

## Bemerkung 2.7

Lineare Abb. zwischen endlich-dim.  
Vektorräumen können durch  
Matrizen beschrieben.

Ziel: Charakterisierung der Abb. ver-  
halten von  $\mathcal{L}$

## Def 2.8 (Operatornorm)

$$\mathcal{L}: (X, \|\cdot\|_X) \rightarrow (Y, \|\cdot\|_Y)$$

$$\|\mathcal{L}\| := \sup_{\|x\|_X = 1} \|\mathcal{L}(x)\|_Y$$

$$= \sup_{x \neq 0} \frac{\| \mathcal{L}(x) \|_Y}{\|x\|_X} =$$

$$\sup_{\|x\|_X \leq 1} \| \mathcal{L}(x) \|_Y$$

“Norm” für die Verformung des  
Einheitskugels unter der Abb.  
 $\mathcal{L}$ .

## Bemerkung 2.9

Ein linearer Operator ist beschränkt genau dann wenn stetig.

$$\|Y(x_2) - Y(x_1)\|_Y =$$

$$\|Y(x_2 - x_1)\|_Y \leq \|Y\| \cdot$$

$$\|x_2 - x_1\|_X \quad (\text{Lipschitz-stetig})$$

$Y$  beschränkt, falls  $\|Y\|$  endlich

Für  $X = \mathbb{R}^n$ ,  $Y = \mathbb{R}^m$  und die  
L-repräsentierende Matrix

$B \in \mathbb{R}^{m \times n}$  berechnet man:

$$\|B\|_{\infty} = \max_{i=1, \dots, m} \sum_{k=1}^n |b_{ik}|$$

(Zeilensummennorm)

$$\|B\|_1 = \max_{k=1, \dots, n} \sum_{i=1}^m |b_{ki}|$$

(Spaltensummennorm)

$$A \in \mathbb{R}^{n \times n}$$

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$$

(Spectralnorm)

Relative Kondition:

$$\text{Ziel: schätze } \frac{\|Y(\bar{x}) - Y(x)\|_p}{\|Y(x)\|_p}$$

$$\text{durch } \frac{\|\bar{x} - x\|}{\|x\|_p}$$

Satz 2.10 (rel. Konditionszahl)

Für injektives  $\mathcal{L}$  gilt

$$\frac{\|\mathcal{L}(\bar{x}) - \mathcal{L}(x)\|_Y}{\|\mathcal{L}(x)\|_Y} \leq \kappa(\mathcal{L}) \cdot \frac{\|\bar{x} - x\|_X}{\|x\|_X}$$

mit  $\kappa(\mathcal{L}) = \frac{\sup_{\|x\|_X=1} \|\mathcal{L}(x)\|_Y}{\inf_{\|x\|_X=1} \|\mathcal{L}(x)\|_Y}$

Ist  $\mathcal{L}$  bijektiv gilt

$$|\mathcal{K}(\mathcal{L})| = \|\mathcal{L}\| \cdot \|\mathcal{L}^{-1}\|$$

Beispiel 2.11 (Kondition einer Basis)

Sei  $V$   $n$ -dimensionaler Vektorraum

und  $\underline{\Phi} = \{\phi_1, \dots, \phi_n\}$  eine Basis

von  $V$ . (z.B. Polynomraum)

$$\mathcal{L}: \mathbb{R}^n \rightarrow V, \quad \mathcal{L}|a| := \sum_{j=1}^n a_j \phi_j$$

Da  $\underline{B}$  Basis ist  $\mathcal{L}$  bijektiv.

$$\mathcal{K}(\underline{B}) := \mathcal{K}(\mathcal{L})$$

(Konditionszahl der Basis)

### Bemerkung 2.12

- a)  $\mathcal{K}(\mathcal{L})$  hängt von der Wahl der Normen ab.
- b) Falls  $\mathcal{L}$  beschränkt, ist  $\mathcal{K}(\mathcal{L})$  auch für injektives (und nicht surjektives  $\mathcal{L}$ ) definiert.

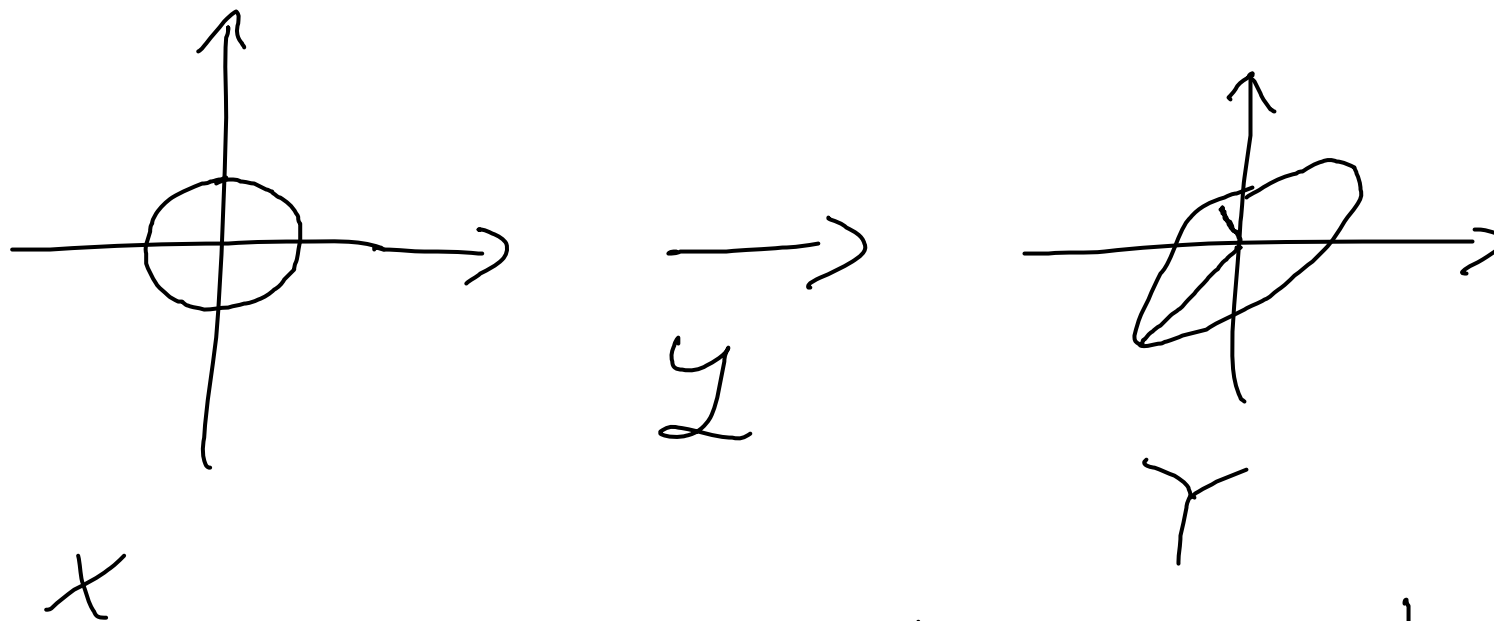


da ist  $\|Y(x)\|_Y \neq 0$ , denn  
 $\|x\|_X = 1$

$Y(x) = 0$  gilt nur für  $x = 0$ .

Anschauliche Bedeutung:

$K(Y)$ : Verhältnis max. Dehnung  
zu max. Stauchung der  
Einheitskugel unter der  
Abb.  $Y$  gemessen in der  
Bildnorm  $\|\cdot\|_Y$



Ass. 2.13 (Condition)

## Beispiel 2.14 (Matrixkondition)

$A \in \mathbb{R}^{n \times n}$  invertierbar

$$\kappa_p(A) := \|A\|_p \|A^{-1}\|_p \quad (2.14)$$

$$p = 1, \dots, \infty$$

Schnittpunkt von Geraden in der Ebene (s. Bsp. 2.1 b)

$$\underbrace{\begin{pmatrix} 3 & 1.001 \\ 6 & 1.997 \end{pmatrix}}_A \underbrace{\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}}_x = \underbrace{\begin{pmatrix} 1.997 \\ 4.003 \end{pmatrix}}_b$$

Lsg. ist  $x = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ . Betrachte

„Datenstörung“ in  $b$ :

$$\bar{b} := \begin{pmatrix} 2.002 \\ 4 \end{pmatrix} \rightarrow \bar{x} := A^{-1} \bar{b}$$

Nun berechnet

$$A^{-1} = -\frac{1}{0.015} \begin{pmatrix} 1.997 & -1.001 \\ -6 & 3 \end{pmatrix}$$

$$\text{und } \bar{x} = \begin{pmatrix} 0.4004 \\ 0.8 \end{pmatrix}$$

in der  $\| \cdot \|_{\infty}$  - Norm gilt:

$$\frac{\|\bar{b} - b\|_{\infty}}{\|b\|_{\infty}} \approx 7.5 \cdot 10^{-4}$$

(Datenstörung)

$$\frac{\|\bar{x} - x\|_{\infty}}{\|x\|_{\infty}} = 1.8$$

(Resultatstörung)

sowie  $\kappa(A) = \|A\|_{\infty} \cdot \|A^{-1}\|_{\infty} = 4738.2$

b) Ist  $X = \mathbb{R}^n$ ,  $Y = \mathbb{R}$  und ein nicht konstantes  $f: X \rightarrow Y$ , das differenzierbar ist gilt:

$$\left| \frac{f(\bar{x}) - f(x)}{f(x)} \right| \leq J_{f, \infty}(x) \cdot$$

$$\sum_{j=1}^n \left| \frac{\bar{x}_j - x_j}{x_j} \right|$$

$$\left\| \frac{\bar{x} - x}{x} \right\|_1$$

$$\|K_{\infty}(x)\| = \max_{j=1, \dots, n} \left| \frac{\partial f}{\partial x_j} \cdot \frac{x_j}{f(x_j)} \right|$$

(das sieht man mit Hilfe von  
Taylorentwicklung 1. Ordnung)

Kondition der Addition:

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto f(x) := x_1 + x_2$$

$$\frac{\partial f}{\partial x_j} \cdot \frac{x_j}{f(x)} = \frac{x_j}{x_1 + x_2}, \quad j = 1, 2$$

$$J_{\infty}(x) = \max \left\{ \left| \frac{x_1}{x_1 + x_2} \right|, \left| \frac{x_2}{x_1 + x_2} \right| \right\} \quad (2.15)$$

haben  $x_1, x_2$  gleiches Vorzeichen:

$$J_{\infty}(x) \leq 1$$

Falls  $x_1 \approx -x_2$ , so kann

$J_{\infty}(x)$  sehr groß werden!

(S. Auslöschung in Kap. 3)



## Kap. 3 Numerische Zahlendarstellung und Gleitpunktarithmetik

Darstellung einer reellen Zahl

$$x = \pm \left( \sum_{j=1}^{\infty} d_j b^{-j} \right) \cdot b^e \quad (3.1)$$

wähle  $e$  so, dass  $d_1 \neq 0$

Rechnerdarstellung ist endlich  
(Maschinenzahlen)

→ normalisierte Gleitpunktendarstellung  
 $x = f \cdot b^e$

$b$ : Grundzahl (Basis) des gewählten  
Zahlensystems (z. B. dezimal, binär)

$r \leq e \leq R$  (Exponent)

$f: \pm 0.d_1d_2 \dots d_m$

$d_j \in \{0, 1, \dots, b-1\}$

$f$ : Mantisse

$m$ : Mantissenlänge

Rechner verarbeitet reelle Zahlen über  
Reduktionsabbildung

$$f_l: \mathbb{R} \rightarrow M(b, m, t, \mathbb{R})$$

(Raschierungszahlen)

Rundung:

$$f_l(x) := \begin{cases} \left( \sum_{j=1}^m d_j b^j \right) \cdot b^e & d_{m+1} < \frac{b}{2} \\ \left( \sum_{j=1}^m d_j b^j + b^{-m} \right) \cdot b^e & d_{m+1} \geq \frac{b}{2} \end{cases}$$

→ kleinstmögliche Zahl

$$x_{\min} = 0.1000\dots \cdot b^t = b^{t-1}$$

→ größtmögliche Zahl

$$x_{\max} = 0.a a \dots a \cdot b^R = \frac{1 - b^{-m}}{b-1} \cdot b^R$$

$a = b - 1$

$$|x| < x_{\min} \rightarrow fl(x) = 0$$

$$|x| > x_{\max} \rightarrow fl(x) = \infty \text{ (overflow)}$$

Rundungsfehler

$$|fl(x) - x| \leq \frac{b^{-m}}{2} b^e \text{ (absolut)}$$

$$(3.2) \left| \frac{fl(x) - x}{x} \right| \leq \frac{\frac{b^{-m}}{2} \cdot b^e}{b^{-n} \cdot b^e} = \frac{b^{n-m}}{2}$$

$$rps := \frac{b^{n-m}}{2} \text{ (Maximalgenauigkeit (relativ))}$$

$$\epsilon \rho S = \inf \left\{ \delta > 0 \mid f(x) (1 + \delta) > 1 \right\}$$

für  $|\epsilon| \leq \epsilon \rho S$  gilt

$$\boxed{|f(x)| = x(1 + \epsilon)} \quad (3.3)$$

### § 3.1 Schrittplan der arithmetik

Algorithmus: „Folge arithmetischer Operationen“. Verküpfung von Maschinenzahlen liefert nicht notwendig wieder eine Maschinenzahl.

Beispiel:  $b=10, m=3$

$$0.346 \cdot 10^2 + 0.785 \cdot 10^2 =$$

$$0.1131 \cdot 10^3 \neq 0.113 \cdot 10^3$$

Ersetze die üblichen arithmet.

Operationen durch Gleitpunktoperationen (Pseudoarithmetik).  $\textcircled{V}$

$$\forall \in \{+, -, \times, \div\}.$$

Führe über Mantissenkette hinaus  
weiter (stehe mit  $t_1$  genaue Rechnung)

nach Exponenten ausgleichen,  
Normalisierung, Ründung, so dass

$$x \otimes y = \text{fl}(x \nabla y) \quad (3.4)$$

wg. (3.3) Annahme:

$$x \otimes y = (x \nabla y) / (1 + \varepsilon)$$

Eigenschaften:

-(3.4) nicht mehr gültig für eine  
Sequenz von Operationen

- Assoziativität und Distributivität gehen verloren

Bsp.  $b=10, m=3$

Nachdem Zahlen  $x = 0.653 \cdot 10^4$

$$y = 0.100 \cdot 10^1$$

$$z = 0.400 \cdot 10^1$$

$(\in M(b=10, m=3))$

$$(x+y) + z = (y+z) + x = 6535 \text{ (exakt)}$$

$$x \oplus y = 0.653 \cdot 10^4$$

$$(x \oplus y) \oplus z = 0.653 \cdot 10^4$$



abz:

$$y \oplus z = 0.500 \cdot 10^2$$

$$(y \oplus z) \oplus x = 0.660 \cdot 10^4 = \text{fl}(x + y + z)$$

(3.5): "lokaler" Fehler einer  
Operation ist im Rahmen  
von  $\epsilon_{PS}$

Frage: Fehlerpropagation im  
Algorithmus

(→) Stabilität des Algorithmus

Algorithmus heißt stabil, falls  
 die Fehler in der „Größenordnung“  
 der unvermeidbaren Fehler liegen  
 ( $\rightarrow$  Kondition des Problems)

Beispiel 3.6 (Ausköschung) /  $b=10$   
 $m=3$

$$x = 0.73563 \quad y = 0.73441 \quad \kappa_{ps} = \frac{1}{2} \cdot 10^{-2}$$

$$x - y = 0.00122$$

$$\bar{x} = fl(x) = 0.736 \quad \bar{y} = fl(y) = 0.734$$

$$|\delta_x| = 0.500 \cdot 10^{-3}$$

$$|\delta_y| = 0.560 \cdot 10^{-3}$$

$$\text{Fehler} \left| \frac{(\bar{x} - \bar{y}) - (x - y)}{x - y} \right|$$

$$= \left| \frac{0.002 - 0.00122}{0.00122} \right| = 0.64 \Rightarrow \delta_x, \delta_y$$

Stimmen die führenden + Ziffern  
bei Subtraktion zweier Zahlen  
überein, steigt der Fehler im  
Resultat mit dem Faktor  $b^t$

⇒ Subtraktion „ähnlich“ großer  
Zahlen ist instabil!