

# **Vorlesung: Numerik 1 für Ingenieure**

Version 28.10.14

Michael Karow

## **5. Vorlesung**

*Thema: Fehlerformeln für lineare Gleichungssysteme*

*Kondition einer Matrix, Matrixnormen*

## Ungenauere Eingangsdaten bei linearen Gleichungssystemen

Aufgabe: Es soll die Gleichung  $Ax = b$  gelöst werden.

Die Ausgangsdaten  $A, b$  können aber z.B. aus folgenden Gründen ungenau sein:

- $A, b$  sind Messwerte und daher nicht genau bekannt.
- $A, b$  sind Ergebnisse früherer (fehlerhafter) Rechnungen
- Die Einträge von  $A, b$  sind keine Maschinenzahlen und können nicht exakt im Rechner abgespeichert werden.

Außerdem machen Algorithmen zur Lösung von  $Ax = b$  Fehler, die man als Fehler in den Daten  $A, b$  interpretieren kann. Macht man z.B. eine  $LR$ -Zerlegung

$$A = LR,$$

dann bekommt man statt  $L, R$  numerisch  $\tilde{L}, \tilde{R}$  heraus. Das Produkt ist

$$\tilde{L}\tilde{R} = A + \Delta A.$$

Beim Vorwärts-Rückwärtseinsetzen löst man also im besten Fall das Gleichungssystem

$$(A + \Delta A)x = b.$$

**Problem:** Wie wirken sich ungenaue Eingangsdaten auf die Lösung eines linearen Gleichungssystems  $Ax = b$  aus?

**Situation:**

Exakte Daten:

Lösung:

$$(A, b) \longmapsto x = A^{-1} b$$

Ungenaue Daten:

Lösung:

$$(\tilde{A}, \tilde{b}) \longmapsto \tilde{x} = \tilde{A}^{-1} \tilde{b}$$

**Frage:** Wie stark unterscheidet sich  $\tilde{x}$  von  $x$  ?

Sei  $Ax = b$  und  $\tilde{A}\tilde{x} = \tilde{b}$ , wobei  $\det(A) \neq 0 \neq \det(\tilde{A})$ .

Dann gelten die folgenden Fehlerabschätzungen.

Für den **absoluten Fehler**:

$$\|\tilde{x} - x\| \leq (\|\tilde{A}^{-1}\| \|x\|) \|\tilde{A} - A\| + \|\tilde{A}^{-1}\| \|\tilde{b} - b\| \quad (*)$$

Für den **relativen Fehler**:

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \|A\| \|\tilde{A}^{-1}\| \left( \frac{\|\tilde{A} - A\|}{\|A\|} + \frac{\|\tilde{b} - b\|}{\|b\|} \right) \quad (**)$$

$$\leq \frac{\text{cond}(A)}{1 - \frac{\|\tilde{A} - A\|}{\|A\|} \text{cond}(A)} \left( \frac{\|\tilde{A} - A\|}{\|A\|} + \frac{\|\tilde{b} - b\|}{\|b\|} \right) \quad (***)$$

Dabei ist  $\text{cond}(A) = \|A\| \|A^{-1}\|$  die **Konditionszahl** von  $A$ .

Die Ungleichungen (\*) und (\*\*) gelten stets, wenn  $A$  und  $\tilde{A}$  invertierbar sind, und wenn für die zugrunde liegenden Normen die Ungleichung  $\|My\| \leq \|M\| \|y\|$  für alle Matrizen  $M$  und alle Vektoren  $y$  erfüllt ist.

Die Ungleichung (\*\*\*) gilt nur unter der zusätzlichen Bedingung, dass  $\frac{\|\tilde{A} - A\|}{\|A\|} \text{cond}(A) < 1$ .

**Ziel der Vorlesung:** Diese Fehlerformeln erklären und 'beweisen'.

# Vektor-Normen

Norm=Maß für die Größe der Einträge in einem Vektor.

## Definition:

Eine Norm auf  $\mathbb{R}^n$  ist eine Funktion  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$  mit den folgenden Eigenschaften:

- 1)  $\|x\| > 0$  für alle  $x \in \mathbb{R}^n$ ,  $x \neq 0$ ,
- 2)  $\|\lambda x\| = |\lambda| \|x\|$  für alle  $x \in \mathbb{R}^n$ ,  $\lambda \in \mathbb{R}$
- 3)  $\|x + y\| \leq \|x\| + \|y\|$  für alle  $x, y \in \mathbb{R}^n$  (Dreiecksungleichung)

## Häufig verwendete Normen:

- a) **Euklidische Norm:**  $\|x\|_2 := \sqrt{\sum_{k=1}^n x_k^2}$
- b) **Summen-Norm:**  $\|x\|_1 := \sum_{k=1}^n |x_k|$
- c) **Maximum-Norm:**  $\|x\|_\infty := \max_{k=1}^n |x_k|$ .

Diese Normen gehören zur Familie der **Hölder- $p$ -Normen**:

$$\|x\|_p := \left( \sum_{k=1}^n |x_k|^p \right)^{1/p}, \quad 1 \leq p < \infty.$$

Es ist

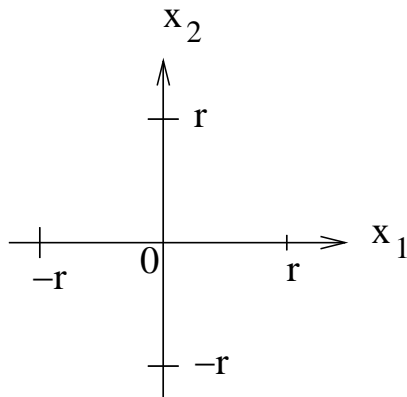
$$\|x\|_\infty = \lim_{p \rightarrow \infty} \|x\|_p.$$

# Die Normsphären

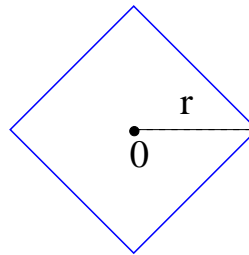
Die Sphäre zur Norm  $\|\cdot\|_\alpha$  und zum Radius  $r > 0$  um den Nullpunkt ist die Menge aller Vektoren  $x \in \mathbb{R}^n$  mit  $\|x\|_\alpha = r$ . Formal:

$$S_\alpha(r) := \{ x \in \mathbb{R}^n \mid \|x\|_\alpha = r \}.$$

**Anschauung:** Für den Fall  $n = 2$  hat man

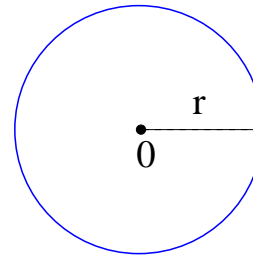


$\alpha = 1$



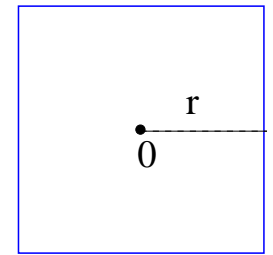
$S_1(r)$

$\alpha = 2$



$S_2(r)$

$\alpha = \infty$



$S_\infty(r)$

Im Fall  $n = 3$  ist

- $S_1(r)$  eine Oktaederoberfläche
- $S_2(r)$  eine Kugeloberfläche (Sphäre im engeren Sinne)
- $S_\infty(r)$  eine Würfeloberfläche.

## Bemerkung: Äquivalenz von Normen

**Definition:** Zwei Normen  $\|\cdot\|$  und  $|\cdot|$  heissen äquivalent, falls es Konstanten  $c_1, c_2 > 0$  gibt, so dass für alle  $x \in \mathbb{R}^n$ ,

$$c_1 \|x\| \leq |x| \leq c_2 \|x\|.$$

Ist dies der Fall, dann folgt für alle  $x \in \mathbb{R}^n$

$$(1/c_2) |x| \leq \|x\| \leq (1/c_1) |x|.$$

## Wichtigste Anwendung der Norm-Äquivalenz:

Sei  $x_k$  eine Folge in  $\mathbb{R}^n$ , die bzgl. der Norm  $\|\cdot\|$  gegen den Punkt  $x_0$  konvergiert, d.h.

$$\lim_{k \rightarrow \infty} \|x_k - x_0\| = 0.$$

Dann konvergiert diese Folge auch bzgl. jeder zu  $\|\cdot\|$  äquivalenten Norm  $|\cdot|$  gegen  $x_0$ , d.h. man hat

$$\lim_{k \rightarrow \infty} |x_k - x_0| = 0.$$

**Satz:** Alle Normen auf  $\mathbb{R}^n$  sind äquivalent.

**Beispiele:** Für alle  $x \in \mathbb{R}^n$  gilt

$$\|x\|_\infty \leq \|x\|_1 \leq n \|x\|_\infty, \quad \|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty.$$

# Induzierte Normen von Matrizen

## Definition:

Sei  $A \in \mathbb{R}^{m \times n}$ , und seien  $\|\cdot\|_\alpha$  und  $\|\cdot\|_\beta$  Normen auf  $\mathbb{R}^n$  bzw.  $\mathbb{R}^m$ . Dann ist die Zahl

$$\|A\|_{\alpha,\beta} := \max_{\|x\|_\alpha=1} \|Ax\|_\beta = \max_{x \neq 0} \frac{\|Ax\|_\beta}{\|x\|_\alpha}$$

die durch  $\|\cdot\|_\alpha$  und  $\|\cdot\|_\beta$  **induzierte Matrixnorm**.

## Alternative Definition:

$\|A\|_{\alpha,\beta}$  ist die kleinste Zahl  $c \geq 0$ , so daß für alle  $x \in \mathbb{R}^n$  gilt  $\|Ax\|_\beta \leq c \|x\|_\alpha$ .

## Interpretation:

$\|A\|_{\alpha,\beta}$  ist der Faktor, um den ein Vektor  $x$  bei Multiplikation mit  $A$  maximal gestreckt werden kann.

## Eigenschaften:

- 1)  $\|A\|_{\alpha,\beta} > 0$  für alle  $A \in \mathbb{R}^{m \times n}$ ,  $A \neq 0$ .
- 2)  $\|\lambda A\|_{\alpha,\beta} = |\lambda| \|A\|_{\alpha,\beta}$  für alle  $A \in \mathbb{R}^{m \times n}$ ,  $\lambda \in \mathbb{R}$
- 3)  $\|A_1 + A_2\|_{\alpha,\beta} \leq \|A_1\|_{\alpha,\beta} + \|A_2\|_{\alpha,\beta}$  für alle  $A_1, A_2 \in \mathbb{R}^{m \times n}$ .
- 4)  $\|Ax\|_\beta \leq \|A\|_{\alpha,\beta} \|x\|_\alpha$  für alle  $x \in \mathbb{R}^n$ . (**Besondere Eigenschaft**)

Ist  $\|\cdot\|_\alpha = \|\cdot\|_\beta$ , dann schreibt man kurz:  $\|A\|_\alpha := \|A\|_{\alpha,\alpha}$ .

Wenn klar ist, welche Norm gemeint ist, dann lässt man den Index  $\alpha$  weg.



## Berechnung von $\|A\|_\infty$

Zu einer Matrix  $A = [a_{ik}] \in \mathbb{R}^{m \times n}$  definieren wir die Zeilensummen:

$$Z_i(A) = \sum_{k=1}^n |a_{ik}|, \quad i = 1, \dots, m.$$

**Satz:** Es gilt stets

$$\|A\|_\infty = \max_{\|x\|_\infty=1} \|Ax\|_\infty = \max_i Z_i(A).$$

**Beweis:** Sei  $y = Ax$ . Dann hat  $y$  die Komponenten

$$y_i = a_{i1} x_1 + a_{i2} x_2 + \dots + a_{ik} x_k + \dots + a_{in} x_n,$$

und es ist

$$\|Ax\|_\infty = \|y\|_\infty = \max\{|y_1|, |y_2|, \dots, |y_m|\}.$$

Wenn  $\|x\|_\infty = 1$ , dann ist  $|x_k| \leq 1$  für alle  $k$  und man hat die Abschätzung:

$$\begin{aligned} |y_i| &= |a_{i1} x_1 + a_{i2} x_2 + \dots + a_{ik} x_k + \dots + a_{in} x_n| \\ &\leq |a_{i1} x_1| + |a_{i2} x_2| + \dots + |a_{ik} x_k| + \dots + |a_{in} x_n| \quad (*) \\ &\leq |a_{i1}| + |a_{i2}| + \dots + |a_{ik}| + \dots + |a_{in}| \\ &= Z_i(A). \end{aligned}$$

Daraus folgt schon mal, dass  $\|A\|_\infty \leq \max_i Z_i(A)$ .

Den maximal möglichen Wert von  $|y_i|$  unter der Nebenbedingung  $\|x\|_\infty = 1$  bekommt man offensichtlich dann, wenn  $x_k \in \{-1, 1\}$  und  $x_k$  dasselbe Vorzeichen hat wie  $a_{ik}$ ,  $k = 1, \dots, n$ . Dann ist  $a_{ik} x_k = |a_{ik}|$  und folglich  $|y_i| = y_i = Z_i(A)$ . Macht man dies für eine Zeile  $i_0$  mit maximaler Zeilensumme, dann bekommt man  $\|Ax\|_\infty = Z_{i_0}(A) = \max_i Z_i(A)$ .

Für die Größe  $\max_{\|x\|_\alpha=1} \|Ax\|_\alpha$  hatten wir die Bezeichnung  $\|A\|_\alpha$  eingeführt:

$$\|A\|_\alpha := \max_{\|x\|_\alpha=1} \|Ax\|_\alpha = \max_{x \neq 0} \frac{\|Ax\|_\alpha}{\|x\|_\alpha}$$

Für die ebenfalls wichtige Größe  $\min_{\|x\|_\alpha=1} \|Ax\|_\alpha$  gibt es keine (allgemein anerkannte) Notation. Dies hat folgenden Grund.

**Satz:** Sei  $A \in \mathbb{R}^{n \times n}$  invertierbar (nicht singulär). Dann gilt:

$$\frac{1}{\|A^{-1}\|_\alpha} = \min_{\|x\|_\alpha=1} \|Ax\|_\alpha = \min_{x \neq 0} \frac{\|Ax\|_\alpha}{\|x\|_\alpha}.$$

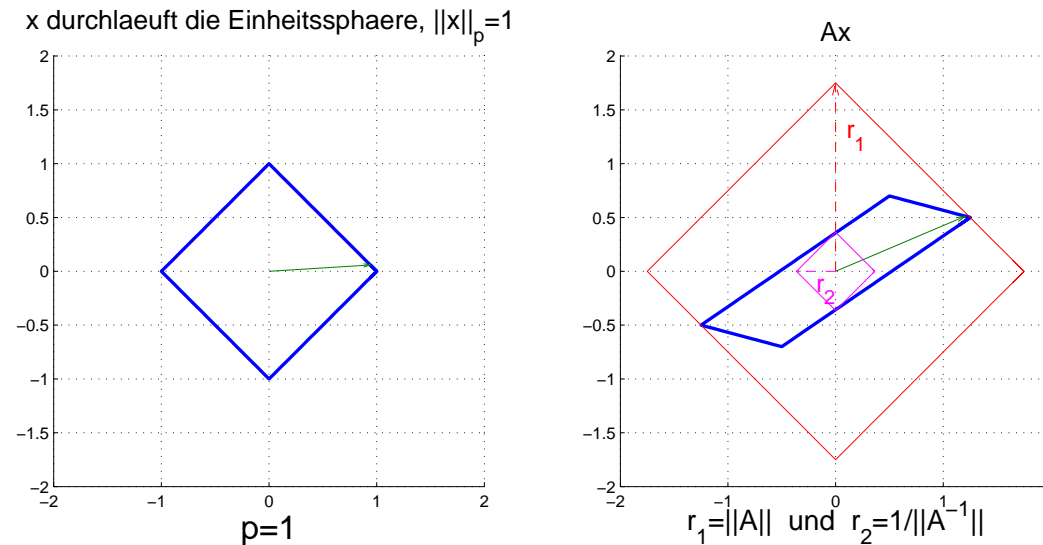
**Beweis:**

$$\begin{aligned} \|A^{-1}\|_\alpha &= \max_{y \neq 0} \frac{\|A^{-1}y\|_\alpha}{\|y\|_\alpha} \\ &= \max_{x \neq 0} \frac{\|A^{-1}(Ax)\|_\alpha}{\|Ax\|_\alpha} && \text{setze } y = Ax \\ &= \max_{x \neq 0} \frac{\|x\|_\alpha}{\|Ax\|_\alpha} \\ &= \frac{1}{\min_{x \neq 0} \frac{\|Ax\|_\alpha}{\|x\|_\alpha}} \end{aligned}$$

Beim letzten Schritt wurde folgende leicht einsehbare Tatsache benutzt:

Sei  $M$  eine Menge positiver Zahlen und  $M^{-1}$  die Menge der Kehrwerte aller Zahlen aus  $M$ . Dann ist  $\max M = \frac{1}{\min M^{-1}}$ .

## Bilder zur Veranschaulichung von $\|A\|$ und $1/\|A^{-1}\|$



Erklärung:

Die dicke blaue Kurve im linken Bild ist die Sphäre

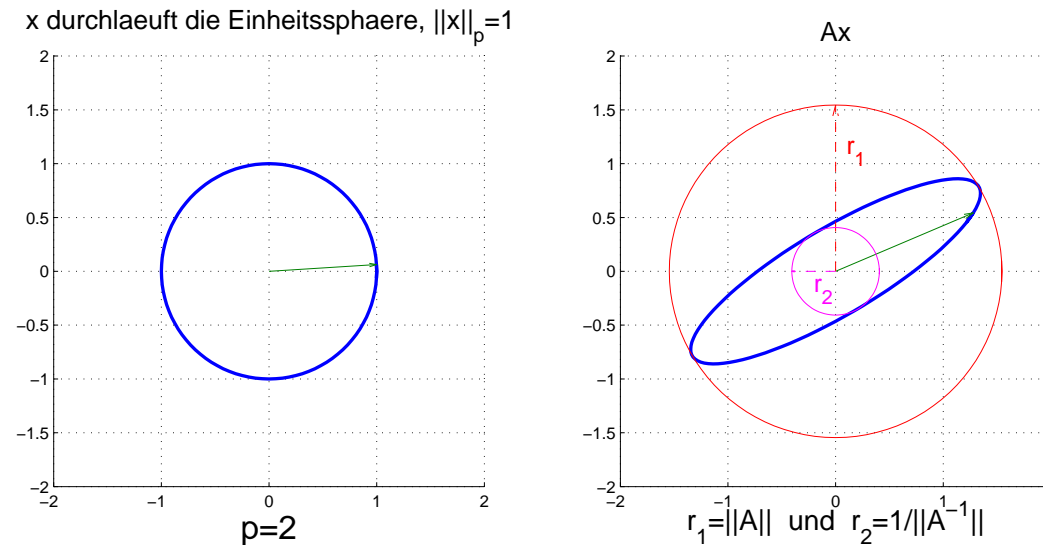
$$S_1(1) = \{ x; \|x\|_1 = 1 \}.$$

Die dicke blaue Kurve im rechten Bild ist das  $A$ -Bild der Sphäre:

$$\{ Ax; \|x\|_1 = 1 \}, \quad \text{wobei } A = \begin{bmatrix} 1.25 & 0.5 \\ 0.5 & 0.7 \end{bmatrix}.$$

Die dünnen Kurven rechts sind die Sphären  $S_1(r_1)$  und  $S_1(r_2)$ .

## Bilder zur Veranschaulichung von $\|A\|$ und $1/\|A^{-1}\|$



Erklärung:

Die dicke blaue Kurve im linken Bild ist die Sphäre

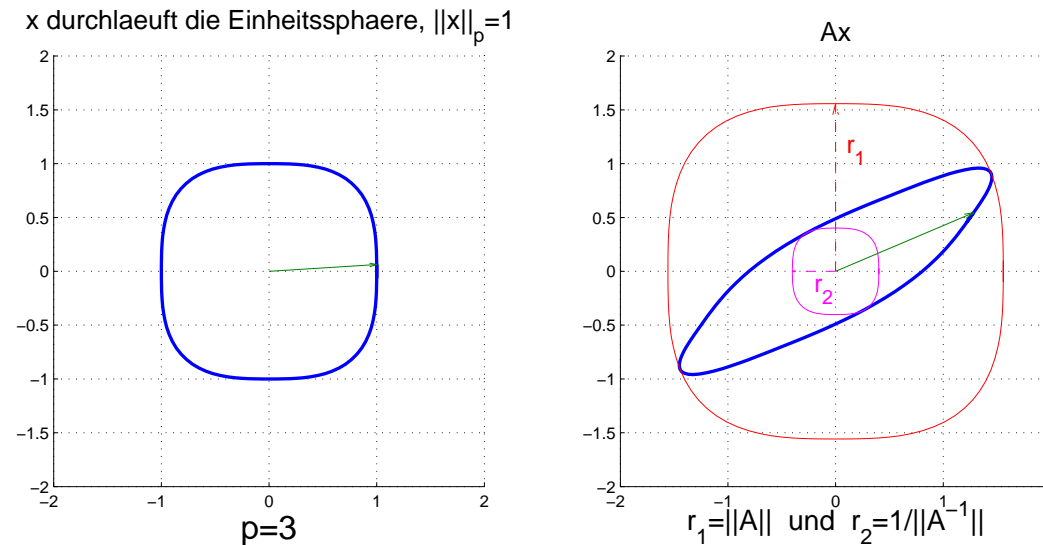
$$S_2(1) = \{ x; \|x\|_2 = 1 \}.$$

Die dicke blaue Kurve im rechten Bild ist das  $A$ -Bild der Sphäre:

$$\{ Ax; \|x\|_2 = 1 \}, \quad \text{wobei } A = \begin{bmatrix} 1.25 & 0.5 \\ 0.5 & 0.7 \end{bmatrix}.$$

Die dünnen Kurven rechts sind die Sphären  $S_2(r_1)$  und  $S_2(r_2)$ .

## Bilder zur Veranschaulichung von $\|A\|$ und $1/\|A^{-1}\|$



Erklärung:

Die dicke blaue Kurve im linken Bild ist die Sphäre

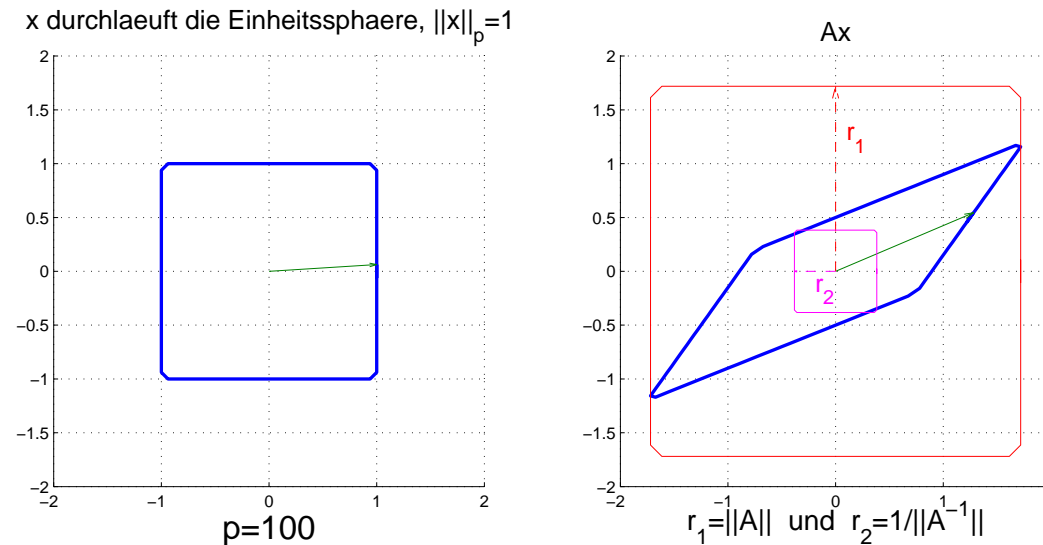
$$S_3(1) = \{ x; \|x\|_3 = 1 \}.$$

Die dicke blaue Kurve im rechten Bild ist das  $A$ -Bild der Sphäre:

$$\{ Ax; \|x\|_3 = 1 \}, \quad \text{wobei } A = \begin{bmatrix} 1.25 & 0.5 \\ 0.5 & 0.7 \end{bmatrix}.$$

Die dünnen Kurven rechts sind die Sphären  $S_3(r_1)$  und  $S_3(r_2)$ .

## Bilder zur Veranschaulichung von $\|A\|$ und $1/\|A^{-1}\|$



Erklärung:

Die dicke blaue Kurve im linken Bild ist die Sphäre

$$S_{100}(1) = \{ x; \|x\|_{100} = 1 \}.$$

Die dicke blaue Kurve im rechten Bild ist das  $A$ -Bild der Sphäre:

$$\{ Ax; \|x\|_{100} = 1 \}, \quad \text{wobei } A = \begin{bmatrix} 1.25 & 0.5 \\ 0.5 & 0.7 \end{bmatrix}.$$

Die dünnen Kurven rechts sind die Sphären  $S_{100}(r_1)$  und  $S_{100}(r_2)$ .

## Konditionszahlen von Matrizen

Für eine invertierbare Matrix  $A \in \mathbb{R}^{n \times n}$  haben wir

$$\|A\|_\alpha = \max_{\|x\|_\alpha=1} \|Ax\|_\alpha, \quad \frac{1}{\|A^{-1}\|_\alpha} = \min_{\|x\|_\alpha=1} \|Ax\|_\alpha.$$

Daraus folgt:

$$\|A\|_\alpha \|A^{-1}\|_\alpha = \frac{\max_{\|x\|_\alpha=1} \|Ax\|_\alpha}{\min_{\|x\|_\alpha=1} \|Ax\|_\alpha}.$$

Diese Größe heißt Konditionszahl von  $A$  bezüglich der Vektornorm  $\|\cdot\|_\alpha$ . Notation:

$$\text{cond}_\alpha(A) := \|A\|_\alpha \|A^{-1}\|_\alpha = \frac{\max_{\|x\|_\alpha=1} \|Ax\|_\alpha}{\min_{\|x\|_\alpha=1} \|Ax\|_\alpha}.$$

Die Konditionszahl ist also der Quotient aus dem maximalen und dem minimalen Streckfaktor, wenn man einen Vektor  $x$  mit der Matrix  $A$  multipliziert. Es gilt stets

$\text{cond}_\alpha(A) \geq 1$  und  $\text{cond}_\alpha(A) = 1$  genau dann, wenn  $\|Ax\|_\alpha = \|x\|_\alpha$  für alle  $x \in \mathbb{R}^n$ .

MATLAB-Anweisung zur Berechnung der Konditionszahl bzgl.  $\|\cdot\|_p$ ,  $p = 1, 2, \infty$ :

`cond(A, p)`

## Extremwerte einer quadratischen Form

**Satz:** Sei  $S \in \mathbb{R}^{n \times n}$  eine symmetrische Matrix. Seien  $\lambda_{\min}, \lambda_{\max} \in \mathbb{R}$  der maximale und der minimale Eigenwert von  $S$  und seien  $\underline{v}, \bar{v} \in \mathbb{R}^n$  zugehörige normierte Eigenvektoren, d.h.:

$$S\underline{v} = \lambda_{\min} \underline{v}, \quad S\bar{v} = \lambda_{\max} \bar{v}, \quad \|\underline{v}\|_2 = \|\bar{v}\|_2 = 1.$$

Dann gilt:

$$\min_{x \neq 0} \frac{x^T S x}{\|x\|_2^2} = \min_{\|x\|_2=1} x^T S x = \underline{v}^T S \underline{v} = \lambda_{\min}$$

$$\max_{x \neq 0} \frac{x^T S x}{\|x\|_2^2} = \max_{\|x\|_2=1} x^T S x = \bar{v}^T S \bar{v} = \lambda_{\max}.$$

**Terminologie:** Der Quotient  $\frac{x^T S x}{\|x\|_2^2}$  heißt Rayleigh-Quotient.

**Beweis:** Seien  $\lambda_{\max} = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n = \lambda_{\min}$  die Eigenwerte von  $S$  und sei  $\bar{v} = v_1, v_2, \dots, v_n = \underline{v}$  eine Orthonormalbasis von Eigenvektoren,  $S v_k = \lambda_k v_k$ . Jeder Vektor  $x \in \mathbb{R}^n$  ist eine Linearkombination der Eigenvektoren:

$$x = x_1 v_1 + x_2 v_2 + \dots + x_n v_n, \quad x_k \in \mathbb{R}.$$

Eine direkte Rechnung ergibt

$$\frac{x^T S x}{\|x\|_2^2} = \frac{\lambda_1 x_1^2 + \lambda_2 x_2^2 + \dots + \lambda_n x_n^2}{x_1^2 + x_2^2 + \dots + x_n^2}.$$

Das Maximum dieses Quotienten wird z.B. angenommen, wenn  $x_1 = 1$  und  $x_2 = \dots = x_n = 0$ . Das Minimum wird z.B. angenommen, wenn  $x_n = 1$  und  $x_1 = \dots = x_{n-1} = 0$ .



## Die 2-Norm einer Matrix

**Satz:** Für jede Matrix  $A \in \mathbb{R}^{m \times n}$  gilt

$$\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sqrt{\lambda_{\max}(A^T A)},$$
$$\min_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sqrt{\lambda_{\min}(A^T A)}.$$

Dabei sind  $\lambda_{\max}$  und  $\lambda_{\min}$  der größte und der kleinste Eigenwert der positiv semidefiniten symmetrischen Matrix  $A^T A$ .

**Beweis:** Man hat  $\|Ax\|^2 = (Ax)^T(Ax) = x^T A^T A x$  und daher

$$\|A\|_2^2 = \max_{x \neq 0} \frac{\|Ax\|_2^2}{\|x\|_2^2} = \max_{x \neq 0} \frac{x^T A^T A x}{\|x\|_2^2} = \lambda_{\max}(A^T A).$$

Bei der letzten Gleichung wurde der Satz über das Maximum des Rayleigh-Quotienten benutzt. Der Beweis für das Minimum ist analog.

## Terminologie:

1. Die Wurzeln aus den Eigenwerten von  $A^T A$  nennt man **Singulärwerte** von  $A$ .

$$\text{Notation: } \sigma_k(A) := \sqrt{\lambda_k(A^T A)},$$

$$\text{Insbesondere: } \sigma_{\max}(A) = \sqrt{\lambda_{\max}(A^T A)}, \quad \sigma_{\min}(A) = \sqrt{\lambda_{\min}(A^T A)}.$$

2. Die 2-Norm  $\|A\|_2$  nennt man auch **Spektralnrm**  
(Spektrum=Menge der Eigenwerte einer Matrix)

Mit den obigen Bezeichnungen hat man

$$\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sigma_{\max}(A),$$

$$\min_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sigma_{\min}(A).$$

Erinnerung: Wenn  $A \in \mathbb{R}^{n \times n}$  invertierbar, dann ist

$$\min_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \frac{1}{\|A^{-1}\|_2}.$$

Daher:

$$\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}.$$

## Überblick: Die wichtigsten Matrixnormen

Sei  $A = [a_{ik}] \in \mathbb{R}^{m \times n}$ .

- $\|A\|_\infty = \max_{i=1}^m \sum_{k=1}^n |a_{ik}|$  (Zeilensummennorm)
- $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)} = \sigma_{\max}(A)$  (Spektralnorm)
- $\|A\|_1 = \max_{k=1}^n \sum_{i=1}^m |a_{ik}|$  (Spaltensummennorm)

**Induzierte Matrixnormen sind submultiplikativ.**

Sei  $\|\cdot\|_\alpha$  irgendeine Vektornorm auf  $\mathbb{R}^n$ . Dann gilt für die induzierte Matrixnorm

$$\|AB\|_\alpha \leq \|A\|_\alpha \|B\|_\alpha, \quad A, B \in \mathbb{R}^{n \times n}.$$

**Beweis:** Nach Definition der induzierten Matrixnorm hat man

$$\|ABx\|_\alpha \leq \|A\|_\alpha \|Bx\|_\alpha.$$

Somit

$$\|AB\|_\alpha = \max_{x \neq 0} \frac{\|ABx\|_\alpha}{\|x\|_\alpha} \leq \max_{x \neq 0} \frac{\|A\|_\alpha \|Bx\|_\alpha}{\|x\|_\alpha} = \|A\|_\alpha \max_{x \neq 0} \frac{\|Bx\|_\alpha}{\|x\|_\alpha} = \|A\|_\alpha \|B\|_\alpha.$$

## **Herleitung der Fehlerformeln**

## Herleitung der Fehlerformeln I

Seien  $A, \tilde{A} \in \mathbb{R}^{n \times n}$  invertierbar,  $Ax = b$ ,  $\tilde{A}\tilde{x} = \tilde{b}$ . Dann folgt

$$\begin{aligned}\tilde{x} - x &= \tilde{A}^{-1}\tilde{b} - x \\ &= \tilde{A}^{-1}b - x + \tilde{A}^{-1}(\tilde{b} - b) \\ &= \tilde{A}^{-1}(A - \tilde{A})x + \tilde{A}^{-1}(\tilde{b} - b)\end{aligned}$$

$\Rightarrow$

$$\begin{aligned}\|\tilde{x} - x\| &= \|\tilde{A}^{-1}(A - \tilde{A})x + \tilde{A}^{-1}(\tilde{b} - b)\| \\ &\leq \|\tilde{A}^{-1}\| \|A - \tilde{A}\| \|x\| + \|\tilde{A}^{-1}\| \|\tilde{b} - b\| \\ &= \|\tilde{A}^{-1}\| (\|A - \tilde{A}\| \|x\| + \|\tilde{b} - b\|)\end{aligned}$$

$\Rightarrow$

$$\begin{aligned}\frac{\|\tilde{x} - x\|}{\|x\|} &\leq \|\tilde{A}^{-1}\| \|A\| \left( \frac{\|A - \tilde{A}\|}{\|A\|} + \frac{\|b\|}{\|A\| \|x\|} \frac{\|\tilde{b} - b\|}{\|b\|} \right) \\ &\leq \|\tilde{A}^{-1}\| \|A\| \left( \frac{\|A - \tilde{A}\|}{\|A\|} + \frac{\|\tilde{b} - b\|}{\|b\|} \right). \quad (\|b\| = \|Ax\| \leq \|A\| \|x\|)\end{aligned}$$

Damit sind die ersten beiden Fehlerformeln (\*) und (\*\*) vom Anfang der Vorlesung bewiesen. Zum Beweis der Fehlerformel (\*\*\*) müssen wir den Faktor  $\|\tilde{A}^{-1}\| \|A\|$  durch die Konditionszahl abschätzen. Dies geschieht auf der nächsten Seite.

## Herleitung der Fehlerformeln II

Für alle  $y \in \mathbb{R}^n$  ist

$$\|Ay\| = \|\tilde{A}y + (A - \tilde{A})y\| \leq \|\tilde{A}y\| + \|(A - \tilde{A})y\|.$$

$\Rightarrow$

$$\|\tilde{A}y\| \geq \|Ay\| - \|(A - \tilde{A})y\|.$$

$\Rightarrow$

$$\frac{\|\tilde{A}y\|}{\|y\|} \geq \frac{\|Ay\|}{\|y\|} - \frac{\|(A - \tilde{A})y\|}{\|y\|} \geq \frac{\|Ay\|}{\|y\|} - \|A - \tilde{A}\|$$

$\Rightarrow$

$$\min_{y \neq 0} \frac{\|\tilde{A}y\|}{\|y\|} \geq \min_{y \neq 0} \frac{\|Ay\|}{\|y\|} - \|A - \tilde{A}\|$$

$\Rightarrow$

$$\frac{1}{\|\tilde{A}^{-1}\|} \geq \frac{1}{\|A^{-1}\|} - \|A - \tilde{A}\|$$

$\Rightarrow$

$$\|\tilde{A}^{-1}\| \leq \frac{1}{\frac{1}{\|A\|} - \|A - \tilde{A}\|} = \frac{\|A^{-1}\|}{1 - \|A - \tilde{A}\| \|A^{-1}\|}$$

$\Rightarrow$

$$\|A\| \|\tilde{A}^{-1}\| \leq \frac{\|A\| \|A^{-1}\|}{1 - \frac{\|A - \tilde{A}\|}{\|A\|} \|A\| \|A^{-1}\|} = \frac{\text{cond}(A)}{1 - \frac{\|A - \tilde{A}\|}{\|A\|} \text{cond}(A)}.$$

Hieraus und aus der letzten Ungleichung auf der vorherigen Seite folgt die Fehlerformel (\* \* \*).

Wir haben eben gezeigt, dass

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \frac{\|\tilde{A} - A\|}{\|A\|} \text{cond}(A)} \left( \frac{\|\tilde{A} - A\|}{\|A\|} + \frac{\|\tilde{b} - b\|}{\|b\|} \right) \quad (***)$$

**Spezialfall:**

$$A = \tilde{A} \quad \Rightarrow \quad \frac{\|\tilde{x} - x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\tilde{b} - b\|}{\|b\|}$$



## Noch einmal: Ungenaue Eingangsdaten bei linearen Gleichungssystemen

Aufgabe: Es soll die Gleichung  $Ax = b$  gelöst werden.

Die Ausgangsdaten  $A, b$  können aber z.B. aus folgenden Gründen fehlerhaft sein:

- $A, b$  sind Messwerte und daher nicht genau bekannt.
- $A, b$  sind Ergebnisse früherer (fehlerhafter) Rechnungen
- Die Einträge von  $A, b$  sind keine Maschinenzahlen und können nicht exakt im Rechner abgespeichert werden.

Außerdem machen Algorithmen zur Lösung von  $Ax = b$  Fehler, die man als Fehler in den Daten  $A, b$  interpretieren kann. Macht man z.B. eine  $LR$ -Zerlegung

$$A = LR,$$

dann bekommt man statt  $L, R$  numerisch  $\tilde{L}, \tilde{R}$  heraus. Das Produkt ist

$$\tilde{L}\tilde{R} = A + \Delta A.$$

Beim Vorwärts-Rückwärtseinsetzen löst man also im besten Fall das Gleichungssystem

$$(A + \Delta A)x = b.$$

## Praktische Konsequenz der Fehlerformeln

Aus der Fehlerformel

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \frac{\|\tilde{A} - A\|}{\|A\|} \text{cond}(A)} \left( \frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right)$$

lässt sich folgende Faustregel ableiten:

**Eine Konditionszahl  $\text{cond}(A) = 10^q$  kostet  $q$  Stellen Genauigkeit bei der Lösung von  $Ax = b$ .**

## Was nützt eine Probe bei schlecht konditionierten Problemen?

Problem:  $Ax = b$

Exakte Lösung:  $x = A^{-1}b$

Numerische Lösung:  $\tilde{x}$

Probe ergibt:  $A\tilde{x} = \tilde{b}$

Angenommen, das Produkt  $A\tilde{x}$  wurde exakt berechnet und der relative Fehler  $\|b - \tilde{b}\|/\|b\|$  ist klein. Kann man daraus schließen, dass auch der relative Fehler  $\|x - \tilde{x}\|/\|x\|$  klein ist?

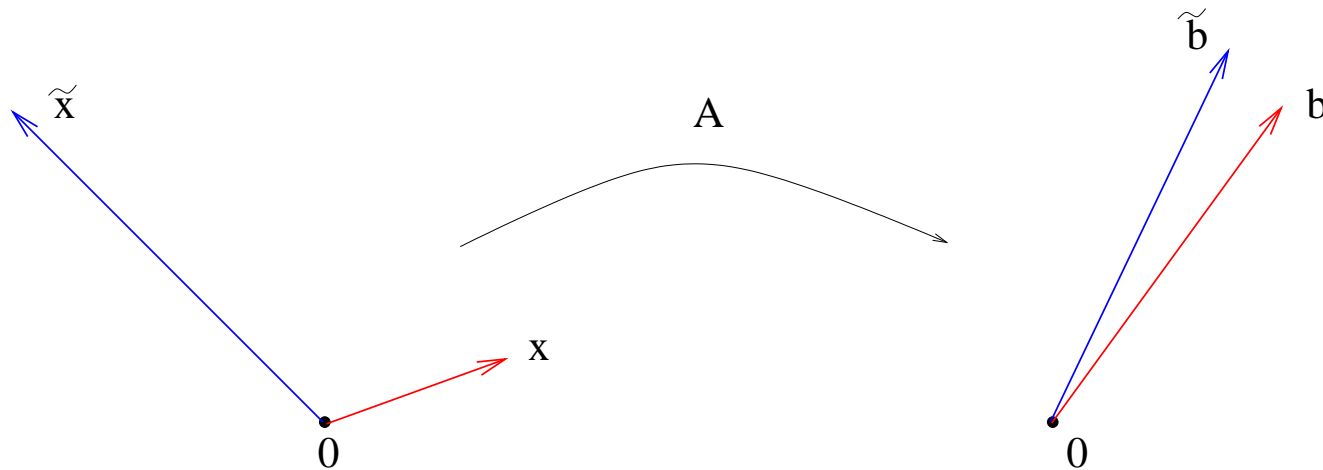
**Antwort:** Das hängt von der Konditionszahl ab. Die (nicht verbesserbare) Fehlerformel

$$\frac{\|x - \tilde{x}\|}{\|x\|} \leq \text{cond}(A) \frac{\|b - \tilde{b}\|}{\|b\|}$$

besagt, dass sich der relative Fehler in  $b$  im ungünstigsten Fall um den Faktor  $\text{cond}(A)$  verstärkt.

## Veranschaulichung der Situation bei einer schlecht konditionierten Matrix:

Relativ weit voneinander entfernte Vektoren  $x, \tilde{x}$  werden auf relativ nahe beieinander liegende Vektoren  $b = Ax, \tilde{b} = A\tilde{x}$  abgebildet.



## Situation bei ungenauer rechter Seite:

Man kennt nur  $\tilde{b}$ , die exakte rechte Seite  $b$  ist unbekannt. Falls das Gleichungssystem exakt gelöst wird, kommt  $\tilde{x}$  heraus. Die Lösung  $x$  zur exakten rechten Seite  $b$  kann aber weit entfernt davon liegen.

## Situation bei einer Proberechnung:

Rechte Seite ist  $b$ . Berechnet wurde die fehlerhafte Lösung  $\tilde{x}$ . Die Probe ergibt  $A\tilde{x} = \tilde{b}$ . Auch wenn  $b$  und  $\tilde{b}$  nahezu übereinstimmen, können die wahre und die berechnete Lösung sich stark unterscheiden.

**Die Konditionszahl einer Matrix ist hoch, wenn die Zeilen und Spalten fast linear abhängig sind**

Beispiel: Sei  $A_\epsilon = \begin{bmatrix} 1 + \epsilon & 3 \\ 2 & 6 \end{bmatrix}$ .

Die Zeilen und Spalten von  $A_0$  sind linear abhängig. Für  $\epsilon \neq 0$ :

$$A_\epsilon^{-1} = \frac{1}{\det(A_\epsilon)} \begin{bmatrix} 6 & -3 \\ -2 & 1 + \epsilon \end{bmatrix} = \frac{1}{6\epsilon} \begin{bmatrix} 6 & -3 \\ -2 & 1 + \epsilon \end{bmatrix}.$$

Die Konditionszahl von  $A_\epsilon$  bzgl. Zeilensummennorm ist für  $\epsilon \in [-6, 4]$ :

$$\text{cond}_\infty(A_\epsilon) = \|A_\epsilon\|_\infty \|A_\epsilon^{-1}\|_\infty = (2 + 6) \frac{6 + 3}{6|\epsilon|} = \frac{12}{|\epsilon|} \rightarrow \infty \quad \text{für } \epsilon \rightarrow 0.$$

**Bemerkung:** In diesem Beispiel ist die hohe Konditionszahl für kleine  $\epsilon$  darauf zurückzuführen, dass  $\det(A_\epsilon)$  klein ist. Eine kleine Determinante impliziert aber nicht notwendig eine kleine Konditionszahl. Beispiel:

$$\text{cond}(\epsilon I) = \|\epsilon I\| \|(\epsilon I)^{-1}\| = 1 \quad \text{für alle } \epsilon > 0.$$

## Verbesserung der Konditionszahl durch Präkonditionierung

Aufgabe: Löse

$$Ax = b. \quad (*)$$

Durch Multiplikation der Gleichung mit einer nicht singulären Matrix  $D \in \mathbb{R}^{n \times n}$  bekommt man die äquivalente Gleichung

$$DAx = Db. \quad (**)$$

Wenn  $A$  große Konditionszahl hat, dann sucht man eine Matrix  $D$  mit

$$\text{cond}(DA) \ll \text{cond}(A)$$

und löst dann  $(**)$  statt  $(*)$ .

### Einfachste Möglichkeit:

Wähle  $D$  als Diagonalmatrix, und zwar so, dass alle Zeilen von  $DA$  die gleiche 1-Norm haben (**Zeilenäquilibration**).