

Kapitel 1: Fließkommazahlen

1
09.10.05

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

unsigned int	\mathbb{N}_0	(char, short, long)
int	\mathbb{Z}	(char, short, long)
float, double	\mathbb{R}	(float)
double	\mathbb{R}	
Complex <double>	\mathbb{C}	(dito float)

int: exakt, aber endlicher Bereich

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

Was hat das für Folgen?

Beispiel 1.1 (Potenzreihe für e^x). Auf Folie

□

1.1 Zahlendarstellung

Zellenwabsystem.

$$x = \left(\begin{array}{c} + \\ - \end{array} \right) \dots m_n \beta^n + \dots + m_1 \beta^1 + m_0 \beta^0 + m_{-1} \beta^{-1} + \dots + m_{-k} \beta^{-k} + \dots$$

\pm ist das Vorzeichen (ein Bit genügt).

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

$m_i \in \{0, 1, 2, \dots, \beta-1\}$ heißen Ziffern.

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

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

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

Basis 10 in Europa ab ca 1585

Reine Power: Jedes $\beta \geq 2$ möglich

Festkommazahlen:

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: $6.6260693 \cdot 10^{-34} \text{ J s}$

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

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

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

Fließkommazahlen erlauben effizientere Darstellung unterschiedlich großer Zahlen. (im Sinne des relativen Fehlers)

Definition 1.2 (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 Fließkomma-Darstellung eindeutig. □

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

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

" $|m(x)| |\beta^{e(x)}|$
 > 0 also $|\beta^{e(x)}| = \beta^{e(x)}$

Beispiel 1.3

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

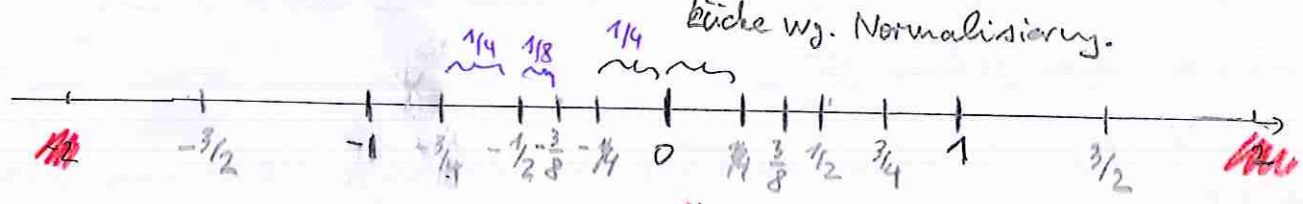
$0, \underline{000.000.000} 0.14 = 0.140 \cdot 10^{-10} \notin F(10, 3, 1)$ da Exponent zu klein.

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

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

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

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



Mögliche Lösung: Nehme für $|x| < \frac{1}{4}$ die nicht normalisierten Zahlen hinzu.

Zahlenbereich:

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

\rightarrow Betrag möglichst groß

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

β^{-r} Wertigkeit der letzten Stelle

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

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

\rightarrow Betrag möglichst klein

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

Damit

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

Abstände zwischen Fließkommazahlen

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)

$$\beta^{e-1} = 0.\overset{m_1}{1}0\dots0 \cdot \beta^e$$

nächste Zahl: $0.10\dots1 \cdot \beta^e$ } Abstand $\beta^{-r} \cdot \beta^e = \beta^{e-r}$

$$e-1 = z \Rightarrow e = z+1$$

nächste Zahl $+ \beta^{e-r}$ $0.\overset{(A-2)}{\beta} \dots \overset{(A-2)}{\beta} \cdot \beta^e$ } Abstand β^{e-r} (da Exponent immer noch gleich)

$$\beta^e = 0.10\dots0 \cdot \beta^{e+1}$$

Abstand β^{e-r} } Abstand springt!

nächste Zahl: $0.10\dots1 \cdot \beta^{e+1}$ } Abstand $\beta^{-r} \cdot \beta^{e+1} = \beta^{e+1-r}$

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

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

a) für $|x| \neq \beta^{-1}$:

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

Da $\beta^{-1} \leq |m(x)| < 1$ folgt $|x| \beta^{-r} < |x-x'| < |x| \beta^{1-r}$
nach (1.1) $x \neq 0$ groß, wenn $|m(x)|$ klein, d.h. β^{-1}

b) $|m(x)| = \beta^{-1}$: $|x-x'| = \beta^{e-1-r} = \beta^{-1} \beta^{e(x)} \beta^{-r} = |x| \cdot \beta^{-r}$

Für den relativen Abstand gilt:

$$\beta^{-r} \leq \frac{|x-x'|}{|x|} \leq \beta^{1-r} (= \beta \beta^{-r})$$

der Fall b

Er schwankt um den Faktor β je nach Größe von $|x|$. Dieser Faktor heißt wobble.

Konsequenz: Kleine β sind besser!

Übung:

$x \in \mathbb{F}$
 β^{1-r} : Die kleinste Zahl, so dass $1+x \neq 1$.
 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_{\max}	+127	≥ 1024	1023	≥ 16384	} nicht sym. andernfalls $F(\beta, r, s)$
e_{\min}	-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.

$$\text{Setze } e = c - 1023 \Rightarrow e \in [-1022, 1023] \quad \rightarrow \text{kein Vorzeichen Extra notwendig!}$$

$\begin{matrix} \text{"} \\ 1-1023 \end{matrix}$
 $\begin{matrix} \text{"} \\ 2046-1023 \end{matrix}$

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

- $c = 0, m = 0$ kodiert die Null
- $c = 0, m \neq 0$ kodiert denormalisierte Darstellung
- $c = 2047, m \neq 0$ kodiert NaN = "not a number" (z.B. Division durch Null)
- $c = 2047, m = 0$ kodiert ∞ (Überlauf).

- $64 - 11 = 53$ 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 $r = 53$. (kleinste Wertigkeit)

- Der Standard definiert 4 Rundungsarten:
nach $+\infty$, nach $-\infty$, nach ϕ , nearest (natürliche Rundung)

1.2 Runden und Rundungsfehler

6
04.10.09

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

$$\text{rd} : D(\beta, r, s) \rightarrow \mathbb{F}(\beta, r, s) \quad (1.2)$$

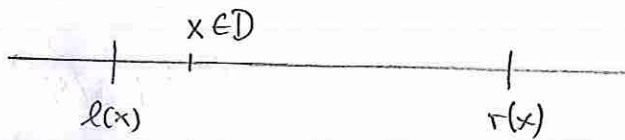
Achtung: rd setzt voraus, dass x im darstellbaren Bereich liegt!

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

Sinnvollerweise soll für rd gelten:

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

○ 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

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

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 $x \in D$,

(alle Ziffern m_i sind korrekt (!), das ist der Unterschied zum TMS) nächste Stelle

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

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

Beispiel 1.4

$$(x)_{10} = 110 \rightarrow x' \in \mathbb{F}(4, 3, 2)$$

1) Exponent

$$\log_4(110) = 3,39$$

Aufrunden damit Mantisse < 1

$$4^x = 110 \Rightarrow x \cdot \log 4 = \log 110$$

$$x = \frac{\log 110}{\log 4} = 3,39$$

$$\Rightarrow \text{Exponent } 4 \cdot (10) \text{ d.h. } x = m \cdot 4^4 \Rightarrow m = \frac{x}{4^4}$$

2) Mantisse

$$\text{teile } x \text{ durch } 4^4 \quad (m)_{10} = \frac{x}{4^4} = \frac{110}{256} = (0,4296875)_{10} \text{ (exakt!)}$$

$$y_1 = \frac{110}{4^4} = 0,4296875 = (m_1 \cdot 4^{-1} + m_2 \cdot 4^{-2} + m_3 \cdot 4^{-3} + \dots)$$

$$y_1 \cdot 4 = 1,71875 \Rightarrow m_1 = 1$$

$$y_2 = (y_1 - m_1) \cdot 4 = 2,875 \Rightarrow m_2 = 2$$

$$y_3 = (y_2 - m_2) \cdot 4 = 3,5 \Rightarrow m_3 = 3$$

$$y_4 = (y_3 - m_3) \cdot 4 = 2,0 \Rightarrow m_4 = 2$$

$$(x)_4 = 0,1 \cdot 2 \cdot 3 \cdot 2 \cdot 4^{10} \text{ (exakt)}$$

3 Stellen ($r=3$)

② Runden mit gerader Rundung

$$\Rightarrow (x')_4 = 1,30 \cdot 4^{10} \text{ (weil } m_3 = 3, \text{ also ungerade.)}$$

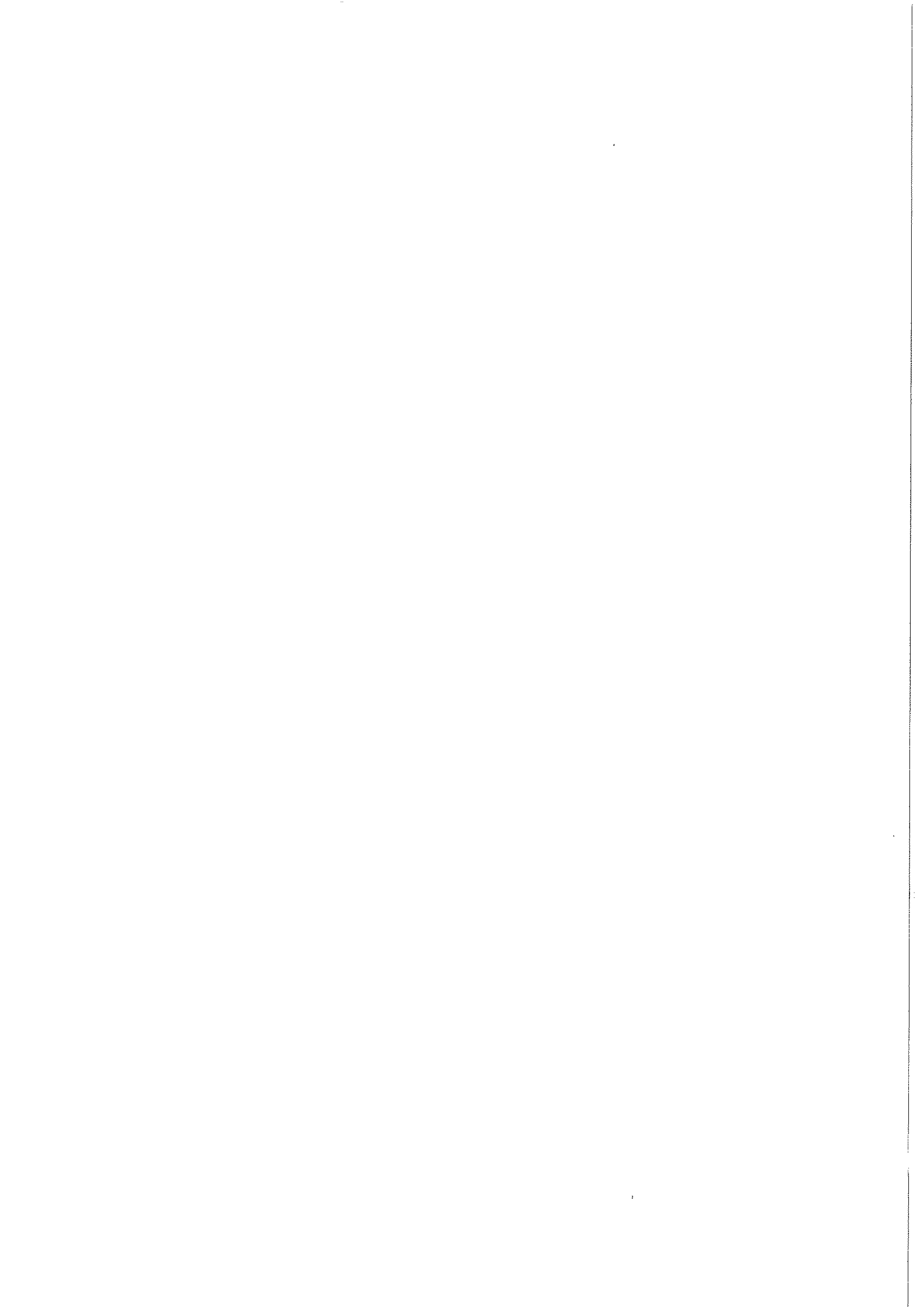
$$\text{Kontrolle: } \left. \begin{array}{l} (1,30 \cdot 4^{10})_4 = (112)_{10} \\ (1,23 \cdot 4^{10})_4 = (108)_{10} \end{array} \right\}$$

das zeigt, dass x genau zwischen den beiden Zahlen ist.

① Natürliche Rundung

$$m_{r+1} = m_4 = 2 = \frac{\beta}{2} \Rightarrow \text{„aufrunden“}$$

$$rd(x) = (0,123 + 0,001) \cdot 4^{10} = 0,130 \cdot 4^{10} \text{ (alles zur Basis 4)}$$



Gerade Rundung (β sei gerade)

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

Damit ist m_r ⁱⁿ $rd(x)$ immer gerade, wenn gerundet wurde (d.h. $\frac{|x-l(x)|}{|x-r(x)|} = 1$)

- 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. $\rightarrow \bar{U}$ -Aufgabe.
Beispiel 1.4 \rightarrow extra Blatt, bis hier vorbereitet

Definition 2.5 (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 $x \neq 0$

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

Umformen liefert:

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

\uparrow
abs. Fehler

Näherungswert =
echter Wert $\times (1 + \epsilon)$

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 $x = \text{Entfernung Leidelberg-Paris} \approx 500 \text{ km}$ ist

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

dagegen relativ groß.

Lemma 1.5 (Rundungsfehler)

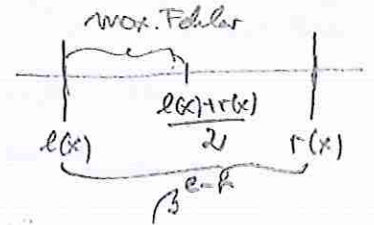
Bei der Rundung in $\mathbb{F}(\beta, r, s)$ gilt für den absoluten Fehler

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

1.3
(Zil)

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 $e(l(x)) = e(r(x)) = e(x)$.
 $l = 0.m_1 \dots m_r \cdot \beta^e$ $m_r < \beta - 1$
 $r = 0.m_1 \dots (m_r + 1) \cdot \beta^e$ $m_r < \beta - 1$
 $|r - l| = \beta^{-r} \cdot \beta^e = \beta^{e-r}$ $m_r < \beta - 1$
nähe S. 4.

wg Fall b)

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

zwei aufeinander folg. 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

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

\uparrow
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 $\epsilon_{ps} := \frac{1}{2} \beta^{1-r}$ heißt Maschinengenauigkeit. (Nimmt Rannach 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.

1.3 Fließkommaarithmetik

Wir benötigen eine Arithmetik auf \mathbb{F} :

$$\otimes : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F} \quad \text{mit } \otimes \in \{+, -, \cdot, / \},$$

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

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

Bsp Multiplikation
→ Verdopplung der
Nachkommastellen.

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

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

1.4
~~(2.3)~~

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

Beispiel 1.7 (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 $x \ominus y = rd(x - y)$.

① $x - y$:

$x = 0.215$	000	000	000	000	0	$\cdot 10^8$
$y = 0.$	000	000	000	000	0.125	$\cdot 10^8 \leftarrow$
	13 Nullen! ... 1 1 1					
$x - y = 0.214$	999	999	999	987	5	$\cdot 10^8$

Schiebe y auf
gleichen Exponenten.
 $y = 0.125 \cdot 10^5 = 0.125 \cdot 10^{-13} \cdot 10^8 = 10^{-5}$

② Runden

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

Bemerkung
zum
Tabellen-
macher-
dilemma.

Problem: Schritt 1 erfordert extrem hohe Stellenzahl $O(\beta^s)$!

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

Dies liefert im übrigen 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 $y = rd(exp(x))$

Problem: $exp(1.626) = 5.0835$

ich verstehe es nicht.

genauer $exp(1.626) = 5.0835000 \dots$
 \downarrow möglichst viele Stellen möglich
 $5.0834999 \dots$
 das ist noch möglich.
 wenn 5 eine gültige Ziffer ist, dann muss man aufrunden

Im allgemeinen ist das gefährlich:

$$\begin{array}{lcl}
 x = 0,101 \cdot 10^1 & \longrightarrow & 0,101 \cdot 10^1 \\
 y = 0,993 \cdot 10^0 & \xrightarrow{\text{schreiben}} & \frac{0,099 \cdot 10^1}{0,002 \cdot 10^1}
 \end{array}$$

Schieben und runden

relativer Fehler im Ergebnis:

$$\frac{(x \oplus y) - (x - y)}{x - y} = \frac{0,02 - 0,017}{0,017} \approx 0,176 \approx 35 \text{ eps}$$

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

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

Mit zwei Stellen erreicht man exakte Rundung!

Die zusätzlichen Stellen nennt man "guard digits".

- IEEE 754: $\oplus, \ominus, \otimes, \oslash$ sind exakt gerundet ebenso \sqrt{x}
- Tabellenmacherdilemma

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!

etwa ~~$(1+\epsilon)-1$ bzw. $(1-1)+\epsilon$~~
 $(\epsilon+1)-1 = 0$ vs. $\epsilon+(1-1) = \epsilon$

- $\exists y \in \mathbb{F}$ so dass $x \oplus y = x$ \leftrightarrow

- Allerdings gilt das Kommutativgesetz!

- Es gelten auch folg. einfache Regeln:

$$(-x) \oplus y = -(x \oplus y), \quad 1 \oplus x = x, \quad x \oplus y = 0 \Rightarrow x = 0 \vee y = 0,$$

$$x \oplus z \leq y \oplus z \quad \text{falls } x \leq y \text{ und } z > 0.$$

7.4 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 (d.h. $+$, \cdot , \ominus , \oslash)

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

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

ii) Jedes φ_i steuert einen (unbekannten) Fehler bei

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

Das machen wir aber nicht so formal.

~~Wie in Kap. 1~~ Nutze Aufspaltung:

$$\underbrace{F(x)}_{\substack{\text{exaktes} \\ \text{Ergebnis}}} - \underbrace{F'(rd(x))}_{\substack{\text{Angabe} \\ \text{Runder} \\ \text{numerische} \\ \text{Auswertung}}} = \underbrace{F(x) - F(rd(x))}_{\text{① Konditionsanalyse von } F} + \underbrace{F(rd(x)) - F'(rd(x))}_{\text{② Rundungsfehleranalyse. Unterschied } F, F' \text{ bei gleicher Eingabe (aus } \mathbb{F} \text{)}}$$

(1.5)
~~(2.4)~~

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 Komponenten 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 ^{1.8} 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, so dass 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:

13
06.10.09

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

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 $x_j \neq 0$

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

1,6
~~(2,5)~~

↑ relativer Eingabefehler. $\leq \epsilon$!

Fasst man die Verstärkungsfaktoren in einer Matrix:

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

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

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

$$K(x) = \sup \left\{ \frac{\|F(x + \Delta x) - F(x)\|}{\|\Delta x\|} \frac{\|x\|}{\|F(x)\|} \right\}$$

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

Definition 2.4 ^{1.9} 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 ^{1.6} gilt

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

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) ^{1.6} gilt gerade

$$f(x) \cdot \frac{\Delta x_j}{x_j} = \frac{\epsilon_j x_j}{x_j} = \epsilon_j$$

Beispiel 2.9 ^{1.10}

U: 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 ^{1.6} (2.5):

$$\frac{F(x_1 + \Delta x_1, x_2 + \Delta x_2) - F(x_1, x_2)}{F(x_1, x_2)} = \underbrace{1 \cdot \frac{x_1}{x_1 + x_2}}_{= k_1} \cdot \underbrace{\frac{\Delta x_1}{x_1}}_{1.1 \leq \epsilon_{ps}} + \underbrace{1 \cdot \frac{x_2}{x_1 + x_2}}_{= k_2} \cdot \underbrace{\frac{\Delta x_2}{x_2}}_{1.1 \leq \epsilon_{ps}}$$

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

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

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

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



② Rundungsfehleranalyse

d.h. Nach ^{1.5}~~(2.4)~~: $F(x) - F'(x)$ mit $x \in \mathbb{F}^m$ Maschinenzahl.

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

Wegen ^{1.4}~~(2.3)~~ (exakt gerundete Arithmetik) und Lemma ^{1.6}~~2.5~~ gilt

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

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

Und damit

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

Analyse, in erster Näherung

Beispiel ^{1.11}~~2.10~~ $F(x_1, x_2) = x_1^2 - x_2^2$ mit zwei Realisierungen

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

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

a) $u = x_1 \otimes x_1 = (x_1 \cdot x_1) (1 + \varepsilon_1)$
 $v = x_2 \otimes x_2 = (x_2 \cdot x_2) (1 + \varepsilon_2)$ $\varepsilon_1 \neq \varepsilon_2$ aber $|\varepsilon_i| \leq \text{eps}!$

$$\begin{aligned} F_a(x_1, x_2) &= u \ominus v = (u - v)(1 + \varepsilon_3) \\ &= (x_1^2(1 + \varepsilon_1) - x_2^2(1 + \varepsilon_2))(1 + \varepsilon_3) \\ &= (x_1^2 + \varepsilon_1 x_1^2 - x_2^2 - \varepsilon_2 x_2^2)(1 + \varepsilon_3) \\ &= \underbrace{x_1^2 - x_2^2}_{= F(x_1, x_2)} + \underbrace{\varepsilon_1 x_1^2 - \varepsilon_2 x_2^2 + \varepsilon_3 x_1^2 - \varepsilon_3 x_2^2}_{\text{erste Ordnung}} + \underbrace{\varepsilon_1 \varepsilon_3 x_1^2 - \varepsilon_2 \varepsilon_3 x_2^2}_{\text{zweite Ordnung}} \end{aligned}$$

relativer Fehler

$$\frac{F_a(x_1, x_2) - F(x_1, x_2)}{F(x_1, x_2)} \stackrel{\text{im ersten Näherung}}{=} \frac{x_1^2}{x_1^2 - x_2^2} (\varepsilon_1 + \varepsilon_3) + \frac{x_2^2}{x_2^2 - x_1^2} (\varepsilon_2 + \varepsilon_3)$$

Bemerkung: gleiche Verstärkungsfaktoren wie im Bsp. ^{1.10}~~2.9~~!

$$b) \quad u = x_1 \ominus x_2 = (x_1 - x_2)(1 + \varepsilon_1)$$

$$v = x_1 \oplus x_2 = (x_1 + x_2)(1 + \varepsilon_2)$$

$$\begin{aligned} F_{\mathbb{R}}(x_1, x_2) = u \odot v &= (u \cdot v)(1 + \varepsilon_3) \\ &= \left((x_1 - x_2)(1 + \varepsilon_1) \cdot (x_1 + x_2)(1 + \varepsilon_2) \right) (1 + \varepsilon_3) \\ &= \underbrace{(x_1 - x_2)(x_1 + x_2)}_{x_1^2 - x_2^2} \underbrace{(1 + \varepsilon_1)(1 + \varepsilon_2)(1 + \varepsilon_3)}_{1 + \varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \varepsilon_1\varepsilon_2 + \dots + \varepsilon_1\varepsilon_2\varepsilon_3} \end{aligned}$$

$$\frac{F_{\mathbb{R}}(x_1, x_2) - F(x_1, x_2)}{F(x_1, x_2)} \stackrel{!}{=} \frac{x_1^2 - x_2^2}{x_1^2 - x_2^2} (\varepsilon_1 + \varepsilon_2 + \varepsilon_3) = \varepsilon_1 + \varepsilon_2 + \varepsilon_3$$

⇒ Verstärkungsfaktor 1, besser als vorher.

Definition 1.12 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 übersteigern.

1) a.w. Verstärkungsfaktoren aus Rundungsfehleranalyse \leq denen aus Konditionsanalyse \Rightarrow „numerisch stabil“

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

1.5 Auslöschung

Obrige Beispiele ^{1.10, 1.11} ~~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 $x_1 \approx x_2$.

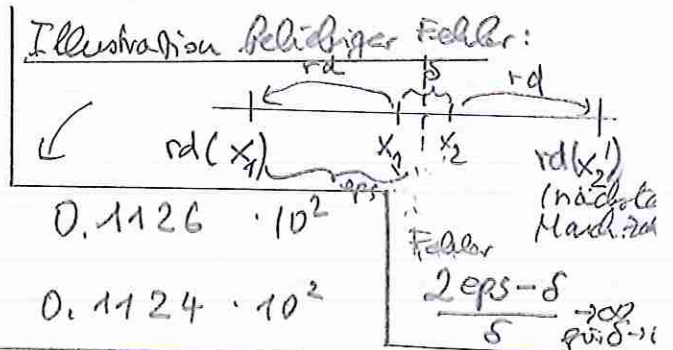
Bemerkung ^{1.13} ~~2.12~~ Bei der Auslöschung werden vor der entsprechenden Addition bzw. Subtraktion eingeführte Fehler extra verstärkt.

Sind $x_1, x_2 \in \mathbb{F}$ Maximalzahlen, 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 ^{1.14} ~~2.13~~ $\mathbb{F} = \mathbb{F}(10, 4, 1)$



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

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

$$x_1 \ominus x_2 = 0.00013871 \cdot 10^2$$

$$= 0.13871 \cdot 10^{-1}$$

$$\text{rd}(x_1) - \text{rd}(x_2) = 0.0002 \cdot 10^2$$

$$= 0.2 \cdot 10^{-1}$$

~~überhaupt kein Fehler~~

relativer Fehler:

keine gültige Ziffer!

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

Hier: Rundung der Eingangsgrößen.

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

Regel ^{1.15} ~~2.14~~ Setze potentiell gefährliche Operationen möglichst früh im Algorithmus ein. (Siehe Beispiel ^{1.19}).

1.6 Die quadratische Gleichung

Die Gleichung

$$y^2 - py + q = 0$$

reell und

hat für $p^2/4 > q \neq 0$ die beiden verschiedenen Lösungen

$$y_{1,2} = f_{\pm}(p,q) = \frac{p}{2} \pm \sqrt{\frac{p^2}{4} - q}, \quad (\text{Definiert zwei } f\text{'s!})$$

Die Konditionsanalyse liefert:

$$\frac{f(p+\Delta p, q+\Delta q) - f(p,q)}{f(p,q)} = \left(1 \pm \frac{p}{2\sqrt{\frac{p^2}{4} - q}} \right) \frac{p}{p \pm 2\sqrt{\frac{p^2}{4} - q}} \frac{\Delta p}{p} + \frac{q}{\sqrt{\frac{p^2}{4} - q} (p \pm 2\sqrt{\frac{p^2}{4} - q})} \frac{\Delta q}{q}$$

- ⇒ - Für $\frac{p^2}{4} \gg q$ und $p < 0$ ist $f_{-}(p,q) = \frac{p}{2} - \sqrt{\frac{p^2}{4} - q}$ gut konditioniert
- Für $\frac{p^2}{4} \gg q$ und $p > 0$ ist $f_{+}(p,q) = \frac{p}{2} + \sqrt{\frac{p^2}{4} - q}$ gut konditioniert
- für $\frac{p^2}{4} \approx q$ sind f_{+} und f_{-} schlecht konditioniert ($@ \rightarrow 0$)

Numerisch stabile Auswertung für den Fall $p^2/4 \gg q$.

$p < 0$: Berechne $y_2 = \frac{p}{2} - \sqrt{\frac{p^2}{4} - q}$ $y_1 = q/y_2$ nach Vieta:
keine Ausl.n. Ver. keine Ausrech. $p < 0$ $p = y_1 + y_2$
 $q = y_1 \cdot y_2$

$p > 0$: Berechne $y_1 = \frac{p}{2} + \sqrt{\frac{p^2}{4} - q}$ $y_2 = q/y_1$.