

Übung 1: Maschinenzahlen

Für $b, t \in \mathbb{N}$, $e_{\min}, e_{\max} \in \mathbb{Z} \cup \{\pm\infty\}$ heißt

$$\mathbb{F}(b, t, e_{\min}, e_{\max}) := \left\{ \overset{\text{Vorzeichen}}{\sigma} \cdot \left(\overset{\text{Mantisse}}{0, z_1 \dots z_t} \right) \cdot \overset{\text{Exponent}}{b^e} \mid \sigma \in \{1, -1\}, \dots \right. \\ \left. z_1 \neq 0, z_1, \dots, z_t \in \{0, \dots, b-1\}, e \in \{e_{\min}, \dots, e_{\max}\} \right\} \cup \{0\}$$

die Menge der Fließkommazahlen zur Basis b und Mantissenlänge t .

Ist b gerade so ist der Rundungsoperator $\text{rd}: \mathbb{R} \rightarrow \mathbb{F}$ für

$\mathbb{R} \ni x = \sigma \cdot (0, z_1 \dots z_t z_{t+1} \dots) \cdot b^e$ durch

$$\text{rd}(x) = \begin{cases} \sigma \cdot (0, z_1 \dots z_t) \cdot b^e & , z_{t+1} \leq \frac{b}{2} - 1 \\ \sigma \cdot (0, z_1 \dots (z_t + 1)) \cdot b^e & , z_{t+1} \geq \frac{b}{2} \end{cases}$$

definiert.

Für die Beispiele betrachten wir im folgenden

$$\mathbb{F} = \mathbb{F}(10, 4, -\infty, \infty).$$

Damit ist z. B.

$$\text{rd}(\pi) = \text{rd}\left(+ (0, 31415 \dots) \cdot 10^{01}\right) = 0,3142 \cdot 10^1$$

und

$$\text{rd}(\sqrt{2}) = \text{rd}\left(0,14142 \dots \cdot 10^1\right) = 0,1414 \cdot 10^1$$

Es heißt

$$\text{macheps} := \varepsilon^* := \text{eps} := \frac{1}{2} b^{(-t+1)}$$

die Maschinengenauigkeit zu $F(b, t, e_{\min}, e_{\max})$.

Satz: Es ist

$$\text{eps} = \inf \{ \delta > 0 \mid \text{rd}(1+\delta) > 1 \}$$

und für $x \in \mathbb{R}$ gilt (falls $e_{\min} = -\infty, e_{\max} = \infty$)

$$\left[\text{eps} \geq \left| \frac{\text{rd}(x) - x}{x} \right| \quad (=: |\varepsilon|) \right] \quad (\Rightarrow)$$

$$\left[\exists \varepsilon \in \mathbb{R}, |\varepsilon| \leq \text{eps}, \text{ so dass } \varepsilon = \frac{\text{rd}(x) - x}{x} \right] \quad (\Leftarrow)$$

$$\left[\exists \varepsilon \in \mathbb{R}, |\varepsilon| \leq \text{eps}, \text{ mit } \text{rd}(x) = x(1+\varepsilon) \right]$$

Beweis: VL \square

Die Maschinen- / Ersatz-Operationen $\tilde{+}, \tilde{-},$

$\tilde{\times}, \tilde{\div} : F \times F \rightarrow F$ sind für $x, y \in F$

durch

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

$$\stackrel{\text{Satz}}{=} (x o y)(1+\varepsilon)$$

$$o \in \{+, -, \times, \div\}$$

gegeben. Weiter bezeichne ε \leftarrow Rundungsfehler

~~folgendes Problem:~~

$$\text{sqrt}(x) := \text{rd}(\sqrt{x}) \stackrel{\text{Satz}}{=} \sqrt{x} (1 + \varepsilon), \quad \text{Rundungsfehler}$$

das maschinelle Wurzelziehen.

Bei verketteten Ausdrücken können sich die Rundungsfehler hochschaukeln / akkumulieren.

Um diesen Effekt zu analysieren betrachten wir die

Fehleranalyse

am Beispiel der äquivalenten Ausdrücke

$$\text{a) } \sqrt{1+x} - 1 = \frac{(\sqrt{1+x} - 1)(\sqrt{1+x} + 1)}{\sqrt{1+x} + 1} =$$

$$\text{b) } \frac{x}{\sqrt{1+x} + 1}, \quad \text{für } |x| \approx 0.$$

Dabei kann man zwei Methoden anwenden:

Methode 1: am Beispiel a)

Ist $f: \mathbb{R} \rightarrow \mathbb{R}$ genügend oft differenzierbar, so heißt

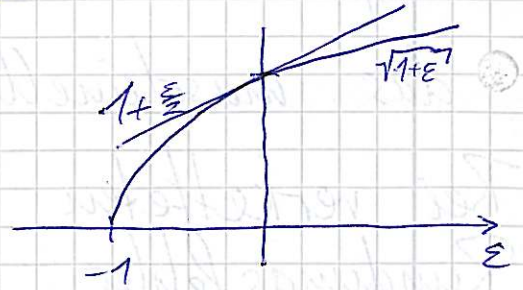
$$f(x+\varepsilon) = f(x) + \varepsilon f'(x) + \frac{\varepsilon^2}{2} f''(x) + \dots$$

die Taylorentwicklung von f . Ist nun $|\varepsilon|$ klein, so ist $|\varepsilon|^2$ sehr klein und daher nennt man

$$f(x+\varepsilon) \approx f(x) + \varepsilon f'(x)$$

Approximation erster Ordnung. Da $\frac{d}{dx} \sqrt{x} = \frac{1}{2\sqrt{x}}$ ist z. B.

$$\sqrt{1+\varepsilon} \approx \sqrt{1} + \varepsilon \frac{1}{2\sqrt{1}} = 1 + \frac{\varepsilon}{2} \quad (*)$$



Nun wird auf der Maschine der Ausdruck a) ausgewertet als

$$\text{sqrt}(1+x) \approx 1$$

$$= \text{sqrt}((1+x)(1+\varepsilon_1)) \approx 1$$

$$= \left[\sqrt{(1+x)(1+\varepsilon_1)} (1+\varepsilon_2) - 1 \right] (1+\varepsilon_3)$$

$$= \left[\sqrt{1+x} \sqrt{1+\varepsilon_1} (1+\varepsilon_2) - 1 \right] (1+\varepsilon_3)$$

$$\stackrel{(*)}{\approx} \left[\sqrt{1+x} \left(1 + \frac{\varepsilon_1}{2} \right) (1+\varepsilon_2) - 1 \right] (1+\varepsilon_3)$$

$$= \left[\sqrt{1+x} \left(1 + \frac{\varepsilon_1}{2} + \varepsilon_2 + \frac{\varepsilon_1 \varepsilon_2}{2} \right) - 1 \right] (1+\varepsilon_3)$$

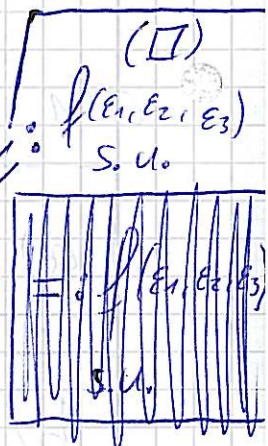
sehr klein

$$\approx \left[\sqrt{1+x} \left(1 + \frac{\varepsilon_1}{2} + \varepsilon_2 \right) - 1 \right] (1+\varepsilon_3)$$

$$= \sqrt{1+x} \left(1 + \frac{\varepsilon_1}{2} + \varepsilon_2 \right) - 1 + \sqrt{1+x} \left(\varepsilon_3 + \frac{\varepsilon_1 \varepsilon_3}{2} + \varepsilon_2 \varepsilon_3 \right) - \varepsilon_3$$

sehr klein

$$\approx \sqrt{1+x} - 1 + \sqrt{1+x} \left(\frac{\varepsilon_1}{2} + \varepsilon_2 + \varepsilon_3 \right) - \varepsilon_3$$



Damit ist der Fehler

$$\left| \sqrt{1+x} - 1 - (\sqrt{1+x} - 1) \right|$$

$$\leq \left| \sqrt{1+x} \right| \cdot \left| \frac{\varepsilon_1}{2} + \varepsilon_2 \right| + \left| \sqrt{1+x} - 1 \right| \cdot \left| \varepsilon_3 \right|$$

$$\leq \sqrt{1+x} \cdot \frac{3}{2} \text{eps} + \left| \sqrt{1+x} - 1 \right| \cdot \text{eps}$$

$$\stackrel{(*)}{\approx} \left(1 + \frac{x}{2}\right) \frac{3}{2} \text{eps} + \left|1 + \frac{x}{2} - 1\right| \cdot \text{eps}$$

x klein

$$\leq \left(\frac{3}{2} + \frac{5}{4}|x|\right) \cdot \text{eps}$$

Methode 2: am Beispiel ~~a)~~ a)

Ist $f: \mathbb{R}^n \rightarrow \mathbb{R}$ genügend oft differenzierbar,
so heißt

$$f(\varepsilon_1, \dots, \varepsilon_n) = f(0, \dots, 0) + Df|_0 \cdot \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} + \frac{1}{2} \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}^T D^2 f|_0 \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} + \dots$$

die mehrdimensionale Taylorentwicklung von f .
Analog heißt

$$f(\underbrace{\varepsilon_1, \dots, \varepsilon_n}_{=: \varepsilon}) \approx f(0, \dots, 0) + Df|_0 \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

die Approximation erster Ordnung.

~~Man wird auf der Maschine der Funktion~~
~~ausgewertet als~~

Mit $f(\varepsilon_1, \varepsilon_2, \varepsilon_3)$ wie oben (\square) ist nun $f(0,0,0)$

der exakte Ausdruck a) und damit ist der Fehler

$$\begin{aligned} & |f(\varepsilon_1, \varepsilon_2, \varepsilon_3) - f(0,0,0)| \\ & \approx \left| Df|_0 \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \end{bmatrix} \right| \quad (\Delta) \\ & \leq \|Df|_0\|_1 \cdot \left\| \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \end{bmatrix} \right\|_\infty \leq \|Df|_0\|_1 \cdot \text{eps} \end{aligned}$$

Mit

$$Df|_{(\varepsilon_1, \varepsilon_2, \varepsilon_3)} = \left[\sqrt{1+x} \frac{1}{2\sqrt{1+\varepsilon_1}} (1+\varepsilon_2)(1+\varepsilon_3), \sqrt{(1+x)(1+\varepsilon_1)} (1+\varepsilon_3), \sqrt{(1+x)(1+\varepsilon_1)} (1+\varepsilon_2) - 1 \right]$$

$$\Rightarrow Df|_{(0,0,0)} = \left[\frac{1}{2} \sqrt{1+x}, \sqrt{1+x}, \sqrt{1+x} - 1 \right]$$

$$\Rightarrow \|Df|_0\|_1 = \frac{3}{2} \sqrt{1+x} + |\sqrt{1+x} - 1|$$

Insgesamt also

$$|\text{sqrt}(1+\tilde{x}) - 1 - (\sqrt{1+x} - 1)|$$

$$\stackrel{(\square)}{=} |f(\varepsilon_1, \varepsilon_2, \varepsilon_3) - f(0,0,0)|$$

(\Delta)

$$\leq \|Df|_0\|_1 \cdot \text{eps} = \sqrt{1+x} \cdot \frac{3}{2} \text{eps} + |\sqrt{1+x} - 1| \cdot \text{eps}$$

(*)
 \tilde{x}
x klein

$$\left(1 + \frac{x}{2}\right) \frac{3}{2} \text{eps} + \left|1 + \frac{x}{2} - 1\right| \text{eps} \leq \left(\frac{3}{2} + \frac{5}{4}|x|\right) \cdot \text{eps}$$

das gleiche Ergebnis wie bei Methode 1.

Methode 2: am Beispiel b)

Fuß der Maschine ist Ausdruck b)

$$x \stackrel{\approx}{=} (\text{sqrt}(1+x) \mp 1)$$

$$= \frac{x(1+\varepsilon_1)}{[\sqrt{(1+x)(1+\varepsilon_2)}(1+\varepsilon_3)+1](1+\varepsilon_4)} =: f(\varepsilon_1, \varepsilon_2, \varepsilon_3, \varepsilon_4)$$

$$\Rightarrow Df|_0 = \left[\frac{x}{\sqrt{1+x}+1}, -\frac{x}{2} \frac{\sqrt{1+x}}{(\sqrt{1+x}+1)^2}, -x \frac{\sqrt{1+x}}{(\sqrt{1+x}+1)^2}, -\frac{x}{\sqrt{1+x}+1} \right]$$

$$\Rightarrow \|Df|_0\|_1 = 2 \frac{|x|}{\sqrt{1+x}+1} + \frac{3}{2} \frac{|x|\sqrt{1+x}}{(\sqrt{1+x}+1)^2}$$

$$\leq 2|x| + \frac{3}{2}|x|\sqrt{1+x} \leq 2 \text{ für } |x| \leq 1$$
$$\leq 5|x|$$

$$\Rightarrow \left| x \stackrel{\approx}{=} (\text{sqrt}(1+x) \mp 1) - \left(\frac{x}{\sqrt{1+x}+1} \right) \right|$$

$$\leq \|Df|_0\|_1 \cdot \text{eps} \leq 5|x| \cdot \text{eps}$$