

# 7. Rundungsfehler

## 7.1 Zahlenherstellung

Menge der Gleitkommazahlen  $M \subset \mathbb{Q}$  :

$$M := \{x \mid x = \sigma \cdot B^e \cdot m, 0 \leq m < B\}$$

mit

$\sigma = \pm 1$  : Vorzeichen

$B$  : Basis, meist  $B=2$

*Taschenrechner und Vorles'g auch  $B=10$*

$e$  : Exponent,  $e_0 \leq e \leq e_1$

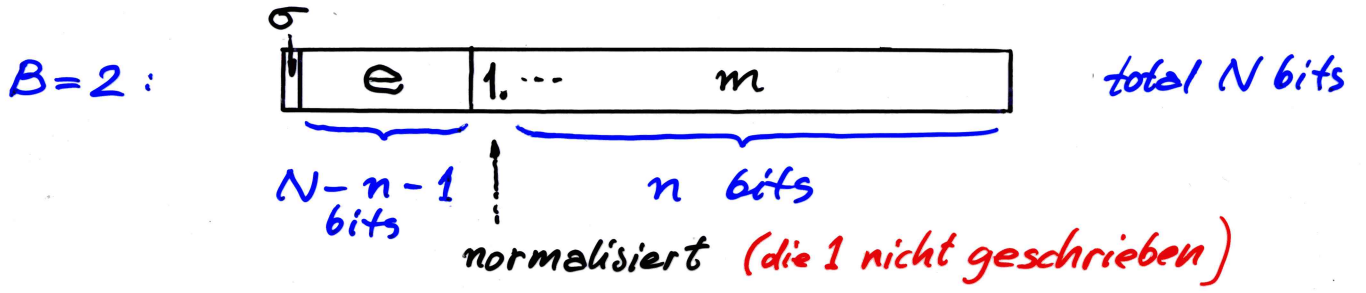
$$m = \sum_{k=0}^{n-1} d_k B^{-k}, d_k \in \{0, 1, \dots, B-1\} : \text{Mantisse}$$

*$d_1 > 0$  : normalisiert,  $d_1 = 0$  "denormal"*

$n$  : Mantissenlänge

### Bemerkungen

- $M$  ist endlich, enthält nur rationale Zahlen mit Nenner  $2^l$  (für  $B=2$ ).
- Gleitkommazahlen werden meist in Speicherbereiche fester Länge  $N$  gepackt, z. B.



Definition : "Maschinenepsilon" :

$$\frac{1}{2} \text{ eps} = \min_{1+\delta > 1} \delta$$

$$B=2: \text{eps} = 2^{-n}$$

eps = eine Einheit an der letzten Stelle der Mantisse !

Definition: "Überflusgrenze":

$$\begin{aligned} \text{realmax} &= \text{grösste darstellbare Zahl} \\ &= 2^{e_1+1} \quad (B=2) \end{aligned}$$

eps,  
realmax,  
realmin  
sind  
Matlab-  
Befehle!

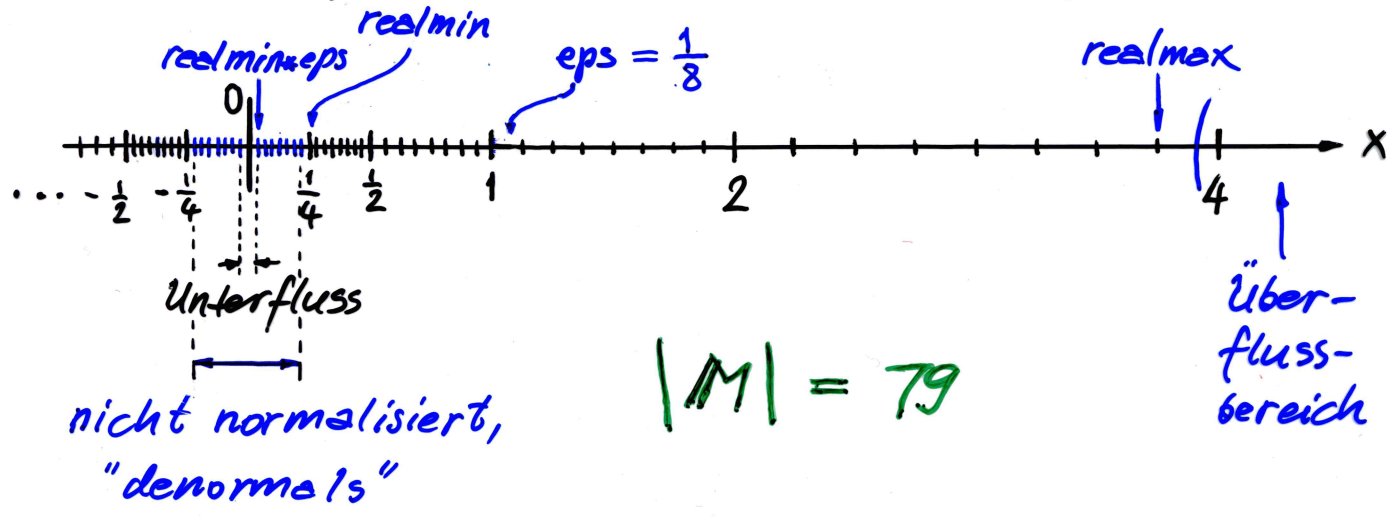
"Unterflusgrenze":

$$\begin{aligned} \text{realmin} &= \text{kleinste darstellbare} \\ &\quad \text{normalisierte positive Zahl} \\ &= 2^{e_0} \end{aligned}$$

Folgerung: Die kleinste positive Zahl  $\in M$  (auch nicht-normalisiert) ist  $\text{realmin} * \text{eps}$

Beispiele

(i)  $B=2, N=7, n=4, -2 \leq e \leq 1$



(ii) Matlab (= IEEE Standard)

$$B=2, N=64, n=52, -1022 \leq e \leq 1023 \quad (11 \text{ bits})$$

$$\text{eps} = 2^{-52} = 2.2204 \cdot 10^{-16}$$

$$\text{realmax} = 2^{1024} = 1.7977 \cdot 10^{308}$$

$$\text{realmin} = 2^{-1022} = 2.2251 \cdot 10^{-308}$$

$$\text{eps} * \text{realmin} = 2^{-1074} = 4.9407 \cdot 10^{-324}$$

# Die Rundungsabbildung $p(x)$

(112)

Sei  $x \in \mathbb{R}$

Def.  $p(x) \in \mathbb{M}$ : eine der (höchstens 2) zu  $x$  nächstgelegenen Gleitkommazahlen.

Bemerkung: In den Grenzfällen ( $x$  genau zwischen 2 Zahlen  $\in \mathbb{M}$ ) sollte gleich häufig auf- bzw. abgerundet werden.

SATZ: Ausserhalb des  $\ddot{U}/\ddot{U}$ -flussbereiches gilt

$$\left| \frac{p(x) - x}{x} \right| \leq \epsilon_{ps}$$

(Schranke für relativen Fehler bei  $\underset{\text{Anwendung von}}{p}$ )

Maschinenoperationen,  $\oplus$ ,  $\otimes$ ,  $--$

Erst exakt, dann  $p$  anwenden, z. B.

Def:  $a, b \in \mathbb{M}$ .  $a \oplus b := p(a + b)$  (z. B.)

Bemerkung: Maschinenoperationen erfüllen die bekannten Rechengesetze manchmal nicht.

Bsp. Basis  $B = 10$ , Mantissenlänge  $n = 3$

$$\begin{array}{l|l} a = 8.02 \\ b = 3.12 \\ c = 0.444 \end{array} \left. \begin{array}{l} \} s_1 := p(a+b) = 11.1 \\ \} p(s_1+c) = 11.5 \end{array} \right| \left. \begin{array}{l} p(s_2+a) = 11.6 \\ \} s_2 := p(c+b) = 3.56 \end{array} \right.$$

Exakt:  $s = a + b + c = 11.584$ ;  $p(s) = 11.6$

Reihenfolge der Terme in einer Summe kann relevant sein. Besser: von kleinen zu grossen Termen summieren.



## Über- und Unterflusseffekte

Vorbemerkung: Ist das Resultat einer Rechnung im Ü/Uflussbereich: **nichts zu machen!**

Ziel: Falls Endresultat  $\in \mathbb{M}$ , sollte die Rechnung funktionieren

Beispiel: Euklidische Norm in  $\mathbb{R}^2$

$$\|x\| = \sqrt{x_1^2 + x_2^2} \quad \text{für Vektor } x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

Schlechter Algorithmus: eingeschränkt auf den Bereich  $(\sqrt{\text{realmin}}, \sqrt{\text{realmax}})$ !

Man kann es besser machen:

$$x_l = |x_1|; \quad x_s = |x_2|;$$

$$\text{if } x_l < x_s, \quad x_l = |x_2|; \quad x_s = |x_1|; \quad \text{end};$$

$$\|x\| = x_l * \text{sqrt}(1 + (x_s/x_l)^2);$$

## 7.2. Fehlerfortpflanzung, Auslöschung

Sei  $a \in \mathbb{R}$ ,  $\tilde{a}$  eine Näherung für  $a$

Def: absoluter Fehler:  $\Delta a := \tilde{a} - a$

$$\text{relativer Fehler: } \delta a := \frac{\Delta a}{a} = \frac{\tilde{a} - a}{a} = \frac{\tilde{a}}{a} - 1, \quad a \neq 0$$

$$\Rightarrow \boxed{\tilde{a} = a + \Delta a = a(1 + \delta a)}$$

Fehlerfortpflanzung:

$$\Delta(a \pm b) = \Delta a \pm \Delta b$$

$$\delta(a \cdot b) \doteq \delta a + \delta b; \quad \delta\left(\frac{a}{b}\right) \doteq \delta a - \delta b, \quad b \neq 0$$

Beweis:

$$\begin{aligned} \tilde{a} &= a(1 + \delta a) \Rightarrow \tilde{a} \cdot \tilde{b} = a \cdot b (1 + \underbrace{\delta a + \delta b}_{\delta(a \cdot b)} + \cancel{\delta a \cdot \delta b}) \\ \tilde{b} &= b(1 + \delta b) \end{aligned}$$

Realistisch: Nur Fehler schranken bekannt: (114)

$$|da| \leq \varepsilon_1, |db| \leq \varepsilon_2 \Rightarrow |\Delta(a \pm b)| \leq \varepsilon_1 + \varepsilon_2$$

(folgt aus Dreiecksungleichung)

$$|da| \leq \varepsilon_1, |db| \leq \varepsilon_2 \Rightarrow |\delta(a \cdot b)| \leq \varepsilon_1 + \varepsilon_2$$

$$|\delta\left(\frac{a}{b}\right)| \leq \varepsilon_1 + \varepsilon_2$$

## Auslöschung

Subtraktion fast gleicher Zahlen

Beispiel:  $a = \frac{22}{7}$ ,  $b = \pi$  in  $M$  mit  $B=10$ ,  $n=5$

$$\Rightarrow \text{eps} = \frac{1}{2} B^{1-n} = 5 \cdot 10^{-5}$$

$$a = 3.1429 \quad da = 1.36 \cdot 10^{-5} < \text{eps}$$

$$b = 3.1416 \quad \delta b = 0.23 \cdot 10^{-5} < \text{eps}$$

$$a - b = 0.0013000$$

$$\frac{22}{7} - \pi = 0.0012645$$

Auffüllen mit Nullen in Basis  $B$

$$\delta(a-b) = 2.8 \cdot 10^{-2} \gg \text{eps}$$

grosser relativer Fehler!

## Messnahmen

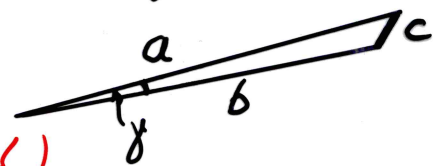
(i) Möglichst lange exakt rechnen

(ii) Mathematische Identitäten verwenden

$$\text{z. B. } \sqrt{x} - \sqrt{y} = \frac{x-y}{\sqrt{x} + \sqrt{y}} \quad \text{oder } 1 - \cos y = 2 \sin^2\left(\frac{y}{2}\right)$$

Bsp: cos-Satz im Dreieck

$$c = \sqrt{a^2 + b^2 - 2ab \cos \gamma} \quad (\text{Auslöschung!})$$
$$= \sqrt{(a-b)^2 + 4ab \sin^2\left(\frac{\gamma}{2}\right)} \quad \text{viel besser!}$$



(iii) Taylorreihen einsetzen

Bsp:  $\log(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} + \dots - \frac{(-x)^N}{N}$

Auslöschung für  $x \rightarrow 0$ !  $\leftarrow$   $\sinh(x) = \frac{e^x - e^{-x}}{2} = x + \frac{x^3}{3!} + \frac{x^5}{5!} + \dots$

Diagnose: Dieselbe Rechnung mit gestörten Daten wiederholen.

### 7.3. Numerische Differentiation

(durch Differenzenquotienten)

Z.B. gebraucht für Quasi-Newton-Verfahren

Sei  $f$  an der Stelle  $x$  differenzierbar, und sei ein numerischer Auswertungsalgorithmus für  $f$  vorhanden:

$$x \longrightarrow \boxed{f} \longrightarrow f(x) \cdot (1 + \delta), \quad |\delta| \leq \mu \cdot \text{eps}$$

$\mu > 0$ , nicht gross, z.B.  $\mu \approx 3$

$$f(x) = \tilde{f} - \delta \cdot f(x) \iff \tilde{f}$$

(\*)

Approximation der Ableitung  $f'(x)$  durch den Differenzenquotienten mit Schritt  $h > 0$ :

$$f'(x) = \frac{f(x+h) - f(x)}{h} - \frac{h}{2} f''(x) - \dots$$

$$\tilde{f}_1 = (1 + \delta_1) f(x+h)$$

$$\tilde{f}_0 = (1 + \delta_0) f(x)$$

$$\uparrow \frac{1}{h} \{ \tilde{f}_1 - \tilde{f}_0 - \delta_1 f(x+h) + \delta_0 f(x) \} - \frac{h}{2} f''(x) - \dots$$

mit (\*)

Fehler:

$$\text{err} := \frac{1}{h} (\tilde{f}_1 - \tilde{f}_0) - f'(x) = \frac{\delta_1}{h} f(x+h) - \frac{\delta_0}{h} f(x) + \frac{h}{2} f''(x) - \dots$$

Abschätzung mit Dreiecks-Ungleichung:

$$|\text{err}| = \left| \frac{1}{h} (\tilde{f}_1 - \tilde{f}_0) - f'(x) \right| \leq \underbrace{\frac{2\mu \cdot \text{eps}}{h} |f(x)| + \frac{h}{2} |f''(x)|}$$

! min durch gute Wahl von  $h$

Resultat:  $h = c \sqrt{\text{eps}}$  mit  $c = 2 \sqrt{\mu |f(x)| / |f''(x)|}$

damit:  $|\text{err}| \leq \sqrt{\text{eps}} \cdot \sqrt{\mu |f(x) \cdot f''(x)|}$

Wähle den Differentiationsschritt  $h$  in der Grössenordnung von  $\sqrt{\text{eps}}$ .  $f'(x)$  hat dann auch die Genauigkeit  $\sqrt{\text{eps}}$ .