

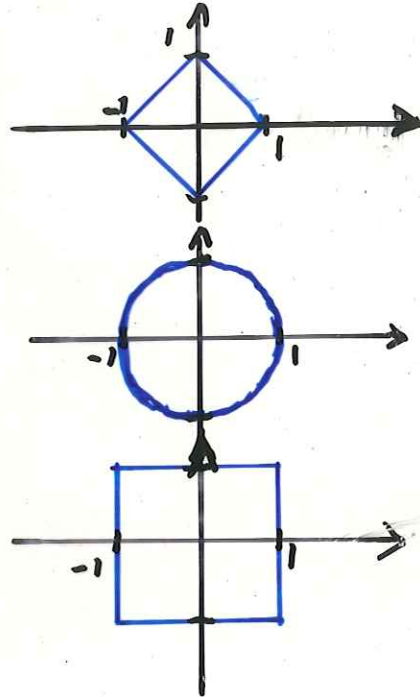
Vektor-Normen : $x \in \mathbb{R}^n$

$$\{x \in \mathbb{R}^2 : \|x\| \leq 1\}$$

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

$$\|x\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2}$$

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$



Matrix-Normen : $A \in \mathbb{R}^{n \times n}$ mit Einträgen a_{ij}

$$\|A\|_1 = \max_{j=1, \dots, n} \sum_{i=1}^n |a_{ij}| \quad \Rightarrow \text{max. Spaltensumme}$$

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)} \quad \Rightarrow \text{Spektralnorm}$$

λ_{\max} : max EW von $A^T A$

$$\|A\|_\infty = \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| \quad \Rightarrow \text{max. Zeilensumme}$$

Landau Symbol \mathcal{O}

Betrachte zwei Funktionen $g, h: \mathbb{R} \rightarrow \mathbb{R}$. Wir verwenden die Notation

$$g(x) = \mathcal{O}(h(x)) \quad (x \rightarrow x_0)$$

wenn es Konstanten $C > 0$ und $\delta > 0$ gibt, so dafs

$$|g(x)| \leq C |h(x)| \quad \forall x \text{ mit } |x - x_0| < \delta$$

gilt.

* Anschauliche Bedeutung

g wächst nicht wesentlich schneller als h

* Mathematische Definition

$$0 \leq \limsup_{x \rightarrow x_0} \left| \frac{g(x)}{h(x)} \right| < \infty$$

Beispiele

I. $g(x) = \mathcal{O}(1)$ g wird durch konstanten Wert beschränkt

II. $g(x) = \mathcal{O}(x)$ g wächst ungefähr auf das Doppelte, wenn sich x verdoppelt

III. $g(x) = \mathcal{O}(x^2)$ g wächst ungefähr auf das Vierfache, wenn sich x verdoppelt

IV. $g(x) = \mathcal{O}(2^x)$ g wächst ungefähr auf das Doppelte, wenn sich x um 1 erhöht.

B2.2

Taylorentwicklung 1. Ordnung für $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$f(\tilde{x}) \doteq f(x) + \nabla f(x)^T (\tilde{x} - x)$$

$$= f(x) + \sum_{j=1}^n \frac{\partial f(x)}{\partial x_j} (\tilde{x}_j - x_j)$$

$$\Rightarrow f(\tilde{x}) - f(x) \doteq \sum_{j=1}^n \frac{\partial f(x)}{\partial x_j} (\tilde{x}_j - x_j)$$

$$| \cdot \frac{1}{f(x)}$$

$$\Rightarrow \frac{f(\tilde{x}) - f(x)}{f(x)} \doteq \sum_{j=1}^n \frac{\partial f(x)}{\partial x_j} \cdot \frac{1}{f(x)} (\tilde{x}_j - x_j) \cdot \frac{x_j}{x_j}$$

$$\Rightarrow \frac{f(\tilde{x}) - f(x)}{f(x)} \doteq \sum_{j=1}^n \frac{\partial f(x)}{\partial x_j} \cdot \frac{x_j}{f(x)} \cdot \frac{(\tilde{x}_j - x_j)}{x_j}$$

$\underbrace{\hspace{10em}}$
rel. Fehler
in der Ausgabe δ_y

$\underbrace{\hspace{10em}}$
Fehler-
verstärkung $\phi_j(x)$

$\underbrace{\hspace{10em}}$
rel. Fehler in der
Eingabe in Eintrag x_j

B2.2

Wir haben

$$\underbrace{\left| \frac{f(\tilde{x}) - f(x)}{f(x)} \right|}_{\delta_y} = \left| \sum_{j=1}^n \phi_j(x) \underbrace{\frac{\tilde{x}_j - x_j}{x_j}}_{\delta_{x_j}} \right|$$

$$\Rightarrow |\delta_y| = \left| \sum_{j=1}^n \phi_j(x) \delta_{x_j} \right|$$

$$\leq \sum_{j=1}^n |\phi_j(x)| \cdot |\delta_{x_j}|$$

$$= |\phi_1(x)| \cdot |\delta_{x_1}| + |\phi_2(x)| \cdot |\delta_{x_2}| + \dots + |\phi_n(x)| \cdot |\delta_{x_n}|$$

$$\leq \max_{j=1, \dots, n} |\phi_j(x)| \cdot |\delta_{x_1}| + \max_{j=1, \dots, n} |\phi_j(x)| \cdot |\delta_{x_2}| + \dots$$

$$\dots + \max_{j=1, \dots, n} |\phi_j(x)| \cdot |\delta_{x_n}|$$

$$= \max_{j=1, \dots, n} |\phi_j(x)| \cdot \sum_{j=1}^n |\delta_{x_j}|$$

$$|\delta_y| \leq \underbrace{\chi_{\infty}^{\infty}(x)}_{\text{max}} \cdot \underbrace{\|\delta x\|_1}_{\text{1-Norm}}$$

rel. Eingabefehler in
der 1-Norm

B2.2

Andere Möglichkeit:

$$|\delta y| \leq \sum_{j=1}^n |\phi_j(x)| \cdot |\delta x_j|$$

$$= |\phi_1(x)| \cdot |\delta x_1| + |\phi_2(x)| \cdot |\delta x_2| + \dots + |\phi_n(x)| \cdot |\delta x_n|$$

$$\leq |\phi_1(x)| \cdot \max_{j=1, \dots, n} |\delta x_j| + |\phi_2(x)| \cdot \max_{j=1, \dots, n} |\delta x_j| + \dots$$

$$\dots + |\phi_n(x)| \cdot \max_{j=1, \dots, n} |\delta x_j|$$

$$= \sum_{j=1}^n |\phi_j(x)| \cdot \max_{j=1, \dots, n} |\delta x_j|$$

$$|\delta y| \leq \underbrace{\sum_{j=1}^n |\phi_j(x)|}_{\chi_{\text{rel}}^1(x)} \cdot \underbrace{\max_{j=1, \dots, n} |\delta x_j|}_{\|\delta x\|_\infty}$$

rel. Eingabefehler in
der ∞ -Norm

Kondition: $f: X \mapsto Y$

I) $f: \mathbb{R} \mapsto \mathbb{R}$ (Eingabe: Skalar, Ausgabe: Skalar)

$$\overbrace{\left| \frac{f(\tilde{x}) - f(x)}{f(x)} \right|}^{\delta_y} \leq \kappa_{\text{rel}}(x) \overbrace{\left| \frac{\tilde{x} - x}{x} \right|}^{\delta_x}$$

$$\text{mit } \kappa_{\text{rel}}(x) = \left| f'(x) \frac{x}{f(x)} \right|$$

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

$$x = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$$

$$\left| \frac{f(\tilde{x}) - f(x)}{f(x)} \right| \leq \kappa_{\text{rel}}(x) \sum_{j=1}^n \left| \frac{\tilde{x}_j - x_j}{x_j} \right|$$

$$\text{mit } \kappa_{\text{rel}}(x) = \kappa_{\text{rel}}^{\infty}(x) = \max_j |\phi_j(x)| \quad \text{und} \quad \phi_j(x) = \frac{\partial f(x)}{\partial x_j} \cdot \frac{x_j}{f(x)}$$

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

$$\text{linear: } y = f(x) = A^{-1} x, \quad A \in \mathbb{R}^{n \times n}$$

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

$\kappa(A)$: Konditionszahl der Matrix A

Rundungsfehler und Gleitpunktarithmetik

Normalisierte Gleitpunkt Darstellung: $x \in \mathbb{R}$

$$x = \pm 0, d_1 d_2 \dots d_m \cdot b^e$$

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

- $b \in \mathbb{N} \setminus \{1\}$ Basis
- Exponent e : $r \leq e \leq R$
- Mantisse $p = \pm 0, d_1 d_2 \dots d_m$
- Mantissenlänge m
- Normalisierung : $d_1 \neq 0$

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

- betragsmäßig kleinste ($\neq 0$) Zahl : $x_{\min} = b^{r-1}$
- betragsmäßig größte Zahl : $x_{\max} = (1 - b^{-m}) b^R$

Reduktionsabbildung

$$fl : \mathbb{R} \mapsto M(b, m, r, R)$$

(Relative) Maschinengenauigkeit ϵ :

$$\epsilon = \frac{b^{1-m}}{2}$$

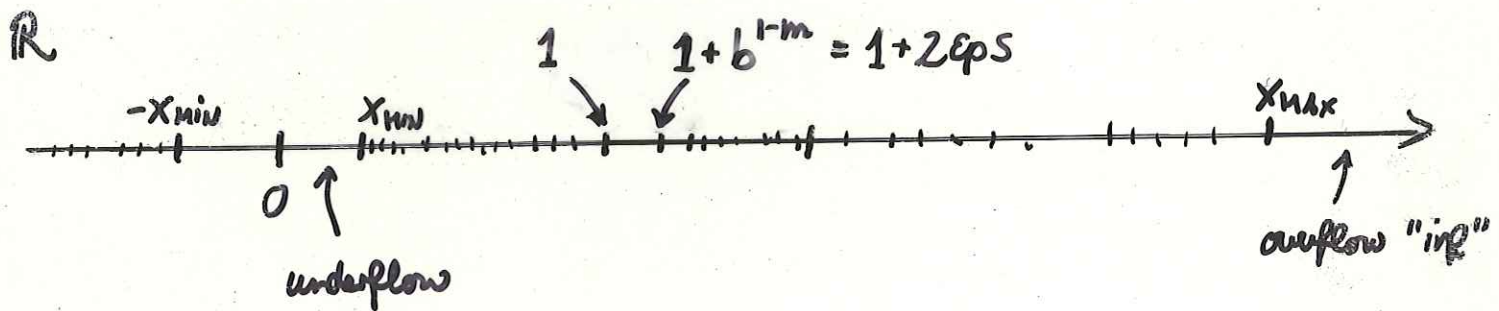
(Relative) Rundungsfehler ϵ :

$$|\epsilon| = \left| \frac{fl(x) - x}{x} \right| \leq \frac{b^{1-m}}{2} = \epsilon$$

Für die Reduktionsabbildung gilt:

$$fl(x) = x(1 + \epsilon)$$

wobei Rundungsfehler $|\epsilon| \leq \epsilon$



Intervall ($b=2$)

$$[1, 2]: 1, 1+b^{1-m}, 1+2 \cdot b^{1-m}, \dots, 2 \hat{=} 1, 1+2\epsilon, 1+4\epsilon, \dots, 2$$

$$[2, 4]: 2, 2+2b^{1-m}, 2+4 \cdot b^{1-m}, \dots, 4$$

$$\vdots$$

$$[2^i, 2^{i+1}]: \text{ " [1, 2] " } \cdot 2^i$$

Anmerkung zu "relativer Maschinengenaugtheit"

Relative Maschinengenaugtheit

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

eps entspricht der halben Distanz von 1 zur nächstgrößeren exakt darstellbaren Zahl nach 1, d.h. die exakt darstellbaren Maschinenzahlen in $M(b, m, r, R)$ sind

$$1, 1 + 2 \text{eps}, 1 + 4 \text{eps}, \dots, 2$$

$$\text{oder } 1, 1 + b^{1-m}, 1 + 2b^{1-m}, \dots, 2$$

ACHTUNG:

- Wenn nicht gerundet sondern abgeschnitten wird, ist die relative Maschinengenaugtheit

$$\text{eps}_c = b^{1-m} \quad (\text{"chopping"})$$

- In manchen Büchern (und in MATLAB) ist die relative Maschinengenaugtheit definiert als

$$\overline{\text{eps}} = b^{1-m}$$

und damit genau dem Abstand von 1 zur nächstgrößeren Maschinenzahl. Um Verwirrung zu vermeiden, definieren wir dieses "MATLAB"-eps als $\overline{\text{eps}}$ mit "Strich".

Float GUI

Parameter (default):

entspricht Notation D & R

$$e_{\min} = -4$$



$$r = e_{\min} + 1 = -3$$

$$e_{\max} = 2$$



$$R = e_{\max} + 1 = 3$$

$$t = 3$$



$$m = t + 1 = 4$$

Damit ergibt sich:

- relative Maschinengenauigkeit

$$\overline{\text{eps}} = b^{1-m} = 2^{-3} = \frac{1}{8}$$

- betragsmäßig kleinste Zahl

$$x_{\min} = b^{r-1} = 2^{-4} = \frac{1}{16}$$

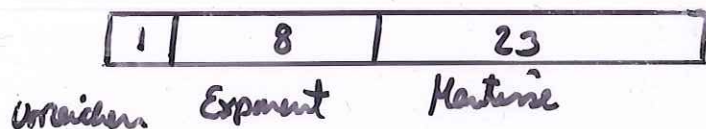
- betragsmäßig größte Zahl

$$x_{\max} = (1 - b^{-m}) b^R = (1 - 2^{-4}) 2^3 = 8 - \frac{1}{2}$$

ACHTUNG: effektive Mantisenlänge um 1 größer als angegeben, da $d_1 = 1$ nicht gespeichert wird.

I) Single precision (32 bit)

Matlab



$$\text{eps} = \frac{2^{(1-24)}}{2} \approx 5.96 \cdot 10^{-8}$$

$$\overline{\text{eps}} = 2^{(1-24)} \approx 1.19 \cdot 10^{-7}$$

$$X_{\min} = 2^{-126} \approx 1.175 \cdot 10^{-38}$$

$$X_{\max} = 2^{127} (2 - 2^{-23}) \approx 3.402 \cdot 10^{38}$$

$\gg \text{eps ('single')}$

$\gg \text{realmin ('single')}$

$\gg \text{realmax ('single')}$

II) double precision (64 bit)



$$\text{eps} = \frac{2^{(1-53)}}{2} \approx 1.11 \cdot 10^{-16}$$

$$\overline{\text{eps}} = 2^{(1-53)} \approx 2.22 \cdot 10^{-16}$$

$$X_{\min} = 2^{-1022} \approx 2.225 \cdot 10^{-308}$$

$$X_{\max} = 2^{1023} (2 - 2^{-52}) \approx 1.798 \cdot 10^{308}$$

$\gg \text{eps}$

$\gg \text{realmin}$

$\gg \text{realmax}$

Gleitpunktoperationen \oplus

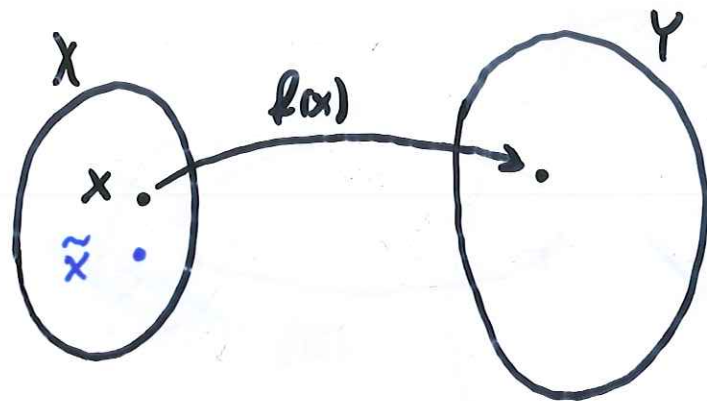
Für $\nabla \in \{+, -, \times, \div\}$ gilt

$$x \oplus y = fl(x \nabla y) \quad \text{für } x, y \in M(b, m, r, R)$$

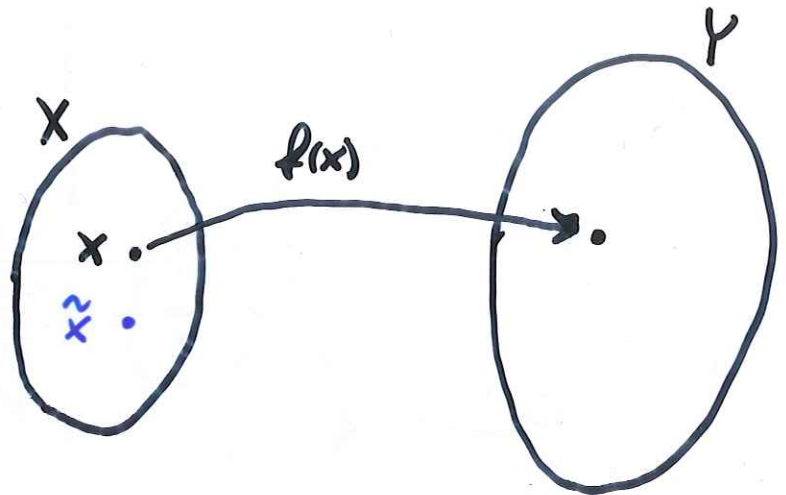
$$\Rightarrow x \oplus y = (x \nabla y)(1 + \varepsilon) \quad \text{für } x, y \in M(b, m, r, R)$$

Aufpassen:

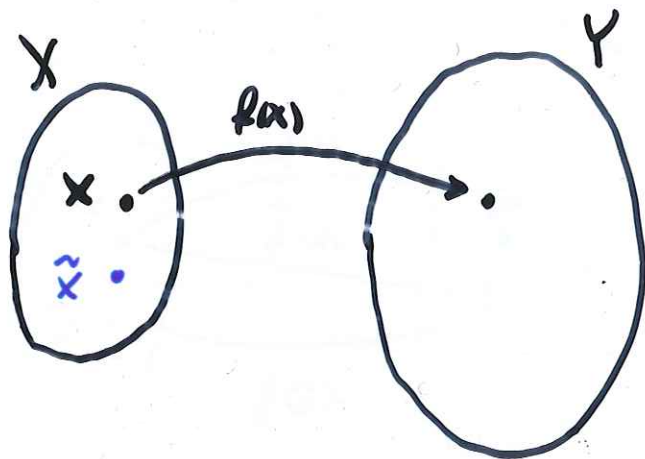
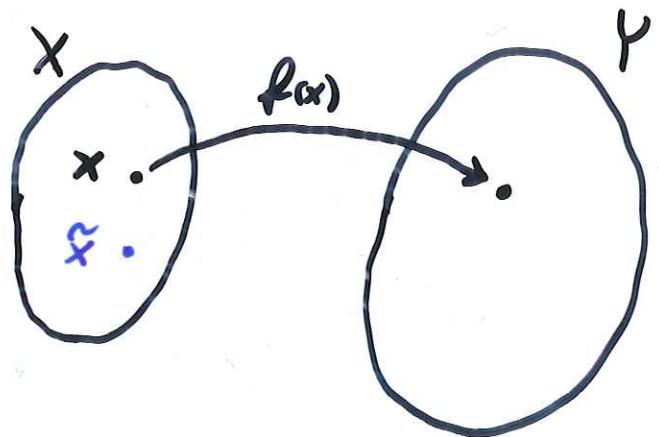
- Assoziativgesetz nicht mehr gültig
- Distributivgesetz nicht mehr gültig
- Gefahr von Auslöschung bei $\nabla \in \{+, -\}$

Skizze Kondition

gut konditioniert



schlecht konditioniert

Skizze Stabilitätrückwärts stabil &
gut konditioniertrückwärts stabil &
schlecht konditioniert