

**Mitschrift**

# **Numerische Methoden der Elektrotechnik**

**TUM EI M.Sc. Wahlpflichtmodul**

Markus Hofbauer

Kevin Meyer

Benedikt Schmidt

SS 2014

**Dozent**

Prof. Dr.-Ing. Ulf Schlichtmann



# Liste der noch zu erledigenden Punkte

Abbildung: plot . . . . .	44
Abbildung: plot . . . . .	44
Abbildung: plot . . . . .	45
Abbildung: plot . . . . .	46
Abbildung: plot . . . . .	46



# Inhaltsverzeichnis

<b>1</b>	<b>Einführung</b>	<b>1</b>
1.1	Zusammenfassung . . . . .	1
1.1.1	Überblick . . . . .	1
1.1.2	Vorgehen . . . . .	1
1.1.3	Aspekte für die Lösung einer mathematischen Aufgabenstellung . . . .	1
1.2	Wiederholung . . . . .	3
1.2.1	Lineares Gleichungssystem . . . . .	3
1.2.2	Matrizen . . . . .	3
1.3	Schaltungsanalyse . . . . .	6
1.3.1	Gerichteter Graph der Schaltung . . . . .	7
1.3.2	Inzidenz Matrix . . . . .	7
1.3.3	Laplace Matrix . . . . .	7
1.3.4	Kirchhoffsches Stromgesetz (KCL) . . . . .	8
1.3.5	Ohmsches Gesetz . . . . .	8
1.3.6	Kirchhoffsches Spannungsgesetz (KVL) . . . . .	8
1.3.7	Alternative Darstellung . . . . .	9
1.4	Maschinendarstellung von Zahlen . . . . .	9
1.4.1	Gleitkomma-Arithmetik: Eigenschaften und Fehlerarten . . . . .	9
<b>2</b>	<b>Nichtlineare Gleichungen</b>	<b>11</b>
2.1	Intervallhalbierung . . . . .	11
2.2	Newton-Raphson-Verfahren . . . . .	12
2.2.1	Variation: Sekanten-Methode . . . . .	14
2.3	Fixpunktiteration . . . . .	14
2.3.1	Fixpunkttheorem . . . . .	15
2.3.2	Konvergenzverhalten iterativer Verfahren . . . . .	15
2.4	Mehrdimensionales Newton-Raphson-Verfahren . . . . .	16
2.4.1	Konvergenz . . . . .	16
2.4.2	Abbruchkriterien . . . . .	16
<b>3</b>	<b>Lineare Gleichungssysteme</b>	<b>19</b>
3.1	Methode zur Bestimmung von $A^{-1}$ . . . . .	19
3.2	Kondition eines Gleichungssystems . . . . .	20
3.3	Einfluss der Pivotisierung auf die Ergebnisgenauigkeit . . . . .	20
3.4	Kondition einer Matrix . . . . .	21
3.5	Allgemeines Iterationsverfahren für lineare Gleichungssysteme . . . . .	22
3.5.1	Iterative Verfahren zum Lösen linearer Gleichungssysteme . . . . .	22
3.5.2	Verschiedene Iterationsverfahren . . . . .	23
3.5.3	Gradientenverfahren . . . . .	25

3.6	Iterativer Ansatz für Minimierungsproblem . . . . .	25
3.6.1	Konvergenz . . . . .	27
3.7	CG - Konjugierte Gradientenmethode . . . . .	27
3.7.1	Schrittweite . . . . .	28
3.7.2	Suchrichtung . . . . .	28
3.7.3	Zusammenfassung des Verfahrens . . . . .	29
3.7.4	Konvergenz . . . . .	29
3.7.5	Vorkonditionierung . . . . .	29
<b>4</b>	<b>Interpolation</b>	<b>31</b>
4.1	Lagrange Interpolationspolynom . . . . .	31
4.2	Dividierte Differenzen . . . . .	32
4.3	Spline - Interpolation . . . . .	33
<b>5</b>	<b>Numerische Infinitesimalrechnung</b>	<b>37</b>
5.1	Numerische Differentiation . . . . .	37
5.1.1	Vorwärtsdifferenz . . . . .	37
5.1.2	Rückwärtsdifferenz . . . . .	37
5.1.3	Zentrierte Differenz . . . . .	37
5.2	Numerische Integration . . . . .	38
5.2.1	Trapezregel . . . . .	38
5.2.2	Simpson Regel . . . . .	39
5.2.3	Genauigkeit . . . . .	41
<b>6</b>	<b>Numerische Lösung von Differentialgleichungen</b>	<b>43</b>
6.1	Analytische Lösung mittels Laplace-Transformation . . . . .	43
6.2	Numerische Lösung mittels explizite Euler-Methode (linear Z-Transformation)	44
6.3	Einfache numerische Integrationsverfahren . . . . .	45
6.4	Eigenschaften numerischer Integrationsverfahren . . . . .	46
6.5	Expliziter Euler Eigenschaften . . . . .	47
6.6	Taylor-Verfahren höherer Ordnung . . . . .	48
6.7	Prädiktor-Korrektor-Ansätze . . . . .	48
6.8	Adams-Bashforth/Adams-Moulton . . . . .	49
6.8.1	Adams-Bashforth (explizit) . . . . .	49
6.8.2	Adams-Moulton (implizit) . . . . .	50
<b>7</b>	<b>Ausgleichsrechnung</b>	<b>51</b>
7.1	Linear Least Squares . . . . .	51
7.2	Polynomiale Least Squares . . . . .	51
7.3	Linear Least Squares - Sichtweise Lineare Algebra . . . . .	52
7.4	QR-Zerlegung . . . . .	53
7.4.1	Givens-Rotation . . . . .	53

<b>Abbildungsverzeichnis</b>	<b>i</b>
------------------------------	----------

# 1 Einführung

## 1.1 Zusammenfassung

### 1.1.1 Überblick

**Definitionsversuch:** Entwicklung und mathematisches Verständnis von numerischen Algorithmen, als von Rechenmethoden zur zahlenmäßigen Lösung mathematischer Probleme.

**Zusammenspiel mit Informatik:**

### 1.1.2 Vorgehen

### 1.1.3 Aspekte für die Lösung einer mathematischen Aufgabenstellung

- Kondition eines Problems (Empfindlichkeit für Störungen)
- Numerische Lösungsverfahren
- Stabilität des Lösungsverfahrens (Empfindlichkeit für Störungen)
- Effizienz des Lösungsverfahrens
- Genauigkeit der Lösung

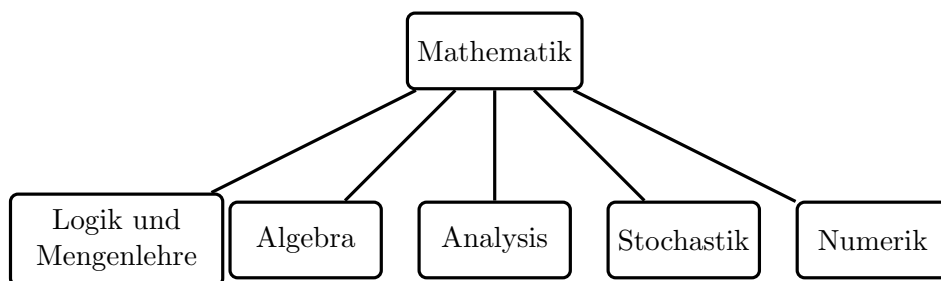


Abbildung 1.1: Teilgebiete der Mathematik

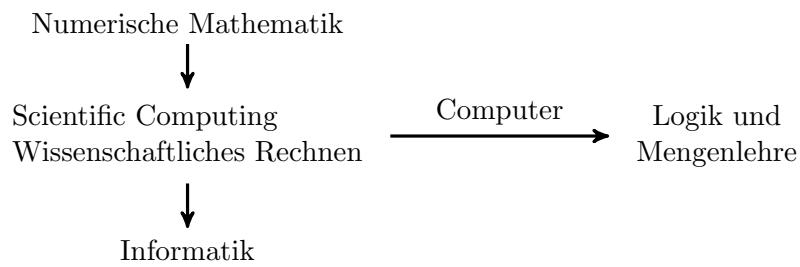


Abbildung 1.2: Zusammenspiel Mathematik und Information in der Numerik

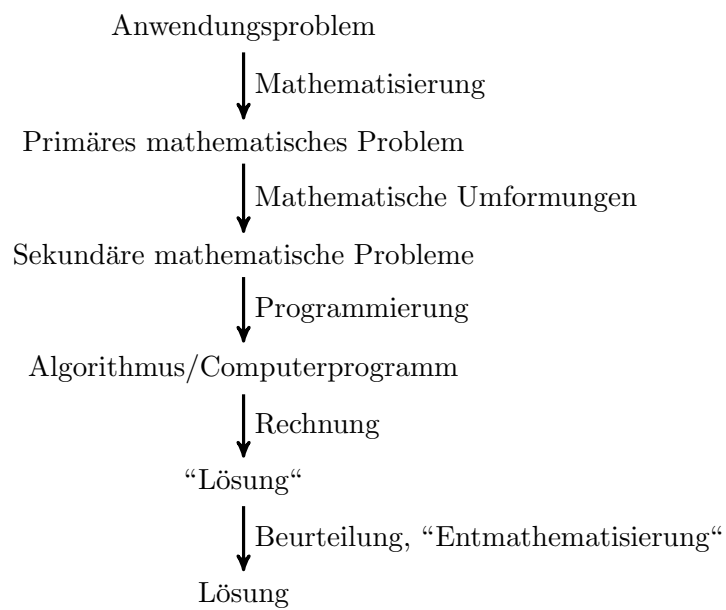


Abbildung 1.3: Vorgehen zur Lösung eines mathematischen Problems



## 1.2 Wiederholung

### 1.2.1 Lineares Gleichungssystem

$$a_{11}x_1 + \dots + a_{1n}x_n = b_1 \quad (1.1)$$

$$a_{21}x_1 + \dots + a_{2n}x_n = b_2 \quad (1.2)$$

$$\vdots \quad (1.3)$$

$$a_{m1}x_1 + \dots + a_{mn}x_n = b_m \quad (1.4)$$

$$\mathbf{A} \cdot \vec{x} = \vec{b} \quad (1.5)$$

$$\mathbf{A} \in \mathbb{R}^{m \times n}, \vec{x} \in \mathbb{R}^n, \vec{b} \in \mathbb{R}^m \quad (1.6)$$

### 1.2.2 Matrizen

#### Elementare Matrix-Operationen

##### Addition

$$\mathbf{A}_{<m \times n>} + \mathbf{B}_{<m \times n>} = \mathbf{C}_{<m \times n>} \quad (1.7)$$

$$c_{\mu\nu} = a_{\mu\nu} + b_{\mu\nu} \quad (1.8)$$

##### Multiplikation mit Skalar

$$k \cdot \mathbf{A} = \mathbf{A} \cdot k = \mathbf{B} \quad (1.9)$$

$$b_{\mu\nu} = k \cdot a_{\mu\nu} \quad (1.10)$$

##### Matrix-Multiplikation

$$\mathbf{A}_{<m \times n>} \cdot \mathbf{B}_{<n \times l>} = \mathbf{C}_{<m \times l>} \quad (1.11)$$

$$c_{\mu\nu} = \sum_{\lambda=1}^n a_{\mu\lambda} \cdot b_{\lambda\nu}; \quad \mu = 1 \dots m, \lambda = 1 \dots n \quad (1.12)$$

##### Multiplikation Matrix-Vektor

$$\text{Spezialfall der Matrix-Multiplikation: } <n \times 1> \text{ bzw. } <1 \times n> \quad (1.13)$$

##### Auffassen als Linearkombination der Spalten der Matrix

$$\mathbf{A} \cdot \vec{x} = \vec{b} \quad (1.14)$$

$$[\vec{a}_1 \ \vec{a}_2 \ \dots \ \vec{a}_n] \cdot \vec{x} = \vec{b} \quad (1.15)$$

$$\vec{a}_1 \vec{x} + \dots + \vec{a}_n \vec{x}_n = \vec{b} \quad (1.16)$$

## Vektoren

$$\text{Spaltenvektor: } \vec{a}_{<n \times 1>} : \vec{a} \quad (1.17)$$

$$\text{Zeilenvektor: } \vec{b}_{<1 \times m>}^T : \vec{b}^T \quad (1.18)$$

## Vektormultiplikation

$$\vec{b}^T \cdot \vec{a} = c \quad (\text{Skalarprodukt } m = n!) \quad (1.19)$$

$$\vec{a} \cdot \vec{b}^T = \mathbf{C} \quad (\text{Vektorprodukt, dyadisches Produkt}) \quad (1.20)$$

## Matrix als Kombination von Zeilen- und Spaltenvektoren

$$\mathbf{A}_{<m \times n>} = \begin{bmatrix} \vec{a}_1^T \\ \vec{a}_2^T \\ \vdots \\ \vec{a}_m^T \end{bmatrix} \quad (1.21)$$

$$\mathbf{B}_{<n \times l>} = \begin{bmatrix} \vec{b}_1 & \vec{b}_2 & \dots & \vec{b}_l \end{bmatrix} \quad (1.22)$$

$$\mathbf{A} \cdot \mathbf{B} = \mathbf{C} \quad (1.23)$$

$$\vec{c}_{\mu\lambda} = \vec{a}_\mu^T \cdot \vec{b}_\lambda \quad (1.24)$$

## Rechenregeln

$$(\mathbf{A} \cdot \mathbf{B}) \cdot \mathbf{C} = \mathbf{A} \cdot (\mathbf{B} \cdot \mathbf{C}) \quad \text{Assoziativitt} \quad (1.25)$$

$$\mathbf{A} \cdot (\mathbf{B} + \mathbf{C}) = \mathbf{A} \cdot \mathbf{B} + \mathbf{A} \cdot \mathbf{C} \quad \text{Distributivitt} \quad (1.26)$$

$$\mathbf{A} \cdot \mathbf{B} \neq \mathbf{B} \cdot \mathbf{A} \quad \text{Kommutativitt gilt im Allgemeinen nicht} \quad (1.27)$$

## Diagonalmatrix

$$\mathbf{D} = \text{diag} [d_1 \quad d_2 \quad \dots \quad d_n] = \begin{bmatrix} d_1 & \dots & 0 \\ 0 & \dots & d_n \end{bmatrix} \quad (1.28)$$

$$\mathbf{D}_1 \cdot \mathbf{D}_2 = \mathbf{D}_2 \cdot \mathbf{D}_1 \quad (1.29)$$

## Einheitsmatrix

$$\mathbf{I} = \mathbf{E} = \mathbf{1} = \text{diag} [1 \quad \dots \quad 1] \quad (1.30)$$

$$\mathbf{A} \cdot \mathbf{I} = \mathbf{I} \cdot \mathbf{A} \quad (1.31)$$

$$\mathbf{I}^n = \mathbf{I} \quad (1.32)$$

$$\mathbf{I}^{-1} = \mathbf{I} \quad (1.33)$$

**Transponierte Matrix**

$$\mathbf{A}_{<m \times n>}^T = \mathbf{B}_{<n \times m>}, \quad b_{\nu\mu} = a_{\mu\nu} \quad (1.34)$$

$$(\mathbf{A} \cdot \mathbf{B})^T = \mathbf{B}^T \cdot \mathbf{A}^T \quad (1.35)$$

$$\mathbf{A} = \mathbf{A}^T \Rightarrow \text{symmetrische Matrix} \quad (1.36)$$

**Inverse Matrix**

$$\mathbf{A} \in \mathbb{R}^{n \times n} \quad (1.37)$$

$$\mathbf{A}^{-1} \cdot \mathbf{A} = \mathbf{A} \cdot \mathbf{A}^{-1} = \mathbf{I} \quad (1.38)$$

$$\mathbf{A}^{-1} \text{ existiert nur f\"ur nicht singul\"are } \mathbf{A} \text{ und ist eindeutig.} \quad (1.39)$$

$$\mathbf{A} = \text{diag} [d_1 \quad \dots \quad d_n] \Rightarrow \mathbf{A}^{-1} = \text{diag} [\frac{1}{d_1} \quad \dots \quad \frac{1}{d_n}] \quad (1.40)$$

$$(\mathbf{A} \cdot \mathbf{B}) = \mathbf{B}^{-1} \cdot \mathbf{A}^{-1} \quad (1.41)$$

**Matrix- und Vektornormen**

$$\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}; \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$$

Bedingungen f\"ur **Vektornormen**:

•

$$\|\vec{x}\| \geq 0, \quad \|\vec{x}\| = 0 \text{ nur f\"ur } \vec{x} = \vec{0}$$

•

$$\|c \cdot \vec{x}\| = |c| \cdot \|\vec{x}\|, \quad c \in \mathbb{R}$$

•

$$\|\vec{x} + \vec{y}\| \leq \|\vec{x}\| + \|\vec{y}\|$$

Dreiecksungleichung

Bedingungen f\"ur **Matrixnormen**:

•

$$\|\mathbf{A}\| \geq 0, \quad \|\mathbf{A}\| = 0 \text{ nur f\"ur } \mathbf{A} = \mathbf{0}$$

•

$$\|c \cdot \mathbf{A}\| = |c| \cdot \|\mathbf{A}\|, \quad c \in \mathbb{R}$$

•

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$$

•

$$\|\mathbf{A} \cdot \mathbf{B}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|$$

Multiplikativit\"atsbedingung

•

$$\|\mathbf{A} \cdot \vec{x}\| \leq \|\mathbf{A}\| \cdot \|\vec{x}\|$$

Kompatibilitätsbedingung

### Vektornormen

$$\begin{aligned} \|\vec{x}\|_1 &= \sum_{i=1}^n |x_i| && \text{Betragssummennorm, } l_1\text{-Norm} \\ \|\vec{x}\|_2 &= \sqrt{\sum_{i=1}^n |x_i|^2} && \text{Euklidnorm, } l_2\text{-Norm, Vektorlänge} \\ \|\vec{x}\|_\infty &= \max_i |x_i| && \text{Maximumsnorm, } l_\infty\text{-Norm, Tschebychefnorm} \\ \|\vec{x}\|_p &= \sqrt[p]{\sum_{i=1}^n |x_i|^p}, p \geq 1 && \text{Höldernormnorm, } l_p\text{-Norm} \end{aligned}$$

### Matrixnormen ( $A \in \mathbb{R}^{n \times n}$ )

$$\begin{aligned} \|\mathbf{A}\|_M &= n \cdot \max_{i,j} |a_{ij}| && \text{Gesamtnorm, Matrixnorm } (\|\mathbf{I}\| = n) \\ (\|\mathbf{A}\|_\infty =) \|\mathbf{A}\|_Z &= \max_i \sum_{j=1}^n |a_{ij}| && \text{Zeilenorm } (\|\mathbf{I}\|_Z = 1) \\ (\|\mathbf{A}\|_1 =) \|\mathbf{A}\|_S &= \max_j \sum_{i=1}^n |a_{ij}| && \text{Spaltennorm } (\|\mathbf{I}\|_S = 1) \\ \|\mathbf{A}\|_E &= \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2} && \text{Euklidnorm, Schurnorm, Frobeniusnorm } (\|\mathbf{I}\|_E = \sqrt{n}) \\ (\mathbf{A}\|_\lambda =) \|\mathbf{A}\|_\lambda &= \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} && \text{Spektralnrm, Hilbertnorm } (\|\mathbf{I}\|_\lambda = 1) \end{aligned}$$

### Kompatibilität zwischen Vektor und Matrixnorm

$$\begin{aligned} l_1 : \quad & \|\mathbf{A}\vec{x}\|_1 \leq \|\mathbf{A}\|_M \cdot \|\vec{x}\|_1 \\ & \|\mathbf{A}\vec{x}\|_1 \leq \|\mathbf{A}\|_S \cdot \|\vec{x}\|_1 \\ l_2 : \quad & M, \lambda, E \\ l_\infty : \quad & M, Z \end{aligned}$$

## 1.3 Schaltungsanalyse

$$\begin{aligned} \vec{i}_{<k>} & \quad \text{Kantenstromvektor} & \vec{u}_{<k>} & \quad \text{Kantenspannungsvektor} \\ \vec{i}_{0<k>} & \quad \text{Kantenquellenstromvektor} & \vec{u}_{0<k>} & \quad \text{Kantenquellenspannungsvektor} \\ \vec{i}_{n<k>} & \quad \text{Knotenquellenstromvektor} & \vec{u}_{n<k>} & \quad \text{Knotenquellenspannungsvektor} \\ \mathbf{A}_{<n \times k>} & \quad \text{Kontenmatrix, Knoteninzidenzmatrix} \\ \mathbf{Y}_{<k \times k>} & \quad \text{Kantenadmittanzmatrix} \\ \mathbf{Y}_{n<n \times n>} & \quad \text{Kontenadmittanzmatrix} \end{aligned}$$

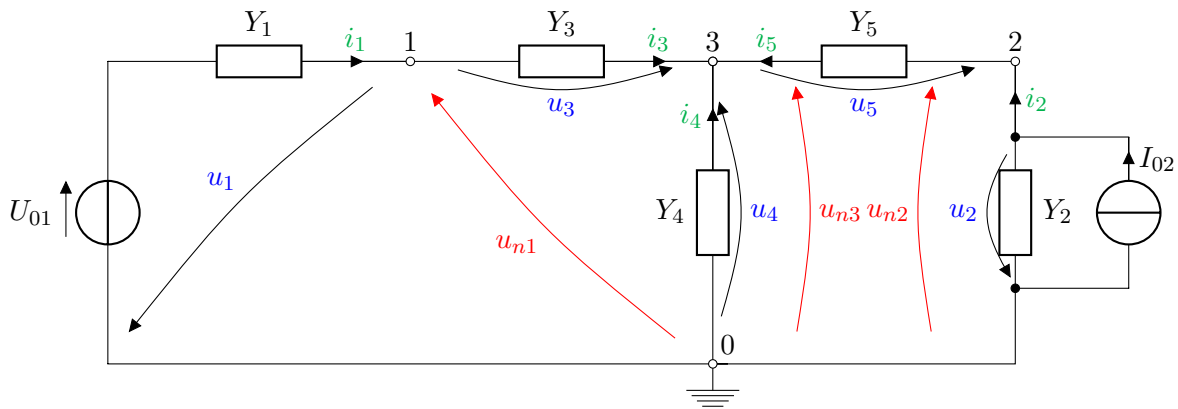


Abbildung 1.4: Beispielschaltung

### 1.3.1 Gerichteter Graph der Schaltung

### 1.3.2 Inzidenz Matrix

$$\mathbf{A}_{<m \times n>} = \begin{bmatrix} 1 & & & -1 \\ & 1 & & -1 \\ -1 & & 1 & \\ & & 1 & -1 \\ & -1 & 1 & \end{bmatrix} \quad (1.42)$$

Summen der Spaltenvektoren =  $\vec{0}$ .  
 $\Rightarrow \mathbf{A}$  hat linear abhängige Spalten.

Rang der Matrix  $\mathbf{A}$ :  $r = \text{rang}(\mathbf{A}) = 3 = n - 1$ .

Dimension des Nullraums: Zahl der Spalten -  $r = 1$ .

Vektor im Nullraum von  $\mathbf{A}$ :

$$\mathbf{A} \cdot \vec{u} = \vec{0} \Rightarrow \vec{u} \in; \quad \vec{u} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad (1.43)$$

### 1.3.3 Laplace Matrix

$$\mathbf{A}^T \cdot \mathbf{A} = \begin{bmatrix} 1 & 0 & -1 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 \\ 0 & 0 & 1 & 1 & 1 \\ -1 & -1 & 0 & -1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & -1 \\ -1 & 0 & 1 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & -1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 & -1 & -1 \\ 0 & 2 & -1 & -1 \\ -1 & -1 & 3 & -1 \\ -1 & -1 & 3 & -1 \end{bmatrix} \quad (1.44)$$

- singular

- $r = 3$
- symmetrisch

$$\mathbf{A}^T \cdot \mathbf{A} = \underset{\text{Grad (degree)}}{\mathbf{D}} - \underset{\text{Adjazenz}}{\mathbf{W}} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 3 \end{bmatrix} - \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix} \quad (1.45)$$

### 1.3.4 Kirchhoffsches Stromgesetz (KCL)

$$\begin{aligned} \mathbf{A}^T \cdot \vec{w} &= \vec{0} & (\text{keine Stromquellen}) & \quad \vec{w}: \text{Kantenströme} \\ \mathbf{A}^T \cdot \vec{w} &= \vec{f} & (\text{mit Stromquellen}) & \quad \vec{f}: \text{Stromquellen} \end{aligned}$$

### 1.3.5 Ohmsches Gesetz

$$\vec{w} = \mathbf{C} \cdot \vec{e} \quad \begin{array}{l} \mathbf{C}: \text{Diagonalmatrix der Kantenleitwerte} \\ \vec{e}: \text{Kantenspannungen} \end{array}$$

### 1.3.6 Kirchhoffsches Spannungsgesetz (KVL)

$$\vec{e} = \vec{b} - \mathbf{A} \cdot \vec{u} \quad \begin{array}{l} \vec{u}: \text{Knotenspannungen - GESUCHT} \\ \vec{b}: \text{Spannungsquellen} \end{array}$$

$$\mathbf{A}^T \cdot \vec{w} = \vec{f} \quad (1.46)$$

$$\mathbf{A}^T \cdot \mathbf{C} \cdot \vec{e} = \vec{f} \quad (1.47)$$

$$\mathbf{A}^T \cdot \mathbf{C} (\vec{b} - \mathbf{A} \cdot \vec{u}) = \vec{f} \quad (1.48)$$

$$\mathbf{A}^T \cdot \mathbf{C} \cdot \mathbf{A} \quad \vec{u} = \mathbf{A}^T \cdot \mathbf{C} \cdot \vec{b} - \vec{f} \quad (1.49)$$

Systemmatrix, gewichtete Laplace Matrix

$$\mathbf{Y}_{<n \times n>} \cdot \vec{u}_{<n>} = \vec{d}_{<n>} \quad (1.50)$$

singulär  $\Rightarrow$  nicht invertierbar

$\Rightarrow$  keine Lösung für Gleichungssystem

Abhilfe: Festlegung eines Bezugspunktes, z.B.  $u_0 = 0$

...

### 1.3.7 Alternative Darstellung

$$C^{-1} \cdot \vec{w} + A \cdot \vec{u} = \vec{b} \quad (1.51)$$

$$A^T \cdot \vec{w} = \vec{f} \quad (1.52)$$

$$\begin{bmatrix} C^{-1} & A \\ A^T & \mathbf{0} \end{bmatrix} \cdot \begin{bmatrix} \vec{w} \\ \vec{u} \end{bmatrix} = \begin{bmatrix} \vec{b} \\ \vec{f} \end{bmatrix} \quad (1.53)$$

## 1.4 Maschinendarstellung von Zahlen

(IEEE 754-2008-“Binary Floating Point Arithmetic Standard“)

Darstellung reeller Zahlen: 64 bit:

s 1 bit Vorzeichen  
c 11 bit Exponent  
f 52 bit Mantisse

$$x = (-1)^s \cdot 2^{c-1023} \cdot (1 + f)$$

z.B.:

$$10 : \underbrace{0}_s \quad \underbrace{100000000010}_{c: 1026-1023=3} \quad \underbrace{0100 \dots 010}_{f: 1,25} \quad (1.54)$$

$$-0,8 : \underbrace{1}_s \quad \underbrace{01111111110}_{c: 1022-1023=-1} \quad \underbrace{10011001 \dots 10011010}_{f: 1,59999} \quad (1.55)$$

### 1.4.1 Gleitkomma-Arithmetik: Eigenschaften und Fehlerarten

$x$ : beliebige reelle Zahl (mathematisch exakt)

$z_1, z_2$ : Gleitkomma-Maschinen Zahlen (endliche Stellenzahl)

**Gleitkomma - Grundoperationen**

gl( $z_1$  op  $z_2$ ), op := +, −, ×, /

- $z_1$  op  $z_2 = x$ ;  $x$  nicht notwendigerweise Gleitkomma-Maschinen Zahl
- +, −: Exponentenangleich erforderlich - Mantissenstellen können verloren gehen.
- Subtraktion nahezu gleichgroßer Zahlen:
  - Ergebnis mit wesentlich kleinerer Mantisse
  - Normalisierung (Linksverschiebung Mantisse, Exponentenangleich): Nachziehen nicht signifikanter Ziffern, Verlust signifikanter Ziffern
- Unterlauf/Überlauf

...





## 2 Nichtlineare Gleichungen

Dies ist nötig beispielsweise in der Schaltungsanalyse inklusiver nichtlinearer Elemente (z.B. Diode, Transistor).

Gesucht sei hierbei das stationäre Verhalten bei konstanter Erregung (*DC-Arbeitspunkt*).

Die Diode kann mithilfe von  $i_D = I_S \cdot \left( e^{\frac{u_D}{U_T}} - 1 \right)$  beschrieben werden. Daraus folgt dann

$$I_0 - \frac{u_D}{R} = I_S \cdot \left( e^{\frac{u_D}{U_T}} - 1 \right)$$

Gesucht ist somit die Nullstelle von

$$F(u_D) = f(u_D) - g(u_D) = 0$$

Nullstellensuchen sind meist iterative Verfahren welche aus einem Startwert  $\vec{x}$  eine Lösung  $\vec{x}^*$  mit  $F(\vec{x}^*) = 0$  ermitteln. Zuerst einmal wollen wir das Problem eindimensional betrachten:

$$F(x) = 0$$

Gegeben ist dabei ein  $F(x)$ , welches stetig auf  $[a, b]$  ist mit  $F(a) \cdot F(b) < 0$  und gesucht ist  $x^*$  mit  $F(x^*) = 0$ ,  $a \leq x^* \leq b$ .

### 2.1 Intervallhalbierung

1. Startintervall  $[a^{(0)}, b^{(0)}] = [a, b]$ ,  $k = 0$
2. Intervallmitte  $m = \frac{a^{(k)} + b^{(k)}}{2}$

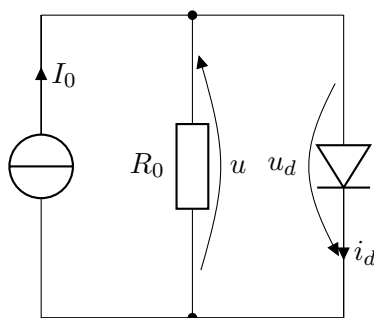


Abbildung 2.1: Beispiel einer Schaltung

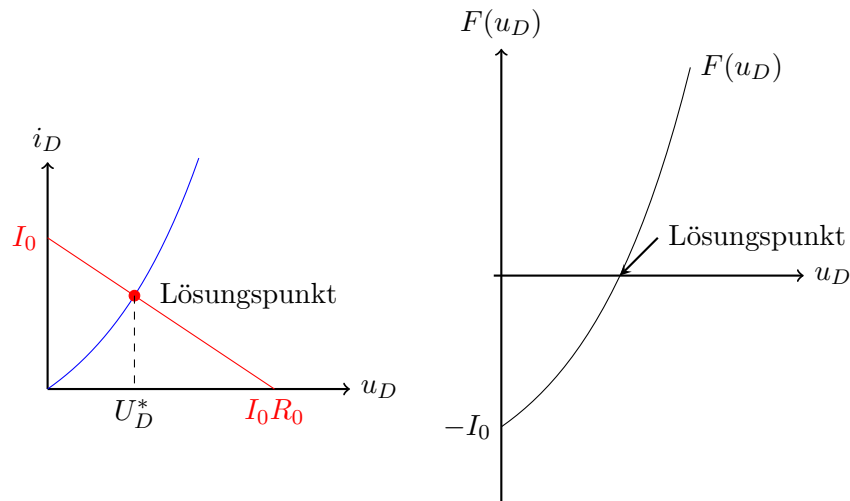


Abbildung 2.2: Arbeitspunktermittlung durch Nullstellensuche

3.  $[a^{(k+1)}, b^{(k+1)}] = \begin{cases} [m, b^{(k)}] & \text{für } F(m) \cdot F(a^{(k)}) > 0 \\ [a^{(k)}, m] & \text{für } F(m) \cdot F(a^{(k)}) < 0 \end{cases}$
4. Falls  $|a^{(k+1)} - b^{(k+1)}| > \epsilon$  dann gehe zu Schritt 2,  $k = k + 1$

Die Konvergenz eines iterativen Verfahrens wird mithilfe der Parameter  $L$ , dem Konvergenzfaktor, und  $p$ , der Konvergenzordnung, beschrieben.

$$\Delta^{(k+1)} = \frac{1}{2} \Delta^{(k)} = \left(\frac{1}{2}\right)^{k+1} \Delta^{(0)}$$

$$\epsilon = x - x^*$$

$$|\epsilon^{(k+1)}| \leq L |\epsilon^{(k)}|^p$$

Die Intervallhalbierung liegt bei  $p = 1$ , also linearer Konvergenz, und  $L = 1/2$ . Die Anzahl der Iterationsschritte bis zu einem Restfehler  $\Delta^R$  ist damit

$$\kappa = \left\lceil \lg \left( \frac{\Delta^{(0)}}{\Delta^R} \right) \right\rceil$$

## 2.2 Newton-Raphson-Verfahren

Die Idee hinter dem Newton-Raphson-Verfahren ist eine Taylorreihe erster Ordnung:

$$F(x) = F(x^{(k)}) + F'(x^{(k)}) \cdot (x - x^{(k)}) + F''(x^{(k)}) \cdot \frac{(x - x^{(k)})^2}{2} + \dots$$

$$x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)})}{F'(x^{(k)})}$$

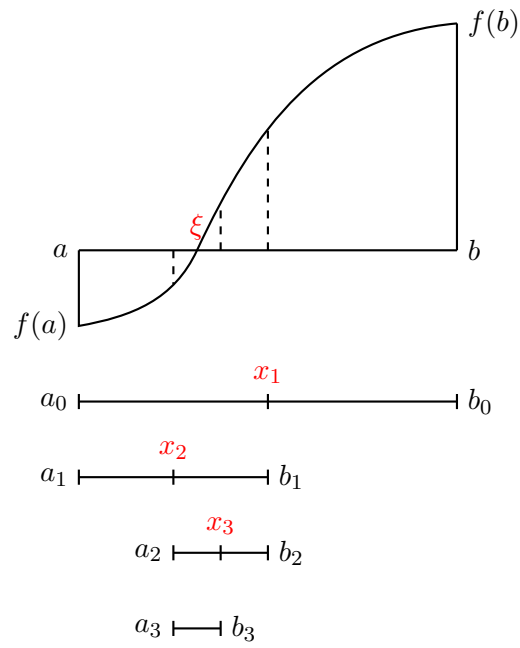


Abbildung 2.3: Intervallhalbierung

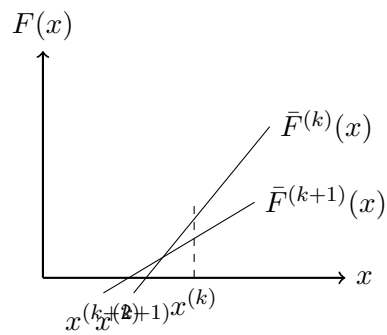


Abbildung 2.4: Newton-Iteration

Das Konvergenzverhalten des Newton-Raphson-Verfahrens erhält man aus einer Taylorreihe:

$$x^{(k+1)} = x^{(k)} - \frac{F(x^*) + F'(x^*)(x^{(k)} - x^*) + \frac{1}{2}F''(x^*) \cdot \epsilon^{k2} + \dots + \frac{1}{n!}F^{(n)}(x^*) \cdot \epsilon^{kn}}{F'(x^*) + F''(x^*) \cdot \epsilon^{(k)} + \frac{1}{(n-1)!} \cdot F^{(n)}\epsilon^{(k)n-1}}$$

$$\epsilon^{(k+1)} = \epsilon^{(k)} \left( 1 - \frac{F' + \frac{1}{2}F'' \cdot \epsilon^{(k)} + \dots}{F' + F'' \cdot \epsilon^{(k)} + \dots} \right) \approx \epsilon^{(k)} \left( \left( 1 + \frac{1}{2}\epsilon^{(k)} \frac{F''}{F'} \right) \cdot \left( 1 - \frac{1}{2}\epsilon^{(k)} \frac{F''}{F'} \right) \right)$$

$$\epsilon^{(k+1)} \approx \epsilon^{(k)} \left( 1 - \left( 1 + \frac{1}{2} \cdot \epsilon^{(k)} \frac{F''}{F'} - \epsilon^{(k)} \frac{F''}{F'} \right) \right) = \frac{1}{2} \cdot \epsilon^{(k)2} \frac{F''}{F'}$$

Probleme beim Newton-Raphson-Verfahren treten bei n-fachen Nullstellen ( $F'(x^*) = \dots = F^{(n-1)}(x^*) = 0$ ,  $F^{(n)}(x^*) \neq 0$ ) auf.

$$\epsilon^{(k+1)} \approx \epsilon^{(k)} \cdot \left( 1 - \frac{(n-1)!}{n!} \right) = \epsilon^{(k)} \cdot \left( 1 - \frac{1}{n} \right) = \epsilon^{(k)} \frac{n-1}{n}$$

Somit erhält man nur mehr lineare Konvergenz. z.B.:  $F(x) = e^x - 1 \Rightarrow x^{(k+1)} = x^{(k)} - \frac{e^{x^{(k)}} - 1}{e^{x^{(k)}}}$

### 2.2.1 Variation: Sekanten-Methode

$$F'(x^{(k)}) = \lim_{x \rightarrow x^{(k)}} \frac{F(x) - F(x^{(k)})}{x - x^{(k)}}$$

$$F'(x^{(k)}) = \frac{F(x^{(k-1)}) - F(x^{(k)})}{x^{(k-1)} - x^{(k)}}$$

Eingesetzt in die Newton-Gleichung erhält man somit

$$x^{(k+1)} = x^{(k)} - \frac{F(x^{(k)}) \cdot (x^{(k)} - x^{(k-1)})}{F(x^{(k)}) - F(x^{(k-1)})}$$

Dadurch ist pro Schritt nur eine Funktionsauswertung nötig.

## 2.3 Fixpunktiteration

Fixpunkt  $x^*$  einer Funktion  $g(x)$

$$x^* = g(x^*) \tag{2.1}$$

Nullstellenproblem und Fixpunktproblem können ineinander überführt werden., z.B.

$$g(x) = x - F(x)$$

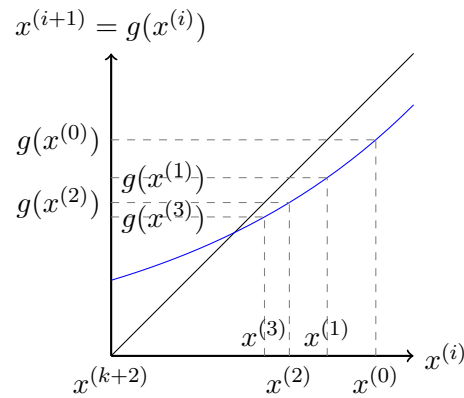


Abbildung 2.5: Fixpunktiteration

### 2.3.1 Fixpunkttheorem

1. Sei  $g(x)$  stetig auf  $[a, b]$ , sowie  $g(x) \in [a, b]$  für alle  $x \in [a, b]$ .  
Dann hat  $g(x)$  mindestens einen Fixpunkt auf  $[a, b]$
2. Existiere zusätzlich  $g'(x)$  auf  $(a, b)$ , sowie eine positive Konstante  $k < 1$  mit  $|g'(x)| \leq k$ , für alle  $x \in (a, b)$ 
  - Dann hat  $g(x)$  exakt einen Fixpunkt auf  $[a, b]$
  - Gilt  $0 < k < 1$ , dann konvergiert für jedes  $x^{(0)} \in [a, b]$  die Folge

$$x^{(n)} = g(x^{(n-1)}) \quad , n \geq 1$$

zum eindeutigen Fixpunkt  $x^* \in [0, b]$ .

### 2.3.2 Konvergenzverhalten iterativer Verfahren

$$x^{(k+1)} = g(x^{(k)}) \quad \text{Iterationsvorschrift} \quad (2.2)$$

$$x^* = g(x^*) \quad \text{Fixpunkt} \quad (2.3)$$

Taylorreihe von  $g(x)$  um  $x^*$

$$\left[ g(x^{(k)}) \right] x^{(k+1)} = \underbrace{g(x^*)}_{x^*} + g'(x^*) \cdot (x^{(k)} - x^*) + \frac{1}{2} g''(x^*) \cdot (x^{(k)} - x^*)^2 + \dots \quad (2.4)$$

$$\varepsilon^{(k+1)} = g'(x^*) \cdot \varepsilon^{(k)} + \frac{1}{2} g''(x^*) \cdot \varepsilon^{(k)2} + \dots \quad (2.5)$$

- $g'(x^*) \neq 0, |g'(x^*)| < 1 \rightarrow$  lineare Konvergenz
- $g'(x^*) = 0, g''(x^*) \neq 0 \Rightarrow$  quadratische Konvergenz ("Konstruktionshinweis")

## 2.4 Mehrdimensionales Newton-Raphson-Verfahren

$$\vec{F}(\vec{x}) = \vec{0} \quad , \vec{F} = (F_1, \dots, F_r)^T \quad (2.6)$$

$$\vec{F}(\vec{x}^{(k+1)}) = \vec{F}(\vec{x}^{(k)}) + \underbrace{\frac{\partial \vec{F}(\vec{x}^{(k)})}{\partial \vec{x}^T}}_{\text{Jacobi-Matrix}} \cdot (\vec{x}^{(k+1)} - \vec{x}^{(k)}) \quad (2.7)$$

$$\begin{array}{l} \text{Jacobi-Matrix} \\ \text{Fundamentalmatrix} \end{array} \quad \mathbf{J}(\vec{x}) = \frac{\partial \vec{F}(\vec{x})}{\partial \vec{x}^T} = \left( \frac{\partial \vec{F}}{\partial \vec{x}_1}, \dots, \frac{\partial \vec{F}}{\partial \vec{x}_r} \right) = \begin{pmatrix} \frac{\partial \vec{F}_1}{\partial \vec{x}_1} & \cdots & \frac{\partial \vec{F}_1}{\partial \vec{x}_r} \\ \vdots & & \vdots \\ \frac{\partial \vec{F}_r}{\partial \vec{x}_1} & \cdots & \frac{\partial \vec{F}_r}{\partial \vec{x}_r} \end{pmatrix} \quad (2.8)$$

$$\vec{F}(\vec{x}^{(k+1)}) = \vec{0} \Rightarrow \mathbf{J}(\vec{x}^{(k)}) \cdot (\vec{x}^{(k+1)} - \vec{x}^{(k)}) = -\vec{F}(\vec{x}^{(k)}) \quad (2.9)$$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} - \mathbf{J}^{-1}(\vec{x}^{(k)}) \cdot \vec{F}(\vec{x}^{(k)})$$

**Praktisch:**

$$\mathbf{J}(\vec{x}^{(k)}) \cdot \vec{x}^{(k+1)} = \mathbf{J}(\vec{x}^{(k)}) \cdot \vec{x}^{(k)} - \vec{F}(\vec{x}^{(k)})$$

### 2.4.1 Konvergenz

$$\vec{x}^{(k+1)} = \vec{\Phi}(\vec{x}^{(k)}) = \vec{\Phi}(\vec{x}^*) + \frac{\partial \vec{\Phi}}{\partial \vec{x}^T} \Big|_{\vec{x}=\vec{x}^*} \cdot (\vec{x}^{(k)} - \vec{x}^*) + \mathcal{O} \left( (\vec{x}^{(k)} - \vec{x}^*) \cdot (\vec{x}^{(k)} - \vec{x}^*)^T \right)$$

$$\vec{\epsilon}^{(k+1)} = \frac{\partial \vec{\Phi}}{\partial \vec{x}^T} \Big|_{\vec{x}=\vec{x}^*} \vec{\epsilon}^{(k)} + \mathcal{O} \left( \vec{\epsilon}^{(k)} \cdot \vec{\epsilon}^{(k)T} \right)$$

$$\begin{aligned} \vec{\Phi}(\vec{x}) &= \vec{x} - \mathbf{J}^{-1} \cdot \vec{F} \\ \frac{\partial \vec{\Phi}}{\partial \vec{x}^T} &= \underbrace{\mathbf{E}}_0 - \underbrace{\mathbf{J}^{-1}}_{\mathbf{J}^{-1}} \underbrace{\frac{\partial \vec{F}}{\partial \vec{x}^T}}_{\mathbf{J}^{-1}} - \underbrace{\frac{\partial \mathbf{J}^{-1}}{\partial \vec{x}^T} \cdot \vec{F}}_0 \end{aligned}$$

### 2.4.2 Abbruchkriterien

Gegeben sei eine Nullstellen-Suche:  $\vec{F}(\vec{x}) = \vec{0} \Rightarrow \vec{x}^*$

Als Abbruchkriterium wäre schön

$$\|\vec{x}^{(k)} - \vec{x}^*\| = \|\vec{\epsilon}^{(k)}\| \leq \|\vec{\epsilon}^{(A)}\|$$

oder

$$\frac{\|\vec{\epsilon}^{(k)}\|}{\|\vec{x}^*\|} \leq \|\vec{\epsilon}^{(R)}\|$$

Problematisch dabei ist, dass  $\vec{x}^*$  unbekannt ist. Daher muss man in der Anwendung ein anderes Kriterium wählen:

- $\|\vec{x}^{(k+1)} - \vec{x}^{(k)}\| \leq \|\vec{\epsilon}^{(A)}\|$  bzw.  $\frac{\|\vec{x}^{(k+1)} - \vec{x}^{(k)}\|}{\|\vec{x}^{(k+1)}\|} \leq \|\vec{\epsilon}^{(R)}\|$
- $\|\vec{F}(\vec{x}^{(k+1)}) - \vec{F}(\vec{x}^{(k)})\| \leq \Delta_F^{(A)}$
- $\|\vec{F}(\vec{x}^{(k)})\| \leq F_{min}$
- $k > k_{max}$





### 3 Lineare Gleichungssysteme

$$\mathbf{A} \cdot \vec{x} = \vec{b}, \quad \mathbf{A} \in \mathbb{R}^{n \times n}; \quad \vec{b}, \vec{x} \in \mathbb{R}^n, \quad \mathbf{A} \neq 0 \quad (3.1)$$

$$\text{Theoretische: } \vec{x} = \mathbf{A}^{-1} \cdot \vec{b} \quad (3.2)$$

$$\text{Alternativ: Lösung durch „Division“-Eliminierung} \quad (3.3)$$

#### 3.1 Methode zur Bestimmung von $\mathbf{A}^{-1}$

$$\mathbf{A} \cdot \mathbf{A}^{-1} = \mathbf{I} \quad (3.4)$$

$$\text{Sei } \mathbf{X} = \mathbf{A}^{-1} \quad (3.5)$$

$$\mathbf{A} \cdot \mathbf{X} = \mathbf{I} \quad (3.6)$$

$$\mathbf{A} \cdot [\vec{x}_1 \quad \vec{x}_2 \quad \dots \quad \vec{x}_n] = [\vec{l}_1 \quad \vec{l}_2 \quad \dots \quad \vec{l}_n]; \quad \vec{l}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix}, \dots, \vec{l}_n = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \quad (3.7)$$

$$\mathbf{A} \cdot \vec{x}_1 = \vec{l}_1; \quad \mathbf{A} \cdot \vec{x}_2 = \vec{l}_2; \quad \dots; \quad \mathbf{A} \cdot \vec{x}_n = \vec{l}_n;$$

Nur rechte Seite ändert sich;  $\mathbf{A} = \mathbf{L} \cdot \mathbf{U}$  nur einmal erforderlich.

Aufwand  $n \cdot 2 \cdot \frac{1}{2}(n^2 + n) = n^3 + n^2 \rightarrow O(n^3)$  für Substitution.

Probleme:

- Pivot  $p = 0 \nleftrightarrow$
- Pivot  $p$  “sehr klein“

Abhilfe: “Pivotisierung“- Zeilen-/Spaltentausch

- “partial pivoting“: Wähle Zeile mit betragsmäßig größtem Element in Pivotspalte
- “complete pivoting“: Zeilen- und Spaltentausch
- Pivotisierung auch zur Aufwandreduktion (Nullelemente erhalten)

Beispiel  $p = 0$

$$\begin{pmatrix} 10 & -7 & 0 \\ -3 & 2,1 & 6 \\ 5 & -1 & 5 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 7 \\ 3,9 \\ 6 \end{pmatrix} \quad \begin{matrix} p_{21} = \frac{5}{10} = \frac{1}{2} \\ p_{31} = \frac{-3}{10} \end{matrix} \quad (3.8)$$

$$\begin{pmatrix} 10 & -7 & 0 \\ 0 & 0 & 6 \\ 0 & 2,5 & 5 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 7 \\ 6 \\ 2,5 \end{pmatrix} \quad p_{32} = \frac{0}{2,5} \quad \nexists \quad (3.9)$$

$$\Rightarrow \text{Zeilentausch} \quad (3.10)$$

$$\begin{pmatrix} 10 & -7 & 0 \\ 0 & 2,5 & 5 \\ 0 & 0 & 6 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 7 \\ 2,5 \\ 6 \end{pmatrix} \quad (3.11)$$

Zeilentausch durch Permutationsmatrix  $P$ .

$$PA\vec{x} = P\vec{b} \Rightarrow \text{dann } PA \rightarrow LU$$

$$\text{hier: } P = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$$

## 3.2 Kondition eines Gleichungssystems

$$\begin{pmatrix} 1 & 1 \\ 1 & 1,0001 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \quad \Rightarrow \vec{x} = \begin{pmatrix} 2 \\ 0 \end{pmatrix} \quad (3.12)$$

$$\begin{pmatrix} 1 & 1 \\ 1 & 1,0001 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 2,0001 \end{pmatrix} \quad \Rightarrow \vec{x} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (3.13)$$

Änderung in der 5. Stelle von  $\vec{b}$  wird zu einer Änderung in der ersten Stelle der Lösung  $\vec{x}$  “verstärkt“. System reagiert sehr sensitiv auf kleine Änderungen der Ausgangsdaten. Kein Lösungsalgorithmus kann etwas daran ändern.

$$\text{cond}(\mathbf{A}) = 4,0002 \cdot 10^4 \quad (3.14)$$

## 3.3 Einfluss der Pivotisierung auf die Ergebnisgenauigkeit

$$B \rightarrow \begin{pmatrix} 0,0001 & 1 \\ 1 & 1 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (3.15)$$

$$\begin{pmatrix} 0,0001 & 1 \\ 0 & -9999 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 1 \\ -9998 \end{pmatrix} \quad \Rightarrow \begin{matrix} x_2 = 0,999\,899\,98 \\ x_1 = 1,000\,100\,01 \end{matrix} \quad (3.16)$$

Annahme: 3 Stellen Genauigkeit

$$\begin{pmatrix} 0,0001 & 1 \\ 0 & -10\,000 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 1 \\ -10\,000 \end{pmatrix} \quad \Rightarrow \begin{matrix} x_2 = 1 \\ x_1 = 0 \end{matrix} \quad \nexists \quad (3.17)$$

Zerlegung  $\mathbf{B} = \mathbf{LDU}$  (ohne Genauigkeitsbeschränkung)

$$\begin{array}{l} \text{("out of} \\ \text{scale} \\ \text{with } \mathbf{B} \text{"}) \end{array} = \begin{pmatrix} 1 & 0 \\ 10\,000 & 1 \end{pmatrix} \cdot \begin{pmatrix} 0,0001 & 0 \\ 0 & -9999 \end{pmatrix} \cdot \begin{pmatrix} 1 & 10\,000 \\ 0 & 1 \end{pmatrix} \quad (3.18)$$

Änderung der Pivotisierungsreihenfolge

$$\begin{pmatrix} 1 & 1 \\ 0,0001 & 1 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \quad (3.19)$$

$$\begin{pmatrix} 1 & 1 \\ 0 & 0,9999 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 2 \\ 0,9998 \end{pmatrix} \Rightarrow \begin{array}{l} x_2 = 0,999\,899\,98 \\ x_1 = 1,000\,100\,01 \end{array} \quad (3.20)$$

Annahme: 3 Stellen Genauigkeit

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \cdot \vec{x} = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \Rightarrow \begin{array}{l} x_2 = 1 \\ x_1 = 0 \end{array} \left. \vphantom{\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}} \right\} \begin{array}{l} \text{nur sehr geringer} \\ \text{Genauigkeitsverlust} \end{array} \quad (3.21)$$

NB:  $\text{cond}(\mathbf{B}) = 2,618\,385\,27$

Uns interessieren also zwei Themen:

- **Kondition** der Gleichungssystems
- **Stabilität** des Lösungsverfahrens

### 3.4 Kondition einer Matrix

Ausgehend von einem linearen Gleichungssystem

$$\mathbf{A} \cdot \vec{x} = \vec{b}$$

kann man auf der rechten Seite einen Fehler  $\Delta\vec{b}$  hinzufügen.

$$\mathbf{A}(\vec{x} + \Delta\vec{x}) = \vec{b} + \Delta\vec{b}$$

Dieser resultiert dann in einem Fehler  $\Delta\vec{x}$  der Lösung in  $\vec{x}$ . Der Zusammenhang zwischen diesen Größen lautet

$$\mathbf{A}\Delta\vec{x} = \Delta\vec{b}$$

und kann umgeformt werden zu

$$\begin{aligned} \Delta\vec{x} &= \mathbf{A}^{-1}\Delta\vec{b} \\ \|\Delta\vec{x}\| &\leq \|\mathbf{A}^{-1}\| \cdot \|\Delta\vec{b}\| \end{aligned}$$

Gleichzeitig kann

$$\vec{b} = \mathbf{A}\vec{x}$$

umgeformt werden zu

$$\|\vec{b}\| \leq \|\mathbf{A}\| \cdot \|\vec{x}\|$$

$$\frac{1}{\|\vec{x}\|} \leq \|\mathbf{A}\| \cdot \frac{1}{\|\vec{b}\|}$$

Diese Gleichungen zusammengefasst ergeben dann

$$\frac{\|\Delta\vec{x}\|}{\|\vec{x}\|} \leq \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\| \cdot \frac{\|\Delta\vec{b}\|}{\|\vec{b}\|}$$

Darin stellen dann  $\frac{\|\Delta\vec{x}\|}{\|\vec{x}\|}$  den relativen Fehler des Ergebnisses,  $\|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\|$  den *Verstärkungsfaktor* für Fehler in  $\vec{b}$  und  $\frac{\|\Delta\vec{b}\|}{\|\vec{b}\|}$  den relativen Fehler in  $\vec{b}$  dar. Wir definieren die Kondition von  $\mathbf{A}$  über

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}^{-1}\| \cdot \|\mathbf{A}\|$$

welche offensichtlich abhängig von der gewählten Matrixnorm ist. Eine  $\text{cond}(\mathbf{A}) \rightarrow \infty$  bedeutet dabei eine schlechte,  $\text{cond}(\mathbf{A}) \rightarrow 0$  eine gute *Kondition*.

### 3.5 Allgemeines Iterationsverfahren für lineare Gleichungssysteme

$$\mathbf{A} \cdot \vec{x} = \vec{b} \quad (3.22)$$

$$\vec{x}^{(k+1)} = \vec{\Phi}(\vec{x}^{(k)}) \quad (3.23)$$

$$\mathbf{B} \text{ beliebig} \quad (3.24)$$

$$\mathbf{B}\vec{x} + (\mathbf{A} - \mathbf{B})\vec{x} = \vec{b} \quad (3.25)$$

$$\Rightarrow \text{Iterationsvorschrift:} \quad (3.26)$$

$$\mathbf{B} \cdot \vec{x}^{(k+1)} + (\mathbf{A} - \mathbf{B}) \cdot \vec{x}^{(k)} = \vec{b} \quad (3.27)$$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} - \mathbf{B}^{-1} \cdot (\mathbf{A} \cdot \vec{x}^{(k)} - \vec{b}) \quad (3.28)$$

$$= \underbrace{(\mathbf{I} - \mathbf{B}^{-1} \cdot \mathbf{A})}_{\text{Konvergenz, wenn } |\lambda| < 1} \cdot \vec{x}^{(k)} + \mathbf{B}^{-1} \cdot \vec{b} \quad (3.29)$$

#### 3.5.1 Iterative Verfahren zum Lösen linearer Gleichungssysteme

$$\mathbf{A}\vec{x} = \vec{b}, \mathbf{A} \in \mathbb{R}^{n \times n}, \vec{x}, \vec{b} \in \mathbb{R}^n, \mathbf{A} \text{ nicht singulär}$$

$$\text{Splitting: } \mathbf{A} = [a_{ij}] = \mathbf{D} + \mathbf{L} + \mathbf{U} \quad (3.30)$$

$$= \begin{bmatrix} a_{11} & & & \\ & a_{22} & & \\ & & \ddots & \\ & & & a_{nn} \end{bmatrix} + \begin{bmatrix} 0 & & & \\ a_{21} & 0 & & \\ \vdots & & \ddots & \\ a_{n1} & \dots & a_{nn-1} & 0 \end{bmatrix} + \begin{bmatrix} 0 & a_{12} & \dots & a_{1n} \\ & 0 & & \vdots \\ & & \ddots & a_{n-1n} \\ & & & 0 \end{bmatrix} \quad (3.31)$$

$$\mathbf{A} = \mathbf{S} - \mathbf{T} \quad (3.32)$$

$$\mathbf{A} \cdot \vec{x} = (\mathbf{S} - \mathbf{T}) \cdot \vec{x} = \vec{b} \quad (3.33)$$

$$\Rightarrow \mathbf{S} \cdot \vec{x} = \mathbf{T} \cdot \vec{x} + \vec{b} \quad (3.34)$$

$$= (\mathbf{S} - \mathbf{A}) \cdot \vec{x} + \vec{b} \quad (3.35)$$

$$\Rightarrow \vec{x} = \mathbf{S}^{-1} \cdot (\mathbf{S} - \mathbf{A}) \cdot \vec{x} + \mathbf{S}^{-1} \cdot \vec{b} \quad (3.36)$$

...

### 3.5.2 Verschiedene Iterationsverfahren

1.  $\mathbf{S} = \mathbf{D}$  : Jacobi (Gesamtschrittverfahren)
2.  $\mathbf{S} = \mathbf{D} + \mathbf{L}$  : Gauß-Seidl (Einzelschrittverfahren)
3.  $\mathbf{S} = \frac{1}{w}(\mathbf{D} + w\mathbf{L})$  : sukzessive Überrelaxation (SOR, successive Over-Relaxation)

1. Jacobi:  $\mathbf{S} = \mathbf{D}$ ;  $\mathbf{D}$  nicht singular (d.h.  $a_{ii} \neq 0$ )

$$\mathbf{A}\vec{x} = \vec{b} \Rightarrow \mathbf{D} \cdot \vec{x}^{(k+1)} = (\mathbf{D} - \mathbf{A}) \cdot \vec{x}^{(k)} + \vec{b} = (-\mathbf{L} - \mathbf{U}) \cdot \vec{x}^{(k)} + \vec{b} \quad (3.37)$$

$$\vec{x}^{(k+1)} = \underbrace{\mathbf{D}^{-1} \cdot (-\mathbf{L} - \mathbf{U})}_{\mathbf{K}_j} \cdot \vec{x}^{(k)} + \underbrace{\mathbf{D}^{-1} \vec{b}}_{\vec{c}_j} \quad (3.38)$$

$$= \mathbf{K}_j \cdot \vec{x}^{(k)} + \vec{c}_j \quad (3.39)$$

2. Gauss-Seidel:  $\mathbf{S} = \mathbf{D} + \mathbf{L}$  (nicht singular)

$$\mathbf{A}\vec{x} + \vec{b} = \vec{0} \Rightarrow (\mathbf{D} + \mathbf{L})\vec{x}^{(k+1)} = (\mathbf{D} + \mathbf{L} - \mathbf{A})\vec{x}^{(k)} + \vec{b} = -\mathbf{U}\vec{x}^{(k)} + \vec{b}$$

$$\vec{x}^{(k+1)} = -\underbrace{(\mathbf{D} + \mathbf{L})^{-1} \cdot \mathbf{U}}_{\mathbf{K}_g} \vec{x}^{(k)} + \underbrace{(\mathbf{D} + \mathbf{L})^{-1} \vec{b}}_{\vec{c}_g}$$

$$\vec{x}^{(k+1)} = \mathbf{K}_g \vec{x}^{(k)} + \vec{c}_g$$

3. Successive Over-Relaxation:  $\mathbf{S} = \frac{1}{\omega}(\mathbf{D} + \omega\mathbf{L})$

$$\mathbf{A}\vec{x} = \vec{b}$$

$$\frac{1}{\omega}(\mathbf{D} + \omega\mathbf{L})\vec{x}^{(k+1)} = \frac{1}{\omega}(\mathbf{D} + \omega\mathbf{L} - \omega\mathbf{A})\vec{x}^{(k)} + \vec{b}$$

$$(\mathbf{D} + \omega\mathbf{L})\vec{x}^{(k+1)} = (\mathbf{D} + \omega\mathbf{L} - \omega\mathbf{A})\vec{x}^{(k)} + \omega\vec{b}$$

$$(\mathbf{D} + \omega\mathbf{L})\vec{x}^{(k+1)} = (\mathbf{D}(1 - \omega) - \omega\mathbf{U})\vec{x}^{(k)} + \omega\vec{b}$$

$$\vec{x}^{(k+1)} = \underbrace{(\mathbf{D} + \omega\mathbf{L})^{-1}(\mathbf{D}(1 - \omega) - \omega\mathbf{U})}_{\mathbf{K}_\omega} \vec{x}^{(k)} + \underbrace{(\mathbf{D} + \omega\mathbf{L})^{-1} \omega\vec{b}}_{\vec{c}_\omega}$$

$$\vec{x}^{(k+1)} = \mathbf{K}_\omega \vec{x}^{(k)} + \vec{c}_\omega$$

Die jeweilige Wahl von  $\omega$  bei SOR beeinflusst dann das Konvergenzverhalten. Für Spezialfälle können Empfehlungen gegeben werden:

- a) Falls  $a_{ii} \neq 0, i = 1, \dots, n \Rightarrow \varrho(\mathbf{K}_\omega) \geq |1 - \omega| \Rightarrow$  Konvergenz für  $0 < \omega < 2$
- b) Falls  $\mathbf{A}$  positiv definit und  $0 < \omega < 2$ , dann konvergiert SOR für jede Initiallösung  $\vec{x}^{(0)}$
- c) Falls  $\mathbf{A}$  positiv definit und tridiagonal ist, dann  $\varrho(\mathbf{K}_g) = \varrho(\mathbf{K}_j)^2 < 1$  und die optimale Wahl für  $\omega$  ist:

$$\omega = \frac{2}{1 + \sqrt{1 - \varrho(\mathbf{K}_j)^2}}$$

Ein Beispiel zum Vergleich zwischen Jacobi und Gauß-Seidel

$$10x_1 - x_2 + 2x_3 = 6 \quad (3.40)$$

$$-x_1 + 11x_2 - x_3 + 3x_4 = 25 \quad (3.41)$$

$$2x_1 - x_2 + 10x_3 - x_4 = -11 \quad (3.42)$$

$$3x_2 - x_3 + 8x_4 = 15 \quad (3.43)$$

Jacobi:

$$x_1^{(k+1)} = \frac{1}{10}x_2^{(k)} - \frac{1}{5}x_3^{(k)} + \frac{3}{5} \quad (3.44)$$

$$x_2^{(k+1)} = \frac{1}{11}x_1^{(k)} + \frac{1}{11}x_3^{(k)} - \frac{3}{11}x_4^{(k)} + \frac{25}{11} \quad (3.45)$$

$$x_3^{(k+1)} = -\frac{1}{5}x_1^{(k)} + \frac{1}{10}x_2^{(k)} + \frac{1}{10}x_4^{(k)} - \frac{11}{10} \quad (3.46)$$

$$x_4^{(k+1)} = -\frac{3}{8}x_2^{(k)} + \frac{1}{8}x_3^{(k)} + \frac{15}{8} \quad (3.47)$$

Gauß-Seidel:

$$x_1^{(k+1)} = \frac{1}{10}x_2^{(k)} - \frac{1}{5}x_3^{(k)} + \frac{3}{5} \quad (3.48)$$

$$x_2^{(k+1)} = \frac{1}{11}x_1^{(k+1)} + \frac{1}{11}x_3^{(k)} - \frac{3}{11}x_4^{(k)} + \frac{25}{11} \quad (3.49)$$

$$x_3^{(k+1)} = -\frac{1}{5}x_1^{(k+1)} + \frac{1}{10}x_2^{(k+1)} + \frac{1}{10}x_4^{(k)} - \frac{11}{10} \quad (3.50)$$

$$x_4^{(k+1)} = -\frac{3}{8}x_2^{(k+1)} + \frac{1}{8}x_3^{(k+1)} + \frac{15}{8} \quad (3.51)$$

Man erkennt bei Gauß-Seidel die Abhängigkeit von Werten, welche erst im aktuellen Iterationsschritt berechnet werden. In der Praxis bietet Gauß-Seidel dadurch in der Regel eine schneller Konvergenz. Vorteil von Jacobi ist jedoch die bessere Parallelisierbarkeit eines Iterationsschrittes.

### Konvergenzbetrachtung

$$\mathbf{S}\vec{x} = \mathbf{T}\vec{x} + \vec{b} \quad (3.52)$$

$$\mathbf{S}\vec{x}^{(x+1)} = \mathbf{T}\vec{x}^{(k)} + \vec{b} \quad (3.53)$$

3.52 - 3.53:

$$\begin{aligned} \mathbf{S}(\vec{x} - \vec{x}^{(k+1)}) &= \vec{T}(\vec{x} - \vec{x}^{(k)}) \\ \mathbf{S}\vec{e}^{(k+1)} &= \mathbf{T}\vec{e}^{(k)} \\ \vec{e}^{(k+1)} &= \mathbf{S}^{-1}\mathbf{T}\vec{e}^{(k)} \\ \vec{e}^{(k+1)} &= \mathbf{G}\vec{e}^{(k)} = \mathbf{G}^k\vec{e}^{(0)} \end{aligned}$$

### 3.5.3 Gradientenverfahren

Sei  $\mathbf{A} \in \mathbb{R}^{n \times n}$  mit

- $\mathbf{A}^T = \mathbf{A}$
- $\vec{x}^T \cdot \mathbf{A} \cdot \vec{x} > 0$  für  $\vec{x} \neq 0$  (positiv definit)
- $\vec{b} \in \mathbb{R}$

Dann existiert genau ein Minimum für die quadratische Form:

$$\Phi(\vec{x}) = 0,5 \cdot \vec{x}^T \cdot \mathbf{A} \cdot \vec{x} - \vec{x}^T \cdot \vec{b} \quad (3.54)$$

$$\nabla \Phi(\vec{x}) = 0,5 (\mathbf{A}^T + \mathbf{A}) \cdot \vec{x} - \vec{b} = \mathbf{A} \cdot \vec{x} - \vec{b} \quad (3.55)$$

$$\text{Minimum: } \nabla \Phi(\vec{x}) = \vec{0} \Rightarrow \mathbf{A} \cdot \vec{x} = \vec{b} \quad (3.56)$$

$\Rightarrow$  Minimierung der quadratischen Form entspricht Lösung des linearen Gleichungssystem.

## 3.6 Iterativer Ansatz für Minimierungsproblem

**Allgemein:**  $\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \cdot \vec{d}^{(k)}$

**Bestimme:**  $\alpha_k, \vec{d} = ?$

**Idee:** Wähle Richtung des steilsten Abstiegs ("steepest descent", siehe Abbildung 3.1)

$$\vec{d}^{(k)} = -\nabla \Phi(\vec{x}^{(k)}) = \vec{b} - \mathbf{A}\vec{x}^{(k)} =: \vec{r}^{(k)} \text{ Residuum} \quad (3.57)$$

$$(3.58)$$

$\alpha_k = ? \Rightarrow$  Minimiere  $\vec{\Phi}$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \vec{r}^{(k)} \quad (3.59)$$

$$\vec{\Phi}(\vec{x}^{(k+1)}) = \frac{1}{2} \left( \vec{x}^{(k)} + \alpha_k \vec{r}^{(k)} \right)^T \cdot \mathbf{A} \cdot \left( \vec{x}^{(k)} + \alpha_k \vec{r}^{(k)} \right) - \left( \vec{x}^{(k)} + \alpha_k \vec{r}^{(k)} \right)^T \cdot \vec{b} \quad (3.60)$$

$$\frac{\partial \vec{\Phi}(\vec{x}^{(k+1)})}{\partial \alpha_k} = \frac{1}{2} \vec{r}^{(k)T} \cdot \mathbf{A} \cdot \vec{x}^{(k)} + \frac{1}{2} \vec{x}^{(k)T} \cdot \mathbf{A} \cdot \vec{r}^{(k)} + \alpha_k \cdot \vec{r}^{(k)T} \cdot \mathbf{A} \cdot \vec{r}^{(k)} - \vec{r}^{(k)T} \cdot \vec{b} \quad (3.61)$$

$$= \vec{r}^{(k)T} \cdot \mathbf{A} \vec{x}^{(k)} - \vec{r}^{(k)T} \cdot \vec{b} + \alpha_k \cdot \vec{r}^{(k)T} \cdot \mathbf{A} \cdot \vec{r}^{(k)} \neq 0 \quad (3.62)$$

$$\dots \quad (3.63)$$

$$\Rightarrow \alpha_k = \frac{\vec{r}^{(k)T} \cdot \vec{r}^{(k)}}{\vec{r}^{(k)T} \cdot \mathbf{A} \cdot \vec{r}^{(k)}} \quad (3.64)$$

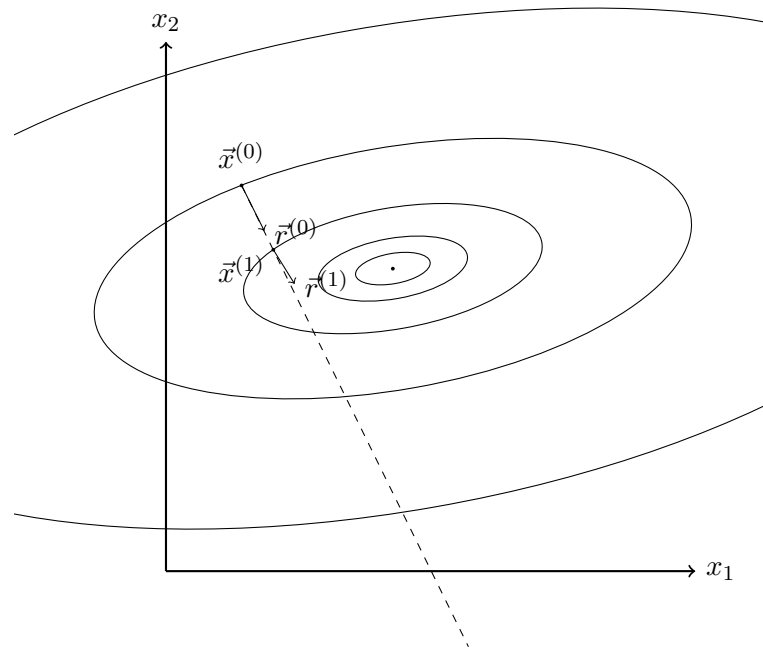


Abbildung 3.1: Gradientenverfahren - steepest descent

**Verfahren:**

Für  $k = 0, 1, 2, \dots$

$$\vec{r}^{(k)} = \vec{b} - \mathbf{A}\vec{x}^{(k)} \quad (3.65)$$

$$\alpha_k = \frac{\vec{r}^{(k)T} \cdot \vec{r}^{(k)}}{\vec{r}^{(k)T} \cdot \mathbf{A} \cdot \vec{r}^{(k)}} \quad (3.66)$$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \cdot \vec{r}^{(k)} \quad (3.67)$$

2 Matrix-Vektormultiplikationen (dominant)

2 Vektor-Skalarprodukte

**Alternative:**

$$-\mathbf{A} \cdot \vec{x}^{(k+1)} + \vec{b} = -\mathbf{A} \cdot \vec{x}^{(k)} - \alpha_k \cdot \mathbf{A} \cdot \vec{r}^{(k)} + \vec{b} \quad (3.68)$$

$$\vec{r}^{(k+1)} = \vec{r}^{(k)} - \alpha_k \cdot \mathbf{A}\vec{r}^{(k)} \quad (3.69)$$

$\Rightarrow$  nur eine Matrix-Vektormultiplikation je Iteration



### 3.6.1 Konvergenz

$$\vec{x}^{(k)} = \vec{x}^* + \vec{e}^{(k)}, \text{ Annahme: } \vec{e} \text{ ist EV von } \mathbf{A} : \quad (3.70)$$

$$\mathbf{A} \cdot \vec{e} = \lambda_e \cdot \vec{e} \quad (3.71)$$

$$\vec{r}^{(k)} = \vec{b} - \mathbf{A} \cdot \vec{x}^{(k)} \quad (3.72)$$

$$= \vec{b} - \mathbf{A} \cdot (\vec{x}^* + \vec{e}^{(k)}) \quad (3.73)$$

$$= \underbrace{\vec{b} - \mathbf{A} \cdot \vec{x}^*}_{\vec{0}} - \mathbf{A} \cdot \vec{e}^{(k)} \quad (3.74)$$

$$\vec{r}^{(k)} = -\mathbf{A} \cdot \vec{e}^{(k)} = -\lambda_e \cdot \vec{e}^{(k)} \Rightarrow \vec{r}^{(k)} \text{ ist ebenfalls EV} \quad (3.75)$$

$$\text{Aus Gleichung 3.67: } \vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \cdot \vec{r}^{(k)} \quad (3.76)$$

$$\vec{e}^{(k+1)} = \vec{e}^{(k)} + \alpha_k \cdot \vec{r}^{(k)} \quad (3.77)$$

$$= \vec{e}^{(k)} + \frac{\vec{r}^{(k)T} \cdot \vec{r}^{(k)}}{\vec{r}^{(k)T} \cdot \mathbf{A} \cdot \vec{r}^{(k)}} \cdot (-\lambda_e \cdot \vec{e}^{(k)}) \quad (3.78)$$

$$= \vec{e}^{(k)} + \frac{\vec{r}^{(k)T} \cdot \vec{r}^{(k)}}{\vec{r}^{(k)T} \cdot \vec{r}^{(k)} \cdot \lambda_e} \cdot (-\lambda_e \cdot \vec{e}^{(k)}) \quad (3.79)$$

$$= \vec{e}^{(k)} - \vec{e}^{(k)} = \vec{0} \Rightarrow \text{Exakte Lösung in einem Schritt} \quad (3.80)$$

### Allgemein

$$\|\vec{x}\|_A = \sqrt{\vec{x}^T \cdot \mathbf{A} \cdot \vec{x}} \text{ Energienorm} \quad (3.81)$$

$$\text{Es gilt: } \|\vec{e}^{(k+1)}\|_A = \frac{\kappa_2(\mathbf{A}) - 1}{\kappa_2(\mathbf{A}) + 1} \cdot \|\vec{e}^{(k)}\|_A \quad (3.82)$$

mit  $\kappa_2$  als die Kondition bezüglich der Spektralnorm

## 3.7 CG - Konjugierte Gradientenmethode

$\vec{p}^{(k)}$  ist eine Menge mit  $\underbrace{\vec{p}^{(i)T}}_{\text{paarweise A-orthogonal bzw. A-konjugiert}} \cdot \mathbf{A} \cdot \vec{p}^{(j)} = 0, \forall i, j = 1, \dots, n \quad i \neq j$   
 $\Rightarrow \vec{p}^{(k)}$  bilden eine Basis des  $\mathbb{R}^n$

$$\text{Daher: } \vec{x}^* = \sum_{i=1}^n \alpha_i \cdot \vec{p}^{(i)} \quad (3.83)$$

$$\Rightarrow \vec{b} = \mathbf{A} \cdot \vec{x}^* = \sum_{i=1}^n \alpha_i \cdot \mathbf{A} \cdot \vec{p}^{(i)} \quad (3.84)$$

### 3.7.1 Schrittweite

$$0 = \frac{\partial \vec{\Phi}}{\partial \alpha_k} \left( \vec{x}^{(k)} + \alpha_k \vec{p}^{(k)} \right) \quad (3.85)$$

$$\Rightarrow \alpha_k = \frac{\vec{p}^{(k)T} \cdot \vec{r}^{(k)}}{\vec{p}^{(k)T} \cdot \mathbf{A} \cdot \vec{p}^{(k)}} \quad (3.86)$$

Wie Gradientenverfahren, aber Suchrichtung  $\vec{p}$

### 3.7.2 Suchrichtung

**Definition:** Lösung  $\vec{x}^{(k)}$  heißt optimal bzgl. einer Richtung  $\vec{p} \neq 0$ , wenn  $\vec{\Phi}(\vec{x}^{(k)}) \leq \vec{\Phi}(\vec{x}^{(k)} + \lambda \vec{p})$ ,  $\forall \lambda \in \mathbb{R}$

$\Leftrightarrow \vec{\Phi}$  besitzt lokales Minimum entlang  $\vec{p}$  für  $\lambda = 0$

$$\Leftrightarrow \frac{\partial \vec{\Phi}}{\partial \lambda} (\vec{x}^{(k)} + \lambda \vec{p}^{(k)})|_{\lambda=0} = \vec{p}^T \cdot \mathbf{A} \cdot \vec{x}^{(k)} - \vec{p}^T \vec{b} + \lambda \vec{p}^T \cdot \mathbf{A} \cdot \vec{p}|_{\lambda=0}$$

$$\Leftrightarrow \vec{p}^T \vec{r}^{(k)}$$

$$\text{Sei } \vec{x}^{(k+1)} = \vec{x}^{(k)} + \vec{q} \Rightarrow \vec{r}^{(k+1)} = \vec{b} \cdot \mathbf{A} \cdot \vec{x}^{(k+1)} \quad (3.87)$$

$$= \vec{b} \cdot \mathbf{A} \cdot \vec{x}^{(k)} - \mathbf{A} \vec{q} \quad (3.88)$$

$$= \vec{r}^{(k)} - \mathbf{A} \vec{q} \quad (3.89)$$

Forderung  $\vec{x}^{(k+1)}$  ebenfalls optimal bzgl.  $\vec{p}$ , also

$$0 = \vec{p}^T \cdot \vec{r}^{(k+1)} = \vec{p}^T \cdot (\vec{r}^{(k)} - \mathbf{A} \vec{q}) \quad (3.90)$$

$$= \underbrace{\vec{p}^T \cdot \vec{r}^{(k)}}_0 - \vec{p}^T \cdot \mathbf{A} \cdot \vec{q} = \underbrace{-\vec{p}^T \cdot \mathbf{A} \cdot \vec{q}}_{\vec{p}, \vec{q}, \mathbf{A} \text{ orthogonal}} = 0 \quad (3.91)$$

### Konstruktion der Suchrichtung

- Start  $\vec{p}^{(0)} = \vec{r}^{(0)}$  (Ausgehend von Initiallösung  $\vec{x}^{(0)}$ , Gradient an  $\vec{x}^{(0)}$ )
- Suche Richtungen  $\vec{p}^{(k+1)} = \vec{r}^{(k+1)} - \beta_k \cdot \vec{p}^{(k)}$ ,  $k = 0, 1, 2, \dots$  mit  $\beta_k \in \mathbb{R}$ , so dass gilt:

$$\vec{p}^{(j)T} \cdot \mathbf{A} \cdot \vec{p}^{(k+1)} = 0 = \left( \mathbf{A} \cdot \vec{p}^{(j)} \right)^T \cdot \vec{p}^{(k+1)} = 0, j = 0, 1, \dots, k \quad (3.92)$$

$$\Rightarrow \left( \mathbf{A} \cdot \vec{p}^{(j)} \right)^T \cdot \left( \vec{r}^{(k+1)} - \beta_k \cdot \vec{p}^{(k)} \right) = 0 \xRightarrow{(j=k)} \beta_k = \frac{\left( \mathbf{A} \cdot \vec{p}^{(k)} \right)^T \cdot \vec{r}^{(k+1)}}{\left( \mathbf{A} \cdot \vec{p}^{(k)} \right)^T \cdot \vec{p}^{(k)}} \quad (3.93)$$

Zu zeigen:  $\vec{p}^{(k+1)}$  ist  $\mathbf{A}$  orthogonal zu  $\vec{p}^{(j)}$   $\forall j < k$

Beweisidee: vollständige Induktion

### 3.7.3 Zusammenfassung des Verfahrens

$$\vec{r}^{(0)} = \vec{b} - \mathbf{A}\vec{x}^{(0)} \quad (3.94)$$

$$\vec{p}^{(0)} = \vec{r}^{(0)} \quad (3.95)$$

Für  $k = 0, 1, \dots$

$$\alpha_k = \frac{\vec{p}^{(k)T} \cdot \vec{r}^{(k)}}{\vec{p}^{(k+1)T} \mathbf{A} \vec{p}^{(k)}} \quad (3.96)$$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \vec{p}^{(k)} \quad (3.97)$$

$$\vec{b} - \mathbf{A}\vec{x}^{(k+1)} = \vec{r}^{(k+1)} = \vec{r}^{(k)} - \alpha_k \mathbf{A} \vec{p}^{(k)} \quad (3.98)$$

$$\beta_k = \frac{(\mathbf{A} \vec{p}^{(k)})^T \cdot \vec{r}^{(k+1)}}{(\mathbf{A} \vec{p}^{(k)})^T \cdot \vec{p}^{(k)}} \quad (3.99)$$

$$\vec{p}^{(k+1)} = \vec{r}^{(k+1)} - \beta_k \vec{p}^{(k)} \quad (3.100)$$

Der Aufwand je Iteration besteht maßgeblich aus

- 1 Matrix-Vektorprodukt
- 2 Vektor-Skalarprodukte (nach einer Umstellung der obigen Formeln)

### 3.7.4 Konvergenz

Sei  $\mathbf{A}$  symmetrisch und positiv definit, dann gilt:

- Das CG-Verfahren bricht nach höchstens  $n$  Schritten mit der exakten Lösung ab
- $\|\vec{e}^{(k)}\|_A \leq \frac{2c^k}{1+c^{2k}} \|\vec{e}^{(0)}\|_A$  mit  $c = \frac{\sqrt{K_2(\mathbf{A})}-1}{\sqrt{K_2(\mathbf{A})}+1}$

In der Praxis mit endlicher Rechengenauigkeit erhält man allerdings nicht die exakte Lösung nach  $n$  Schritten, anstelle dessen bricht man nach einem gewissen Kriterium (zum Beispiel Betrag des Residiums) ab.

Verbessert werden kann das Konvergenzverhalten durch eine Vorkonditionierung.

### 3.7.5 Vorkonditionierung

Für eine Linksvorkonditionierung wird  $\mathbf{A}\vec{x} = \vec{b}$  transformiert in

$$\mathbf{P}^{-1} \mathbf{A} \vec{x} = \mathbf{P}^{-1} \vec{b} \quad (3.101)$$

Um eine Rechtsvorkonditionierung vorzunehmen (welche seltener benutzt wird) wird nach

$$(\mathbf{A} \mathbf{P}^{-1}) \mathbf{P} \vec{x} = \vec{b} \quad (3.102)$$

transformiert und in zwei Schritten gelöst, zuerst

$$\mathbf{A} \mathbf{P}^{-1} \vec{y} = \vec{b} \quad (3.103)$$

und danach

$$\mathbf{P}\vec{x} = \vec{y} \quad (3.104)$$

Gilt  $K(\mathbf{P}^{-1}\mathbf{A}) < K(\mathbf{A})$ , dann erhält man ein besseres Konvergenzverhalten. Ideal wäre  $\mathbf{P}^{-1} = \mathbf{A}^{-1}$ , allerdings ist die Berechnung der Inversen aufwändiger als die Lösung des ursprünglichen Problems (die Lösung des Gleichungssystems). Möglichst einfach wäre der Einsatz von  $\mathbf{P} = \mathbf{I}$ , allerdings erhält man damit keine Gewinnung durch die Vorkonditionierung. Praktischerweise wählt man eine möglichst einfach ermittelbare Matrix  $\mathbf{P}$  welche zwischen diesen beiden Extrema liegt.

Eine mögliche Lösung für dieses Problem ist der Einsatz von *Preconditioned Conjugate Gradient* (PCG): Sei  $\mathbf{P}$  symmetrisch und positiv definit  $\Rightarrow \exists$  Zerlegung

$$\mathbf{P} = \mathbf{P}^{\frac{1}{2}} \mathbf{P}^{\frac{1}{2}T} \quad (3.105)$$

CG wird dann angewandt auf

$$\tilde{\mathbf{A}}\tilde{\vec{x}} = \tilde{\vec{b}} \quad (3.106)$$

$$\tilde{\mathbf{A}} = \mathbf{P}^{-\frac{1}{2}} \mathbf{A} \mathbf{P}^{-\frac{1}{2}T} \quad (3.107)$$

$$\tilde{\vec{x}} = \mathbf{P}^{\frac{1}{2}T} \vec{x} \quad (3.108)$$

$$\tilde{\vec{b}} = \mathbf{P}^{-\frac{1}{2}} \vec{b} \quad (3.109)$$

Das Verfahren sieht dann wie folgt aus

$$\vec{r}^{(0)} = \vec{b} - \mathbf{A}\vec{x}^{(0)} \quad (3.110)$$

$$\mathbf{P}\vec{z}^{(0)} = \vec{r}^{(0)} \Rightarrow \vec{z}^{(0)} \quad (3.111)$$

$$\vec{p}^{(0)} = \vec{z}^{(0)} \quad (3.112)$$

Für  $k = 0, 1, 2, \dots$

$$\alpha_k = \frac{\vec{p}^{(k)T} \cdot \vec{r}^{(k)}}{\vec{p}^{(k)T} \cdot \mathbf{A}\vec{p}^{(k)}} \quad (3.113)$$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \alpha_k \vec{p}^{(k)} \quad (3.114)$$

$$\vec{r}^{(k+1)} = \vec{r}^{(k)} - \alpha_k \mathbf{A}\vec{p}^{(k)} \quad (3.115)$$

$$\mathbf{P}\vec{z}^{(k+1)} = \vec{r}^{(k+1)} \Rightarrow \vec{z}^{(k+1)} \quad (3.116)$$

$$\beta_k = \frac{(\mathbf{A}\vec{p}^{(k)})^T \cdot \vec{z}^{(k+1)}}{(\mathbf{A}\vec{p}^{(k)})^T \cdot \vec{p}^{(k)}} \quad (3.117)$$

$$\vec{p}^{(k+1)} = \vec{z}^{(k+1)} - \beta_k \vec{p}^{(k)} \quad (3.118)$$

## 4 Interpolation

Aufgabe: Eine (stetige) Funktion zu bestimmen, welche gegebene Datenpunkte *bestmöglichst* abbildet.

Eine möglicher Herangehensweise an diese Problem ist eine Entwicklung in eine Taylorreihe, welche allerdings gewisse Nachteile aufweist:

- konzentriert sich auf einen Punkt
- keine Aussage über Genauigkeit an anderen Punkten
- Approximation verbessert sich nicht mit höheren Graden des Approximationspolynoms: z.B.:  $e^x$ ,  $\frac{1}{x}$

Eine Interpolation hingegen zielt auf eine exakte Abbildung an den gegebenen Stützstellen eine Näherung dazwischen ab.

Allgemein formuliert lautet das Interpolationsproblem: Gegeben  $n+1$  Paare von reellen oder komplexen Zahlen  $((x_0, f(x_0)), (x_1, f(x_1)), \dots, (x_n, f(x_n)))$ . Mit Stützstellen bezeichnet man  $x_i$ , mit Stützwert  $f(x_i)$  und mit Stützpunkt  $(x_i, f(x_i))$ .

Ziel ist die Bestimmung von  $\Phi(x, a_0, \dots, a_n)$ ,  $a_0, a_1, \dots, a_n$ , so dass  $\forall_{i=0}^n f(x_i) = \Phi(x_i)$

### 4.1 Lagrange Interpolationspolynom

$$L_{n,k}(x) = \frac{(x - x_0) \cdot (x - x_1) \cdot \dots \cdot (x - x_{k-1}) \cdot (x - x_{k+1}) \cdot \dots \cdot (x - x_n)}{(x_k - x_0) \cdot (x_k - x_1) \cdot \dots \cdot (x_k - x_{k-1}) \cdot (x_k - x_{k+1}) \cdot \dots \cdot (x_k - x_n)} \quad (4.1)$$

$$L_k(x) = L_{n,k} = \prod_{i=0, i \neq k}^n \frac{(x - x_i)}{(x_k - x_i)} \quad (4.2)$$

$$L_{n,k}(x_k) = 1, L_{n,k}(x_i) = 0 \text{ für } x_i \neq x_k$$

$$P(x) = \sum_{k=0}^n f(x_k) \cdot L_{n,k}(x) \quad (4.3)$$

Theorem: Seien  $n+1$  paarweise unterschiedliche Stützstellen  $x_i$  sowie die zugehörigen Stützwerte  $f(x_i)$  gegeben. Dann existiert ein eindeutiges Interpolationspolynom  $P(x)$  vom Grad höchstens  $n$  mit  $P(x_i) = f(x_i)$

Fehlerabschätzung: Gegeben seien  $n + 1$  paarweise unterschiedliche Stützstellen  $x_i \in [a, b]$  sowie eine Funktion  $f$  die auf diesem Intervall  $[a, b]$   $(n + 1)$ -mal stetig differenzierbar ist. Dann existiert für  $x \in [a, b]$  ein  $\xi(x) \in (a, b)$  mit der Eigenschaft

$$f(x) = P(x) + \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \prod_{i=0}^n (x - x_i) \quad (4.4)$$

## 4.2 Dividierte Differenzen

Sei  $P(x)$  das Lagrange-Polynom  $n$ -ter Ordnung. Wähle

$$P_n(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \dots + a_n(x - x_0)(x - x_1) \cdot \dots \cdot (x - x_{n-1}) \quad (4.5)$$

Setzt man nun kontinuierlich die Stützstellen  $x_i$  in  $P_n(x)$  erhält man

$$P_n(x_0) = a_0 = f(x_0) \quad (4.6)$$

und somit den Wert für  $a_0$ . Aus

$$P_n(x_1) = f(x_0) + a_1(x_1 - x_0) = f(x_1) \quad (4.7)$$

folgt dann die Bestimmung von  $a_1$

$$a_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0} \quad (4.8)$$

*Dividierte Differenzen-Notation:*

- 0. dividierte Differenz

$$f[x_i] = f(x_i) \quad (4.9)$$

- 1. dividierte Differenz

$$f[x_i, x_{i+1}] = \frac{f[x_{i+1}] - f[x_i]}{x_{i+1} - x_i} = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} \quad (4.10)$$

- k. dividierte Differenz

$$f[x_i, x_{i+1}, \dots, x_{i+k}] = \frac{f[x_{i+1}, x_{i+2}, \dots, x_{i+k}] - f[x_i, x_{i+1}, x_{i+2}, \dots, x_{i+k-1}]}{x_{i+k} - x_i} \quad (4.11)$$

Mithilfe dieser Notation erhält man  $P_n(x)$  als sogenannte Newtons dividierte Differenz

$$P_n(x) = f[x_0] + \sum_{k=1}^n f[x_0, x_1, \dots, x_k] \cdot (x - x_0) \cdot (x - x_{k-1}) \quad (4.12)$$

Durch eine tabellarische Darstellung erhält man eine systematische Vorgehensweise zur Berechnung der dividierten Differenzen.

$x$	$f(x)$	1. div. Diff.	2. div. Diff.
$x_0$	$f[x_0]$	$f[x_0, x_1] = \frac{f[x_1] - f[x_0]}{x_1 - x_0}$	$f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_1, x_0]}{x_2 - x_0}$
$x_1$	$f[x_1]$	$f[x_1, x_2] = \frac{f[x_2] - f[x_1]}{x_2 - x_1}$	$\vdots$
$x_2$	$f[x_2]$	$f[x_2, x_3] = \frac{f[x_3] - f[x_2]}{x_3 - x_2}$	
$\vdots$	$\vdots$	$\vdots$	

Bei äquidistanten Stützstellen bietet sich eine vereinfachte Notation an

$$h = x_{i+1} - x_i \quad (4.13)$$

$$x = x_0 + sh \quad (4.14)$$

$$x - x_i = (s - i) \cdot h \quad (4.15)$$

Damit ist dann

$$P_n(x) = P_n(x_0 + s \cdot h) = s \cdot h \cdot f[x_0, x_1] + s \cdot (s-1) \cdot h^2 f[x_0, x_1, x_2] + \dots + (s-n+1) h^n f[x_0, x_1, \dots, x_n] \quad (4.16)$$

Mithilfe von Binomialkoeffizienten

$$\binom{s}{k} = \frac{s(s-1) \cdot \dots \cdot (s-k+1)}{k!} \quad (4.17)$$

$$P_n(x) = f[x_0] + \sum_{k=1}^n \binom{s}{k} k! h^k f[x_0, x_1, \dots, x_k] \quad (4.18)$$

$\Delta$ -Notation:  $\Delta f(x_i) = f(x_{i+1}) - f(x_i)$

$$f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0} = \frac{\Delta f(x_0)}{h} \quad (4.19)$$

$$f[x_0, x_1, x_2] = \frac{1}{2h} \frac{\Delta f(x_1) - \Delta f(x_0)}{h} = \frac{1}{2} \Delta^2 f(x_0) \quad (4.20)$$

$$f[x_0, x_1, \dots, x_k] = \frac{1}{k! h!} \Delta^k f(x_0) \quad (4.21)$$

Damit erhält man dann die Newton'sche Vorwärtsdifferenz

$$P_n(x) = f(x_0) + \sum_{k=1}^n \binom{s}{k} \Delta^k f(x_0) \quad (4.22)$$

## 4.3 Spline - Interpolation

**Grundidee:** Interpolation durch stückweise stetige Polynome.

Einfachste Variante: stückweise lineare Interpolation (siehe Abbildung 4.1)

Nachteil: keine Differenzierbarkeit an Endpunkten der einzelnen linearen Polynome gewährleistet.

Alternative: stückweise quadratische Polynome

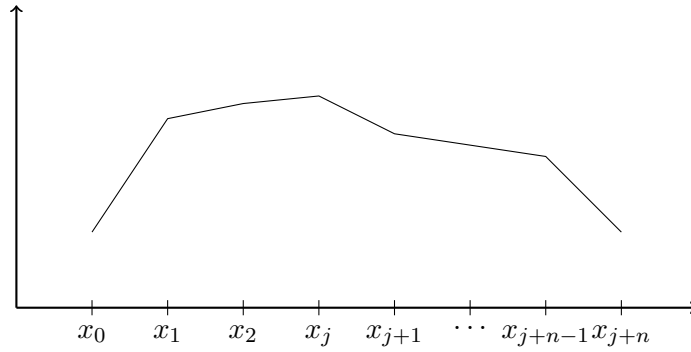


Abbildung 4.1: Beispielhafte Darstellung einer stückweise linearen Interpolation

- je 3 Koeffizienten  
 $\Rightarrow$  nur 2 für Bestimmung der Endpunkte erforderlich  
 $\Rightarrow$  3. Koeffizient kann stetige Differenzierbarkeit auf  $[x_0, x_n]$  sicherstellen
- ABER: Anforderung an Ableitung in  $[x_0, x_n]$  können nicht berücksichtigt werden.

Weitere Alternative: stückweise kubische Polynome ("cubic splines")

- je 4 Koeffizienten  
 $\Rightarrow$  stetige Differenzierbarkeit auf  $[x_0, x_n]$   
 $\Rightarrow$  stetige 2. Ableitung

**Definition:** Gegeben sei eine auf  $[a, b]$  definierte Funktion  $f$  sowie Stützstellen  $a = x_0 < x_1 < \dots < x_n = b$ . Ein kubischer Spline-Interpolant  $S$  erfüllt folgende Bedingungen:

1.  $S(x)$  ist ein kubisches Polynom, auf  $[x_j, x_{j+1}]$  als  $S_j(x)$  bezeichnet, für  $j = 0, 1, \dots, n-1$
2.  $S_j(x_j) = f(x_j)$ ,  $S_j(x_{j+1}) = f(x_{j+1})$ , für  $j = 0, 1, \dots, n-1$
3.  $S_j(x_{j+1}) = S_{j+1}(x_{j+1})$ , für  $j = 0, 1, \dots, n-2$
4.  $S'_{j+1}(x_{j+1}) = S'_j(x_{j+1})$ , für  $j = 0, 1, \dots, n-2$
5.  $S''_{j+1}(x_{j+1}) = S''_j(x_{j+1})$ , für  $j = 0, 1, \dots, n-2$
6. Eine der folgenden Randbedingungen ist erfüllt
  - $S''(x_0) = S''(x_n) = 0$  (freier bzw. natürlicher Rand)
  - $S'(x_0) = f'(x_0)$ ,  $S'(x_n) = f'(x_n)$  (eingespannter Rand)

$$S_j(x) = a_j + b_j(x - x_j) + c_j(x - x_j)^2 + d_j(x - x_j)^3, \quad j = 0, 1, \dots, n-1 \quad (4.23)$$

$$S_j(x_j) = a_j = f(x_j), \quad j = 0, 1, \dots, n-1 \quad (4.24)$$

$$a_{j+1} = S_{j+1}(x_{j+1}) = S_j(x_{j+1}) = a_j + b_j(x_{j+1} - x_j) + c_j(x_{j+1} - x_j)^2 + d_j(x_{j+1} - x_j)^3, \quad j = 0, 1, \dots, n-1 \quad (4.25)$$

$$\text{mit } h_j = x_{j+1} - x_j: \quad a_{j+1} = a_j + b_j \cdot h_j + c_j \cdot h_j^2 + d_j \cdot h_j^3, \quad j = 0, 1, \dots, n-1 \quad (4.26)$$



$$a_n := f(x_n) \text{ (Hilfsdefinition)} \quad (4.27)$$

$$b_n := S'(x_n) \quad (4.28)$$

$$S'_j(x) = b_j + 2c_j \cdot (x - x_j) + 3d_j \cdot (x - x_j)^2, \quad j = 0, 1, \dots, n-1 \quad (4.29)$$

$$S'_j(x_j) = b_j \quad (4.30)$$

$$b_{j+1} = b_j + 2c_j \cdot h_j + 3d_j \cdot h_j^2 \quad (4.31)$$

$$c := \frac{S''(x_n)}{2} \quad (4.32)$$

$$S''_j(x) = 2c_j + 6d_j \cdot (x - x_j) \quad (4.33)$$

$$S''_j(x) = 2c_j \quad (4.34)$$

$$S''_{j+1}(x_{j+1}) = 2c_{j+1} \quad (4.35)$$

$$S''_{j+1}(x_{j+1}) = S''_j(x_{j+1}) \quad (4.36)$$

$$c_{j+1} = c_j + 3d_j \cdot h_j \quad (4.37)$$

$$d_j = \frac{1}{3h_j} \cdot (c_{j+1} - c_j), \quad j = 0, 1, \dots, n-1 \quad (4.38)$$

Gleichung 4.26 in Gleichung 4.31:

$$a_{j+1} = a_j + b_j \cdot h_j + \frac{h_j^2}{3} \cdot (2c_j + c_{j+1}) \quad (4.39)$$

$$b_{j+1} = b_j + h_j \cdot (c_j + c_{j+1}) \quad (4.40)$$

$$b_j = b_{j-1} + h_j \cdot (c_{j+1} + c_j) \quad (4.41)$$

Aus Gleichung 4.39:

$$b_j = \frac{1}{h_j} \cdot (a_{j+1} - a_j) - \frac{h_j}{3} \cdot (2c_j + c_{j+1}) \quad (4.42)$$

$$b_{j-1} = \frac{1}{h_{j-1}} \cdot (a_j - a_{j-1}) - \frac{h_{j-1}}{3} \cdot (2c_{j-1} + c_j) \quad (4.43)$$

In Gleichung 4.42:

$$h_{j-1} c_{j-1} + 2 \cdot (h_{j-1} + h_j) \cdot c_j + h_j c_{j+1} = \frac{3}{h_j} \cdot (a_{j+1} - a_j) - \frac{3}{h_{j-1}} \cdot (a_j - a_{j-1}), \quad j = 0, 1, \dots, n-1 \quad (4.44)$$

**Natürliche Splines:**  $S''(x_0) = S''(x_n) = 0$

$$\Rightarrow c_n = \frac{S'(x_n)}{2} = 0 \quad (4.45)$$

$$0 = S''(x_0) = 2c_0 + 6d_0 \cdot (x_0 - x_0) \Rightarrow c_0 = 0 \quad (4.46)$$

Mit Gleichung 4.44:

$$\begin{bmatrix} 1 & 0 & 0 & 0 & \dots & \\ h_0 & 2 \cdot (h_0 + h_1) & h_1 & 0 & & \\ 0 & h_1 & 2 \cdot (h_1 + h_2) & h_2 & 0 & \\ 0 & \dots & & & & \\ 0 & \dots & & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} 0 \\ \frac{3}{h_j} \cdot (a_2 - a_1) - \frac{3}{h_0} \cdot (a_1 - a_0) \\ \vdots \\ \vdots \\ 0 \end{bmatrix} \quad (4.47)$$

Matrix streng diagonaldominant (Diagonalelemente betragsmäßig größer als Summe aller anderen Elemente einer Zeile)

$\Rightarrow$  eindeutige Lösung

# 5 Numerische Infinitesimalrechnung

## 5.1 Numerische Differentiation

$$f'(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h} \quad (5.1)$$

Taylor-Polynom:

$$f(x_0 + h) = f(x_0) + h \cdot f'(x_0) + \frac{h^2}{2} \cdot f''(x_0) + \frac{h^3}{3} \cdot f^{(3)}(x_0) + \dots$$

Abbruch nach linearem Term:

$$f(x_0 + h) = f(x_0) + h \cdot f'(x_0) + R(h^2)$$

### 5.1.1 Vorwärtsdifferenz

$$f'(x_0) \approx \frac{f(x_0 + h) - f(x_0)}{h} + R(h) \quad (5.2)$$

### 5.1.2 Rückwärtsdifferenz

$$f'(x_0) \approx \frac{f(x_0) - f(x_0 - h)}{h} \quad (5.3)$$

### 5.1.3 Zentrierte Differenz

$$f'(x_0) \approx \frac{f(x_0 + h) - f(x_0 - h)}{2h}$$

## Fehleranalyse

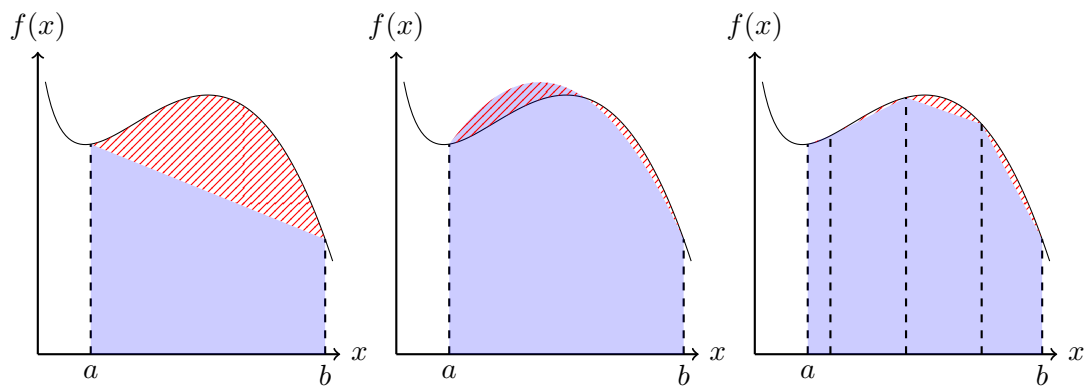
$$\begin{aligned} f(x_0 + h) &= \underset{\text{gerundeter Wert}}{\tilde{f}(x_0 + h)} + e_{-1} \\ f(x_0 - h) &= \tilde{f}(x_0 - h) + e_{+1} \end{aligned} \quad (5.4)$$

$$f'(x_0) = \frac{\tilde{f}(x_0 + h) - \tilde{f}(x_0 - h)}{2h} + \frac{e_{+1} - e_{-1}}{2h} - \frac{f^{(3)}(\xi)}{6} h^2 \quad (5.5)$$

Annahmen:  $|e_{-1}|, |e_{+1}| \leq \varepsilon$  und  $|f^{(3)}(\xi)| \leq M$

Damit  $\left| f'(x_0) - \frac{\tilde{f}(x_0 + h) - \tilde{f}(x_0 - h)}{2h} \right| \leq \frac{\varepsilon}{h} + \frac{h^2 M}{6}$

$$\Rightarrow h_{\text{opt}} = \sqrt[3]{\frac{3\varepsilon}{M}}$$



(a) Polynom 1. Ordnung

(b) Polynom 2. Ordnung

(c) mehrere Polynome

Abbildung 5.1: Unterschiedliche Methoden eine Funktion bei einer numerischen Integration anzunähern

## 5.2 Numerische Integration

Alternativer Begriff: „Numerische Quadratur“ (bestimmte Integrale)

$$J = \int_a^b f(x) \, dx \cong \int_a^b P_h(x) \, dx \quad (5.6)$$

$$P_n(x) = a_0 + a_1x + \dots + a_nx^n$$

### 5.2.1 Trapezregel

Lineares Lagrange Polynom

$$P_1(x) = \frac{x - x_1}{x_0 - x_1} \cdot f(x_0) + \frac{x - x_0}{x_1 - x_0} \cdot f(x_1)$$

$$\int_a^b = \int_{x_0}^{x_1} P_1(x) \, dx + \underbrace{\frac{1}{2} \int_{x_0}^{x_1} f''(\xi(x)) \cdot (x - x_0) \cdot (x - x_1) \, dx}_{\text{Fehlerterm } \epsilon}$$

Mittelwert Integralrechnung

$$\epsilon = \frac{1}{2} f''(\xi) \int_{x_0}^{x_1} (x - x_0) \cdot (x - x_1) \, dx = \frac{1}{2} f''(\xi) \left[ \frac{x^3}{3} - \frac{(x_1 + x_0)}{2} x^2 + x_0 x_1 x \right]_{x_0}^{x_1} = \frac{-h^3}{12} f''(\xi) \quad (5.7)$$

mit  $h = b - a = x_1 - x_0$

Also:

$$\int_a^b f(x) \, dx = \left[ \frac{(x-x_1)^2}{2(x_0-x_1)} \cdot f(x_0) + \frac{(x-x_0)^2}{2(x_1-x_0)} \cdot f(x_1) \right]_{x_0}^{x_1} - \frac{h^3}{12} f''(\xi) = \frac{h}{2} [f(x_0) + f(x_1)] - \frac{h^3}{12} f''(\xi) \quad (5.8)$$

Exakt, sofern  $f(x)$  höchstens linear ist

### 5.2.2 Simpson Regel

Quadratisches Lagrange Polynom

$$\begin{aligned} \int_a^b f(x) \, dx &= \int_{x_0}^{x_2} \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} f(x_0) + \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} f(x_1) + \frac{(x-x_0)(x-x_1)}{(x_2-x_1)(x_2-x_1)} f(x_2) \, dx \\ &\quad + \int_{x_0}^{x_2} \frac{(x-x_0)(x-x_1)(x-x_2)}{6} f^{(3)}(\xi(x)) \, dx \\ &= \frac{h}{3} [f(x_0) + 4f(x_1) + f(x_2)] - \frac{h^5}{90} f^{(4)}(\xi) \end{aligned} \quad (5.9)$$

mit  $h = x_1 - x_0 = x_2 - x_1$

### Simpson $\frac{3}{8}$ Regel

$$\int_a^b f(x) \, dx = \frac{3h}{8} [f(x_0) + 3f(x_1) + 3f(x_2) + f(x_3)] - \frac{h^5}{6480} f^{(4)}(\xi) \quad (5.10)$$

- benötigt  $3m$  Segmente bzw.  $3m + 1$  Stützpunkte
- Grad der Genauigkeit/„Degree of precision“  
Größte ganze Zahl  $n$ , so dass eine Integrationsformel für  $x^k$ ,  $k = 0, 1, \dots, n$  exakte Resultate liefert.

**Composite Simpson  $\frac{1}{3}$** 

Teile  $[a, b]$  in  $n$  Teilintervalle auf. Wende die Simpson Regel auf jedes Paar nacheinanderfolgender Teilintervalle an.

$$h = \frac{b-a}{n}; x_j = a + j \cdot h, j = 0, 1, \dots, n$$

$$\begin{aligned} \int_a^b f(x) dx &= \sum_{j=1}^{\frac{n}{2}} \int_{x_{2j-2}}^{x_{2j}} f(x) dx \\ &= \sum_{j=1}^{\frac{n}{2}} \frac{h}{3} [f(x_{2j-2}) + 4f(x_{2j-1}) + f(x_{2j})] \frac{h^5}{90} f^{(4)}(\xi_j) \\ &= \frac{h}{3} \left[ f(x_0) + 2 \sum_{j=1}^{\frac{n}{2}-1} f(x_{2j}) + 4 \sum_{j=1}^{\frac{n}{2}} f(x_{2j-1}) + f(x_n) \right] - \frac{h^5}{90} \sum_{j=1}^{\frac{n}{2}} f^{(4)}(\xi_0) \end{aligned} \quad (5.11)$$

mit  $x_{2j-2} < \xi_j < x_{2j}$ ,  $f$  auf  $[a, b]$  4-mal stetig differenzierbar.

**Fehler**

$$E(f) = \frac{-h^5}{90} \sum_{j=1}^{\frac{n}{2}} f^{(4)}(\xi_j), \text{ mit } x_{2j-2} < \xi_j < x_{2j}, j = 1, 2, \dots, \frac{n}{2}$$

**Extremwertsatz**

:  $f^{(4)}$  nimmt Max und Min auf  $[a, b]$  an, also

$$\begin{aligned} \min_{x \in [a, b]} f^{(4)}(x) &\leq f^{(4)}(\xi_j) \leq \max_{x \in [a, b]} f^{(4)}(x) \\ \frac{n}{2} \min_{x \in [a, b]} f^{(4)}(x) &\leq \sum_{j=1}^{\frac{n}{2}} f^{(4)}(\xi_j) \leq \frac{n}{2} \max_{x \in [a, b]} f^{(4)}(x) \\ \min_{x \in [a, b]} f^{(4)}(x) &\leq \frac{n}{2} \sum_{j=1}^{\frac{n}{2}} f^{(4)}(\xi_j) \leq \max_{x \in [a, b]} f^{(4)}(x) \end{aligned} \quad (5.12)$$

**Zwischenwertsatz**

Es existiert ein  $\mu \in (a, b)$ , so dass

$$f^{(4)}(\mu) = \frac{2}{n} \sum_{j=1}^{\frac{n}{2}} f^{(4)}(\xi_j) \quad (5.13)$$

Damit

$$E(f) = \frac{-h^5}{90} \sum_{j=1}^{\frac{n}{2}} f^{(4)}(\xi_j) = \frac{-h^5}{180} n \cdot f^{(4)}(\mu) \stackrel{h=\frac{b-a}{n}}{=} \frac{-(b-a)}{180} h^4 \cdot f^{(4)}(\mu) \quad (5.14)$$

Also, mit  $\mu \in (a, b)$ :

$$\int_a^b f(x) dx = \frac{h}{3} \left[ f(a) + 2 \sum_{j=1}^{n/2-1} f(x_{2j}) + 4 \sum_{j=1}^{n/2} f(x_{2j-1}) + f(b) \right] - \frac{b-a}{180} h^4 f^{(4)}(\mu) \quad (5.15)$$

Wie man sich nun vielleicht vorstellen kann gibt es unendlich viele mögliche numerische Integrationsverfahren. Das Composite Simpson-Verfahren ist allerdings ein häufig eingesetztes allgemeines Verfahren.

### 5.2.3 Genauigkeit

Zur Bestimmung der Genauigkeit gibt es zwei Ansätze:

1. Gegeben sei ein Verfahren, eine Funktion  $f$ , ein Intervall  $[a, b]$  und eine Stützstellenanzahl  $n$  und gesucht sei der Fehler
2. Gegeben sei eine Funktion  $f$ , ein Intervall  $[a, b]$  und ein gewünschter Fehler und gesucht sei ein Verfahren und eine Stützstellenanzahl  $n$

Der Rundungsfehler wird definiert über

$$f(x_i) = \underbrace{\tilde{f}(x_i)}_{\text{gerundeter Wert}} + \underbrace{e_i}_{\text{Rundungsfehler}} \quad (5.16)$$

Damit erhält man dann einen akkumulierten Fehler  $e(h)$  bei Composite-Simpson von

$$e(h) = \left| \frac{h}{3} \left[ e_0 + 2 \sum_{j=1}^{n/2-1} e_{2j} + 4 \sum_{j=1}^{n/2} e_{2j-1} + e_n \right] \right| \leq \frac{h}{3} \left[ |e_0| + 2 \sum_{j=1}^{n/2-1} |e_{2j}| + 4 \sum_{j=1}^{n/2} |e_{2j-1}| + |e_n| \right] \quad (5.17)$$

Seine alle Rundungsfehler  $|e_i| \leq \varepsilon$ , dann gilt

$$e(h) \leq \frac{h}{3} \left[ \varepsilon + 2\left(\frac{n}{2} - 1\right)\varepsilon + 4\frac{n}{2}\varepsilon + \varepsilon \right] = nh\varepsilon = (b-a)\varepsilon \neq f(h, n) \quad (5.18)$$

Erstaunlich ist somit, dass der Rundungsfehler nicht wächst, wenn der zu integrierende Bereich in mehr Teilintervalle unterteilt wird.





## 6 Numerische Lösung von Differentialgleichungen

Bsp: Einschwinganalyse (TRANSIENTE Analyse)

### 6.1 Analytische Lösung mittels Laplace-Transformation

$$\begin{aligned}
 C \cdot \dot{y}(t) + \frac{1}{R} \cdot y(t) &= i_0(t) \\
 \dot{y}(t) + \underbrace{\frac{1}{RC}}_{-p_\infty} \cdot y(t) &= \frac{1}{RC} R \cdot i_0(t) \\
 \dot{y}(t) - p_\infty y(t) &= -p_\infty R \cdot i_0(t) \\
 \circ \xrightarrow{\mathcal{L}} \bullet \text{ Laplace} & \\
 p \cdot Y(p) - \underbrace{y(+0)}_{=0} - p_\infty Y(p) &= -p_\infty R \frac{1}{p} \cdot I_0 \\
 Y(p) &= \frac{-p_\infty R \cdot I_0}{(p - p_\infty)} = \frac{A}{p} + \frac{B}{p - p_\infty} \\
 A = \lim_{p \rightarrow 0} p \cdot Y(p) &= R \cdot I_0 \\
 B = \lim_{p \rightarrow p_\infty} (p - p_\infty) \cdot Y(p) &= -R \cdot I_0
 \end{aligned}
 \tag{6.1}$$

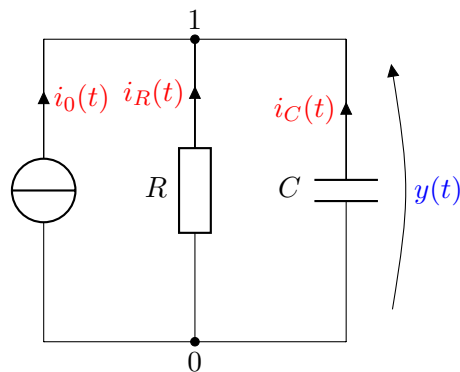


Abbildung 6.1: Beispielschaltung Transiente Analyse

$$y(t) = A + Be^{p_{\infty}t} = R \cdot I_0 \cdot (1 - e^{p_{\infty}t}) \quad (6.4)$$

## 6.2 Numerische Lösung mittels explizite Euler-Methode (linear Z-Transformation)



Diskretisierung der Zeit

$$y(\nu) \hat{=} y(\nu \cdot \Delta t) = y(t) \quad (6.5)$$

Diskretisierung der DGL

$$\dot{y}(\nu) \approx p_{\infty} \cdot y(\nu) - p_{\infty} - p_{\infty} R \cdot i_0(\nu) \quad (6.6)$$

Differenzengleichung (explizite Euler-Methode)

$$y(\nu + 1) \approx y(\nu) + \Delta t \cdot \dot{y}(\nu) \quad (6.7)$$



$$y(\nu + 1) \approx y(\nu) + \Delta t \cdot [p_{\infty} y(\nu) - p_{\infty} R \cdot i_0(\nu)] \quad , \nu = 0, 1, 2, \dots \quad (6.8)$$

$$\hat{y}(\nu + 1) = (1 + p_{\infty} \Delta t) \cdot y(\nu) - \Delta t p_{\infty} R i_0(\nu) \quad (6.9)$$

sukzessive Berechnung ausgehend von  $y(0)$ ,  $\nu = 1, 2, 3, \dots$

Bei linearen Schaltungen (linearen Differenzengleichungen) ist geschlossene numerische Lösung möglich durch  $\circ \xrightarrow{z} \bullet$  von Gleichung 6.9

$$z \cdot \dot{Y}(z) - \underbrace{z \cdot y(0)}_0 - (1 - p_{\infty} \Delta t) \cdot Y(z) = -p_{\infty} R \Delta t \frac{z}{z - 1} I_0 \quad (6.10)$$

$$Y(z) = \frac{-p_\infty \Delta t R}{z - (1 + p_\infty \Delta t)} \frac{z I_0}{z - 1} = \frac{z R I_0}{z - 1} \frac{z R I_0}{z - (1 + p_\infty \Delta t)}$$

$$\circ \xrightarrow{z} \bullet$$

$$\hat{y}(\nu) = [1 - (1 + p_\infty \Delta t)^\nu] \cdot R I_0 \quad , \nu = 0, 1, 2, \dots$$
(6.11)

Zahlenbeispiel:  $I_0 R = 1$ ;  $p_\infty = -1$ ,  $t = 1$

exakte Lösung:  $y(1) = 1 - e^{-1} = 0,632\,121$

numerische Lösung:  $y(\nu) = 1 - (1 - \Delta t)^\nu$  ,  $\nu \cdot \Delta t = 1$

$\Delta t$	$\nu$	$\hat{y}(1)$	$\varepsilon^{(A)} =  \hat{y}(1) - y(1) $	Fehler $\varepsilon^{(A)} \sim \Delta t$
0,1	10	0,651 322	0,019 201	
0,05	20	0,641 514	0,009 393	
0,025	40	0,636 768	0,004 647	

#### Stabilität:

Für  $|1 + p_\infty \cdot \Delta t| = |1 - \Delta t| < 1 \quad \Rightarrow \lim_{x \rightarrow \infty} \hat{y}(\nu) = 1$

$0 < \Delta t < 2$  Algorithmus stabil

Für  $|1 - \Delta t| > 1 \quad \Rightarrow \lim_{\nu \rightarrow \infty} \hat{y}(\nu) \rightarrow \infty$

$\Delta t > 2$  Algorithmus instabil

## 6.3 Einfache numerische Integrationsverfahren

$$y(0) = 0, \quad y(t) = \int_0^t \dot{y}(\tau) \, d\tau$$

$$y(t + \Delta t) = \int_0^t \dot{y}(\tau) \, d\tau + \int_t^{t+\Delta t} \dot{y}(\tau) \, d\tau = y(t) + \int_t^{t+\Delta t} \dot{y}(\tau) \, d\tau$$

$$y(\nu + 1) \approx y(\nu) + \Delta y = y(\nu) + \frac{\Delta t \cdot \tilde{\Delta y}}{\Delta t}$$
(6.12)

**Expliziter Euler:**  $y(\nu + 1) \approx y(\nu) + \Delta t \cdot \dot{y}(\nu)$  (FE, Forward Euler)



**Impliziter Euler:**  $y(\nu + 1) \approx y(\nu) + \Delta t \cdot \dot{y}(\nu + 1)$  (BE, Backward Euler)



**Trapez-Methode:**  $y(\nu + 1) \approx y(\nu) + \frac{\Delta t}{2} \cdot (\dot{y}(\nu) + \dot{y}(\nu + 1))$  (TR, Trapezvidal)



## 6.4 Eigenschaften numerischer Integrationsverfahren

- explizit/implizit: wird Wert aus aktuellem Zeitschritt  $(\nu + 1)$  verwendet?
- Schrittzahl: Anzahl der verwendeten vergangenen Zeitschritte
- Genauigkeit:  $\varepsilon^{(A)}(\nu \cdot \Delta t) = \hat{x}(\nu) - x(\nu) \approx (\Delta t)^k$   
(für  $|p_\infty \Delta t| \ll 1$ ) Fehler der Ordnung  $k$
- Stabilität:
  - für  $\text{Re}\{p_\infty\} < 0$
  - für  $\Delta t \rightarrow \infty$  (“asymptotische Stabilität“)

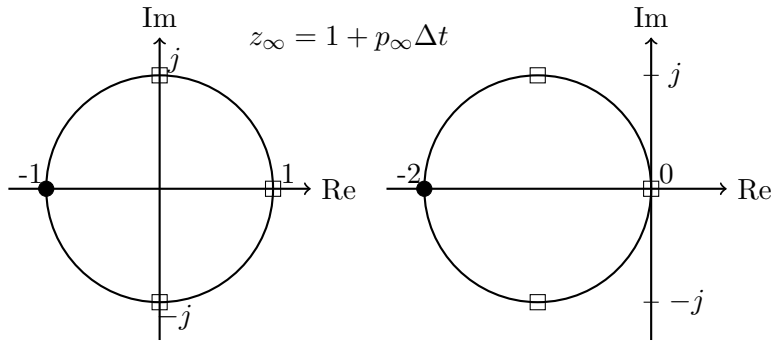


Abbildung 6.2: Stabilitätsbereich für den expliziten Euler

<b>Test-DGL:</b>	$\dot{x}(t) = p_\infty \cdot x(t)$	$p_\infty$ : Eigenwerte des Systems
exakte Lösung:	$x(t) = x(0) \cdot e^{p_\infty t}$	
	$x(\nu) = x(0) \cdot e^{p_\infty \cdot \nu \cdot \Delta t}, x(0) \neq 0$	
numerische Lösung:	$\hat{x}(\nu) = \int \{\hat{f}(\nu)\}$	
Fehler:	$\varepsilon^{(A)}(\Delta t) = \hat{x}(\nu) - x(\nu) = \hat{x}(\nu) - x(0) \cdot e^{p_\infty \nu \Delta t}$	

## 6.5 Expliziter Euler Eigenschaften

**Differenzengleichung:**  $\hat{x}(\nu + 1) = \hat{x}(\nu) + \Delta t \cdot \hat{f}(\nu), f(t) = \dot{x}(t)$   
 Einschnitt-Verfahren

<b>Test-DGL:</b>	$f(t) = \dot{x}(t) = p_\infty \cdot x(t)$
	$x(0) \neq 0$
Differenzengleichung:	$\hat{x}(\nu + 1) = \hat{x}(\nu) + \Delta t p_\infty \cdot \hat{x}(\nu) = (1 + p_\infty \cdot \Delta t) \cdot \hat{x}(\nu)$
numerische Lösung:	$\hat{x}(\nu) = (1 + p_\infty \Delta t) \cdot x(\nu)$
Genauigkeit:	$\varepsilon^{(A)} = \hat{x}(\nu) - x(\nu) = \Delta t \cdot \overbrace{\frac{\nu \Delta t}{2}}^{=t} \cdot p_\infty^2 \cdot x(0) \cdot e^{p_\infty \overbrace{\nu \Delta t}^{=t}}$
Stabilität:	$\lim_{t \rightarrow \infty} x(t) \rightarrow 0$ für $\text{Re}\{p_\infty\} < 0$
	$\lim_{\nu \rightarrow \infty} \hat{x}(\nu) = \lim_{\nu \rightarrow \infty} (1 + p_\infty \Delta t)^\nu x(0) \rightarrow 0$ falls $ 1 + p_\infty \Delta t  < 1$
Ordnung	$k = 1$
Asymptotische Stabilität ( $\Delta t \rightarrow \infty$ ):	$\hat{x}(\nu + 1) \approx p_\infty \Delta t \hat{x}(\nu) \Rightarrow$ instabiles Verhalten!

Somit ist die Wahl von  $\Delta t$  aus Stabilitätsgründen eingeschränkt. Bei zu großem  $\Delta t$  liegt das Produkt  $p_\infty \Delta t$  nicht mehr im Stabilitätsbereich. Praktisch beim expliziten Euler ist hingegen, dass es keine Stabilität in der rechten Halbebene gibt, es werden also instabile System als instabil detektiert.

Im folgenden noch einmal die Eigenschaften der Verfahren kurz zusammengefasst:

	explizit/implizit	Schritte	Ordnung	stabil für $\operatorname{Re}\{p_\infty\} < 0$	asymptotisch stabil
expliziter Euler	explizit	1	1	nein	nein
impliziter Euler	implizit	1	1	ja	ja
Trapez	implizit	1	2	ja	nein
Gear	implizit	2	2	ja	ja

Zudem haben Einschrittverfahren Vorteile beim Starten, mit Mehrschrittverfahren kann nicht direkt gestartet werden.

## 6.6 Taylor-Verfahren höherer Ordnung

Die eher im technischen Bereich gängige Schreibweise

$$y(\nu + 1) = y(\nu) + \Delta t T^{(n)}(\nu, y(\nu)) \quad (6.13)$$

ist äquivalent zur in der Mathematik üblichen Variante

$$y_{i+1} = y_i + h T^{(n)}(t_i, y_i) \quad (6.14)$$

mit

$$T^{(n)}(t_i, y_i) = f(t_i, y_i) + \frac{h}{2} f'(t_i, y_i) + \cdots + \frac{h^{n-1}}{n!} f^{(n-1)}(t_i, y_i) \quad (6.15)$$

Dadurch erhält man eine bessere Genauigkeit, genauer gesagt:  $\epsilon O(h^n)$ . Nachteilig kann sich allerdings auswirken, dass höhere Ableitungen der Funktion benötigt werden.

## 6.7 Prädiktor-Korrektor-Ansätze

Die Grundidee dahinter ist ein explizites Verfahren zur Vorhersage des Wertes  $y_{i+1}$  zu nutzen, und diesen Wert nachträglich mithilfe eines weiteren Verfahrens zu verbessern. Dieser Ansatz kann auch iterativ verwendet werden.

Ein Beispiel für so einen Prädiktor-Korrektor-Ansatz ist das Verfahren von Heun, welches auf einem expliziten Euler basiert. Der explizite Euler liefert

$$y_{i+1}^0 = y_i + h \cdot f(t_i, y_i) \quad (6.16)$$

und daraus berechnet Heun ein verbessertes

$$y_{i+1}^1 = f(t_{i+1}, y_{i+1}^0). \quad (6.17)$$

Der iterative Ansatz dazu wäre dann

$$y_{i+1}^0 = y_i^m + h \cdot f(t_i, y_i) \quad (6.18)$$

$$y_{i+1}^j = y_i^m + h \cdot \frac{f(t, y_i^m) + f(t_{i+1}, y_{i+1}^{j-1})}{2} \quad (6.19)$$

## 6.8 Adams-Bashforth/Adams-Moulton

Die Grundidee hinter diesen Verfahren beruht darauf  $f$  durch ein Interpolationspolynom zu ersetzen.

Ausgangspunkt ist die Differentialgleichung

$$y' = f(t, y) \quad (6.20)$$

wobei  $a \leq t \leq b$  und  $y(a) = \alpha$ .

Adams-Bashforth und Adams-Moulton sind  $m$ -schrittige Mehrschrittverfahren, in denen  $y_{i+1}$  ersetzt wird durch

$$y_{i+1} = a_{m-1}y_i + a_{m-2}y_{i-1} + \cdots + a_0y_{i+1-m} + h[b_m f(t_{i+1}, y_{i+1}) + b_{m-1}f(t_i, y_i) + \cdots + b_0f(t_{i+1-m}, y_{i+1-m})] \quad (6.21)$$

mit  $i = m-1, m, \dots, N-1$ ,  $h = \frac{b-a}{N}$ . Falls  $b_n = 0$  spricht man von einem expliziten Verfahren, ansonsten von einem impliziten.

### Herleitung

$$y(t_{i+1}) = y_{i+1} \quad (6.22)$$

$$y_{i+1} - y_i = \int_{t_i}^{t_{i+1}} y'(t) dt = \int_{t_i}^{t_{i+1}} f(t, y(t)) dt \quad (6.23)$$

$$y_{i+1} = y_i + \int_{t_i}^{t_{i+1}} f(t, y(t)) dt \quad (6.24)$$

Ersetze  $f(t, y(t))$  durch Interpolationspolynom  $P(t)$

$$y_{i+1} \approx y_i + \int_{t_i}^{t_{i+1}} P(t) dt \quad (6.25)$$

Für das Interpolationspolynom bieten sich an

- Taylor
- Lagrange
- Newton'sche Rückwärtsdifferenzen

### 6.8.1 Adams-Bashforth (explizit)

Die explizite Einschnittvariante ist äquivalent zum expliziten Euler:

$$y_{i+1} = y_i + hf(t_i, y_i) \quad (6.26)$$

Die Variante mit zwei Schritten lautet

$$y_{i+1} = y_i + \frac{h}{2} [3f(t_i, y_i) - f(t_{i-1}, y_{i-1})] \quad (6.27)$$

und hat eine quadratisches Fehlerverhalten.

### 6.8.2 Adams-Moulton (implizit)

Die Einschnittvariante hiervon ist für  $m = 0$  der implizite Euler, für  $m = 1$  das Trapezverfahren. Wenn zwei Schritten verwendet werden sieht das Verfahren wie folgt aus

$$y_{i+1} = y_i + \frac{h}{12} [5f(t_{i+1}, y_{i+1}) + 8f(t_i, y_i) - f(t_{i-1}, y_{i-1})] \quad (6.28)$$

und besitzt ein kubisches Fehlerverhalten.



## 7 Ausgleichsrechnung

Auch bekannt unter der Bezeichnung Approximation oder Least Square **Ziele der Approximation:** („curve fitting“)

- Bestgeeignete Funktion eines bestimmten Typs, um gegebenen Daten anzunähern
- Für gegebene Funktionen „bestmögliche“, einfachere Funktionen zu finden.

**„Bestmögliche“ Funktion zur Annäherung gegebener Stützpunkte**

- Gegeben: Wertepaare  $(x_i, y_i)$ ,  $i = 1, \dots, m$
- Annahme: Gerade beschreibt Wertepaare  $y = f(x) = ax + b$

**Ansätze:**

- $E = \sum_{i=1}^m l_i = \sum_{i=1}^m [y_i - (ax_i + b)]$
- $E = \sum_{i=1}^m |l_i| = \sum_{i=1}^m |[y_i - (ax_i + b)]|$
- Minmax:  $\min E_\infty = \min_{i=1, \dots, m} \max \{|y_i - (ax_i + b)|\}$

### 7.1 Linear Least Squares

Minimiere:  $E_2(a, b) = \sum_{i=1}^m (y_i - (ax_i + b))^2$

Bestimmung von  $a, b$ :  $\frac{\partial E_2}{\partial a} = 0$ ,  $\frac{\partial E_2}{\partial b} = 0$

$\Rightarrow$  Normalengleichung:

- $a = \dots$
- $b = \dots$

### 7.2 Polynomiale Least Squares

Approximation durch  $P_n(x) = a_n x^n + \dots + a_1 x + a_0$

$$E_2(a_0, a_1, \dots, a_n) = \sum_{i=1}^m (y_i - P(x_i))^2 \quad (7.1)$$

Minimierung  $\frac{\partial E_2}{\partial a_i} = 0 \Rightarrow n + 1$  Normalengleichungen in den  $n + 1$  Unbekannten  $a_0, \dots, a_n$

### 7.3 Linear Least Squares - Sichtweise Lineare Algebra

$$\mathbf{A}_{<m \times n} \cdot \vec{x}_n = \vec{b}_n, m > n \quad (7.2)$$

Beispiel:

$$\begin{aligned} \bullet \quad m=2, n=1 & \quad \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \cdot x = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \\ \bullet \quad m=3, n=2 & \quad \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \end{aligned}$$

Fehler je Zeile:  $e_i = b_i - \vec{a}_i \cdot \vec{x}$

Vektor der Fehler:  $\vec{e} = \vec{b} - \mathbf{A} \cdot \vec{x}$

Summe der Fehlerquadrate:  $E_2 = \|\vec{e}\|^2 = \|\vec{b} - \mathbf{A}\vec{x}\|^2$

Gesucht  $\hat{\vec{x}}$ , so dass  $E_2$  minimiert wird.

$$\mathbf{A}\hat{\vec{x}} = \vec{p} \quad ; \quad \vec{e} = \vec{b} - \vec{p} \quad (7.3)$$

$\vec{p}$  ist die Projektion von  $\vec{b}$  in den Spaltenraum von  $\text{col}(\mathbf{A})$

$\vec{e}$  steht senkrecht auf den Spaltenraum  $\text{col}(\mathbf{A})$

$$\begin{aligned} E(x) &= \|\vec{b} - \mathbf{A}\vec{x}\|^2 = \|\mathbf{A}\vec{x} - \vec{b}\|^2 \\ &= \vec{x}^T \mathbf{A}^T \mathbf{A} \vec{x} - (\mathbf{A}\vec{x})^T \vec{b} - \vec{b}^T \mathbf{A} \vec{x} + \vec{b}^T \vec{b} \\ &= \vec{x}^T \underbrace{\mathbf{A}^T \mathbf{A}}_{\mathbf{K}} \vec{x} - 2\vec{x}^T \underbrace{\mathbf{A}^T \vec{b}}_{\vec{f}} + \vec{b}^T \vec{b} \\ &= \vec{x}^T \mathbf{K} \vec{x} - 2\vec{x}^T \vec{f} + \vec{b}^T \vec{b} \quad \text{zu minimieren} \\ &= (\vec{x} - \mathbf{K}^{-1} \vec{f})^T \mathbf{K} (\vec{x} - \mathbf{K}^{-1} \vec{f}) - \vec{f}^T \mathbf{K}^{-1} \vec{f} + \vec{b}^T \vec{b} \end{aligned} \quad (7.4)$$

Ausdruck kann minimal 0 werden für  $\hat{\vec{x}} = \mathbf{K}^{-1} \vec{f} \Rightarrow$  dort also Minimum vom  $E(\vec{x})$

$$\begin{aligned} E_{\min} &= E(\hat{\vec{x}}) = (\vec{b} - \mathbf{A}\hat{\vec{x}})^T \cdot (\vec{b} - \mathbf{A}\hat{\vec{x}}) \\ &= \vec{b}^T \vec{b} - \vec{f}^T \mathbf{K}^{-1} \vec{f} \\ &= \vec{b}^T \vec{b} - \vec{b}^T \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \vec{b} \end{aligned} \quad (7.5)$$

$\vec{e} = \vec{b} - \mathbf{A}\vec{x}$  steht senkrecht auf  $\text{col}(\mathbf{A}) \Rightarrow \mathbf{A}^T \vec{e} = \vec{0}$

$$\Rightarrow \mathbf{A}^T (\vec{b} - \mathbf{A}\hat{\vec{x}}) = \vec{0} \Rightarrow \mathbf{A}^T \mathbf{A} \hat{\vec{x}} = \mathbf{A}^T \vec{b} \quad (7.6)$$

$\underbrace{\mathbf{A}^T \mathbf{A} \hat{\vec{x}} = \mathbf{A}^T \vec{b}}_{\text{quadratisch, symmetrisch}} \quad \text{LGS} \Rightarrow \text{Bestimmung von } \hat{\vec{x}}$

**Lösungsverfahren:**

- Gauß LU: OK, aber  $\text{cond}(\mathbf{A}^T \mathbf{A}) = \text{cond}(\mathbf{A})^2 \Rightarrow$  Problematisch, wenn  $\mathbf{A}$  schlecht konditioniert
- Orthogonalzerlegung

## 7.4 QR-Zerlegung

$$A = QR \text{ mit } \begin{cases} A & \in \mathbb{R}^{n \times n} \\ Q & \in \mathbb{R}^{n \times n} \text{ orthogonale Matrix } (Q^T = Q^{-1}) \\ R & \in \mathbb{R}^{n \times n} \text{ rechte obere Dreiecksmatrix} \end{cases}$$

Damit:

$$\begin{aligned} A\vec{x} &= \vec{b} \\ A^T A\vec{x} &= A^T \vec{b} \\ (QR)^T QR\vec{x} &= (QR)^T \vec{b} \\ \underbrace{R^T Q^T Q}_{I} R\vec{x} &= R^T Q^T \vec{b} \\ R^T R\vec{x} &= R^T Q^T \vec{b} \quad | \cdot R^{-1} \\ R\vec{x} &= Q^T \vec{b} \Rightarrow \text{simple Rücksubstitution} \end{aligned} \tag{7.7}$$

OK, wie bestimme ich nun  $Q, R$ ?

### 3 wesentliche Ansätze:

- Gram-Schmidt
- Householder-Transformation
- Givens-Rotation

#### 7.4.1 Givens-Rotation

auch bekannt unter der Bezeichnung Jacobi-Rotation

Eine Rotation ist gegeben durch

$$G(l, k, \Theta) = \begin{bmatrix} 1 & 0 & \dots & \dots & \dots & \dots & \dots & 0 \\ 0 & 1 & & & & & & 0 \\ \vdots & & c & & & -s & & \vdots \\ \vdots & & & & & & & \vdots \\ \vdots & & -s & & & c & & \vdots \\ \vdots & & & & & & 1 & 0 \\ 0 & \dots & \dots & \dots & \dots & \dots & 0 & 1 \end{bmatrix} \tag{7.8}$$

wobei  $c = \cos \Theta$ ,  $s = \sin \Theta$ . Komponentenweise  $G(l, k, \Theta) = (g_{i,j}(l, k, \Theta))$  dargestellt ergibt das

$$g_{i,j} = \begin{cases} \cos \Theta & \text{für } i = l, j = l \vee i = k, j = k \\ \sin \Theta & \text{für } i = l, j = k \\ -\sin \Theta & \text{für } i = k, j = l \\ 1 & \text{für } i = j \text{ außer } i, j = l, k \\ 0 & \text{sonst} \end{cases} \tag{7.9}$$

Damit bedeutet  $\mathbf{G} \cdot \vec{x}$  bzw.  $\mathbf{G}^T \cdot \vec{x}$  eine Drehung des Vektors  $\vec{x}$  um  $\pm\Theta$  in der  $(l, k)$ -Ebene. Die Hauptanwendung hierfür ist ein iteratives Vorgehen um Nulleinträge in Matrizen oder Vektoren zu erreichen. Desweiteren ist  $\mathbf{G}\mathbf{G}^T = \mathbf{I}$  und somit die Rotationsmatrix  $\mathbf{G}$  orthogonal.

$$\begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} c & -s \\ s & c \end{pmatrix} = \begin{pmatrix} c^2 + s^2 & 0 \\ 0 & c^2 + s^2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad (7.10)$$

Damit gilt

$$\mathbf{G}_r \mathbf{G}_{r-1} \cdot \dots \cdot \mathbf{G}_2 \mathbf{G}_1 \mathbf{A} = \mathbf{R} \quad (7.11)$$

$$\mathbf{A} = \mathbf{G}^T \mathbf{R} = \mathbf{Q}^T \mathbf{R} \quad (7.12)$$

Die Frage ist nun, wie  $\Theta$  gewählt werden muss. Hierfür genügt die Betrachtung einer zweidimensionalen Struktur

$$\begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} x_l \\ x_k \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix} \quad (7.13)$$

Daraus kann gefolgert werden, dass

$$s = c \frac{x_k}{x_l} \quad (7.14)$$

und aus

$$c^2 + s^2 = 1 \quad (7.15)$$

$$c = \sqrt{1 - s^2} = \sqrt{\frac{x_l^2 - c^2 x_k^2}{x_l^2}} \quad (7.16)$$

$$c^2 x_l^2 = x_l^2 - c^2 - x_k^2 \quad (7.17)$$

$$c = \pm \frac{x_l}{\sqrt{x_l^2 + x_k^2}} \quad (7.18)$$

Das ganze an dem Beispiel

$$\mathbf{A} = \begin{pmatrix} 3 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 3 \end{pmatrix} \quad (7.19)$$

führt über

$$\mathbf{G} = \begin{pmatrix} c & s & 0 \\ -s & c & 1 \\ 0 & 0 & 1 \end{pmatrix} \quad (7.20)$$

zu

$$\mathbf{GA} = \begin{pmatrix} \sqrt{10} & 6/\sqrt{10} & 1/\sqrt{10} \\ 0 & 3/\sqrt{10} & 3/\sqrt{10} \\ 0 & 1 & 3 \end{pmatrix} \quad (7.21)$$

# Abbildungsverzeichnis

1.1	Teilgebiete der Mathematik . . . . .	1
1.2	Zusammenspiel Mathematik und Information in der Numerik . . . . .	2
1.3	Vorgehen zur Lösung eines mathematischen Problems . . . . .	2
1.4	Beispielschaltung . . . . .	7
2.1	Beispiel einer Schaltung . . . . .	11
2.2	Arbeitspunktermittlung durch Nullstellensuche . . . . .	12
2.3	Intervallhalbierung . . . . .	13
2.4	Newton-Iteration . . . . .	13
2.5	Fixpunktiteration . . . . .	15
3.1	Gradientenverfahren - steepest descent . . . . .	26
4.1	Beispielhafte Darstellung einer stückweise linearen Interpolation . . . . .	34
5.1	Unterschiedliche Methoden eine Funktion bei einer numerischen Integration anzunähern . . . . .	38
6.1	Beispielschaltung Transiente Analyse . . . . .	43
6.2	Stabilitätsbereich für den expliziten Euler . . . . .	47