdot B^{e_{\max}}$
- Kleinste darstellbare positive Zahl: $x_{\min} = B^{e_{\min}-1}$

Approximations- und Rundungsfehler

Fehlerarten Sei \tilde{x} eine Näherung des exakten Wertes x :

Absoluter Fehler: **Relativer Fehler:**

$$|\tilde{x} - x| \qquad \left| \frac{\tilde{x} - x}{x} \right| \text{ bzw. } \frac{|\tilde{x} - x|}{|x|} \text{ für } x \neq 0$$

Maschinengenauigkeit eps ist die kleinste positive Zahl, für die gilt:

Allgemein:

Dezimal:

$$\text{eps} := \frac{B}{2} \cdot B^{-n} \qquad \text{eps}_{10} := 5 \cdot 10^{-n}$$

Sie begrenzt den maximalen relativen Rundungsfehler:

$$\left| \frac{rd(x) - x}{x} \right| \leq \text{eps}$$

Rundungseigenschaften Für alle $x \in \mathbb{R}$ mit $|x| \geq x_{\min}$ gilt:

Absoluter Fehler:

Relativer Fehler:

$$|rd(x) - x| \leq \frac{B}{2} \cdot B^{e-n-1} \qquad \left| \frac{rd(x) - x}{x} \right| \leq \text{eps}$$

Fehlerfortpflanzung

Konditionierung

Die Konditionszahl K beschreibt die relative Fehlervergrößerung bei Funktionsauswertungen:

$$K := \frac{|f'(x)| \cdot |x|}{|f(x)|}$$

Interpretation

- $K \leq 1$: gut konditioniert
- $K > 1$: schlecht konditioniert
- $K \gg 1$: sehr schlecht konditioniert

Fehlerfortpflanzung

Für eine differenzierbare Funktion f gilt näherungsweise:

Absoluter Fehler:

Relativer Fehler:

$$|f(\tilde{x}) - f(x)| \approx |f'(x)| \cdot |\tilde{x} - x| \qquad \frac{|f(\tilde{x}) - f(x)|}{|f(x)|} \approx K \cdot \frac{|\tilde{x} - x|}{|x|}$$

Fehleranalyse einer Funktion

So analysieren Sie die Fehlerfortpflanzung einer Funktion:

- Berechnen Sie $f'(x)$
- Bestimmen Sie die Konditionszahl K
- Schätzen Sie den absoluten Fehler ab
- Schätzen Sie den relativen Fehler ab
- Beurteilen Sie die Konditionierung anhand von K

$$\underbrace{|f(\tilde{x}) - f(x)|}_{\text{absoluter Fehler von } f(x)} \approx |f'(x)| \cdot \underbrace{|\tilde{x} - x|}_{\text{absoluter Fehler von } x}$$

$$\underbrace{\frac{|f(\tilde{x}) - f(x)|}{|f(x)|}}_{\text{relativer Fehler von } f(x)} \approx \underbrace{\frac{|f'(x)| \cdot |x|}{|f(x)|}}_{\text{Konditionszahl } K} \cdot \underbrace{\frac{|\tilde{x} - x|}{|x|}}_{\text{relativer Fehler von } x}$$

Fehleranalyse Beispiel: Fehleranalyse von $f(x) = \sin(x)$

- $f'(x) = \cos(x)$
- $K = \frac{|x \cos(x)|}{|\sin(x)|}$
- Für $x \rightarrow 0$: $K \rightarrow 1$ (gut konditioniert)
- Für $x \rightarrow \pi$: $K \rightarrow \infty$ (schlecht konditioniert)
- Der absolute Fehler wird nicht vergrößert, da $|\cos(x)| \leq 1$

Auslöschung Besonders problematisch: Auslöschung
Bei der Subtraktion fast gleich großer Zahlen können signifikante Stellen verloren gehen. Beispiel:

- $1.234567 - 1.234566 = 0.000001$
- Aus 7 signifikanten Stellen wird 1 signifikante Stelle

Praktische Fehlerquellen

Kritische Operationen

Die häufigsten Quellen für numerische Fehler sind:

- Auslöschung bei Subtraktion ähnlich großer Zahlen
- Überlauf (overflow) bei zu großen Zahlen
- Unterlauf (underflow) bei zu kleinen Zahlen
- Verlust signifikanter Stellen durch Rundung

Auslöschung bei der Berechnung von $\sqrt{x^2 + 1} - 1$:
Für kleine x führt die direkte Berechnung zu Auslöschung:

- Für $x = 10^{-8}$:
- $\sqrt{10^{-16} + 1} - 1 \approx 1.000000000 - 1 = 0$
- Korrekte Lösung durch Umformung:
- $\sqrt{x^2 + 1} - 1 = \frac{x^2}{\sqrt{x^2 + 1} + 1}$

Vermeidung von Auslöschung

So vermeiden Sie Auslöschungseffekte:

- Identifizieren Sie Subtraktionen ähnlich großer Zahlen
- Suchen Sie nach algebraischen Umformungen
- Prüfen Sie alternative Berechnungswege
- Verwenden Sie Taylorentwicklungen für kleine Werte

Analyse von Algorithmen

Fehlerakkumulation

Bei n aufeinanderfolgenden Operationen mit relativen Fehlern $\leq \varepsilon$ gilt für den Gesamtfehler:

- Best case: $\mathcal{O}(n\varepsilon)$ bei gleichverteilten Fehlern
- Worst case: $\mathcal{O}(2^n\varepsilon)$ bei systematischen Fehlern

Numerische Stabilität

Ein Algorithmus heißt numerisch stabil, wenn:

- Kleine Eingabefehler zu kleinen Ausgabefehlern führen
- Rundungsfehler sich nicht übermäßig akkumulieren
- Die Konditionszahl des Problems nicht künstlich verschlechtert wird

```
1 def fib(n):
2     if n <= 1:
3         return n
4     return fib(n-1) + fib(n-2)
```

1. Bestimmen Sie kritische Operationen
2. Schätzen Sie Rundungsfehler pro Operation ab
3. Analysieren Sie die Fehlerfortpflanzung
4. Berechnen Sie die worst-case Fehlerschranke
5. Vergleichen Sie alternative Implementierungen

```

1 def quadratic_stable(a, b, c):
2     #  $ax^2 + bx + c = 0$ 
3     if a == 0:
4         return [-c/b] if b != 0 else []
5
6     # Calculate discriminant
7     disc = b*b - 4*a*c
8     if disc < 0:
9         return []
10
11     # Choose numerically stable formula
12     if b >= 0:
13         q = -0.5*(b + sqrt(disc))
14     else:
15         q = -0.5*(b - sqrt(disc))
16
17     x1 = q/a
18     x2 = c/(q)
19     return sorted([x1, x2])

```

```
1 from decimal import Decimal, getcontext
2
3 # Set precision
4 getcontext().prec = 40
5
6 # Precise calculation
7 x = Decimal('1.0') / Decimal('7.0')
8 print(x)    # 0.1428571428571428571428571428571428571428
```

- F hat genau einen Fixpunkt \bar{x} in $[a, b]$
- Die Fixpunktiteration konvergiert gegen \bar{x} für alle $x_0 \in [a, b]$
- Fehlerabschätzungen:
 - a-priori: $|x_n - \bar{x}| \leq \frac{\alpha^n}{1-\alpha} \cdot |x_1 - x_0|$
 - a-posteriori: $|x_n - \bar{x}| \leq \frac{\alpha}{1-\alpha} \cdot |x_n - x_{n-1}|$

Benötigt zwei Startwerte x_0 und x_1 .

Konvergenzverhalten

Konvergenzordnung

Sei (x_n) eine gegen \bar{x} konvergierende Folge. Die Konvergenzordnung $q \geq 1$ ist definiert durch:

|x_{n+1} - \bar{x}| \leq c \cdot |x_n - \bar{x}|^q

wobei $c > 0$ eine Konstante ist. Für $q = 1$ muss zusätzlich $c < 1$ gelten.

Konvergenzordnungen der Verfahren

Die verschiedenen Verfahren zeigen folgende Konvergenzgeschwindigkeiten:

Newton-Verfahren:	Vereinfachtes Newton:	Sekantenverfahren:
Quadratische Konvergenz	Lineare Konvergenz $q = 1$	Superlineare Konvergenz $q = \frac{1+\sqrt{5}}{2} \approx 1.618$

Konvergenzgeschwindigkeit Vergleich der Verfahren Startwert $x_0 = 1$, Funktion $f(x) = x^2 - 2$, Ziel: $\sqrt{2}$

n	Newton	Vereinfacht	Sekanten
1	1.5000000	1.5000000	1.5000000
2	1.4166667	1.4500000	1.4545455
3	1.4142157	1.4250000	1.4142857
4	1.4142136	1.4125000	1.4142136

Fehlerabschätzung

Nullstellensatz von Bolzano

Sei $f : [a, b] \rightarrow \mathbb{R}$ stetig. Falls $f(a) \cdot f(b) < 0$, dann existiert mindestens eine Nullstelle $\xi \in (a, b)$.

Fehlerabschätzung für Nullstellen

So schätzen Sie den Fehler einer Näherungslösung ab:

- 1. Sei x_n der aktuelle Näherungswert
- 2. Wähle Toleranz $\epsilon > 0$
- 3. Prüfe Vorzeichenwechsel: $f(x_n - \epsilon) \cdot f(x_n + \epsilon) < 0$
- 4. Falls ja: Nullstelle liegt in $(x_n - \epsilon, x_n + \epsilon)$
- 5. Damit gilt: $|x_n - \xi| < \epsilon$

Praktische Fehlerabschätzung Fehlerbestimmung bei $f(x) = x^2 - 2$

- 1. Näherungswert: $x_3 = 1.4142157$
- 2. Mit $\epsilon = 10^{-5}$:
- 3. $f(x_3 - \epsilon) = 1.4142057^2 - 2 < 0$
- 4. $f(x_3 + \epsilon) = 1.4142257^2 - 2 > 0$
- 5. Also: $|x_3 - \sqrt{2}| < 10^{-5}$

Abbruchkriterien Praktische Implementierung

In der Praxis verwendet man meist mehrere Abbruchkriterien:

- Absolute Änderung: $|x_n - x_{n-1}| < \epsilon_1$
- Funktionswert: $|f(x_n)| < \epsilon_2$
- Maximale Iterationszahl: $n < n_{max}$
- Kombination dieser Kriterien

Numerische Lösung linearer Gleichungssysteme

Grundlagen

Lineares Gleichungssystem

Ein lineares Gleichungssystem der Form $Ax = b$ besteht aus:

A = [a11 ... a1n; ...; an1 ... ann] in R^{n x n}, x = (x1; ...; xn) in R^n, b = (b1; ...; bn)

Der Gauss-Algorithmus

Grundidee Gauss-Elimination

Transformation des Gleichungssystems $Ax = b$ in ein äquivalentes System $\tilde{A}x = \tilde{b}$, wobei \tilde{A} eine obere Dreiecksmatrix ist.

Erlaubte Operationen:

- $z_j := z_j - \lambda z_i$ für $i < j$ und $\lambda \in \mathbb{R}$
- $z_i \rightarrow z_j$ (Vertauschen von Zeilen)

Gauss-Algorithmus

So lösen Sie ein lineares Gleichungssystem mit dem Gauss-Algorithmus:

1. Elimination (Vorwärts)

- 1. Für $i = 1, \dots, n - 1$:
- 2. Für $j = i + 1, \dots, n$:
- 3. Berechne $\lambda_{ji} = a_{ji} / a_{ii}$
- 4. $z_j := z_j - \lambda_{ji} z_i$

2. Rückwärtseinsetzen

x_i = (b_i - sum_{j=i+1}^n a_ij x_j) / a_ii, i = n, n-1, ..., 1

Gauss-Elimination Lösung eines 3x3 Systems Gegebenes System:

A = [2 1 -1; 4 -1 3; -2 2 1], b = (3; 1; 4)

1. Elimination der ersten Spalte:

[2 1 -1 | 3; 0 -3 5 | -5; 0 3 -1 | 7]

2. Elimination der zweiten Spalte:

[2 1 -1 | 3; 0 -3 5 | -5; 0 0 4 | 2]

3. Rückwärtseinsetzen:

x3 = 2/4 = 1/2; x2 = (-5 - 5(1/2)) / -3 = -2; x1 = (3 - 1(-2) - (-1)(1/2)) / 2 = 1

Pivotisierung

Spaltenpivotisierung

Strategie zur numerischen Stabilisierung des Gauss-Algorithmus durch Auswahl des betragsmäßig größten Elements als Pivotelement. Vor jedem Eliminationsschritt in Spalte i:

- Suche k mit $|a_{ki}| = \max\{|a_{ji}| \mid j = i, \dots, n\}$
- Falls $a_{ki} \neq 0$: Vertausche Zeilen i und k
- Falls $a_{ki} = 0$: Matrix ist singulär

Gauss mit Pivotisierung

Erweiterter Gauss-Algorithmus mit Spaltenpivotisierung:

- 1. Für $i = 1, \dots, n - 1$:
- 2. Finde $k \geq i$ mit $|a_{ki}| = \max\{|a_{ji}| \mid j = i, \dots, n\}$
- 3. Falls $a_{ki} = 0$: Stop (Matrix singulär)
- 4. Vertausche Zeilen i und k
- 5. Für $j = i + 1, \dots, n$:
- 6. $z_j := z_j - \frac{a_{ji}}{a_{ii}} z_i$

Pivotisierung Gauss-Elimination mit Pivotisierung System:

$$\begin{pmatrix} 1 & 2 & -1 \\ 4 & 8 & -3 \\ 9 & 18 & -8 \end{pmatrix} x = \begin{pmatrix} 1 \\ 4 \\ 9 \end{pmatrix}$$

1. Erste Spalte: Pivot $|9| \rightarrow$ Tausche Zeilen 1 und 3
2. Nach Elimination der ersten Spalte:

$$\left(\begin{array}{ccc|c} 9 & 18 & -8 & 9 \\ 0 & 0 & 0.89 & 0 \\ 0 & 0 & 0.89 & 0 \end{array} \right)$$

3. System ist schlecht konditioniert (identische Zeilen)

Matrix-Zerlegungen

Dreieckszerlegung

Eine Matrix $A \in \mathbb{R}^{n \times n}$ kann zerlegt werden in:

Untere Dreiecksmatrix L: **Obere Dreiecksmatrix R:**

$l_{ij} = 0$ für $j > i$ $r_{ij} = 0$ für $i > j$

Diagonale meist normiert Diagonalelemente $\neq 0$

($l_{ii} = 1$)

LR-Zerlegung

Jede reguläre Matrix A , für die der Gauss-Algorithmus ohne Zeilenvertauschungen durchführbar ist, lässt sich zerlegen in:

$$A = LR$$

wobei L eine normierte untere und R eine obere Dreiecksmatrix ist.

Berechnung der LR-Zerlegung

So berechnen Sie die LR-Zerlegung:

1. Führen Sie Gauss-Elimination durch
2. R ist die resultierende obere Dreiecksmatrix
3. Die Eliminationsfaktoren $-\frac{a_{ji}}{a_{ii}}$ bilden L
4. Lösen Sie dann nacheinander:
 - $Ly = b$ (Vorwärtseinsetzen)
 - $Rx = y$ (Rückwärtseinsetzen)

LR-Zerlegung Berechnung einer LR-Zerlegung Matrix:

$$A = \begin{pmatrix} 2 & 1 & 1 \\ 4 & -1 & 0 \\ -2 & 3 & 1 \end{pmatrix}$$

Schrittweise Elimination führt zu:

$$L = \begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 2 & 1 \end{pmatrix}, \quad R = \begin{pmatrix} 2 & 1 & 1 \\ 0 & -3 & -2 \\ 0 & 0 & -2 \end{pmatrix}$$

