

§5 Zahlendarstellung auf Digitalrechnern und Fehleranalyse

Bemerkung 5.1, Zahlendarstellung:

Bei den Digitalrechnern werden reelle Zahlen durch endlich viele physikalische Zustände dargestellt. Weil nur endlich viele reelle Zahlen dargestellt werden können, müssen die anderen approximiert werden.

Sei $g \in \mathbb{N}$ fest mit $g \geq 2$. Zu $r \in \mathbb{R}$ gibt es ein eindeutig bestimmtes $j \in \{0, 1\}$, ein eindeutig festgelegtes $l \in \mathbb{Z}$ und zu $k \in \mathbb{Z}$ mit $k \leq l$ eindeutig bestimmte $a_k \in \mathbb{Z}$ mit:

$$(5.1) \quad r = (-1)^j \cdot \sum_{k=-\infty}^l a_k \cdot g^k$$

$$(5.1a) \quad a_l \neq 0 \text{ für } r \neq 0, \quad j = l = 0 \text{ für } r = 0,$$

$$(5.1b) \quad a_k \in S_g := \{0, 1, \dots, g-1\} \text{ für alle } k \leq l, \\ a_k < g-1 \text{ für unendlich viele } k < l.$$

g hängt dabei von der Maschine ab;
in der Regel ist $g = 10$ (Decimaldarstellung)
oder $g = 2$ (Dualdarstellung)

Bemerkung 5.2, Die Festpunktdarstellung:

Bei der Festpunktdarstellung sind die Stellen s und t der Stellen vor bzw. nach dem Komma durch die Maschine fixiert. Die darstellbaren

Zahlen r haben die Gestalt

$$(5.2) \quad r = \pm \sum_{k=-t}^{s-1} b_k \cdot g^k \quad \text{mit } b_k \in S_g \text{ für } -t \leq k < s.$$

Darstellung auf der Maschine

$$\boxed{\pm b_{s-1} \dots b_0 \mid b_{-1} \dots b_{-t}}$$

Die darstellbaren Zahlen sind in dem offenen Intervall $(-g^s, g^s)$ äquidistant verteilt mit dem Abstand g^{-t} .

Die Festpunktdarstellung ist beispielsweise angemessen für kaufmännische Rechnungen.

Bemerkung 5.3, Die Gleitpunktdarstellung:

Zu $r \in \mathbb{R} \setminus \{0\}$ gibt es eindeutig bestimmte $a \in \mathbb{R}$ mit $\frac{1}{g} \leq |a| < 1$ und $b \in \mathbb{Z}$ mit

$$(5.3) \quad r = a \cdot g^b.$$

a heißt die Mantisse und b der Exponent von r .

Für gewisse, durch die Maschine fixierte, $s, t \in \mathbb{N}$ wird a bzw. b durch t bzw. s Stellen und das Vorzeichen dargestellt.

Ist r exakt darstellbar, so ist

$$-(g^t - 1) \leq b \leq g^s - 1$$

und folglich

$$-g^s (g^t - 1) < r < g^s (g^t - 1)$$

Beispiele:

Für $g=10$, $s=2$ und $t=8$ erhält man folgende Darstellungen:

$$\begin{array}{rcl}
 5138 & \mapsto & \boxed{0.51380000} \quad \boxed{04} \\
 -18,23197 & \mapsto & \boxed{-0.18231970} \quad \boxed{02} \\
 0,03164 & \mapsto & \boxed{0.31640000} \quad \boxed{-01} \\
 -0,03164 & \mapsto & \boxed{-0.31640000} \quad \boxed{-01}
 \end{array}$$

Bemerkung 5.4, Rundung:

Es sei M die -endliche- Menge der in der Maschine darstellbaren Zahlen - auch Maschinenzahlen - genannt.

Eine Zahl $r \in \mathbb{R}$ zu runden heißt:

Suche ein $\tilde{r} \in M$ mit:

$$|r - \tilde{r}| \leq |r - s| \quad \text{für alle } s \in M.$$

Im Falle der Festpunktdarstellung setzen wir für r wie in (5.1) und $l \leq s-1$:

$$(5.4) \quad rd(r) := \begin{cases} (-1)^j \cdot \sum_{k=-t}^l a_k \cdot g^k & \text{falls } a_{-t-1} < \frac{1}{2} \cdot g \\ (-1)^j \cdot \left(\sum_{k=-t}^l a_k \cdot g^k + g^{-t} \right) & \text{falls } a_{-t-1} \geq \frac{1}{2} \cdot g. \end{cases}$$

Wir erhalten dann die -von r unabhängige- Abschätzung

$$(5.4a) \quad |r - rd(r)| \leq \frac{1}{2} \cdot g^{-t}.$$

Im Falle der Gleitpunktdarstellung setzen wir

$$gl(0) := 0,$$

$$(5.4) \quad gl(a \cdot g^b) := rd(a) \cdot g^b \quad \text{für } \frac{1}{g} \leq |a| < 1, b \in \mathbb{Z}.$$

Für den absoluten Fehler $|r - gl(r)|$ erhalten wir die
- von $r = a \cdot g^b$ abhängige - Abschätzung

$$(5.4a) \quad |r - gl(r)| = |a - rd(a)| \cdot g^b \leq \frac{1}{2} \cdot g^{-t} \cdot g^b.$$

Für den relativen Fehler $\frac{|r - gl(r)|}{|r|}$, der den
absoluten Fehler in Relation zur betrachteten
Zahl $r \neq 0$ setzt, erhalten wir die - von $r = a \cdot g^b$
unabhängige - Abschätzung

$$(5.4b) \quad \frac{|r - gl(r)|}{|r|} \leq \frac{1}{2} \cdot g^{-t+b} \cdot \frac{g}{g^b} = \frac{1}{2} \cdot g^{-t+1} =: eps.$$

eps heißt die Maschinengenauigkeit.

Zu jedem $r \in \mathbb{R} \setminus \{0\}$ gibt es also ein $\varepsilon \in [-eps, eps]$ mit:

$$(5.4c) \quad gl(r) = r \cdot (1 + \varepsilon).$$

Beispiele für $g=10, t=8, s=2$

$$gl(246813795) = 246813800 = 0,2468138 \cdot 10^9,$$

$$gl(-\frac{1}{15}) = gl(-0,0\bar{6}) = -0,66666667 \cdot 10^{-1},$$

$$gl(e) = gl(2,7182818) = 0,27182818 \cdot 10^1.$$

Bemerkung 5.5. Rechenoperationen in der Gleitpunktdarstellung

Für $v_1, v_2 \in M$ setzen wir

$$(5.5 a) \quad v_1 \oplus v_2 := \text{gl}(v_1 + v_2),$$

$$(5.5 b) \quad v_1 \ominus v_2 := \text{gl}(v_1 - v_2),$$

$$(5.5 c) \quad v_1 \odot v_2 := \text{gl}(v_1 \cdot v_2),$$

$$(5.5 d) \quad v_1 \oslash v_2 := \text{gl}(v_1 : v_2) \text{ für } v_2 \neq 0.$$

Das seien jeweils die Ergebnisse, die von der Maschine berechnet werden.

Bemerkung 5.6:

Die Gleitpunktoperationen genügen nicht den üblichen Gesetzen der arithmetischen Operationen.

Zum Beispiel gilt:

$$(5.6) \quad \text{eps} = \frac{1}{2} \cdot g^{-A+1} = \min \{x \in M \mid 1 \oplus x > 1\}.$$

Für $0 \leq x < \text{eps}$, $x \in M$, ist also $1 \oplus x = 1$.

Gegenbeispiel zum Assoziativgesetz:

Sei $g = 10$, $A = 4$, $s = 2$ sowie

$$a = 0,1349 \cdot 10^{-1} = 0,01349;$$

$$b = 0,3368 \cdot 10^2 = 33,68;$$

$$c = -0,3367 \cdot 10^2 = -33,67.$$

Dann folgt:

$$a + b + c = 0,02349;$$

$$(a \oplus b) \oplus c = 33,69 \oplus (-33,67) = 0,02000;$$

$$a \oplus (b \oplus c) = 0,01349 \oplus 0,01 = 0,02349.$$

Konvention 5.6:

Für $v \in \mathbb{R} \setminus \{0\}$ mit gerundetem $v' \in M$ bezeichne

$$(5.6) \quad \varepsilon_v := \frac{|v - v'|}{|v|}$$

den relativen Fehler.

Satz 5.7, Fehlerfortpflanzung

- exakte Rechnung mit falschen Daten:

Es seien $v_1, v_2 \in \mathbb{R} \setminus \{0\}$, sind $v_1', v_2' \in M$ die gerundeten Werte, so gilt also:

$$(5.7) \quad \varepsilon_{v_1} \cdot |v_1| = |v_1 - v_1'|, \quad \varepsilon_{v_2} \cdot |v_2| = |v_2 - v_2'|.$$

Ferner nehmen wir an, dass $v_1' + v_2'$, $v_1' - v_2'$, $v_1' \cdot v_2'$ und $\frac{v_1'}{v_2'}$ exakt berechnet werden.

Dann gelten folgende Abschätzungen für die relativen Fehler $\varepsilon_{v_1 \circ v_2}$, wobei \circ für die 4 Grundrechenarten steht:

$$(5.7a) \quad \varepsilon_{v_1 + v_2} \leq \frac{1}{|v_1 + v_2|} \cdot (\varepsilon_{v_1} \cdot |v_1| + \varepsilon_{v_2} \cdot |v_2|),$$

falls $v_1 + v_2 \neq 0$;

$$(5.7b) \quad \varepsilon_{v_1 - v_2} \leq \frac{1}{|v_1 - v_2|} \cdot (\varepsilon_{v_1} \cdot |v_1| + \varepsilon_{v_2} \cdot |v_2|),$$

falls $v_1 \neq v_2$;

$$(5.7c) \quad \varepsilon_{v_1 \cdot v_2} \leq \varepsilon_{v_1} + \varepsilon_{v_2} + \varepsilon_{v_1} \cdot \varepsilon_{v_2} \approx \varepsilon_{v_1} + \varepsilon_{v_2} ;$$

$$(5.7d) \quad \varepsilon_{v_1 : v_2} \leq \varepsilon_{v_1} + (\varepsilon_{v_2} + \varepsilon_{v_1} \cdot \varepsilon_{v_2}) \cdot \frac{|v_2|}{|v_2'|} \approx \varepsilon_{v_1} + \varepsilon_{v_2} .$$

Beweis:

Verifikation von (5.7a):

Wir erhalten im Falle $v_1 + v_2 \neq 0$:

$$\begin{aligned} \varepsilon_{v_1+v_2} &= \frac{|(v_1+v_2) - (v_1'+v_2')|}{|v_1+v_2|} \leq \frac{|v_1-v_1'| + |v_2-v_2'|}{|v_1+v_2|} \\ &= \frac{1}{|v_1+v_2|} \cdot (\varepsilon_{v_1} \cdot |v_1| + \varepsilon_{v_2} \cdot |v_2|) \end{aligned}$$

Verifikation von (5.7b):

Ist $v_1 \neq v_2$, so folgt:

$$\begin{aligned} \varepsilon_{v_1-v_2} &= \frac{|(v_1-v_2) - (v_1'-v_2')|}{|v_1-v_2|} \leq \frac{|v_1-v_1'| + |v_2'-v_2|}{|v_1-v_2|} \\ &= \frac{1}{|v_1-v_2|} \cdot (\varepsilon_{v_1} \cdot |v_1| + \varepsilon_{v_2} \cdot |v_2|) \end{aligned}$$

Verifikation von (5.7c):

Wir erhalten:

$$\begin{aligned} \varepsilon_{v_1 \cdot v_2} &= \frac{|v_1 \cdot v_2 - v_1' \cdot v_2'|}{|v_1 \cdot v_2|} = \frac{|v_1 \cdot (v_2 - v_2') + (v_1 - v_1') \cdot v_2'|}{|v_1 \cdot v_2|} \\ &= \frac{|v_1 \cdot (v_2 - v_2') + (v_1 - v_1') \cdot v_2 - (v_1 - v_1') \cdot (v_2 - v_2')|}{|v_1 \cdot v_2|} \\ &\leq \frac{|v_1| \cdot |v_2 - v_2'|}{|v_1| \cdot |v_2|} + \frac{|v_1 - v_1'| \cdot |v_2|}{|v_1| \cdot |v_2|} + \frac{|v_1 - v_1'| \cdot |v_2 - v_2'|}{|v_1| \cdot |v_2|} \\ &= \varepsilon_{v_2} + \varepsilon_{v_1} + \varepsilon_{v_1} \cdot \varepsilon_{v_2} \end{aligned}$$

Verifikation von (5.7d):

Es folgt:

$$\begin{aligned} \varepsilon_{v_1: v_2} \cdot \left| \frac{v_1}{v_2} \right| &= \left| \frac{v_1}{v_2} - \frac{v_1'}{v_2'} \right| = \frac{|v_2' \cdot (v_1 - v_1') - v_1' \cdot (v_2 - v_2')|}{|v_2 \cdot v_2'|} \\ &= \frac{|v_2' \cdot (v_1 - v_1') - v_1' \cdot (v_2 - v_2') + (v_1 - v_1') \cdot (v_2 - v_2')|}{|v_2 \cdot v_2'|} \\ &\leq \frac{|v_2'| \cdot |v_1 - v_1'|}{|v_2| \cdot |v_2'|} + \frac{|v_1| \cdot |v_2 - v_2'|}{|v_2| \cdot |v_2'|} + \frac{|v_1 - v_1'| \cdot |v_2 - v_2'|}{|v_2| \cdot |v_2'|} \end{aligned}$$

Multiplikation mit $\frac{|v_2|}{|v_1|}$ liefert:

$$\begin{aligned} \varepsilon_{v_1: v_2} &\leq \frac{|v_1 - v_1'|}{|v_1|} + \frac{|v_2 - v_2'|}{|v_2'|} + \frac{|v_1 - v_1'|}{|v_1|} \cdot \frac{|v_2 - v_2'|}{|v_2'|} \\ &= \varepsilon_{v_1} + (\varepsilon_{v_2} + \varepsilon_{v_1} \cdot \varepsilon_{v_2}) \cdot \frac{|v_2|}{|v_2'|} \end{aligned}$$

□

Bemerkung 5.8, Analyse der Fehlerfortpflanzung - Auslöschung

Bei der Multiplikation und der Division verhält sich der relative Fehler unproblematisch. Bei Vernachlässigung sehr kleiner Produkte wie $\varepsilon_{v_1} \cdot \varepsilon_{v_2}$ addieren sich die relativen Fehler im ungünstigsten Fall.

Die Addition ist ebenfalls unproblematisch, wenn v_1 und v_2 gleiches Vorzeichen haben; gegebenenfalls ist $\varepsilon_{v_1 + v_2} \leq \varepsilon_{v_1} + \varepsilon_{v_2}$.

Auch die Subtraktion ist unproblematisch, wenn $v_1 - v_2$ nicht ungefähr 0 ist.

Gilt jedoch $v_1 \approx v_2$, so kann $\epsilon_{v_1 - v_2}$ gemäß (5.7b) groß werden!

Dieses Verhalten heißt Fehlerverstärkung durch Auslöschung.

Beispiel 5.9:

Wir betrachten das folgende lineare Gleichungssystem:

$$3x = 10 \quad \wedge \quad 5x - \frac{1}{1000} \cdot y = 16,66.$$

Die exakte Lösung berechnet sich wie folgt:

$$x = \frac{10}{3} = 3, \bar{3} ;$$

$$y = 1000 \cdot (16, \bar{6} - 16,66) = 1000 \cdot 0,00\bar{6} = 6, \bar{6} = 6 \frac{2}{3} .$$

Ein Rechner mit vierstelliger Gleitpunkt-darstellung liefert jedoch folgende gerundete Werte:

$$x' = 3,333 ;$$

$$y' = 1000 \cdot (5 \cdot 3,333 - 16,66) = 1000 \cdot (16,67 - 16,66) = 10 .$$

Während der Wert x' noch möglichst genau ist (bei Berücksichtigung von nur 4 Stellen), ergibt sich aber für y' ein mehr als nur ungenauer Wert!

Beispiel 5.10, Der Horner-Algorithmus

Sei $m \in \mathbb{N}$, seien $a_0, \dots, a_m \in \mathbb{R}$, und definiere die Polynom-Funktion $f: \mathbb{R} \rightarrow \mathbb{R}$ durch

$$\begin{aligned} f(x) &:= \sum_{j=0}^m a_j \cdot x^j = a_m \cdot x^m + \dots + a_1 \cdot x + a_0 \\ &= (\dots (a_m \cdot x + a_{m-1}) \cdot x + \dots + a_1) \cdot x + a_0. \end{aligned}$$

Für $x \in \mathbb{R}$ ist der Horner-Algorithmus zur Berechnung von $f(x)$ gegeben durch:

$$r_0 := a_m$$

for $j=1$ until m do

$$r_j := r_{j-1} \cdot x + a_{m-j}.$$

Für $m=2$ werden -bei exakten Anfangsdaten- folgende Werte geliefert:

Exakte Werte	Werte im Rechner
$r_0 := a_2$	$r_0 := a_2$
$r_1 := r_0 \cdot x + a_1$	$\tilde{r}_1 := r_0 \odot x \oplus a_1$
$r_2 := r_1 \cdot x + a_0$	$\tilde{r}_2 := \tilde{r}_1 \odot x \oplus a_0$

Aus (5.4c) folgt -für passende $\varepsilon_1, \varepsilon_2, \varepsilon_3, \varepsilon_4 \in [-\text{eps}, \text{eps}]$:

$$r_0 = a_2,$$

$$\begin{aligned} \tilde{r}_1 &= ((r_0 \cdot x) \cdot (1 + \varepsilon_1) + a_1) \cdot (1 + \varepsilon_2) \\ &= a_2 \cdot (1 + \varepsilon_1) \cdot (1 + \varepsilon_2) \cdot x + a_1 \cdot (1 + \varepsilon_2), \end{aligned}$$

$$\begin{aligned} \tilde{r}_2 &= ((\tilde{r}_1 \cdot x) \cdot (1 + \varepsilon_3) + a_0) \cdot (1 + \varepsilon_4) \\ &= a_2 \cdot (1 + \varepsilon_1) \cdot (1 + \varepsilon_2) \cdot (1 + \varepsilon_3) \cdot (1 + \varepsilon_4) \cdot x^2 \\ &\quad + a_1 \cdot (1 + \varepsilon_2) \cdot (1 + \varepsilon_3) \cdot (1 + \varepsilon_4) \cdot x + a_0 \cdot (1 + \varepsilon_4). \end{aligned}$$