

# Einleitung

## Numerische Mathematik:

Zahlenmäßige Lösung eines mathem. Problems mit Hilfe des Computers

## Historische Bemerkungen:

### a) Mechanische Rechner:

Blaise Pascal (1623 - 1662)

Gottfried W. Leibniz (1646 - 1716)

### b) Elektronische Rechner:

Alan M. Turing (1912 - 1954)

Johann v. Neumann (1903 - 1957)

### c) Rasante Entwicklung 1950 - heute

- floating-point representation
- Fehleranalyse numerischer Verfahren  
Lanczos, Givens, Wilkinson, ...

## Problem (vereinfacht)

Eine vorgegebene Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$   
 ist an einer vorgegebenen Stelle  $x \in \mathbb{R}$   
 auszuwerten:  $y = f(x)$ .

$x$ : **Eingangsgröße, Eingabe**

$y$ : **Ergebnis, Resultat**

### Prinzipieller Ablauf:

- Eingabe in den Rechner:

$$x \mapsto \tilde{x} \quad \text{Maschinenzahl}$$

- Ersetzung von  $f$  durch einen **Algorithmus**  
 (= endliche Zahl elementarer Rechenoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit festgelegter Reihenfolge)

$$f \mapsto \tilde{f}$$

- **Numerisches Resultat:**  $\tilde{y} = \tilde{f}(\tilde{x})$
- Weitere Verfälschung: **Rundungsfehler**

## Aufgaben der Numerischen Mathem.:

- Konstruktion von Algorithmen
- Analyse von Algorithmen bezüglich
  - Genauigkeit
  - Rundungsfehler einfluss
  - Effizienz (Aufwand)
  - Zuverlässigkeit



# 1. Fehleranalyse

## 1.1 Maschinenzahlen M

Auf einem Computer sind nur endlich viele reelle Zahlen darstellbar, diese heißen **Maschinenzahlen**.

Darstellung:

$$\begin{aligned}x &= \pm (x_1 g^{-1} + x_2 g^{-2} + \dots + x_L g^{-L}) * g^e \\ &= \pm (0.x_1 x_2 \dots x_L)_g * g^e\end{aligned}\quad \underline{(1.1)}$$

$g$ : **Basis**,  $g \in \mathbb{N}$ ,  $g > 1$   
(zumeist  $g = 2, 10, 16$ )

$e$ : **Exponent**,  $e \in \mathbb{Z}$ ,  $e_{\min} \leq e \leq e_{\max}$

$x_i$ : **Ziffern**,  $x_i \in \{0, 1, \dots, g-1\}$

Die Ziffernfolge  $x_1 x_2 \dots x_L$  heißt **Mantisse**

Zumeist wird  $x_1 \neq 0$  verlangt (für  $x \neq 0$ )

**(Normalisierte Gleitpunkt-Darstellung)**

$L$ : **Mantissenlänge**



Beispiel : (Intel, IEC/IEEE - Norm)

einfach genau :  $g=2, L=24,$   
 $e_{\min} = -125, e_{\max} = 128$

doppelt genau :  $g=2, L=53,$   
 $e_{\min} = -1021, e_{\max} = 1024$

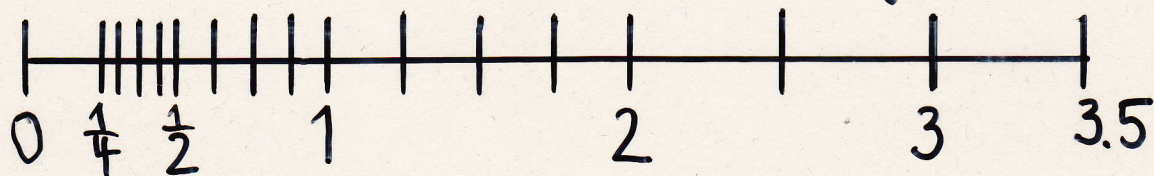
Folgerung :

a) Es gibt nur endlich viele Maschinenzahlen,  
 nämlich  $2(g-1)g^{L-1}(e_{\max} - e_{\min} + 1) + 1$

b) Es gibt keine beliebig großen und keine  
 beliebig kleinen Maschinenzahlen,  
 $x_{\max} = g^{e_{\max}}(1 - g^{-L})$  ist die größte,  
 $x_{\min} = g^{e_{\min} - 1}$  die kleinste pos. Masch.zahl

c)  $x, y \in \mathbb{M} \Rightarrow$  I. Allg.  $x \pm y, x * y, x/y \notin \mathbb{M}$

d) Die Maschinenzahlen sind nicht gleichverteilt:



( $g=2, L=3, e_{\min} = -1, e_{\max} = 2$ )



## 1.2 Absolute u. relative Fehler

Ist  $\tilde{x}$  Näherung für  $x \in \mathbb{R}$ , so heißt

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

$$\varepsilon_x := \frac{\tilde{x} - x}{x} \quad \text{relativer Fehler}$$

(für  $x \neq 0$ )

(1.2)

### Beispiel:

$$\begin{array}{l} \underline{a)} \\ \tilde{x} = 0.1238 \cdot 10^8 \end{array} \left. \begin{array}{l} x = 0.1237 \cdot 10^8 \\ \tilde{x} = 0.1238 \cdot 10^8 \end{array} \right\} \Rightarrow \begin{array}{l} \Delta x = 10^4 \\ \varepsilon_x = 0.8 \cdot 10^{-3} \end{array}$$

$$\begin{array}{l} \underline{b)} \\ \tilde{x} = 0.7921 \cdot 10^{-5} \end{array} \left. \begin{array}{l} x = 0.7321 \cdot 10^{-5} \\ \tilde{x} = 0.7921 \cdot 10^{-5} \end{array} \right\} \Rightarrow \begin{array}{l} \Delta x = 6 \cdot 10^{-7} \\ \varepsilon_x = 0.8 \cdot 10^{-1} \end{array}$$

### Merkregel: (1.3)

Gilt für den relativen Fehler  $\varepsilon_x = \alpha \cdot 10^{-k}$  mit  $0.1 \leq |\alpha| < 1$ ,  $k \in \mathbb{N}$ , so besitzt  $\tilde{x}$  etwa  $k$  korrekte Dezimalziffern.



## Erläuterung:

$$x = 0.x_1 \dots x_k x_{k+1} \dots \cdot 10^e, \quad x_1 \neq 0$$

$$\tilde{x} = 0.x_1 \dots x_k \tilde{x}_{k+1} \dots \cdot 10^e, \quad \tilde{x}_{k+1} \neq x_{k+1}$$

⇒

$$\begin{aligned} \varepsilon_x &= \frac{\tilde{x} - x}{x} \\ &= \frac{0.\tilde{x}_{k+1} \dots - 0.x_{k+1} \dots}{\underbrace{0.x_1 x_2 \dots}_{=: \alpha}} \cdot 10^{-k} \end{aligned}$$

wobei

$$0.1 \leq |\alpha| \leq 10.$$

## 1.3 Rundungsfehler

Neben den **Datenfehlern** ( $x \mapsto \tilde{x}$ ) und den **Verfahrensfehlern** ( $f \mapsto \tilde{f}$ ) spielen Rundungsfehler eine wichtige Rolle für die Beurteilung eines Algorithmus.

**Rundungsfehler:** Approximation einer reellen Zahl  $x$  durch eine Maschinenzahl  $\tilde{x} \in M$ . Dies tritt i. Allg. bei jeder element. Rechenoperation auf.

Gegeben:  $x = \pm (0.x_1x_2 \dots x_L x_{L+1} \dots)_g \cdot g^e$   
 mit  $x_1 \neq 0$ ,  $e_{\min} \leq e \leq e_{\max}$

Gesucht: Eine zu  $x$  benachbarte Maschinenzahl  $\tilde{x} \in \mathbb{M}$ .

### a) Abschneiden:

$$\tilde{x} := \text{rd}_1(x) := \pm (0.x_1 \dots x_L)_g \cdot g^e$$

$$\Rightarrow \tilde{x} = x(1 + \varepsilon), \quad |\varepsilon| \leq g^{1-L} \quad \underline{(1.4)}$$

denn:

$$\begin{aligned} \varepsilon_x &= \frac{\tilde{x} - x}{x} = - \frac{(0.0 \dots 0 x_{L+1} \dots)_g \cdot g^e}{(0.x_1 x_2 \dots)_g \cdot g^e} \\ &= - \underbrace{\frac{(0.x_{L+1} x_{L+2} \dots)_g}{(0.x_1 x_2 \dots)_g}}_{|\dots| \leq g} \cdot g^{-L} \quad \blacktriangle \end{aligned}$$

### b) Runden:

$$\tilde{x} := \text{rd}_2(x) :=$$

$$:= \begin{cases} (0.x_1 \dots x_L)_g \cdot g^e, & \text{falls } x_{L+1} < g/2 \\ (0.x_1 \dots \underbrace{x_L}_{\uparrow})_g \cdot g^e + g^{e-L}, & \text{sonst} \end{cases}$$



$$\Rightarrow \boxed{\tilde{x} = x(1 + \varepsilon), \quad |\varepsilon| \leq \frac{1}{2} g^{1-L}}$$

Zusatz: Falls  $x_j = g^{-1} (\forall j=1, \dots, L)$  und  $x_{L+1} \geq \frac{g}{2}$ , so setze  $rd_2(x) := (0.10\dots0) g^{e+1}$  für  $e < e_{max}$  und "overflow" für  $e = e_{max}$ .

Definition (1.5) Die Zahl  $eps := \frac{1}{2} g^{1-L}$  heißt relative Maschinengenauigkeit.

**Beispiel** (IEC/IEEE Norm)

- single precision :  $g=2, L=24, eps \doteq 0.6 \cdot 10^{-7}$
- double precision :  $g=2, L=53, eps \doteq 1.1 \cdot 10^{-16}$
- extended :  $g=2, L=64, eps \doteq 0.5 \cdot 10^{-19}$

Gleitpunkt-Rechnung:

Die Grundoperationen  $+, -, *, /$  werden in höherer Genauigkeit ausgeführt, normalisiert und auf  $L$  Stellen gerundet.

Bezeichnung :  $fl(a \pm b), fl(a * b), fl(a / b)$

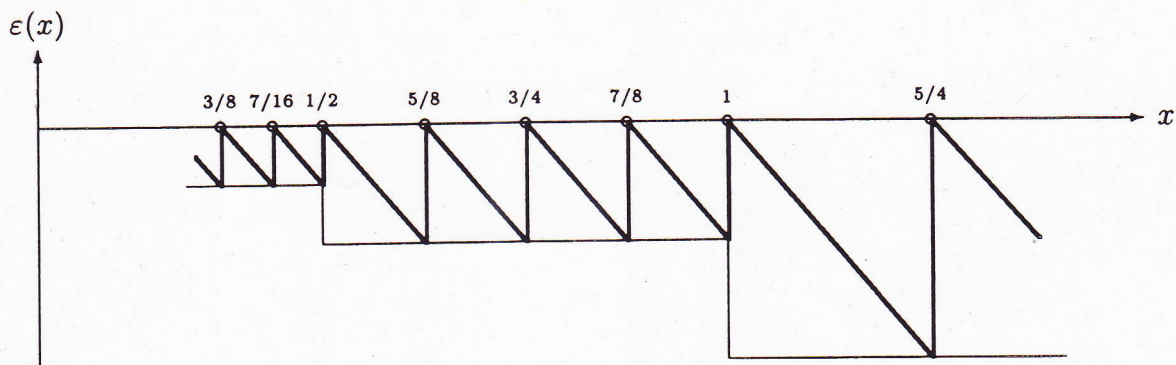


Abb. 4.6: Absoluter Rundungsfehler bei Rundung durch Abschneiden

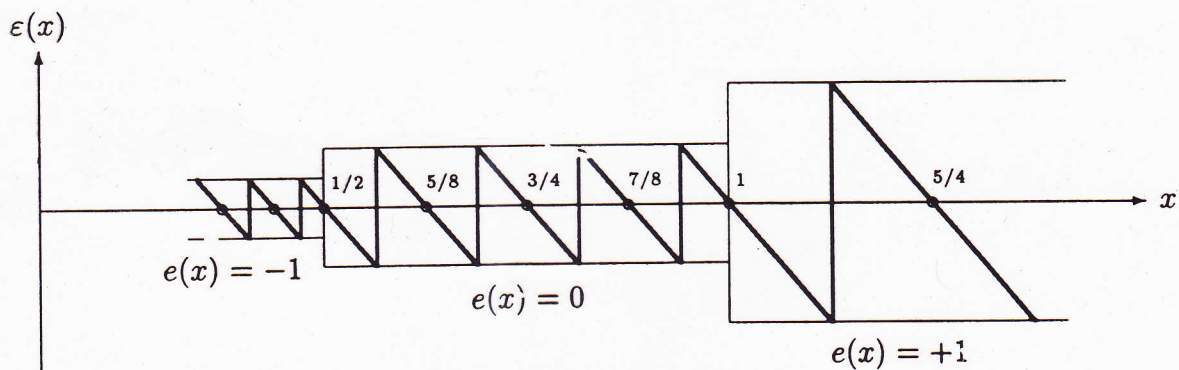


Abb. 4.7: Absoluter Rundungsfehler bei optimaler Rundung

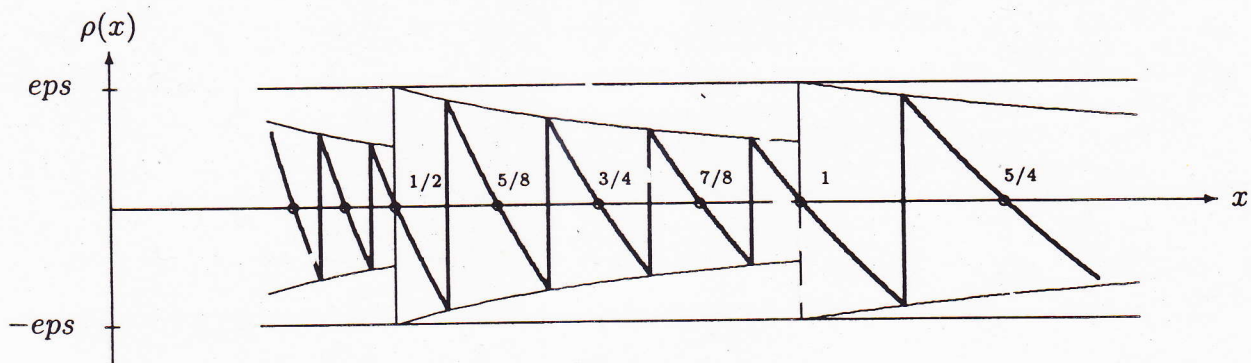


Abb. 4.8: Relativer Rundungsfehler bei optimaler Rundung



## Genauigkeit:

$$fl(a \pm b) = (a \pm b) (1 + \epsilon)$$

$$fl(a * b) = (a * b) (1 + \mu) \quad \underline{(1.6)}$$

$$fl(a/b) = (a/b) (1 + \delta)$$

mit  $|\epsilon|, |\mu|, |\delta| \leq \epsilon_{ps}$

## 1.4 Kondition

Ein mathem. Problem - z.B.  $y = f(x)$  - heißt **gut konditioniert**, falls kleine Änderungen der Eingangsdaten nur zu kleinen Änderungen im Resultat führt (bei exakter Rechnung!)

Andernfalls:

**Schlecht konditioniert ... chaotisch**

Quantitativ:  $(x, f(x) \neq 0)$

$$\Delta y = f(\tilde{x}) - f(x) \approx f'(x) \Delta x$$

$$\epsilon_y = \frac{\Delta y}{y} \approx \frac{f'(x)}{f(x)} \Delta x = \frac{x \cdot f'(x)}{f(x)} \cdot \epsilon_x$$

Definition (1.7) Die Größen

$$K_{\text{abs}} := f'(x), \quad K_{\text{rel}} := \frac{x f'(x)}{f(x)}$$

heißen absolute u. relative Konditionszahlen.

Bei mehreren Eingangsdaten  $y = f(x_1, \dots, x_n)$  addieren sich die Fehler, die durch die Verfälschung der einzelnen  $x_i$  hervorgerufen werden ( $\rightarrow$  Analysis II)

Man erhält:

$$\Delta y \approx \sum_{i=1}^n \frac{\partial f}{\partial x_i}(x) \Delta x_i \quad (\text{absoluter Fehler})$$

$$\varepsilon_y \approx \sum_{i=1}^n \frac{x_i}{f(x)} \frac{\partial f}{\partial x_i}(x) \cdot \varepsilon_{x_i} \quad (\text{rel. Fehler}) \quad (1.8)$$

Die  $\frac{\partial f}{\partial x_i}(x)$  heißen partielle Ableitungen von  $f$ .

Dabei wird jeweils nur  $x_i$  als Variable angesehen, die anderen  $x_j$  für  $j \neq i$  sind fest!



Die Zahlen  $K_{i,abs} := \frac{\partial f}{\partial x_i}(x)$  und

$K_{i,rel} := \frac{x_i}{f(x)} \frac{\partial f}{\partial x_i}(x)$ ,  $i=1, \dots, n$  heißen

**absolute u. relative Konditionzahlen**

bzgl. der Eingangsgröße  $x_i$ .

**Beispiel (1.9)** (Elementare Operationen)

$$\varepsilon_{x_1 \pm x_2} \approx \frac{x_1}{x_1 \pm x_2} \varepsilon_{x_1} \pm \frac{x_2}{x_1 \pm x_2} \varepsilon_{x_2}$$

$$\varepsilon_{x_1 \cdot x_2} \approx \varepsilon_{x_1} + \varepsilon_{x_2}$$

$$\varepsilon_{x_1/x_2} \approx \varepsilon_{x_1} - \varepsilon_{x_2}$$

⇒ Multiplikation u. Division sind gut konditionierte Operationen, ebenso die Addition von Zahlen mit gleichen Vorzeichen.

**Aber:** Subtraktion von annähernd gleichen Zahlen führt zu ev. erheblicher Fehler = Verstärkung. **Auslöschung!**

**(Cancellation!)**

Beispiel (1.10)

$$y = x_1 - x_2$$

exakt :

$$x_1 = 0.10005482410 * 10^5$$

$$x_2 = 0.09997342213 * 10^5$$

$$\begin{aligned}
 y &= 0.10005482410 * 10^5 \\
 &\quad - 0.09997342213 * 10^5 \\
 &\quad \hline
 &= 0.00008140197 * 10^5 \\
 &= \underline{0.8140197 * 10^1}
 \end{aligned}$$

Rechner mit  $g=10$  und  $L=5$  :

$$\tilde{x}_1 = 0.10005 * 10^5$$

$$\tilde{x}_2 = 0.99973 * 10^4$$

$$\begin{aligned}
 \tilde{y} &= 0.1000500 * 10^5 \\
 &\quad - 0.0999730 * 10^5 \\
 &\quad \hline
 &= 0.0000770 * 10^5 \\
 &= \underline{0.77000 * 10^1}
 \end{aligned}$$



## Beispiel (1.11) (Skalarprodukt)

$$x = (x_1, \dots, x_n)^T, \quad y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$$

$$\langle x, y \rangle := \sum_{i=1}^n x_i y_i = x_1 y_1 + \dots + x_n y_n$$

Nach Formel (1.8) erhält man:

$$\varepsilon_{\langle x, y \rangle} \approx \sum_{i=1}^n \frac{x_i}{\langle x, y \rangle} y_i \varepsilon_{x_i} + \sum_{i=1}^n \frac{y_i}{\langle x, y \rangle} x_i \varepsilon_{y_i}$$

Annahme:  $|\varepsilon_{x_i}|, |\varepsilon_{y_i}| \leq \text{eps}$ ,

Vorzeichen stimmt mit denen der Vorfaktoren überein

$$\Rightarrow |\varepsilon_{\langle x, y \rangle}| \leq \underbrace{2 \frac{\sum |x_i y_i|}{|\sum x_i y_i|}}_{\text{Konditionszahl}} \cdot \text{eps}$$

Ergebnis: Die Skalarprodukt-Berechnung kann schlecht konditioniert sein, wenn nämlich  $|\sum x_i y_i|$  klein, aber  $\sum |x_i y_i|$  groß ist. Der Grund ist wiederum:

Auslöschung!

Konkretes Beispiel : (Kulisch, Miranker, '86)

$$X := \begin{bmatrix} 2.718281828 \\ -3.141592654 \\ 1.414213562 \\ 0.5772156649 \\ 0.3010299957 \end{bmatrix}, \quad y := \begin{bmatrix} 1486.2497 \\ 878366.9879 \\ -22.37492 \\ 4773714.647 \\ 0.000185049 \end{bmatrix}$$

Numerische Ergebnisse:

FORTRAN single prec. :  $\langle x, y \rangle \doteq -0.49994\dots$

FORTRAN double prec. :  $\langle x, y \rangle \doteq 0.10251\dots \cdot 10^{-9}$

MATLAB :  $\langle x, y \rangle \doteq -0.52750\dots \cdot 10^{-10}$

Exaktes Ergebnis:

$$\underline{\underline{\langle x, y \rangle = -0.100657107 \cdot 10^{-10}}}$$



# 1.5 Fehleranalyse, Stabilität

In einfachen Fällen lässt sich der Einfluss von Rundungsfehler folgendermaßen untersuchen

- Simulation der Gleitpunkt-Rechnung
- Linearisierung der Fehler

**Beispiel (1.12)**  $y = a^2 - b^2 = a \cdot a - b \cdot b$

Gleitpunkt-Rechnung:

$$\begin{aligned}\tilde{y} &= \text{fl}(a^2 - b^2) \\ &= \left[ a^2 (1 + \varepsilon_a)^2 (1 + \mu_1) - b^2 (1 + \varepsilon_b)^2 (1 + \mu_2) \right] (1 + \sigma)\end{aligned}$$

$\varepsilon_a, \varepsilon_b$  : Datenfehler

$\mu_1, \mu_2$  : Multiplikationsfehler

$\sigma$  : Subtraktionsfehler

Linearisierung:

$$(1 + \alpha)(1 + \beta) \approx 1 + \alpha + \beta$$

$$(1 + \alpha)^k \approx 1 + k\alpha$$

$$\sqrt{1 + \alpha} \approx 1 + \alpha/2$$

⇒

$$\begin{aligned}\tilde{y} &\approx a^2(1+2\varepsilon_a+\mu_1+\sigma) - b^2(1+2\varepsilon_b+\mu_2+\sigma) \\ &= (a^2-b^2) \left\{ 1 + \frac{a^2}{a^2-b^2}(2\varepsilon_a+\mu_1) - \frac{b^2}{a^2-b^2}(2\varepsilon_b+\mu_2) + \right. \\ &\quad \left. + \sigma \right\} \\ &= y \cdot \{1 + \varepsilon_y\}\end{aligned}$$

⇒

$$\begin{aligned}\varepsilon_y &\approx \left( \frac{2a^2}{a^2-b^2} \varepsilon_a - \frac{2b^2}{a^2-b^2} \varepsilon_b \right) + \\ &\quad + \frac{a^2}{a^2-b^2} \mu_1 - \frac{b^2}{a^2-b^2} \mu_2 + \sigma\end{aligned}$$

(...)  $\hat{=}$  Kondition bzw. Eingangsfehler

Rest  $\hat{=}$  Rundungsfehler

### Definition (1.13)

Ein Algorithmus heißt stabil, wenn der Einfluss der Rundungsfehler nicht (wesentlich) größer ist als der Einfluss der Eingabefehler.

In obigem Beispiel :

$$|\varepsilon_{\text{Eingabe}}| \approx 2 \left| \frac{a^2+b^2}{a^2-b^2} \right| \cdot \text{eps}$$

$$|\varepsilon_{\text{Rundung}}| \approx \left( \left| \frac{a^2+b^2}{a^2-b^2} \right| + 1 \right) \cdot \text{eps}$$



⇒ Algorithmus ist stabil !

Für ein „vernünftiges“ numerisches Verfahren hat man zwei Anforderungen zu stellen :

- 1.) Das Problem sollte gut konditioniert sein.
- 2.) Der numerische Algorithmus sollte stabil sein.