# Numerische Mathematik für Elektrotechniker Zahlendarstellung, Rundungsfehler, Stabilität eines Algorithmus

Benjamin Berkels

Karl-Heinz Brakhage, Thomas Jankuhn, Christian Löbbert

Institut für Geometrie und Praktische Mathematik RWTH Aachen

Wintersemester 2019/2020

# Zusammenfassung letzte Vorlesung

#### Kondition

- ist unabhängig von einem speziellen Lösungweg (Algorithmus)
- gibt an, welche Genauigkeit man bestenfalls (bei exakter Rechnung) bei gestörten Eingangsdaten erwarten kann.

#### Frage:

- ▶ Was ist die (relative) Kondition eines Problems?
- ⇒ Die relative Kondition eines Problems bezeichnet das Verhältnis des relativen Ausgabefehlers zum relativen Eingabefehler, d.h. die Sensitivität des Problems unter Störungen der Eingabedaten.

### Zusammenfassung letzte Vorlesung

 $f:\mathbb{R}^n \to \mathbb{R}$  (Eingabe: Vektor, Ausgabe: Skalar)

$$\left| rac{f( ilde{x}) - f(x)}{f(x)} 
ight| \stackrel{.}{\leq} \kappa_{ ext{rel}}(x) \left| \sum_{j=1}^{n} \left| rac{ ilde{x}_{j} - x_{j}}{x_{j}} 
ight|$$

mit

$$\kappa_{ ext{rel}}(x) = \kappa_{ ext{rel}}^{\infty}(x) := \max_{j} \left| rac{\partial f(x)}{\partial x_{j}} \cdot rac{x_{j}}{f(x)} 
ight| \left( rac{\cancel{1}^{j} (\cancel{x})}{\cancel{1}^{j} (\cancel{x})} \stackrel{\raisebox{3pt}{$\chi$}}{\downarrow} 
ight)$$

 $f:\mathbb{R}^n o\mathbb{R}^n$  (Eingabe: Vektor, Ausgabe: Vektor) Für f(x)=Ax linear, invertierbar gilt

$$\frac{\|f(\tilde{x}) - f(x)\|}{\|f(x)\|} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{\kappa(A)} \frac{\|\tilde{x} - x\|}{\|x\|}$$

IGPM, RWTH Aachen

# Heute in der Vorlesung

#### Themen:

Dahmen & Reusken Kap 2.2-2.3

Zusammenfassung

- Zahlendarstellung und Rundungsfehler
- Gleitpunktarithmetik
- Stabilität eines Algorithmus

#### Was Sie mitnehmen sollten:

- Wie werden Zahlen im Computer dargestellt
- Wichtige Eigenschaften der Gleitpunktarithmetik
- Stabilität vs. Kondition

### Motivation

Rundungsfehler und Gleitpunktarithmetik

### Warum betrachten wir Gleitpunktdarstellung?

Aufgrund der Art und Weise, wie Zahlen im Computer dargestellt werden, können überraschende Ergebnisse auftreten

>> 
$$u = 0.3/0.1$$
;  
>>  $3 - u$   
ans =  $4.4409e-16$ 

#### Ein paar schlechte Beispiele:

- 1. D.N. Arnold, Some disasters attributable to bad numerical computing, 1998. https://www.ima.umn.edu/~arnold/disasters/
- 2. T. Huckle, Collection of Software Bugs, 2011. https://www5.in.tum.de/~huckle/bugse.html
- 3. K. Vuik, Some disasters caused by numerical errors. http://ta.twi.tudelft.nl/users/vuik/wi211/disasters.html

# Beispiel 2.31.

Wir betrachten als Beispiel die Zahl 123.75:

Dezimalsystem (Basis 10)

Rundungsfehler und Gleitpunktarithmetik

$$= 1 \cdot 10^{2} + 2 \cdot 10^{1} + 3 \cdot 10^{0} + 7 \cdot 10^{-1} + 5 \cdot 10^{-2}$$
$$= 10^{3} (1 \cdot 10^{-1} + 2 \cdot 10^{-2} + 3 \cdot 10^{-3} + 7 \cdot 10^{-4} + 5 \cdot 10^{-5})$$

Binärsystem (Basis 2)

#### 123.75

$$= 1 \cdot 2^{6} + 1 \cdot 2^{5} + 1 \cdot 2^{4} + 1 \cdot 2^{3} + 0 \cdot 2^{2} + 1 \cdot 2^{1} + 1 \cdot 2^{0}$$

$$+ 1 \cdot 2^{-1} + 1 \cdot 2^{-2}$$

$$= 2^{7} (1 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3} + 1 \cdot 2^{-4} + 0 \cdot 2^{-5} + 1 \cdot 2^{-6}$$

$$+ 1 \cdot 2^{-7} + 1 \cdot 2^{-8} + 1 \cdot 2^{-9})$$

# Zahlendarstellung

0000000000000000000

Rundungsfehler und Gleitpunktarithmetik

Seien  $b \in \mathbb{N}$ , b > 1, fest gewählt. Jedes  $x \in \mathbb{R}$ ,  $x \neq 0$ , lässt sich in der Form

$$x = \pm \left(\sum_{j=1}^{\infty} d_j \, b^{-j} 
ight) \cdot b^e$$

darstellen, mit  $d_i \in \{0, 1, \ldots, b-1\}, d_1 \neq 0$ , und e eine ganze Zahl.

# Zahlendarstellung

0000000000000000000

Rundungsfehler und Gleitpunktarithmetik

Seien  $b \in \mathbb{N}$ , b > 1, fest gewählt. Jedes  $x \in \mathbb{R}$ ,  $x \neq 0$ , lässt sich in der Form

$$x = \pm \left(\sum_{j=1}^{\infty} d_j b^{-j}\right) \cdot b^e$$

darstellen, mit  $d_i \in \{0, 1, \dots, b-1\}, d_1 \neq 0$ , und e eine ganze Zahl.

- $\triangleright$  Dezimalsystem (Basis b=10)  $123.75 \Rightarrow 0.12375 \cdot 10^3$
- Binärsystem (Basis b = 2)  $123.75 \Rightarrow 0.1111011111 \cdot 2^{111}$

# Normalisierte Gleitpunktdarstellung

### Gleitpunktdarstellung:

Rundungsfehler und Gleitpunktarithmetik

$$x = \pm 0.d_1d_2...d_m \times b^e$$
$$= \pm \left(\sum_{j=1}^m d_j b^{-j}\right) \times b^e$$

#### wobei

- ightharpoonup Basis  $b \in \mathbb{N} \setminus \{1\}$
- ightharpoonup Exponent  $e \in \mathbb{Z}$  mit r < e < R
- Mantisse  $f = \pm 0.d_1d_2...d_m, d_i \in \{0, 1, ..., b-1\}$
- Mantissenlänge m
- Normalisierung:  $d_1 \neq 0$  für  $x \neq 0$

$$M(10, 1, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

$$(11, 0, 0)$$

### Maschinenzahlen

00000000000000000000000

Rundungsfehler und Gleitpunktarithmetik

Nur endliche Anzahl von Zahlen darstellbar m, r, R vs.  $\infty$ :

$$x = \pm \left(\sum_{j=1}^{m} d_j b^{-j}\right) imes b^e, \quad r \leq e \leq R$$

 $\Rightarrow$  Maschinenzahlen  $\mathbb{M}(b, m, r, R)$ ,

Kardinalität von  $\mathbb{M}(b, m, r, R)$  ist endlich und hat

$$2 imes (b-1) imes b^{m-1} imes (R-r+1)$$
 Elemente.

### Definition (Rundung)

Reduktionsabbildung fl:  $\mathbb{R} \to \mathbb{M}(b, m, r, R)$  definiert durch

$$\mathrm{fl}(x) := \pm \begin{cases} \left(\sum_{j=1}^m d_j \, b^{-j}\right) \times b^e & \text{falls } d_{m+1} < \frac{b}{2}, \\ \left(\sum_{j=1}^m d_j \, b^{-j} + b^{-m}\right) \times b^e & \text{falls } d_{m+1} \geq \frac{b}{2}, \end{cases}$$

### Bildbereich

Rundungsfehler und Gleitpunktarithmetik

00000000000000000000000

 $\triangleright$  Betragsmäßig kleinste ( $\neq 0$ ) Zahl:

$$d_1 = 1, \ d_2 = \cdots = d_m = 0; \ e = r: \quad x_{\text{min}} = b^{r-1}$$

Betragsmäßig größte Zahl:

$$d_1 = \cdots = d_m = b-1; \ e = R: \quad x_{\text{MAX}} = (1 - b^{-m}) \times b^R$$

Erinnerung: Geomet. Reihe  $(b-1)\sum_{i=1}^m b^{-i} = 1 - b^{-m}$ 

ightharpoonup Bildbereich  $\mathbb{D} := [-x_{\text{MAY}}, -x_{\text{MIN}}] \cup [x_{\text{MIN}}, x_{\text{MAY}}]$ 

### **Achtung:** Die Endlichkeit von e beschränkt den **Bildbereich!**

- ▶ Unterlauf, wenn  $0 \neq |x| < |x_{\text{MIN}}|$ ;
- ightharpoonup Überlauf, wenn  $|x|>|x_{\text{MAX}}|$ .

# Bildbereich und Genauigkeit

Rundungsfehler und Gleitpunktarithmetik

### Maschinenzahlen $\mathbb{M}(b, m, r, R)$ :

Es gibt einen begrenzten Bereich von Zahlen, die dargestellt werden können

Die Endlichkeit von e beschränkt den Bildbereich

Es gibt nur eine endliche Anzahl von Zahlen, die innerhalb des Bildbereichs dargestellt werden können

Die Endlichkeit von f beschränkt die Genauigkeit.

# Maschinengenauigkeit – Beispiel

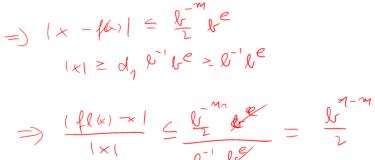
### Gleitpunktdarstellung: b=10, m=6

x	fl(x)	$\left \frac{fl(x)-x}{x}\right $
$\frac{1}{3} = 0.33333333$	$0.333333 * 10^{0}$	$1.0*10^{-6}$
$\sqrt{2} = 1.41421356$	$0.141421 * 10^{1}$	$2.5 * 10^{-6}$
$e^{-10} = 0.000045399927$	$0.453999 * 10^{-4}$	$6.6 * 10^{-7}$
$e^{10} = 22026.46579$	$0.220265 * 10^5$	$1.6 * 10^{-6}$
$\frac{1}{10} = 0.1$	$0.100000 * 10^{0}$	0.0

### Gleitpunktdarstellung: b=2, m=10

x	fl(x)	$\left \frac{fl(x)-x}{x}\right $
$\frac{1}{3}$	$0.1010101011 * 2^{-1}$	$4.9 * 10^{-4}$
$\sqrt{2}$	$0.1011010100 * 2^{1}$	$1.1 * 10^{-4}$
$e^{-10}$	$0.10111111010 * 2^{-111}$	$3.3*10^{-4}$
$e^{10}$	$0.1010110000 * 2^{1111}$	$4.8 * 10^{-4}$
$\frac{1}{10}$	$0.1100110011 * 2^{-11}$	$2.4 * 10^{-4}$

Dahmen & Reusken



# Maschinengenauigkeit

Rundungsfehler und Gleitpunktarithmetik

lacktriangle Für den relativen Rundungsfehler erhält man für  $x\in\mathbb{D}$ 

$$\left| \frac{\mathrm{fl}(x) - x}{x} \right| \le \frac{\frac{b^{-m}}{2} \cdot b^e}{b^{-1} \cdot b^e} = \frac{b^{1-m}}{2}.$$

Die (relative) Maschinengenauigkeit

$$\mathsf{eps} := \frac{b^{1-m}}{2}$$

charakterisiert das Auflösungsvermögen des Rechners, d.h.

$$\mathsf{eps} = \inf\{\delta > 0 \mid \mathrm{fl}(1+\delta) > 1\}$$

▶ Der Rundungsfehler  $\varepsilon$  erfüllt  $|\varepsilon| \le eps$  und es gilt

$$f(x) = x (1 + \varepsilon).$$

# Maschinengenauigkeit – Beispiel

# Gleitpunktdarstellung: b=10, $m=6 ightarrow ext{eps} = rac{1}{2} imes 10^{-5}$

x	fl(x)	$\left \frac{fl(x)-x}{x}\right $
$\frac{1}{3} = 0.33333333$	$0.333333 * 10^{0}$	$1.0*10^{-6}$
$\sqrt{2} = 1.41421356$	$0.141421 * 10^{1}$	$2.5 * 10^{-6}$
$e^{-10} = 0.000045399927$	$0.453999 * 10^{-4}$	$6.6 * 10^{-7}$
$e^{10} = 22026.46579$	$0.220265 * 10^5$	$1.6 * 10^{-6}$
$\frac{1}{10} = 0.1$	$0.100000 * 10^{0}$	0.0

### Gleitpunktdarstellung: b=2, $m=10 ightarrow {\rm eps} = 9.8 imes 10^{-4}$

x	fl(x)	$\left \frac{fl(x)-x}{x}\right $
1/3	$0.1010101011 * 2^{-1}$	$4.9 * 10^{-4}$
$\sqrt{2}$	$0.1011010100 * 2^{1}$	$1.1 * 10^{-4}$
$e^{-10}$	$0.10111111010 * 2^{-111}$	$3.3*10^{-4}$
$e^{10}$	$0.1010110000 * 2^{1111}$	$4.8 * 10^{-4}$
$\frac{1}{10}$	$0.1100110011 * 2^{-11}$	$2.4 * 10^{-4}$

Dahmen & Reusken

### Historie

#### Vor dem Jahr 1985

Rundungsfehler und Gleitpunktarithmetik

- Kein einheitlicher Standard
- Jeder Computer hatte seine eigene Gleitpunktdarstellung
- Manche binär (Basis 2, 8, 16), manche dezimal, sogar trinär!
- Gleitpunktarithmetik hat sich auf unterschiedlichen Computern unterschiedlich verhalten!

#### Im Jahr 1985

- ANSI/IEEE Standard 754-1985 for Binary Floating-Point Arithmetic
- ► ANSI American National Standards Institute
- ► IEEE Institute of Electrical and Electronics Engineers
- Alle Computer seit 1985 benutzen diesen Standard
- Maschinen-unabhängiges Modell, wie sich Gleitpunktarithmetik verhält

### IEEE Standard

Rundungsfehler und Gleitpunktarithmetik

Double-precision floating-point

64-bit Wort mit

52 bits für f

11 hits für e

1 bit für das Vorzeichen

Der Exponent e ist eine ganze Zahl im Intervall

$$-1022 \le e \le 1023$$

 $\triangleright$  Effektive Mantisse "1 + f"  $(d_1 = 1 \text{ durch Normalisierung, muss nicht gespeichert werden})$ 

### IEEE Standard

Rundungsfehler und Gleitpunktarithmetik

- $x_{MIN}$ : f = 0 und e = -1022
- $x_{MAX}$ : f = 1 eps und e = 1023
- ightharpoonup Überlauf: e=1024 und f=0
  - Schreibweise: infinity oder Inf
  - Frfüllt: 1/Inf = 0 und Inf+Inf = Inf
- Not-a-Number oder NaN: e = 1024 und  $f \neq 0$ 
  - Undefinierte Zahl, z.B. 0/0
- ▶ Unterlauf: e = -1023 (Sonderfall: denormalisierte Zahlen)
- ▶ In MATLAB:

	Binary	Decimal
eps	2^(-52)	2.2204e-16
realmin	2^(-1022)	2.2251e-308
realmax	(2-eps)*2^1023	1.7977e+308

### Pseudoarithmetik

Rundungsfehler und Gleitpunktarithmetik

Exakte elementare arithmetische Operation von Maschinenzahlen ⇒ Maschinenzahl

### **Beispiel**

$$b = 10, m = 3$$
:

$$0.346 \times 10^2 + 0.785 \times 10^2 = 0.1131 \times 10^3 \neq 0.113 \times 10^3$$

Ähnliches passiert bei Multiplikation und Division.

Exakte Arithmetik \ighthapprox Pseudoarithmetik (Gleitpunktarithmetik),

z.B.: 
$$+ \rightsquigarrow \bigoplus$$
.

### Pseudoarithmetik

Rundungsfehler und Gleitpunktarithmetik

### **Forderung**

Für 
$$\nabla \in \{+,-,\cdot,\div\}$$
 gelte  $x \ \, \mathop{\bigtriangledown} y = \mathrm{fl}(x \nabla y) \quad \text{für } x,y \in \mathbb{M}(b,m,r,R).$  Da  $\mathrm{fl}(x) = x \ (1+\varepsilon)$ , folgt somit, dass für  $\nabla \in \{+,-,\cdot,\div\}$   $x \ \, \mathop{\bigtriangledown} y = (x \nabla y)(1+\varepsilon) \quad \text{für } x,y \in \mathbb{M}(b,m,r,R)$  und ein  $\varepsilon$  mit  $|\varepsilon| \leq \mathrm{eps}$  gilt (falls  $|x \nabla y| \leq x_{\mathrm{MAX}} \wedge x \ \, \mathop{\bigtriangledown} y \neq 0$ ).

#### Vorsicht bei Pseudoarithmetik:

- Grundlegende Regeln der Algebra, die bei exakter Arithmetik gelten, sind nicht mehr gültig.
- Reihenfolge der Verknüpfung spielt eine Rolle (Assoziativität der Addition geht verloren).

### Assoziativgesetz

#### Beispiel 2.36.

Rundungsfehler und Gleitpunktarithmetik

Zahlensystem mit b = 10, m = 3. Maschinenzahlen

$$x = 6590 = 0.659 \times 10^4$$
  
 $y = 1 = 0.100 \times 10^1$   
 $z = 4 = 0.400 \times 10^1$ 

Exakte Rechnung:

$$(x + y) + z = (y + z) + x = 6595.$$

Pseudoarithmetik:

$$x\oplus y=0.659 imes 10^4$$
 und  $(x\oplus y)\oplus z=0.659 imes 10^4$ , aber

$$y\oplus z=0.500\times 10^1$$
 und  $(y\oplus z)\oplus x=0.660\times 10^4$ .

# Distributivgesetz

0000000000000000000

Rundungsfehler und Gleitpunktarithmetik

#### Beispiel 2.37.

Für 
$$b=10,\; m=3, x=0.156\cdot 10^2$$
 und  $y=0.157\cdot 10^2$   $(x-y)\cdot (x-y)\;=\;0.01$   $(x\ominus y)\odot (x\ominus y)\;=\;0.100\times 10^{-1}$ 

aber

$$(x\odot x)\ominus (x\odot y)\ominus (y\odot x)\oplus (y\odot y)=-0.100\times 10^1.$$

# Auslöschung

# Beispiel 2.38.

00000000000000000000

Rundungsfehler und Gleitpunktarithmetik

Betrachte

$$x = 0.73563, \quad y = 0.73441, \quad x - y = 0.00122.$$

Bei 3-stelliger Rechnung ( $b=10,\ m=3,\ {\sf eps}=\frac{1}{2}\times 10^{-2}$ ):

$$\tilde{x} = \text{fl}(x) = 0.736, \quad |\delta_x| = 0.50 \cdot 10^{-3}$$

$$ilde{y} \ = \ \mathrm{fl}(y) = 0.734, \quad |\delta_y| \ = \ 0.56 \cdot 10^{-3}$$

Die relative Störung im Resultat:

$$\left| \frac{(\tilde{x} - \tilde{y}) - (x - y)}{x - y} \right| = \left| \frac{0.002 - 0.00122}{0.00122} \right| = 0.64$$

also sehr groß im Vergleich zu  $\delta_x$ ,  $\delta_y$ .

# Auslöschung

**Betrachte** 

$$a = 1.23456$$
 und  $b = 1.23567$ .

Sind a und b Ergebnisse aus Rechnungen mit Gleitpunktarithmetik, so sind die niedrigen Stellen durch Rundungsfehler verfälscht.

Angenommen, die letzten drei Stellen sind verfälscht.

Differenz der korrekten Stellen

$$1.23 - 1.23 = 0$$

Allerdings hat

$$b - a = 0.00111 = 0.111 \cdot 10^{-2}$$

keine einzige korrekte Stelle!

Die korrekten Stellen in a und b löschen sich in a-b aus!

Führende Nullen sind nicht Teil der resultierenden Maschinenzahl.

Rundungsfehler und Gleitpunktarithmetik

00000000000000000000

# Zusammenfassung Gleitpunktarithmetik

$$\left|\frac{(x {\textstyle \bigtriangledown} y) - (x {\textstyle \nabla} y)}{(x {\textstyle \nabla} y)}\right| \leq \mathsf{eps}, \ \ x,y \in \mathbb{M}, \ \ \nabla \in \{+,-,\cdot, \rdots\}$$

Die relativen Rundungsfehler bei den elementaren Gleitpunktoperationen sind  $\leq$  eps, wenn die Eingangsdaten x,yMaschinenzahlen sind.

Sei  $f(x,y)=x
abla y,\; x,y\in\mathbb{R},\; 
abla\in\{+,-,\cdot,\div\}$  und  $\kappa_{\mathrm{rel}}$  die relative Konditionszahl von f. Es gilt

$$abla\in \{\cdot\,,\div\}: \;\; \kappa_{
m rel}\le 1 \quad ext{ für alle } x,y, \ 
abla\in \{+,-\}: \;\; \kappa_{
m rel}\gg 1 \quad ext{ wenn } |x
abla y|\ll \max\{|x|,|y|\}$$

Sehr große Fehlerverstärkung bei +, - möglich (Auslöschung).

# Beispiel: Polynom 7. Grades

Rundungsfehler und Gleitpunktarithmetik

0000000000000000000

$$p(x) = (x-1)^{7}$$

$$= x^{7} - 7x^{6} + 21x^{5} - 35x^{4} + 35x^{3} - 21x^{2} + 7x - 1$$

Matlab-Demo

### Stabilität

#### Definition

Rundungsfehler und Gleitpunktarithmetik

Ein Algorithmus heißt gutartig oder stabil, wenn die durch ihn im Laufe der Rechnung erzeugten Fehler in der Größenordnung des durch die Kondition des Problems bedingten unvermeidbaren Fehlers bleiben.

- Kondition ist Eigenschaft des Problems
- Stabilität ist Eigenschaft des Verfahrens/Algorithmus
- ⇒ Wenn ein Problem schlecht konditioniert ist, kann man nicht erwarten, dass die Numerische Methode (ein stabiler Algorithmus) gute Ergebnisse liefert.

Ziel: Numerische Methode soll Fehlerverstärkung nicht noch weiter vergrößern

$$\Phi_{1}(x) = \frac{\partial f(x)}{\partial x_{1}} \cdot \frac{x_{1}}{f(x)} = \frac{-x_{1}}{|x_{1}|^{2} - x_{2}}$$

$$\Phi_{2}(x) = \frac{1}{2} - \frac{1}{2} \epsilon_{1}(x)$$

$$\Phi_{3}(6.000227, 0.01) = -1,00$$

$$A_{2}(x) = 1_{100}$$
=)  $K_{rd} = min | E_{1}, A_{2} | = 1_{100}$ 

### Beispiel 2.39.

Rundungsfehler und Gleitpunktarithmetik

Bestimmung der kleineren Lösung  $u^*$  von

$$y^2 - 2a_1y + a_2 = 0$$

für  $a_1 = 6.000227$ ,  $a_2 = 0.01$ .

Algorithmus I

$$egin{align} u^* &= f(a_1,a_2) = -rac{-2\cdot a_1}{2} - \sqrt{\left(rac{-2\cdot a_1}{2}
ight)^2 - a_2} \ &= a_1 - \sqrt{a_1^2 - a_2}. \end{split}$$

$$egin{array}{lll} y_1 &=& u_1\,u_1 \ \longrightarrow & y_2 &=& y_1-a_2 \ \longrightarrow & y_3 &=& \sqrt{y_2} \ \longrightarrow & u^* &=& a_1-y_2 \end{array}$$

### Beispiel 2.39.

Algorithmus I

Rundungsfehler und Gleitpunktarithmetik

$$u^* = f(a_1, a_2) = a_1 - \sqrt{a_1^2 - a_2}$$
.

In Gleitpunktarithmetik mit b=10, m=5 (eps  $=\frac{1}{2}\times 10^{-4}$ ):

$$\tilde{u}^* = 0.90000 \times 10^{-3}$$

Exakte Lösung:

$$u^* = 0.83336 \cdot 10^{-3}$$

- ▶ Problem für diese Eingangsdaten  $a_1$ ,  $a_2$  gut konditioniert.
- Durch Algorithmus erzeugte Fehler sind sehr viel größer als der unvermeidbare Fehler.
- ⇒ Algorithmus I ist nicht stabil

# Beispiel 2.39.

Rundungsfehler und Gleitpunktarithmetik

Bestimmung der Lösung  $u^*$  von

$$y^2 - 2a_1y + a_2 = 0$$

für  $a_1 = 6.000227$ ,  $a_2 = 0.01$ .

Algorithmus II (Alternative)

$$u^* = rac{a_2}{a_1 + \sqrt{a_1^2 - a_2}} \ y_1 = a_1 a_1 \ 
ightarrow y_2 = y_1 - a_2 \ 
ightarrow y_3 = \sqrt{y_2} \ 
ightarrow y_4 = a_1 + y_3 \ 
ightarrow u^* = rac{a_2}{y_4}$$

# Beispiel 2.39.

Algorithmus II

Rundungsfehler und Gleitpunktarithmetik

$$u^* = \frac{a_2}{a_1 + \sqrt{a_1^2 - a_2}}$$

In Gleitpunktarithmetik mit b=10, m=5 (eps  $=\frac{1}{2}\times 10^{-4}$ ):

$$\tilde{u}^* = 0.83333 \times 10^{-3}$$

Exakte Lösung:

$$u^* = 0.83336 \cdot 10^{-3}$$

- Gesamtfehler bleibt im Rahmen der Maschinengenauigkeit.
- Auslöschung tritt nicht auf.
- ⇒ Algorithmus II ist stabil

### Rückwärtsstabilität

Rundungsfehler und Gleitpunktarithmetik

Ein Verfahren zur Berechnung von f(x) liefert als Ergebnis  $\tilde{f}(x)$ .

#### Definition

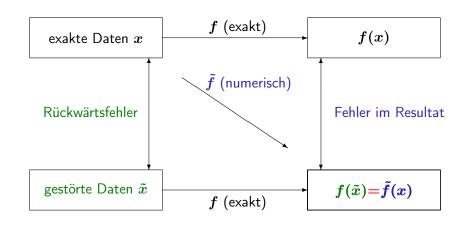
Das Verfahren heißt rückwärts stabil, wenn es für alle  $x \in X$  ein  $\tilde{x} \in X$  gibt, so dass

$$ilde{f}(x) = f( ilde{x})$$
 und  $frac{\|x- ilde{x}\|}{\|x\|} = \mathcal{O}(\mathsf{eps}).$ 

⇒ Ein rückwärts stabiler Algorithmus gibt die exakte Lösung des nahezu richtigen Problems (Daten, d.h.  $x \to \tilde{x} = x + \Delta x$ ).

# Rückwärtsanalyse

Rundungsfehler und Gleitpunktarithmetik



### Rückwärtsstabilität

Rundungsfehler und Gleitpunktarithmetik

#### Satz

Wird ein rückwärts stabiler Algorithmus zur Lösung des Problems  $m{f}$ mit Kondition  $\kappa(x)$  angewendet, so gilt

$$rac{\| ilde{f}(x)-f(x)\|}{\|f(x)\|}=\mathcal{O}(\kappa(x)\operatorname{eps}).$$

Beweis:

$$rac{\| ilde{f}(x)-f(x)\|}{\|f(x)\|} = rac{\|f( ilde{x})-f(x)\|}{\|f(x)\|} \lesssim \kappa(x) \ \underbrace{rac{\| ilde{x}-x\|}{\|x\|}}_{\mathcal{O}(\mathsf{eps})}.$$

### Rückwärtsstabilität

Rundungsfehler und Gleitpunktarithmetik

Was haben wir gemacht?

Fehler im Algorithmus  $\tilde{f}$  wurden

zurückgespiegelt auf Fehler in den Daten  $\tilde{x}$ .

Zusammenfassung

1

Vorteil: Auswertung von  $f(\tilde{x})$  ist Frage nach Kondition von f.

# Beispiel 2.40.: Summation ist rückwärts stabil

Geg.: Maschinenzahlen  $x_1, x_2, x_3$ , Maschinengenauigkeit eps.

Ges.: Summe  $S = (x_1 + x_2) + x_3$ .

Man erhält

Rundungsfehler und Gleitpunktarithmetik

$$\widetilde{S} = ((x_1 + x_2) (1 + \varepsilon_2) + x_3) (1 + \varepsilon_3)$$

mit  $|\varepsilon_i| < \text{eps}, i = 2, 3$ .

Daraus folgt

$$egin{aligned} \widetilde{S} &= x_1 \left( 1 + arepsilon_2 
ight) \left( 1 + arepsilon_3 
ight) + x_2 \left( 1 + arepsilon_2 
ight) \left( 1 + arepsilon_3 
ight) + x_3 \left( 1 + arepsilon_3 
ight) \\ &\doteq x_1 \left( 1 + arepsilon_2 + arepsilon_3 
ight) + x_2 \left( 1 + arepsilon_2 + arepsilon_3 
ight) + x_3 \left( 1 + arepsilon_3 
ight) \\ &= x_1 \left( 1 + \delta_1 
ight) + x_2 \left( 1 + \delta_2 
ight) + x_3 \left( 1 + \delta_3 
ight) \end{aligned}$$

wobei

$$|\delta_1| = |\delta_2| = |\varepsilon_2 + \varepsilon_3| < 2 \cdot \mathsf{eps}, \quad |\delta_3| = |\varepsilon_3| < \mathsf{eps}$$

# Beispiel 2.40.: Summation ist rückwärts stabil

Es gilt

Rundungsfehler und Gleitpunktarithmetik

$$\widetilde{S} = x_1 (1 + \delta_1) + x_2 (1 + \delta_2) + x_3 (1 + \delta_3)$$
  
=:  $\tilde{x}_1 + \tilde{x}_2 + \tilde{x}_3$ ,

wobei

$$|\delta_1| = |\delta_2| = |arepsilon_2 + arepsilon_3| \leq 2 \cdot \mathsf{eps}, \quad |\delta_3| = |arepsilon_3| \leq \mathsf{eps}$$

 $\Rightarrow$  Fehlerbehaftetes Resultat  $\widetilde{S}$  als exaktes Ergebnis zu gestörten Eingabedaten  $\widetilde{x}_i = x_i (1 + \delta_i)$ .

Der durch Rechnung bedingte Fehler ist höchstens

$$egin{array}{ccc} \left| rac{f( ilde{x}) - f(x)}{f(x)} 
ight| & \leq & \kappa_{ ext{rel}}(x) \cdot \sum\limits_{j=1}^{3} \left| rac{ ilde{x}_{j} - x_{j}}{x_{j}} 
ight| \ & \leq & \kappa_{ ext{rel}}(x) \cdot \sum\limits_{j=1}^{3} \left| \delta_{j} 
ight| \leq \kappa_{ ext{rel}}(x) \cdot 5 \cdot ext{eps.} \end{array}$$

# Beispiel 2.40.

Rundungsfehler und Gleitpunktarithmetik

Der für die Summation  $f(x) = f(x_1, x_2, x_3) = x_1 + x_2 + x_3$ unvermeidbare Fehler ist

$$\left| rac{f( ilde{x}) - f(x)}{f(x)} 
ight| \leq \kappa_{ ext{rel}}(x) \cdot \sum_{j=1}^{3} \left| rac{ ilde{x}_{j} - x_{j}}{x_{j}} 
ight| \leq \kappa_{ ext{rel}}(x) \cdot 3 \cdot \mathsf{eps},$$

wenn Daten höchstens mit Maschinengenauigkeit gestört werden  $(\tilde{x}_i = x_i(1+\varepsilon), |\varepsilon| < \text{eps}).$ 

Die Größenordnung der Fehler ist identisch

 $\Rightarrow$  Berechnung von S ist ein stabiler Algorithmus.

### Summenbildung

Rundungsfehler und Gleitpunktarithmetik

Summenbildung tritt in vielen Problemen z. B. Skalarprodukte, Matrix/Vektor-Multiplikation, ... auf.

Wir betrachten: 
$$S_n = \sum\limits_{j=1}^n x_j$$

Analog zum Fall  $S = (x_1 + x_2) + x_3$  kann man zeigen, dass

$$(x_1 \oplus x_2 \oplus \cdots \oplus x_n) - (x_1 + x_2 + \cdots + x_n)$$

$$\stackrel{\cdot}{=} x_1(\varepsilon_1 + \varepsilon_2 + \cdots + \varepsilon_n)$$

$$+ x_2(\varepsilon_2 + \cdots + \varepsilon_n) + \cdots + x_n \varepsilon_n$$

mit  $|\varepsilon_i| < \text{eps}, \ i = 1, \dots, n$ .

- Der erste Summand wird mit größtem Fehler multipliziert.
- Reihenfolge bei der Summation wichtig

⇒ der relative Fehler wird am kleinsten, wenn die betragsgrößten Summanden zuletzt aufsummiert werden (vgl. Beispiel 2.36.).

### Zusammenfassung

Was Sie mitnehmen sollten:

Wie werden Zahlen im Computer dargestellt

ightharpoonup Maschinenzahlen  $\mathbb{M}(b,m,r,R)$  $\Rightarrow x_{\text{min}}, x_{\text{max}}, \text{ eps}, |\varepsilon| < \text{eps}$ 

Welche Probleme können dabei/deswegen auftreten?

- Assoziativ- und Distributivgesetz nicht mehr gültig
- ▶ Gefahr der Auslöschung bei  $\nabla \in \{+, -\}$

### Zusammenfassung

#### Stabilität vs. Kondition

- Bei einem stabilen Lösungsverfahren bleiben die im Laufe der Rechnung erzeugten Rundungsfehler in der Größenordnung der durch die Kondition des Problems bedingten unvermeidbaren Fehler.
- Kenntnisse über die Kondition eines Problems sind oft für die Interpretation oder Bewertung der Ergebnisse von entscheidender Bedeutung
  - "Schlechtes Ergebnis" bedeutet nicht unbedingt gleich "instabiler Algorithmus", sondern deutet evtl. auf eine schlechte Kondition des Problems hin.
- ► In einem Algorithmus sollen (wegen Stabilität) Auslöschungseffekte vermieden werden.

Rundungsfehler und Gleitpunktarithmetik

Es seien  $x_{\min}$  bzw.  $x_{\max}$  die kleinste bzw. größte (strikt) positive Zahl sowie eps die relative Maschinengenauigkeit in der Menge  $\mathbb{M}(b,m,r,R)$  der Maschinenzahlen und  $\mathbb{D}:=[-x_{\text{MAX}},-x_{\text{MIN}}]$  $\cup [x_{ ext{ iny MIN}}, x_{ ext{ iny MAX}}]$ . Ferner beschreibe  $\mathrm{fl}: \mathbb{D} o \mathbb{M}(b,m,r,R)$  die Standardrundung. Alle Zahlen sind im Dezimalsystem angegeben.

Berechnen Sie  $x_{\text{MAX}}$  für  $\mathbb{M}(3,2,-1,3)$ 



### Verständnisfragen

Es seien  $x_{\min}$  bzw.  $x_{\max}$  die kleinste bzw. größte (strikt) positive Zahl sowie eps die relative Maschinengenauigkeit in der Menge  $\mathbb{M}(b,m,r,R)$  der Maschinenzahlen und  $\mathbb{D}:=[-x_{\text{MAX}},-x_{\text{MIN}}]$  $\cup [x_{\text{MIN}}, x_{\text{MAX}}]$ . Ferner beschreibe  $\mathrm{fl}: \mathbb{D} \to \mathbb{M}(b, m, r, R)$  die Standardrundung. Alle Zahlen sind im Dezimalsystem angegeben.

Berechnen Sie  $x_{\text{MAX}}$  für  $\mathbb{M}(3,2,-1,3)$  24

- $ig| \mathbf{W} ig|$  Es gilt  $\Big| rac{\mathrm{fl}(x-y) (x-y)}{(x-y)} \Big| \leq \mathrm{eps}$  für alle  $x,y \in \mathbb{M}(b,m,r,R)$ mit  $x \neq y$  und  $|x - y| \leq x_{\text{MAX}}$ .
- $igg| ext{f}$  Es gilt  $\Big|rac{ ext{fl}(x-y)-(x-y)}{(x-y)}\Big| \leq ext{eps}$  für alle  $x,y \in \mathbb{D}$  mit x 
  eq yund  $|x-y| \leq x_{\text{MAX}}$ .
- Bei einem stabilen Algorithmus ist der Ausgabefehler nicht viel größer als der Eingabefehler.