

**Definition 2.1: Maschinenzahlen / Gleitpunktzahlen**

- Unter der zusätzlichen Normierungs-Bedingung  $m_1 \neq 0$  (falls  $x \neq 0$ ) ergibt sich eine eindeutige Darstellung der sogenannten **maschinendarstellbaren Zahlen  $M$**  zur Basis  $B$ :

$$M = \{x \in \mathbb{R} \mid x = \pm 0.m_1m_2m_3\dots m_n \cdot B^{\pm e_1e_2\dots e_l}\} \cup \{0\}$$

Dabei gilt  $m_i, e_i \in \{0, 1, \dots, B-1\}$  für  $i \neq 0$  und  $B \in \mathbb{N}$ ,  $B > 1$ )

- Der **Wert** einer solchen Zahl ist definiert als

$$\sum_{i=1}^n m_i B^{e-i}$$

und ergibt gerade die (nicht normierte) Darstellung der Zahl im Dezimalsystem. Dabei ist  $e$  ebenfalls im Dezimalsystem zu nehmen, also  $e = \sum_{i=1}^l e_i B^{l-i}$  und es gilt  $e \in \mathbb{Z}$ , d.h.  $e$  kann natürlich auch negativ sein. Weiter gibt es eine obere und untere Schranke:  $e_{\min} \leq e \leq e_{\max}$ .

- Man redet dann auch von einer  **$n$ -stelligen Gleitpunktzahl zur Basis  $B$**  (engl: floating point). Zahlen, die nicht in dieser Menge  $M$  liegen, müssen durch Rundung in eine maschinendarstellbare Zahl umgewandelt werden.

**Definition 2.2: Absoluter / Relativer Fehler**

- Hat man eine Näherung  $\tilde{x}$  zu einem exakten Wert  $x$ , so ist der Betrag der Differenz  $|\tilde{x} - x|$  der **absolute Fehler**.
- Falls  $x \neq 0$ , so ist  $|\frac{\tilde{x}-x}{x}|$  bzw.  $\frac{|\tilde{x}-x|}{|x|}$  der **relative Fehler** dieser Näherung.

Bemerkung: In der Numerik ist der relative Fehler der wichtigere. Weshalb?

**Definition 2.3: Maschinengenauigkeit**

- Die Zahl  $\text{eps} := 5 \cdot 10^{-n}$  heisst **Maschinengenauigkeit**. Bei allgemeiner Basis  $B$  gilt  $\text{eps} := \frac{B}{2} \cdot B^{-n}$ . Sie gibt den maximalen relativen Fehler, der durch Rundung entstehen kann.
- Alternative Definition: Die Maschinengenauigkeit ist die kleinste positive Maschinenzahl, für die auf dem Rechner  $1 + \text{eps} \neq 1$  gilt.

$$\left| \frac{rd(x) - x}{x} \right| \leq 5 \cdot 10^{-n} \text{ (da } x \geq 10^{e-1}).$$

- Näherung für den **absoluten Fehler bei Funktionsauswertungen**:

$$|f(\tilde{x}) - f(x)| \approx |f'(x)| \cdot |\tilde{x} - x|$$

- Näherung für den **relativen Fehler bei Funktionsauswertungen**:

$$\frac{|f(\tilde{x}) - f(x)|}{|f(x)|} \approx \frac{|f'(x)| \cdot |x|}{|f(x)|} \cdot \frac{|\tilde{x} - x|}{|x|}$$

#### Definition 2.4: Konditionszahl

- Den Faktor

$$K := \frac{|f'(x)| \cdot |x|}{|f(x)|}$$

nennt man **Konditionszahl**.

- Man unterscheidet **gut konditionierte Probleme**, d.h. die Konditionszahl ist klein, und **schlecht konditionierte Probleme** (ill posed problems) mit grosser Konditionszahl. Bei gut konditionierten Problemen wird der relative Fehler durch die Auswertung der Funktion nicht grösser.

#### Satz 3.1: Nullstellensatz von Bolzano

- Sei  $f : [a, b] \rightarrow \mathbb{R}$  stetig mit  $f(a) \leq 0 \leq f(b)$  oder  $f(a) \geq 0 \geq f(b)$ . Dann muss  $f$  in  $[a, b]$  eine Nullstelle besitzen.

#### Satz 3.2: Bisektionsverfahren [1]

- Gegeben sei eine stetige Funktion  $f : [a, b] \rightarrow \mathbb{R}$  mit  $f(a) \cdot f(b) < 0$ . In jedem der über die Rekursion für  $i = 0, 1, \dots$  erzeugten Intervalle

$$\begin{aligned} [a_0, b_0] &= [a, b]; \\ [a_{i+1}, b_{i+1}] &= \begin{cases} [a_i, \frac{a_i+b_i}{2}] & \text{falls } f(\frac{a_i+b_i}{2}) \cdot f(a_i) \leq 0 \\ [\frac{a_i+b_i}{2}, b_i] & \text{sonst} \end{cases} \end{aligned}$$

befindet sich eine Nullstelle von  $f$  und es gilt

$$b_i - a_i = \frac{b-a}{2^i}, \text{ insbesondere also } \lim_{i \rightarrow \infty} (b_i - a_i) = 0$$

**Definition 3.1: Fixpunktgleichung / Fixpunkt [1]**

- Eine Gleichung der Form  $F(x) = x$  heisst **Fixpunktgleichung**.
- Ihre Lösungen  $\bar{x}$ , für die  $F(\bar{x}) = \bar{x}$  erfüllt ist, heissen **Fixpunkte** (da die Funktion  $F$  die Punkte  $\bar{x}$  auf sich selbst abbildet).

Anstelle eines Nullstellenproblems kann man also ein dazu äquivalentes Fixpunktproblem betrachten. Dazu muss aber  $f(x) = 0$  in die Fixpunktform  $F(x) = x$  gebracht werden, wozu es viele Möglichkeiten gibt. Bei dieser Überführung muss unbedingt auf Äquivalenz geachtet werden, d.h. die Lösungsmenge darf nicht verändert werden.

**Beispiel 3.2:**

- Die Gleichung  $p(x) = x^3 - x + 0.3$  soll in Fixpunktform gebracht werden.  
Lösung: Die einfachste Möglichkeit ist  $p(x) = 0 \iff F(x) \equiv x^3 + 0.3 = x$   
Aber auch  $F(x) \equiv \sqrt[3]{x - 0.3} = x$  ist möglich.
- Die Gleichung  $x = \cos(x)$ , die wir weiter oben graphisch gelöst haben, ist bereits in der Fixpunktform.

**Definition 3.2: Fixpunktiteration [1]**

- Gegeben sei  $F : [a, b] \rightarrow \mathbb{R}$ , mit  $x_0 \in [a, b]$ . Die rekursive Folge

$$x_{n+1} \equiv F(x_n), \quad n = 0, 1, 2, \dots$$

heisst Fixpunktiteration von  $F$  zum Startwert  $x_0$ .

**Satz 3.2 zur Fixpunktiteration [1]:**

- Sei  $F : [a, b] \rightarrow \mathbb{R}$  mit stetiger Ableitung  $F'$  und  $\bar{x} \in [a, b]$  ein Fixpunkt von  $F$ . Dann gilt für die Fixpunktiteration  $x_{n+1} = F(x_n)$ :
  - Ist  $|F'(\bar{x})| < 1$ , so konvergiert  $x_n$  gegen  $\bar{x}$ , falls der Startwert  $x_0$  nahe genug bei  $\bar{x}$  liegt. Der Punkt  $\bar{x}$  heisst dann **anziehender Fixpunkt**.
  - Ist  $|F'(\bar{x})| > 1$ , so konvergiert  $x_n$  für keinen Startwert  $x_0 \neq \bar{x}$ . Der Punkt  $\bar{x}$  heisst dann **abstossender Fixpunkt**.

**Satz 3.3: Banachscher Fixpunktsatz [1]**

- Sei  $F : [a, b] \rightarrow [a, b]$  (d.h.  $F$  bildet  $[a, b]$  auf sich selber ab) und es existiere eine Konstante  $\alpha$  mit  $0 < \alpha < 1$  und

$$|F(x) - F(y)| \leq \alpha |x - y| \quad \text{für alle } x, y \in [a, b]$$

(d.h.  $F$  ist “Lipschitz-stetig” und “kontraktiv”,  $\alpha$  nennt man auch Lipschitz-Konstante). Dann gilt:

- $F$  hat genau einen Fixpunkt  $\bar{x}$  in  $[a, b]$
- Die Fixpunktiteration  $x_{n+1} = F(x_n)$  konvergiert gegen  $\bar{x}$  für alle Startwerte  $x_0 \in [a, b]$
- Es gelten die Fehlerabschätzungen

$$|x_n - \bar{x}| \leq \frac{\alpha^n}{1 - \alpha} |x_1 - x_0| \quad \text{a-priori Abschätzung}$$

$$|x_n - \bar{x}| \leq \frac{\alpha}{1 - \alpha} |x_n - x_{n-1}| \quad \text{a-posteriori Abschätzung}$$

**Bemerkungen:**

- Aus  $|F(x) - F(y)| \leq \alpha |x - y|$  für alle  $x, y \in [a, b]$  folgt

$$\frac{|F(x) - F(y)|}{|x - y|} \leq \alpha,$$

wobei die linke Seite sämtliche möglichen Steigungen der Sekanten durch die beiden Punkte  $(x, F(x))$  und  $(y, F(y))$  für alle  $x, y \in [a, b]$  darstellt. Aus diesem Grund kann man  $\alpha$  als die grösstmögliche Steigung von  $F(x)$  auf dem Intervall  $[a, b]$  interpretieren, bzw.

$$\alpha = \max_{x_0 \in [a, b]} |F'(x_0)|$$

- Wählt man das Intervall  $[a, b]$  sehr nahe um einen anziehenden Fixpunkt  $\bar{x}$ , so ist also  $\alpha \approx |F'(\bar{x})|$ .
- In der Praxis gestaltet es sich meist schwierig, ein Intervall  $[a, b]$  zu finden, dass unter  $F$  auf sich selbst abgebildet wird. Hat man ein solches Intervall gefunden, dann sind die Fehlerabschätzungen aber recht nützlich. Wir werden diesen Satz nochmals im Zusammenhang mit der iterativen Lösung von linearen Gleichungssystemen in Kap. 4 aufgreifen.

Auflösen nach  $x_{n+1}$  liefert

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad (n = 0, 1, 2, 3, \dots).$$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_0)} \quad (n = 0, 1, 2, 3, \dots).$$

**Definition 3.3: Konvergenzordnung [1]**

- Sei  $(x_n)$  eine gegen  $\bar{x}$  konvergierende Folge. Dann hat das Verfahren die **Konvergenzordnung**  $q \geq 1$  wenn es eine Konstante  $c > 0$  gibt mit

$$|x_{n+1} - \bar{x}| \leq c \cdot |x_n - \bar{x}|^q$$

für alle  $n$ . Falls  $q = 1$  verlangt man noch  $c < 1$ . Im Fall  $q = 1$  spricht man von linearer, im Fall  $q = 2$  von quadratischer Konvergenz.

**Definition 4.1: Untere Dreiecksmatrix / Obere Dreiecksmatrix [6]**

- Eine  $n \times n$  Matrix  $L = (l_{ij})$  heisst **untere Dreiecksmatrix**, wenn  $l_{ij} = 0$  für  $j > i$  gilt; sie heisst **normierte untere Dreiecksmatrix**, wenn ausserdem  $l_{ii} = 1$  für alle  $i$  gilt.
- Eine  $n \times n$  Matrix  $R = (r_{ij})$  heisst **obere Dreiecksmatrix**, wenn  $r_{ij} = 0$  für  $i > j$  gilt; sie heisst **normierte obere Dreiecksmatrix**, wenn ausserdem  $r_{ii} = 1$  für alle  $i$  gilt.

- Untere normierte Dreiecksmatrix:

$$L = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ l_{21} & 1 & 0 & \cdots & 0 \\ l_{31} & l_{32} & 1 & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & 0 \\ l_{n1} & l_{n2} & \cdots & l_{nn-1} & 1 \end{pmatrix}$$

- Obere Dreiecksmatrix:

$$R = \begin{pmatrix} r_{11} & r_{12} & r_{13} & \cdots & r_{1n} \\ 0 & r_{22} & r_{23} & \cdots & r_{2n} \\ 0 & 0 & r_{33} & \cdots & r_{3n} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & r_{nn} \end{pmatrix}$$

**Gauss-Algorithmus zur Transformation von  $Ax = b$  auf ein oberes Dreieckssystem [1]:**

- für  $i = 1, \dots, n-1$ :

erzeuge Nullen unterhalb des Diagonalelementes in der  $i$ -ten Spalte

- Falls nötig und möglich, Sorge durch Zeilenvertauschung für  $a_{ii} \neq 0$ :  
falls  $a_{ii} \neq 0$ : tue nichts

$$\text{falls } a_{ii} = 0 : \begin{cases} \text{falls } a_{ji} = 0 \text{ für alle } j = i+1, \dots, n : \\ \quad A \text{ ist nicht regulär; stop;} \\ \text{wenn } a_{ji} \neq 0 \text{ für ein } j = i+1, \dots, n : \\ \quad \text{sei } j \geq i+1 \text{ der kleinste Index mit } a_{ji} \neq 0 \\ \quad z_i \longleftrightarrow z_j \end{cases}$$

- Eliminationsschritt:

für  $j = i+1, \dots, n$  eliminiere das Element  $a_{ji}$  durch:

$$z_j := z_j - \frac{a_{ji}}{a_{ii}} \cdot z_i$$

### Gauss-Algorithmus zur Transformation von $Ax = b$ mit Spaltenpivotisierung [1]:

- für  $i = 1, \dots, n - 1$  :
  - erzeuge Nullen unterhalb des Diagonalelementes in der  $i$ -ten Spalte
    - Suche das betragsgrösste Element unterhalb der Diagonalen in der  $i$ -ten Spalte:  
Wähle  $k$  so, dass  $|a_{ki}| = \max\{|a_{ji}| \mid j = i, \dots, n\}$   
 $\begin{cases} \text{falls } a_{ki} = 0 : A \text{ ist nicht regulär; stop;} \\ \text{falls } a_{ki} \neq 0 : z_k \longleftrightarrow z_i; \end{cases}$
    - Eliminationsschritt:  
für  $j = i + 1, \dots, n$  eliminiere das Element  $a_{ji}$  durch:

$$z_j := z_j - \frac{a_{ji}}{a_{ii}} \cdot z_i$$

### Satz 4.1: $LR$ -Zerlegung [1]

Zu jeder regulären  $n \times n$  Matrix  $A$ , für die der Gauss-Algorithmus ohne Zeilenvertauschung durchführbar ist, gibt es  $n \times n$  Matrizen  $L$  und  $R$  mit den folgenden Eigenschaften:

- $L$  ist eine normierte untere Dreiecksmatrix (also mit  $l_{ii} = 1$  für  $i = 1, \dots, n$ )
- $R$  ist eine obere Dreiecksmatrix mit  $r_{ii} \neq 0$  für  $i = 1, \dots, n$
- $A = L \cdot R$  ist die  $LR$ -Zerlegung von  $A$ .

Aufwand: Die Berechnung der  $LR$ -Zerlegung mit dem Gauss-Algorithmus benötigt  $\frac{1}{3}(n^3 - n)$  Punktoperationen

## 4.5.2 Die Cholesky-Zerlegung

Im folgenden lernen wir ein weiteres Verfahren zur Dreieckszerlegung von Matrizen kennen, einen Spezialfall der  $LR$ -Zerlegung. Dieses Verfahren ist nach seinem Entdecker André-Louis Cholesky (1875 -1918) benannt, einem französischen Mathematiker. Er entwickelte es für Anwendungen in der Geodäsie. <sup>9</sup> Die Cholesky-Zerlegung funktioniert nicht für allgemeine Matrizen sondern nur für *symmetrische, positiv definite* Matrizen. Falls anwendbar, ist es aber etwa um einen Faktor zwei effizienter als die allgemeine  $LR$ -Zerlegung.

### Definition 4.2: Symmetrische / positiv definite Matrizen [1]

- Eine Matrix  $A \in \mathbb{R}^{n \times n}$  heisst symmetrisch, falls  $A^T = A$  gilt ( $A^T$  ist die transponierte Matrix).
- Eine Matrix  $A \in \mathbb{R}^{n \times n}$  heisst positiv definit, falls für alle  $x \in \mathbb{R}^n$ ,  $x \neq 0$  gilt  $x^T A x > 0$ .

**Satz 4.2: Cholesky Zerlegung [1]**

Für jede positiv definite  $n \times n$  Matrix  $\mathbf{A}$  gibt es genau eine rechts-obere Dreiecksmatrix  $\mathbf{R}$  mit  $r_{ii} > 0$  für  $i = 1, \dots, n$  und  $\mathbf{A} = \mathbf{R}^T \mathbf{R}$ . Diese Zerlegung heisst **Cholesky-Zerlegung** von  $\mathbf{A}$ .

Die Berechnung der Cholesky-Zerlegung geschieht anhand des folgenden Algorithmus, der uns die Koeffizienten  $r_{ij}$  der oberen Dreiecksmatrix  $\mathbf{R}$  berechnet und gleichzeitig überprüft, ob  $\mathbf{A}$  positiv definit ist:

**Cholesky-Algorithmus [1]:**

Gegeben sei eine symmetrische  $n \times n$  Matrix  $\mathbf{A}$ . Für  $i = 1, \dots, n$  berechne:

- $S = a_{ii} - \sum_{k=1}^{i-1} r_{ki}^2$  (für  $i = 1$  ist also  $S = a_{ii}$ )
- falls  $S \leq 0$ , dann ist  $\mathbf{A}$  nicht positiv definit  $\rightarrow$ stopp.
- falls  $S > 0$ :
  - $r_{ii} = \sqrt{S}$
  - für  $j = i + 1, \dots, n$  :  $r_{ij} = \frac{1}{r_{ii}} \left( a_{ij} - \sum_{k=1}^{i-1} r_{ki} r_{kj} \right)$

**Definition 4.3: Vektornorm [1]**

Eine Abbildung  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$  heisst Vektornorm, wenn die folgenden Bedingungen für alle  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,  $\lambda \in \mathbb{R}$  erfüllt sind:

- $\|\mathbf{x}\| \geq 0$  und  $\|\mathbf{x}\| = 0 \iff \mathbf{x} = \mathbf{0}$
- $\|\lambda \mathbf{x}\| = |\lambda| \|\mathbf{x}\|$
- $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$  "Dreiecksungleichung"

**Definition 4.4: Vektornormen / Matrixnormen [1]**

- Für Vektoren  $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$  gibt es die folgenden Vektornormen:

$$\begin{aligned} \text{1-Norm, Summennorm} & : \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i| \\ \text{2-Norm, euklidische Norm} & : \|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \\ \infty\text{-Norm, Maximumnorm} & : \|\mathbf{x}\|_\infty = \max_{i=1, \dots, n} |x_i| \end{aligned}$$

- Für eine  $n \times n$  Matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  sind mit den Vektornormen die folgenden Matrixnormen verbunden, welche die Eigenschaften der Definition 4.3 ebenfalls erfüllen:

$$\begin{aligned} \text{1-Norm, Spaltensummennorm} & : \|\mathbf{A}\|_1 = \max_{j=1, \dots, n} \sum_{i=1}^n |a_{ij}| \\ \text{2-Norm, Spektralnrm} & : \|\mathbf{A}\|_2 = \sqrt{\rho(\mathbf{A}^T \mathbf{A})} \\ \infty\text{-Norm, Zeilensummennorm} & : \|\mathbf{A}\|_\infty = \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| \end{aligned}$$

- Berechnen Sie die 1-, 2-, und  $\infty$ - Norm des Vektors  $\begin{pmatrix} -1 \\ 2 \\ 3 \end{pmatrix}$  sowie die 1- und  $\infty$ - Norm von  $\begin{pmatrix} 1 & 2 & 3 \\ 3 & 4 & -2 \\ 7 & -3 & 5 \end{pmatrix}$ .

Lösung:

$$\left\| \begin{pmatrix} -1 \\ 2 \\ 3 \end{pmatrix} \right\|_1 = 1 + 2 + 3 = 6, \quad \left\| \begin{pmatrix} -1 \\ 2 \\ 3 \end{pmatrix} \right\|_2 = \sqrt{1 + 2^2 + 3^2} = \sqrt{14},$$

$$\left\| \begin{pmatrix} -1 \\ 2 \\ 3 \end{pmatrix} \right\|_\infty = \max\{1, 2, 3\} = 3$$

$$\left\| \begin{pmatrix} 1 & 2 & 3 \\ 3 & 4 & -2 \\ 7 & -3 & 5 \end{pmatrix} \right\|_1 = \max\{1 + 3 + 7, 2 + 4 + 3, 3 + 2 + 5\} = 11$$

$$\left\| \begin{pmatrix} 1 & 2 & 3 \\ 3 & 4 & -2 \\ 7 & -3 & 5 \end{pmatrix} \right\|_\infty = \max\{1 + 2 + 3, 3 + 4 + 2, 7 + 3 + 5\} = 15.$$

#### Satz 4.3: Abschätzung für fehlerbehaftete Vektoren [1]

- Sei  $\|\cdot\|$  eine Norm,  $A \in \mathbb{R}^{n \times n}$  eine reguläre  $n \times n$  Matrix und  $x, \tilde{x}, b, \tilde{b} \in \mathbb{R}^n$  mit  $Ax = b$  und  $A\tilde{x} = \tilde{b}$ . Dann gilt für den absoluten und den relativen Fehler in  $x$ :

$$- \|x - \tilde{x}\| \leq \|A^{-1}\| \cdot \|b - \tilde{b}\|$$

$$- \frac{\|x - \tilde{x}\|}{\|x\|} \leq \|A\| \cdot \|A^{-1}\| \cdot \frac{\|b - \tilde{b}\|}{\|b\|} \text{ falls } \|b\| \neq 0$$

- Die Zahl  $\text{cond}(A) = \|A\| \cdot \|A^{-1}\|$  nennt man Konditionszahl der Matrix  $A$  bzgl. der verwendeten Norm.

#### Satz 4.4: Abschätzung für fehlerbehaftete Matrix [1]

Sei  $\|\cdot\|$  eine Norm,  $A, \tilde{A} \in \mathbb{R}^{n \times n}$  reguläre  $n \times n$  Matrizen und  $x, \tilde{x}, b, \tilde{b} \in \mathbb{R}^n$  mit  $Ax = b$  und  $A\tilde{x} = \tilde{b}$ . Falls

$$\text{cond}(A) \cdot \frac{\|A - \tilde{A}\|}{\|A\|} < 1$$

dann gilt:

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



**Bemerkung:**

- Für den Fall, dass  $\mathbf{A}$  exakt gegeben ist, gilt  $\frac{\|\mathbf{A} - \tilde{\mathbf{A}}\|}{\|\mathbf{A}\|} = 0$  und der relative Fehler für  $\mathbf{x}$  aus Satz 4.4 reduziert sich auf den relativen Fehler in Satz 4.3.

**Beispiel 4.11 [1]:**

- Nehmen Sie noch einmal das Beispiel 4.9 und untersuchen Sie die Fehlerfortpflanzung unter der zusätzlichen Annahme, dass die Matrix  $\mathbf{A}$  um maximal 0.003 elementweise gestört ist.
- Lösung: Wir hatten bereits die folgenden Größen berechnet

$$\|\mathbf{A}\|_{\infty} = 12.1, \quad \text{cond}(\mathbf{A}) = 732.05, \quad \|\mathbf{b}\|_{\infty} = 1.5, \quad \|\mathbf{b} - \tilde{\mathbf{b}}\|_{\infty} \leq 0.1$$

Wenn nun jedes Element von  $\mathbf{A}$  um maximal 0.003 gestört wird, summiert sich diese Störung in der  $\infty$ -Norm auf und wir erhalten  $\|\mathbf{A} - \tilde{\mathbf{A}}\|_{\infty} \leq 0.006$  und damit

$$\text{cond}(\mathbf{A}) \cdot \frac{\|\mathbf{A} - \tilde{\mathbf{A}}\|_{\infty}}{\|\mathbf{A}\|_{\infty}} \leq 0.363 < 1.$$

Wir können also die Abschätzung aus Satz 4.4 anwenden und erhalten

$$\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|_{\infty}}{\|\mathbf{x}\|_{\infty}} \leq \frac{732.05}{1 - 0.363} \left( \frac{0.006}{12.1} + \frac{0.1}{1.5} \right) \leq 77.2$$

**Definition 4.6: Jacobi- bzw. Gesamtschrittverfahren [1]**

- Zu lösen sei  $\mathbf{Ax} = \mathbf{b}$ . Die Matrix  $\mathbf{A} = (a_{ij})$  sei zerlegt in der Form

$$\mathbf{A} = \underbrace{\begin{pmatrix} 0 & 0 & 0 & \cdots & 0 \\ a_{21} & 0 & 0 & \cdots & 0 \\ a_{31} & a_{32} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn-1} & 0 \end{pmatrix}}_{=: \mathbf{L}} + \underbrace{\begin{pmatrix} a_{11} & 0 & 0 & \cdots & 0 \\ 0 & a_{22} & 0 & \cdots & 0 \\ 0 & 0 & a_{33} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & a_{nn} \end{pmatrix}}_{=: \mathbf{D}} + \underbrace{\begin{pmatrix} 0 & a_{12} & a_{13} & \cdots & a_{1n} \\ 0 & 0 & a_{23} & \cdots & a_{2n} \\ 0 & 0 & 0 & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & a_{n-1,n} \\ 0 & 0 & \cdots & 0 & 0 \end{pmatrix}}_{=: \mathbf{R}}$$

Dann heisst die Fixpunktiteration

$$\begin{aligned} \mathbf{D}\mathbf{x}^{(k+1)} &= -(\mathbf{L} + \mathbf{R})\mathbf{x}^{(k)} + \mathbf{b} \quad \text{bzw.} \\ \mathbf{x}^{(k+1)} &= -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{R})\mathbf{x}^{(k)} + \mathbf{D}^{-1}\mathbf{b} \end{aligned}$$

**Gesamtschrittverfahren** oder **Jacobi-Verfahren**.

**Definition 4.7: Gauss-Seidel bzw. Einzelschrittverfahren [1]**

- Zu lösen sei  $Ax = b$ . Die Matrix  $A = (a_{ij})$  sei zerlegt in der Form

$$A = \underbrace{\begin{pmatrix} 0 & 0 & 0 & \cdots & 0 \\ a_{21} & 0 & 0 & \cdots & 0 \\ a_{31} & a_{32} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn-1} & 0 \end{pmatrix}}_{=: L} + \underbrace{\begin{pmatrix} a_{11} & 0 & 0 & \cdots & 0 \\ 0 & a_{22} & 0 & \cdots & 0 \\ 0 & 0 & a_{33} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & a_{nn} \end{pmatrix}}_{=: D} + \underbrace{\begin{pmatrix} 0 & a_{12} & a_{13} & \cdots & a_{1n} \\ 0 & 0 & a_{23} & \cdots & a_{2n} \\ 0 & 0 & 0 & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & 0 \end{pmatrix}}_{=: R}$$

Dann heisst die Fixpunktiteration

$$\begin{aligned} (D + L)x^{(k+1)} &= -Rx^{(k)} + b \text{ bzw.} \\ x^{(k+1)} &= -(D + L)^{-1}Rx^{(k)} + (D + L)^{-1}b \end{aligned}$$

**Einzelschrittverfahren oder Gauss-Seidel-Verfahren.**

**Definition 4.8: anziehender / abstossender Fixpunkt [1]**

- Gegeben sei eine Fixpunktiteration

$$x^{(n+1)} = Bx^{(n)} + b =: F(x^{(n)})$$

wobei  $B$  eine  $n \times n$  Matrix ist und  $b \in \mathbb{R}^n$ . Weiter sei  $\|\cdot\|$  eine der in Kap. 4.6.1 eingeführten Normen und  $\bar{x} \in \mathbb{R}^n$  erfülle  $\bar{x} = B\bar{x} + b = F(\bar{x})$ . Dann heisst

- $\bar{x}$  anziehender Fixpunkt, falls  $\|B\| < 1$  gilt
- $\bar{x}$  abstossender Fixpunkt, falls  $\|B\| > 1$  gilt.

**Satz 4.5: Abschätzungen [1]**

- Gegeben sei wie in obiger Definition eine Fixpunktiteration

$$x^{(n+1)} = Bx^{(n)} + b =: F(x^{(n)})$$

und  $\bar{x} \in \mathbb{R}^n$  sei ein bezüglich der Norm  $\|\cdot\|$  anziehender Fixpunkt. Dann konvergiert die Fixpunktiteration für alle Startvektoren  $x^{(0)} \in \mathbb{R}^n$  gegen  $\bar{x}$  und es gelten die Abschätzungen

$$\begin{aligned} \|x^{(n)} - \bar{x}\| &\leq \frac{\|B\|^n}{1 - \|B\|} \|x^{(1)} - x^{(0)}\| \quad \text{a-priori Abschätzung} \\ \|x^{(n)} - \bar{x}\| &\leq \frac{\|B\|}{1 - \|B\|} \|x^{(n)} - x^{(n-1)}\| \quad \text{a-posteriori Abschätzung} \end{aligned}$$

Bemerkung: Der Vergleich mit den Definitionen für das Gesamt- und Einzelschrittverfahren liefert die Matrix  $B$ :

- für das Gesamtschrittverfahren (Jacobi) ist

$$B = -D^{-1}(L + R),$$

- für das Einzelschrittverfahren (Gauss-Seidel) ist

$$B = -(D + L)^{-1}R.$$

**Definition 4.8: Diagonaldominanz [1]**

- $A$  ist eine **diagonaldominante Matrix**, falls eines der beiden folgenden Kriterien gilt:
  - für alle  $i = 1, \dots, n$ :  $|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}|$  (Zeilensummenkriterium)
  - für alle  $j = 1, \dots, n$ :  $|a_{jj}| > \sum_{i=1, i \neq j}^n |a_{ij}|$  (Spaltensummenkriterium)

**Satz 4.6: Konvergenz [1]**

- Falls  $A$  diagonaldominant ist, konvergiert das Gesamtschrittverfahren (Jacobi) und auch das Einzelschrittverfahren (Gauss-Seidel) für  $Ax = b$ .

**Bemerkungen:**

- Die Bedingung  $\|B\| < 1$  für einen anziehenden Fixpunkt  $\bar{x}$  impliziert, dass  $A$  diagonaldominant ist.
- Diagonaldominanz ist nur ein hinreichendes Kriterium. Es gibt durchaus nicht diagonaldominante Matrizen, für die die Verfahren trotzdem konvergieren kann. Ein notwendiges und hinreichendes Kriterium ist, dass der Spektralradius  $\rho(B) < 1$ .

**Definition 5.1: Funktionen mit mehreren Variablen**

- Unter einer Funktion mit  $n$  unabhängigen Variablen  $x_1, \dots, x_n$  und einer abhängigen Variablen  $y$  versteht man eine Vorschrift, die jedem geordneten Zahlentupel  $(x_1, x_2, \dots, x_n)$  aus einer Definitionsmenge  $D \subset \mathbb{R}^n$  genau ein Element  $y$  aus einer Wertemenge  $W \subset \mathbb{R}$  zuordnet. Symbolische Schreibweise:

$$\begin{aligned} f : D \subset \mathbb{R}^n &\longrightarrow W \subset \mathbb{R} \\ (x_1, x_2, \dots, x_n) &\mapsto y = f(x_1, x_2, \dots, x_n) \end{aligned}$$

- Da das Ergebnis  $y \in \mathbb{R}$  ein Skalar (eine Zahl) ist, redet man auch von einer **skalarwertigen** Funktion.

**Bemerkungen:**

- Die obige Definition lässt sich einfach erweitern auf beliebige **vektorwertige** Funktionen, die nicht einen Skalar, sondern einen Vektor als Wert zurückgeben:

$$f : \mathbb{R}^n \longrightarrow \mathbb{R}^m,$$

mit

$$f(x_1, \dots, x_n) = \begin{pmatrix} y_1 = f_1(x_1, x_2, \dots, x_n) \\ y_2 = f_2(x_1, x_2, \dots, x_n) \\ \vdots \\ y_m = f_m(x_1, x_2, \dots, x_n) \end{pmatrix},$$

wobei die  $m$  Komponenten  $f_i : \mathbb{R}^n \longrightarrow \mathbb{R}$  ( $i = 1, \dots, m$ ) von  $f$  wieder skalarwertige Funktionen sind, entsprechend Def. 5.1.

- Wie bei einem Vektor  $x$  stellen wir zur besseren Unterscheidbarkeit vektorwertige Funktionen  $f$  fett gedruckt dar, im Gegensatz zu einem Skalar  $x$  und einer skalarwertigen Funktion  $f$ .
- Wir werden uns bei der Lösung nichtlinearer Gleichungssysteme auf vektorwertige Funktionen  $f : \mathbb{R}^n \longrightarrow \mathbb{R}^n$  konzentrieren.

**Definition 5.2 [8]: Partielle Ableitungen 1. Ordnung**

Unter den partiellen Ableitungen 1. Ordnung einer Funktion  $z = f(x, y)$  und der Stelle  $(x, y)$  werden die folgenden Grenzwerte verstanden (falls sie vorhanden sind):

- Partielle Ableitung 1. Ordnung nach  $x$ :

$$\frac{\partial f}{\partial x}(x, y) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x, y) - f(x, y)}{\Delta x}$$

- Partielle Ableitung 1. Ordnung nach  $y$ :

$$\frac{\partial f}{\partial y}(x, y) = \lim_{\Delta y \rightarrow 0} \frac{f(x, y + \Delta y) - f(x, y)}{\Delta y}$$

**Definition 5.3: Jacobi-Matrix / Linearisierung / Tangentialebene**

- Sei  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  mit  $\mathbf{y} = \mathbf{f}(\mathbf{x}) = \begin{pmatrix} y_1 = f_1(\mathbf{x}) \\ y_2 = f_2(\mathbf{x}) \\ \vdots \\ y_m = f_m(\mathbf{x}) \end{pmatrix}$  und  $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$ . Die **Jacobi-Matrix** enthält sämtliche partiellen Ableitung 1. Ordnung von  $\mathbf{f}$  und ist definiert als

$$D\mathbf{f}(\mathbf{x}) := \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \frac{\partial f_1}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \frac{\partial f_2}{\partial x_1}(\mathbf{x}) & \frac{\partial f_2}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_2}{\partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \frac{\partial f_m}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{pmatrix}$$

- Die “verallgemeinerte Tangentengleichung”

$$\mathbf{g}(\mathbf{x}) = \mathbf{f}(\mathbf{x}^{(0)}) + D\mathbf{f}(\mathbf{x}^{(0)}) \cdot (\mathbf{x} - \mathbf{x}^{(0)})$$

beschreibt eine lineare Funktion und es gilt  $\mathbf{f}(\mathbf{x}) \approx \mathbf{g}(\mathbf{x})$  in einer Umgebung eines gegebenen Vektors  $\mathbf{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})^T \in \mathbb{R}^n$ . Man spricht deshalb auch von der **Linearisierung** der Funktion  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  in einer Umgebung von  $\mathbf{x}^{(0)}$  (ein hochgestellter Index in Klammern  $\mathbf{x}^{(k)}$  bezeichnet wie bisher einen Vektor aus  $\mathbb{R}^n$  nach der  $k$ -ten Iteration).

- Für den speziellen Fall  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  mit  $y = f(x_1, x_2)$  liefert die linearisierte Funktion

$$g(x_1, x_2) = f(x_1^{(0)}, x_2^{(0)}) + \frac{\partial f}{\partial x_1}(x_1^{(0)}, x_2^{(0)}) \cdot (x_1 - x_1^{(0)}) + \frac{\partial f}{\partial x_2}(x_1^{(0)}, x_2^{(0)}) \cdot (x_2 - x_2^{(0)})$$

die Gleichung der **Tangentialebene**. Sie enthält sämtliche im Flächenpunkt  $P = (x_1^{(0)}, x_2^{(0)}, f(x_1^{(0)}, x_2^{(0)}))$  an die Bildfläche von  $y = f(x_1, x_2)$  angelegten Tangenten.

**Newton-Verfahren für Systeme [1]:**

Gesucht sind Nullstellen von  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . Sei  $\mathbf{x}^{(0)}$  ein Startvektor in der Nähe einer Nullstelle. Das Newton-Verfahren zur näherungsweisen Bestimmung dieser Nullstelle lautet:

- für  $n = 0, 1, \dots$ :

- Berechne  $\delta^{(n)}$  als Lösung des linearen Gleichungssystems

$$D\mathbf{f}(\mathbf{x}^{(n)})\delta^{(n)} = -\mathbf{f}(\mathbf{x}^{(n)})$$

- Setze

$$\mathbf{x}^{(n+1)} := \mathbf{x}^{(n)} + \delta^{(n)}$$

**Vereinfachtes Newton-Verfahren für Systeme [1]:**

Gesucht sind Nullstellen von  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . Sei  $\mathbf{x}^{(0)}$  ein Startvektor in der Nähe einer Nullstelle. Das vereinfachte Newton-Verfahren zur näherungsweisen Bestimmung dieser Nullstelle lautet:

- für  $n = 0, 1, \dots$  :

- Berechne  $\boldsymbol{\delta}^{(n)}$  als Lösung des linearen Gleichungssystems

$$D\mathbf{f}(\mathbf{x}^{(0)})\boldsymbol{\delta}^{(n)} = -\mathbf{f}(\mathbf{x}^{(n)})$$

- Setze

$$\mathbf{x}^{(n+1)} := \mathbf{x}^{(n)} + \boldsymbol{\delta}^{(n)}$$

**Gedämpftes Newton-Verfahren für Systeme [6]:**

Gesucht sind Nullstellen von  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . Sei  $\mathbf{x}^{(0)}$  ein Startvektor in der Nähe einer Nullstelle,  $k_{max} \in \mathbb{N}$  sei vorgegeben. Das gedämpfte Newton-Verfahren zur näherungsweisen Bestimmung dieser Nullstelle lautet:

- für  $n = 0, 1, \dots$  :

- Berechne  $\boldsymbol{\delta}^{(n)}$  als Lösung des linearen Gleichungssystems

$$D\mathbf{f}(\mathbf{x}^{(n)})\boldsymbol{\delta}^{(n)} = -\mathbf{f}(\mathbf{x}^{(n)})$$

- Finde das minimale  $k \in \{0, 1, \dots, k_{max}\}$  mit

$$\left\| \mathbf{f} \left( \mathbf{x}^{(n)} + \frac{\boldsymbol{\delta}^{(n)}}{2^k} \right) \right\|_2 < \left\| \mathbf{f} \left( \mathbf{x}^{(n)} \right) \right\|_2$$

- Falls kein minimales  $k$  gefunden werden kann, rechne mit  $k = 0$  weiter
- Setze

$$\mathbf{x}^{(n+1)} := \mathbf{x}^{(n)} + \frac{\boldsymbol{\delta}^{(n)}}{2^k}$$

