

## Inhaltsverzeichnis

<b>1 Einführung in die gewöhnliche Differentialgleichungen</b>	<b>1</b>
1.1 Beispiel 1: „stetige“ Verzinsung . . . . .	1
1.1.1 Taylor-Entwicklung . . . . .	1
1.2 ‘Lineare Dgl‘ mit konstanten Koeffizienten . . . . .	2
1.3 Beispiel 2: Newton’sche Teilchenmechanik . . . . .	2
1.3.1 Gravitation . . . . .	2
1.3.2 Impulserhaltung . . . . .	2
1.3.3 Energieerhaltung . . . . .	3
1.4 Beispiel 3: Molekulardynamik . . . . .	3
1.5 Beispiel 4: Reaktionskinetik . . . . .	3
1.6 Beispiel 4: DNA-Replikation . . . . .	3
1.7 Beispiel 5: Populationsdynamik . . . . .	4
1.8 Beispiel 6: Räuber-Beute Modell(Lotke-Voltena Modell) . . . . .	4
1.9 Beispiel 7: Minimalkurven . . . . .	4
<b>2 Existenz und Eindeutigkeit des AWP</b>	<b>4</b>
2.1 Lokale Existenz . . . . .	4
2.2 Satz: Satz von Peno . . . . .	5
2.3 Schaudersche Fixpunktsatz . . . . .	5
2.3.1 Anwendung auf DGL . . . . .	5
2.4 Satz von Auzela Ascoli . . . . .	6
2.5 Satz von Picard-Lindelöf (globale Version) . . . . .	6
<b>3 Stabilität von DGL</b>	<b>8</b>
3.1 Kondition . . . . .	8
3.2 Einschrittverfahren für gewöhnliche DGL . . . . .	9
3.2.1 Eulersches Polygonzugverfahren(1726) . . . . .	9
3.3 Konsistenz . . . . .	10
3.3.1 Definition: Konsistenz . . . . .	10
3.3.2 Definition: Gitterfehler . . . . .	10
3.3.3 Definition(Konvergenz) . . . . .	10
3.4 Lemma von Gronwall . . . . .	12
3.5 Steifheit bzgl. Anfangswert . . . . .	13
3.5.1 Beispiel: Implizites Eulerverfahren . . . . .	14
3.6 4.2 Explizite Runge-Kutta Verfahren . . . . .	14
3.7 7.1.1 Konsistenz . . . . .	16
3.8 Definition 7.6: . . . . .	16
3.9 7.8 Lemma . . . . .	16
3.10 7.1.2 Stabilität . . . . .	18
3.10.1 Beispiel: 7.1.1: . . . . .	18
3.11 7.12 Definition . . . . .	18
3.12 Satz . . . . .	19
3.13 Definition: 7.14 . . . . .	20

3.14 Satz 7.15	20
3.15 Satz 7.23	21
3.16 Satz 7.16:	21
3.17 Konstruktion stabiler Mehrschrittverfahren	21
3.17.1 Adams-Verfahren	21
3.17.2 Adams-Moulton-Verfahren der Ordnung K	21
3.17.3 Adams-Basforth Verfahren	22
<b>4 §8 Numerische Lösung von Randwertproblemen</b>	<b>22</b>
4.1 (8.1) Randwertprobleme für lineare gewöhnliche Dgl. zweiter Ordnung	22
4.1.1 Randwertprobleme	22
4.2 8.1: Existenz und Eindeutigkeit von Lösungen	22
4.2.1 lineare DGL mit konstanten Koeffizienten	25
4.3 FREDHOLM-ALTERNATIVE	25
4.3.1 Adjungierten Integraloperator	25
4.4 8.2 Numerische Lösung von AWP mit finiten Differenzen	26
4.5 Diskretisierung der Randbedingung	27
4.5.1 Maximumsprinzip	29
4.6 Stabilitätsabschätzung für RWP	29
4.7 Finite Differenzen	32
4.8 M-Matrix Eigenschaft	32
4.9 Upwind-Verfahren:	33
4.10 Anfangswertproblem	33
4.10.1 Numerische Lösung:	33
4.11 Freie Randwertprobleme	34
<b>5 Buch Teil 1, Kapitel 8</b>	<b>35</b>
5.1 8.1 Klassische Iterationsverfahren	35
5.2 Satz 8.0:	36
5.3 Satz 8.1:	36
5.4 Iterative Verfahren:	36
5.4.1 RICHARDSON-Verfahren:	36
5.4.2 Jacobi-Verfahren:	37
5.5 Satz 8.2:	37
5.6 Lemma 8.3	38
5.7 Satz: 8.4:	38
5.8 Relaxierung von FP-Iterationen	39
5.9 Def. 8.5	39
5.10 Bsp: 8.6:	39
5.11 Lemma 8.7:	39
5.12 8.3 Verfahren der konjugierten Gradienten (CG)	40
5.12.1 Lemma 8.15:	41
5.12.2 Algorithmus 8.16 (CG-Verfahren)	42
5.12.3 Satz: 8.17:	42

5.13 8.4 Vorkonditionierung . . . . .	42
5.14 Satz 8.22: Für CG-Verfahren mit Vorkonditionierung (PCG) . . . . .	43
5.15 Lemma: 8.23: . . . . .	43
5.16 Iterative Lösung symmetrisch indefiniter Probleme . . . . .	45
5.17 Lemma: . . . . .	46
5.18 Satz: . . . . .	46
5.19 Satz: . . . . .	47
5.20 Inexaktes Uzawa-Verfahren . . . . .	47
5.21 Lemma: . . . . .	49

13.10.2009

# 1 Einführung in die gewöhnliche Differentialgleichungen

GG: in der wir eine Funktion  $x : I \rightarrow R$  suchen, in der Gleichung kommen auch Ableitungen von  $x$  vor:  
 (Partielle DGL: Gleichung, in der wir Funktionen  $u : \Omega \rightarrow R$  suchen,  $\Omega \subset R^d, d \geq 2$ , in der Gleichung kommen partielle Ableitungen nach mindestens zwei Variablen vor.)  
 System: GDGL, selbe Definition, n Gleichungen für  $x : I \rightarrow R^n$

$x' = 0 \rightarrow$  gelöst von jeder Konstanten Funktion  
 $x(t_0) = x_0$  Anfangswert

$$\left. \begin{array}{l} x' = f(x(t), t) \in (t_0, T) \\ x(t_0) = x_0 \end{array} \right\} \text{Anfangswertproblem}$$

$$x' = f(x(t), t) \in (a, b)$$

## 1.1 Beispiel 1: „stetige“ Verzinsung

Tägliche Verzinsung mit Zinssatz  $r$

$$t = 0, \text{ Vermögen } V_0, \delta t = \frac{1}{365}$$

$$V(\delta) = V_0(1 + r\delta t)$$

$$V(2 * \delta) = V(\delta)(1 + r\delta t)$$

...

$$V(t + \delta) = V(t)(1 + r\delta t)$$

$$\frac{V(t+\delta) - V(t)}{\delta t} = V(t)r$$

### 1.1.1 Taylor-Entwicklung

$$V(t + \delta t) = V(t) + V'(t)\delta t + 1/2V''(\tau)\delta t^2, \tau \in (t, t + \delta t]$$

$$V(t + \delta t) = V(t) + \sum_{j=1}^k \frac{1}{j!} V^{(j)}(\delta t)^j + \underbrace{\frac{1}{(k+1)!} V^{(k+1)}(\tau)(\delta t)^{k+1}}_{\leq \frac{(\delta t)^{k+1}}{(k+1)!} ||V^{(k+1)}||_\infty}$$

$$V'(t) + 1/2V''(\tau)\delta t = rV(t)$$

$$\delta t \rightarrow 0, V'(t) = rV(t)$$

$$\log(V)' = V'/V = r$$

$$\log(V) = rt + c$$

$$t = 0 : \log(V_0) = c, \log(V(t)) = rt + \log(V_0), V(t) = V_0 e^{rt}$$

Spezialfälle von  $x' = f(x, t)$

Separierte Dgl.:  $x' = g(x)h(t)$

$$\frac{x'}{g(x)} = h(t)$$

Suchen Stammfunktion von  $1/g$  und von  $h$

$$G'(x) = 1/g(x), H'(x) = h(t)$$

$$\frac{d}{dt} G(x(t)) = G'(x(t))x'(t) = \frac{x'(t)}{g(x(t))}$$

$$G(x(t)) = H(t) + c$$

$$x(0) = x_0 \rightarrow c = G(x_0) - H(0)$$

$$G(x(t)) = G(x_0) + H(t) - H(0)$$

$$x(t) = G^{-1}(G(x_0) + H(t) - H(0))$$

## 1.2 ‘Lineare Dgl‘ mit konstanten Koeffizienten

$$x' = Ax, A \in R^{n \times n}$$

$$x(t) = \underbrace{e^{At}}_{\in R^n} \underbrace{x_0}_{\in R^n}, e^{At} = \sum_{k=0}^{\infty} \frac{(At)^k}{k!}$$

einfacher für  $A$  diagonalisierbar

$A = BDB^{-1}$ ,  $D$  Diagonalisierbar,  $B$  reguläre Matrix

$$A^2 = BDB^{-1}BDB^{-1} = BD^2B^{-1}$$

$$A^k = BD^kB^{-1}$$

$$x(t) = B(e^{d_i^k t})B^{-1}x_0$$

$$x' = Ax, x = By, y = B^{-1}x, x' = By'$$

$$By' = ABY \rightarrow y'B^{-1}ABY \rightarrow y' = DY \rightarrow y'_k = d_k y_k$$

## 1.3 Beispiel 2: Newton'sche Teilchenmechanik

$x(t)$ : Raumkoordinaten zur Zeit  $t$

$v = x'$

Massenbeobachtung  $mv' = F(x, v, t)$

$$\left. \begin{array}{l} x' = v \\ mx'' = mv' = F(x, v, t) \end{array} \right\} \text{g. Dgl}$$

$N$  Teilchen,  $i = 1, \dots, N$ ,  $x_i(t), v_i(t)$

$$x'_i = v_i$$

$$m_{v'_i} = F_i^{ext}(x_i, v_i, t) + \sum_{j \neq i} F_{ij}(x_i, x_j)$$

$$F_{ij} = -F_{ji}$$

6N Differentialgleichungen

### 1.3.1 Gravitation

$$F_{ij} = -G \frac{m_i m_j}{|x_i - x_j|^3} (x_i - x_j) = -\nabla(G \frac{m_i m_j}{|x_i - x_j|})$$

$N = 2$ :

$$x'_1 = v_1, x'_2 = v_2$$

$$m_1 v'_2 - G \frac{m_1 m_2}{|x_1 - x_2|^3} (x_1 - x_2)$$

$$m_2 v_2 = +G \frac{m_1 m_2}{|x_1 - x_2|^3} (x_1 - x_2)$$

### 1.3.2 Impulserhaltung

$$m_1 v'_1 + m_2 v'_2 = 0$$

$$\Rightarrow m_1 v_1 + m_2 v_2 = const$$

$$\Rightarrow m_1 x'_1 + m_2 x'_2 = A$$

$$\Rightarrow m_1 x_1 + m_2 v_2 = At + B$$

### 1.3.3 Energieerhaltung

$$\frac{d}{dt} \left( -G \frac{m_1 m_2}{|x_1 - x_2|} + \frac{1}{2} m_1 v_1^2 + \frac{1}{2} m_2 v_2^2 \right) = 0$$

16.10.2009

## 1.4 Beispiel 3: Molekularodynamik

$N$  Moleküle,  $x_i \in R^2$ ,  $v_i \in R^3$ ,  $(p_i - m_i v_i \in R^3)$

$$x'_i = v_i$$

$$m'_{v_i} = -\nabla U_{ext}(x_i, t) - \sum_{j \neq i} \nabla U_{int}(x_i - x_j)$$

$$\text{Pot Energie } \sum_i U_{ext}(x_i) + 1/2 \sum_{i,j} U_{int}(x_i - x_j)$$

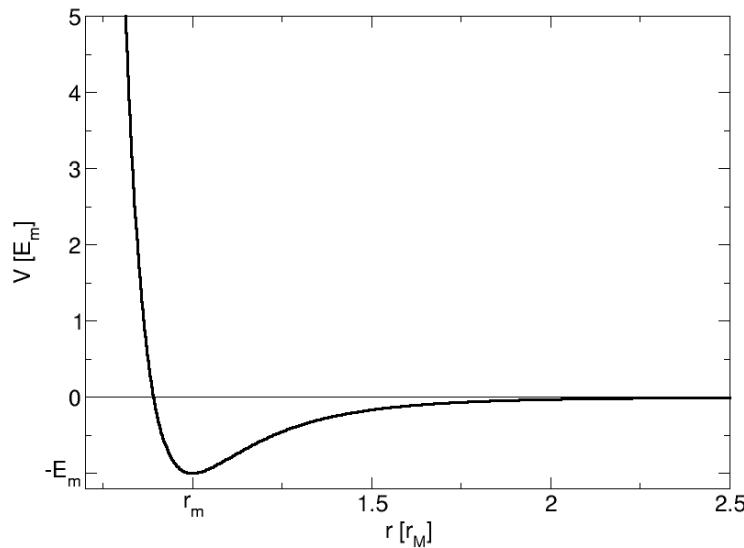


Abbildung 1: Lennard-Jones-Potential  $V$  vs. Abstand  $r$  zwischen den Atomen

Problem: extrem kleine Zeitschritte nötig bei der numerischen Lösung, um Durchdringung zu vermeiden.

## 1.5 Beispiel 4: Reaktionskinetik

Monomolekulare Reaktion: (z.B. Radioaktiver Zerfall)

Wahrscheinlichkeit einer Reaktion von  $A$  zu  $B$  in Zeitintervall  $(t, t + \delta t) \approx k\delta t$

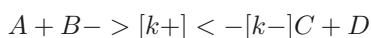
$n_A(t)$  = Anzahl von Molekülen  $A$  zum Zeitpunkt  $t$ ;  $4n_A(t + \delta t) \approx n_A(t) - k\delta t n_A(t)$

$$\frac{n_A(t + \delta t) - n_A(t)}{\delta t} = -kn_A(t)$$

$$\delta t > 0, n'_A(t) = -kn_A(t)$$

$$n_A(t) = n_A(0)e^{-kt}$$

Bimolekulare Reaktion:



## 1.6 Beispiel 4: DNA-Replikation

$$n'_A = n'_B = -k_+ \underbrace{n_A(t)n_B(t)}_{= \text{Anzahl der } (A,B)\text{-Paare}} + k_- n_C(t)n_D(t)$$

$$n'_C = n'_D = k_+ n_A(t)n_B(t) - k_- n_C(t)n_D(t)$$

$$\text{Einhaltung der Gesamtanzahl: } (n_A + n_B + n_C + n_D)' = 0$$

## 1.7 Beispiel 5: Populationsdynamik

$$\begin{aligned} n(t) &= \text{Anzahl der Individuen zur Zeit } t \\ n'(t) &= gn(t) - sn(t) = (g - s)n(t) \\ n'(t) &= (g - s)n(t) + E(t) - A(t) \end{aligned}$$

## 1.8 Beispiel 6: Räuber-Beute Modell(Lotke-Voltena Modell)

$$\begin{aligned} R(t) &= \text{erw. Anzahl an Räubern} \\ B(t) &= \text{erw. Anzahl an Beutetieren} \\ B'(t) &= aB(t) - bB(t)R(t) \\ R'(t) &= -dR(t) + cB(t)R(t) \\ B(0) &= B_0, R(0) = R_0 \end{aligned}$$

## 1.9 Beispiel 7: Minimalkurven

Was ist die kürzeste Verbindung zwischen zwei Punkte?

$$E(y) = \int_a^b \sqrt{1+y'(x)^2} dx + \int_a^b F * y dx$$

Suche das Minimum über alle Funktionen  $y$  mit  $y(a) = c$

$$y(b) = d$$

Störungen  $\phi(x)$  mit  $\phi(a) = \phi(b) = 0$

$y + \epsilon\phi$  ist ebenfalls Kurve zwischen  $(0, c)$  und  $(b, d)$ .

Für Minimum  $y$  gilt:  $E(y) \leq E(y + \epsilon\phi), \forall \epsilon > 0$

$$f(\epsilon) := E(y + \epsilon\phi) \text{ hat Minimum bei } \epsilon = 0.$$

$$\Rightarrow f'(0) = 0$$

$$f(\epsilon) = \int \sqrt{1+(y'+\epsilon\phi')^2} dx + \int F(y+\epsilon\phi) dx$$

$$f'(\epsilon) = \int_a^b \frac{2(y'+\epsilon\phi')\phi'}{1+\sqrt{1+(y'+\epsilon\phi')^2}} dx + \int_a^b F\phi dx$$

$$f'(0) = \int_a^b \left( \frac{y'\phi'}{\sqrt{1+y'^2}} F\phi \right) dx$$

$$= \frac{y'}{\sqrt{1+y'^2}} \phi|_a^b - \int_a^b \phi \left( \left( \frac{y'}{\sqrt{1-y'^2}} \right)' + F \right) dx$$

$$f'(0) = 0 \Rightarrow \int_a^b \phi(..) dx = 0, \forall \phi \text{ mit } \phi(a) = 0, \phi(b) = 0$$

$$\Rightarrow \left( \frac{y'}{\sqrt{1+y'^2}} \right)' = F, \text{ für } x \in (a, b)$$

$$y(a) = c, y(b) = d$$

$$F = 0, \frac{y'}{\sqrt{1+y'^2}} = const$$

$$y'^2 = c^2(1+y'^2) \Rightarrow y'^2 = \frac{c^2}{1-c^2}$$

$$y' = \frac{c}{\sqrt{1-c^2}} = const \Rightarrow y \text{ linear}$$

$$\left( \frac{y'}{\sqrt{1+y'^2}} \right)' = \frac{y''(1+y'^2)-y'^2y''}{\sqrt{1+y'^2}^3} = \frac{y''}{\sqrt{1+y'^2}^3} \Rightarrow y'' = 0$$

## 2 Existenz und Eindeutigkeit des AWP

$$x'(t) = f(x(t), t), x : [t_0, T] \rightarrow \mathbb{R}^n, f : \mathbb{R}^n \times [t_0, T] \rightarrow \mathbb{R}^n$$

$$x(t_0) = x_0$$

$$\text{Suchen } x \in C^1((t_0, T))^n$$

### 2.1 Lokale Existenz

Es existiert  $T > t_0$ , sodass AWP eine Lösung  $x \in C^1((t_0, T))^n$

Globale Existenz: Für alle  $T > t_0$  hat AWP eine Lösung  $x \in C^1((t_0, T))^n$

Beispiel:  $x' = x - > x(t) = x_0 e^t$ , globaler Existenz

Beispiel:  $x' = -e^{-x} \Rightarrow e^x x' = -1$

integriert:  $e^x = -t + c, t = 0 : e^{x_0} = c$

$$x = \log(e^{x_0} - t)$$

Nur Existenz für  $t < e^{x_0}$

$$t_{max} = e^{x_0}$$

Minimalanforderung für Existenz:  $f$  stetig

## 2.2 Satz: Satz von Peano

**Satz von Peano:**  $f$  stetig  $\Rightarrow$  lokale Existenz

---

20.30.2009

$$x'(t) = f(x(t), t), t \in (t_0, T)$$

$$x(t_0) = x_0$$

$$f \in C(R^n \times [t_0, T]), x_0 \in R^n$$

Suche Lösung:  $x \in C((t_0, T))^n$

$$x(t) - x_0 = \int_{t_0}^T f(x(s), s) ds$$

$$x(t) = x_0 + \int_{t_0}^T f(x(s), s) ds$$

$x = F(x)$ , Fixpunktgleichung

$$f : C((t_0, T))^n \rightarrow C((t_0, T))^n$$

$$x - > x_0 + \int_{t_0}^t f(x(s), s) ds$$

$\Rightarrow$  wohldefiniert

## 2.3 Schaudersche Fixpunktsatz

$F : X \rightarrow X$ , sei stetig ( $x_1 \rightarrow x \Rightarrow F(x_1) \rightarrow F(x)$ ) und kompakt ( $M$  beschränkt  $\Rightarrow F(M)$  präkompakt  $\Rightarrow \overline{F(M)}$  kompakt)  $X$  Banachraum

$D$  abgeschlossen und beschränkt mit  $F(D) \subset D \Rightarrow$  Dann existiert ein Fixpunkt von  $F \in D$ , d.h.  $\bar{x} \in D$  mit  $\bar{x} = F(\bar{x})$

### 2.3.1 Anwendung auf DGL

$$D = B_R(x_0) = x \in C((t_0, T))^n \quad |||x - x_0|||_\infty \leq R$$

$$||x||_\infty = \sup_{t \in [t_0, T]} ||x(t)||, \text{ (euklidische Norm, da } x(t) \in R^n)$$

$$\begin{aligned} ||F(x) - x_0||_\infty &= \sup_{t \in [t_0, T]} \left| \left| \int_{t_0}^t f(x(s), s) ds \right| \right| \leq \sup_t ((t - t_0) \sup_{s \in (t_0, t)} ||f(x(s), s)||) \leq \sum_t ((T - t_0) \sup_{s \in [t_0, T]} ||f(x(s), s)||) = \\ &(T - t_0) \sup_{s \in [t_0, T]} ||f(x(s), s)|| \end{aligned}$$

$$||x - x_0|| \leq R$$

$$\sup_s ||x(s) - x_0|| \leq R$$

$$\forall s : ||x(s) - x_0|| \leq R$$

$$x(s) \in B_R^{R^n}(x_0)$$

$$\leq (T - t_0) \sup_{s \in [t_0, T], z} ||f(z, s)||$$

$$z \in B_R^{R^n}(x_0)$$

$f$  ist stetig,  $B_R^{R^n}(x_0) \times [t_0, T]$  ist Kompakt in  $R^{n+1} \Rightarrow f$  ist auf dieser Menge beschränkt,  $\sup_{s, z} ||f(z, s)||$

$$\leq C$$

$$\Rightarrow ||F(x) - x_0|| \leq (T - t_0)C(R, T) \leq R$$

$$\lim_{T \rightarrow \infty} (T - t_0)C(R, T) = 0$$

$C(R, T) = \sup_{s \in [t_0, T], z \in B_R^{R^n}(x_0)} \|f(z, s)\| \rightarrow \sup_{z \in B_R^{R^n}(x_0)} \|f(z, t_0)\|$  beschränkt  
 $\Rightarrow \exists T > t_0 \text{ mit } (T - t_0)C(R, T) \leq R$   
 $\Rightarrow \|F(x) - x_0\|_\infty \leq R \Rightarrow F(x) \in B_R^{R^n}(x_0)$

**Kompaktheit:**  $x_n \in C([t_0, T])^n$  beschränkt  
 $\Rightarrow F(x_n) \leq C([t_0, T])^n$  hat konvergente Teilfolge

## 2.4 Satz von Auzela Ascoli

Die Funktionenfolge  $x_n$  hat auf einem Intervall  $[t_0, T]$  eine konvergente Teilfolge in der  $\infty$ -Norm ( $\Rightarrow$  gleichmäßige Konvergenz), wenn sie gleichgradig stetig ist.  $\forall x_n \forall \delta > 0 \exists \epsilon > 0 \forall x_n : |t - s| < \epsilon \Rightarrow |x_0(t) - x_n(s)| < \delta$  = gleichmäßig stetig

$x_n$  beschränkt in  $C([t_0, T])^n$

$$\|x_n\|_\infty \leq C$$

$$t, s \in (t_0, T]$$

$$\begin{aligned} \|F(x_n)(s) - F(x_n)(t)\| &= \left\| \int_{t_0}^s f(x_n(\tau), \tau) d\tau - \int_{t_0}^t f(x_n(\tau), \tau) d\tau \right\| = \left\| \int_t^s f(x_n(\tau), \tau) d\tau \right\| \leq |t - s| \sum_{\tau \in [t, s]} \|f(x_n(\tau), \tau)\| \leq \\ &\leq |t - s| \sup_{\tau \in [t_0, T]} \|f(\underbrace{x_n(\tau)}_{\in B_c^{R^n}(0)}, \tau)\| \end{aligned}$$

$$\leq |t - s| \sup_{\tau \in [t_0, T], z \in B_c^{R^n}(0)} \|f(z, \tau)\| \leq \tilde{C}|t - s|$$

$\tilde{C}$  unabhängig von  $n$ .

$$\epsilon(\delta) = \frac{\delta}{\tilde{C}}$$
 unabhängig von  $n(x_n)$

**Satz von Anzelo-Ascoli**  $\Rightarrow F(x_n)$  hat konvergente Teilfolge in  $C([t_0, T])^n$

$\Rightarrow F$  ist kompakt,  $F(B_R(x_0)) \subset B_R(x_0)$ , falls  $T - t_0$  hinreichend klein ist.

**Schouders Fixpunktsatz**  $\Rightarrow x = F(x)$

$x(t) = x_0 + \int_{t_0}^t f(x(s), s) ds$ , hat einen Fixpunkt in  $B_R(x_0)$ , falls  $T$  hinreichend klein

Fixpunkt  $\bar{x}$

$$\bar{x}(t) = x_0 + \int_{t_0}^t f(\bar{x}(s), s) ds \in C^1([t_0, T])^n$$

$$\bar{x}'(t) = f(\bar{x}(t), t)$$

$$\bar{x}(t_0) = x_0$$

$\Rightarrow \bar{x}$  löst das AWP

**Satz von Peano:**  $f$  stetig,  $x_0 \in R^n$ . Dann existiert  $T > t_0$ , sodass das AWP eine Lösung in  $(t_0, T)$  hat.

Beispiel:  $x' = x^\alpha$ ,  $t > 0$  und  $x(0) = 0$

$0 < \alpha < 1$ ,  $x^\alpha$  stetig, aber keine eindeutige Lösung.

## 2.5 Satz von Picard-Lindelöf (globale Version)

$f$  stetig,  $x_0 \in R^n$ ,  $f$  Lipschitzstetig im ersten Argument, d.h.  $\forall x, y \in R^n \forall t > t_0 : \|f(x, t) - f(y, t)\| \leq L\|x - y\|$   
Dann existiert für alle  $T > t_0$  eine eindeutige Lösung des AWP.

Beweis über Banach'schen Fixpunktsatz

Banach'scher Fixpunktsatz:  $F : X \rightarrow X$ ,  $D$  beschränkt, abgeschlossen, konvex mit  $F(D) \subset D$  und  $F$  sei kontraktiv, d.h.  $\exists \gamma < 1 \forall x, y \in X \|F(x) - F(y)\| \leq \gamma \|x - y\|$

**Lipschitz-stetigkeit von  $F$**

$$\begin{aligned} \|F(x) - F(y)\|_\infty &= \sup_{t \in [t_0, T]} \left\| \int_{t_0}^t (f(x(s), s) - f(y(s), s)) ds \right\| \leq \sup_t ((t - t_0) \sup_{s \in [t_0, t]} \|f(x(s), s) - f(y(s), s)\|) \leq \\ &\leq (T - t_0) \sup_{s \in [t_0, T]} \|f(x(s), s) - f(y(s), s)\| \leq [\text{Lipschitzstetig}] (T - t_0) \sup_s L \|x(s) - y(s)\| \leq \underbrace{(T - t_0)L}_{\gamma} \|x - y\|_\infty \end{aligned}$$

Wir brauchen  $\gamma < 1$ , OK für  $T < t_0 + 1/L$

23.10.2009

$$x' = f(x, t) \text{ in } (t_0, T)$$

$$x(0) = x_0$$

$$F(x)(t) = x_0 + \int_0^t f(x(s), s) ds$$

$x' = Lx \Rightarrow x = x_0 e^{Lt}$ ,  $|L|$  ist Lipschitz-Konstante von  $f = Lx$

$e^{-Lt}x$  ist beschränkt

$$\|x\|_c = \sup_{t \in [t_0, T]} |e^{-ct}x(t)|$$

$\|\cdot\|_c$  ist äquivalent zur  $\infty$ -Norm

$$x^{-ct_0}\|x(t)\| \geq \|e^{-ct}x(t)\| \geq e^{-cT}\|x(t)\|, c \geq 0$$

$$\Rightarrow e^{-ct_0}\|x\|_\infty \geq \|x\|_c \geq e^{-cT}\|x\|_\infty$$

$C([t_0, T])$  mit  $\|\cdot\|_c$  ist Banachraum

$$\begin{aligned} \|F(x_1) - F(x_2)\|_c &= \sup_{t \in [t_0, T]} (e^{-ct} \left\| \int_{t_0}^t (f(x_1(s), s) - f(x_2(s), s)) ds \right\|) \\ &\leq \sup_t (e^{-ct} \int_{t_0}^t \|f(x_1(s), s) - f(x_2(s), s)\|) \\ &\leq \sup_t (e^{-ct} \int_{t_0}^t L \|x_1(s) - x_2(s)\| ds) \\ &= \sup_t (e^{-ct} \int_{t_0}^t L e^{cs} \underbrace{(e^{-c\tau} \|x_1(\tau) - x_2(\tau)\|)}_{\leq \sup_{\tau \in [t_0, T]} (e^{-c\tau} \|x_1(\tau) - x_2(\tau)\|) = \|x_1 - x_2\|_c}) \\ &\leq \sup_t (e^{-ct} \int_{t_0}^t L e^{cs} \|x_1 - x_2\|_c ds) \\ &= \sup_{t(e^{-ct} L 1/c(e^{ct} - e^{-ct_0})) \|x_1 - x_2\|_c} \\ &= L/c \|x_1 - x_2\|_c \sup_t (1 - e^{c(t_0 - t)}) \\ &= L/c \|x_1 - x_2\|_c (1 - e^{c(t_0 - T)}) \end{aligned}$$

$$\gamma = L/c(1 - e^{c(t_0 - T)})$$

< 1 für  $c \geq L$

$$\gamma = 1 - e^{L(t_0 - T)} < 1 \text{ für } c = L$$

$\Rightarrow F$  ist kontraktiv

$F$  bildet  $B_R(x_0) = \{x \in C([t_0, T])^n : \|x - y_0\|_c \in R\}$  in sich selbst ab

$$\|F(x) - x_0\|_c = \sup_{t \in [t_0, T]} e^{-ct} \left\| \int_{t_0}^t f(x(s), s) ds \right\|$$

$$= \sup_t (e^{-ct} \left\| \int_{t_0}^t (f(x(s), s) - f(x_0, s)) ds + \int_{t_0}^t f(x_0, s) ds \right\|)$$

$$\leq \sup_t (e^{-ct} (\int_{t_0}^t \|f(x(s), s) - f(x_0, s)\| ds + \int_{t_0}^t \|f(x_0, s)\| ds))$$

$$\leq \sup_t e^{-ct} \int_{t_0}^t \|f(x(s), s) - f(x_0, s)\| ds + e^{-ct_0} \underbrace{\int_{t_0}^T \|f(x_0, s)\| ds}_{=: A}$$

$$\leq \gamma \|x - x_0\|_c + A$$

$$\leq \gamma R + A \leq R$$

erfüllt für  $R \geq \frac{A}{1-\gamma}$

Voraussetzung des Banach'schen Fixpunktsatzes erfüllt auf  $B_R(x_0) \Rightarrow F$  hat einen eindeutigen Fixpunkt  $\bar{x}$  und  $\bar{x} \in B_R(x_0)$

$$\Rightarrow \bar{x}' = f(\bar{x}, t), t \in [t_0, T]$$

$$\bar{x}(t_0) = x_0 \bar{x} \in C^1([t_0, T])$$

Was passiert, wenn  $f$  nun lokal Lipschitz-stetig ist?  $f$  stetig diffbar bzgl.  $x$

$$L_M = \sup_{||\partial_x f(x, \dots)||}$$

lokale Lipschitz-Konstante in  $M$

$$f(x_1), f(x_2) = \partial_x f(\xi)(x_1 - x_2), \xi \in [x_1, x_2]$$

$$\|f(x_1) - f(x_2)\| \leq \|\partial_x f(\xi)\| \|x_1 - x_2\| \leq L_M \|x_1 - x_2\| \text{ für alle } x_1, x_2 \in M$$

$\partial_x f$  stetig,  $f$  stetig

$\Rightarrow$  (i) es gibt  $T > t_0$ , sodass die Lösung in  $[t_0, T]$  existiert.

(ii) Es gibt höchstens eine Lösung.

Beweis von (ii) Angenommen,  $x_1$  und  $x_2$  sind Lösungen mit  $x_1 \neq x_2, x_1(t_0) = x_0 = x_2(t_0)$

$x_1 \neq x_2 \Rightarrow \exists \tau \geq t_0$ , so dass  $x_1(t) \neq x_2(t)$  in  $(\tau, \tau + \epsilon)$  und  $x_1(t) = x_2(t)$  in  $[t_0, \tau]$

$$t > \tau : x_1(t) = x_1(\tau) + \int_{\tau}^t f(x_1(s), s) ds$$

$$x_2(t) = x_2(\tau) + \int_{\tau}^t f(x_2(s), s) ds$$

$$\|x_1(t) - x_2(t)\| = \left\| \int_{\tau}^t f(x_1(s), s) - f(x_2(s), s) ds \right\|$$

$$\leq \int_{\tau}^t \|f(x_1(s), s) - f(x_2(s), s)\| ds$$

$$= \int_{\tau}^t \|\partial_x f(\xi(s), s)(x_1(s) - x_2(s))\| ds$$

$$\|x_1(t) - x_2(t)\| \leq \int_{\tau}^t \|\partial_x f(\xi(s), s)\| - \sup_{\sigma \in (\tau, t)} \|x_1(\sigma) - x_2(\sigma)\| ds$$

$$= \int_{\tau}^{\tau+\epsilon} \|\partial_x f(\xi(s), s)\| ds \sup_{\sigma \in [\tau, \tau+\epsilon]} \|x_1(\sigma), x_2(\sigma)\|$$

$$\|x_1 - x_2\|_{\infty} = \underbrace{\int_{\tau}^{\tau+\epsilon} \|\partial_x f(\xi(s), s)\| ds}_{\rightarrow 0 \text{ für } \epsilon \rightarrow 0, < 1 \text{ für } \epsilon \text{ hinreichend klein}} \|x_1 - x_2\|_{\infty}$$

$\epsilon$  hinreichend klein  $\Rightarrow \|x_1 - x_2\|_{\infty} < \eta \|x_1 - x_2\|_{\infty} \Rightarrow x_1 = x_2$  auf  $(\tau, \tau + \epsilon)$  Widerspruch, da  $\eta < 1$

$\Rightarrow$  es gibt kein solches  $\tau$

$\Rightarrow x_1 = x_2$  für alle  $t$ .

**Beweis:** (i)  $M(R, T) := \sup_{z \in B_R^{R^n}(x_0), t \in [t_0, T]} \|\partial_x f(z, t)\|$

$$F(x)(t) = x_0 + \int_{t_0}^t f(x(s), s) ds$$

$$\|F(x_1) - F(x_2)\| = \sup_t \left\| \int_{t_0}^t (f(x_1(s), s) - f(x_2(s), s)) ds \right\| \leq \sup_t \int_{t_0}^t \|\partial_x f(\xi(s), s)(x_1(s) - x_2(s))\| ds$$

$$\leq \sup_t \int_{t_0}^t \underbrace{\|\partial_x f(s, s)\|}_{\leq \sup_{z \in B_R} \|\partial_x f(z, s)\| \leq M(R, T)} ds \|x_1 - x_2\|_{\infty}$$

$x_1, x_2 \in B_R(x_0)$

$$\xi(s) \in [x_1, x_2] \Rightarrow \xi(s) \in B_R^{R^n}(x_0)$$

$$\leq \sup_t \int_{t_0}^t M(R, T) ds \|x_1 - x_2\|_{\infty}$$

$$\leq (T - t_0) M(R, T) * \|x_1 - x_2\|_{\infty}$$

### 3 Stabilität von DGL

Verschiedene Begriffe

#### 3.1 Kondition

Änderung von  $x$  bezüglich der Eingabe

$$-x \in R^n$$

$$-f : R^n \times [t_0, T] \rightarrow R^n$$

vgl.  $\|x - \tilde{x}\|_{zul}((x_0, f) - (\tilde{x}_0, \tilde{f}))$ ?

$$x' = f(t)$$

$$x(t) = x_0 + \int_{t_0}^t f(s) ds$$

$$x(t) - \tilde{x}_0 + \int_{t_0}^t (f(s) - \tilde{f}(s)) ds$$

$\infty$ -Norm für  $x$  und  $f$

$$\|x - \tilde{x}\|_{\infty} \leq \|x_0 - \tilde{x}_0\| + (T - t_0) \|f - \tilde{f}\|_{\infty}$$

$\infty$ -Norm für  $x, x'$ ,  $\infty$ -Norm

$$\|x\|_{c_1} = \max\|x\|_{\infty}, \|x'\|_{\infty}$$

$$\|x' - \tilde{x}'\|_{\infty} = \|f - \tilde{f}\|_{\infty}$$

Fehler in der Krümmung ( $x''$ ),  $\|x\|_{C^2} = \max\|x\|_{\infty}, \|x'\|_{\infty}, \|x''\|_{\infty}$

nur sinnvoll für  $f \in C^1$

$$x'' = f'$$

$$\|x'' - \tilde{x}''\|_{\infty} \leq \|f' - \tilde{f}'\|_{\infty} \rightarrow \text{Stabilität } f \text{ gemessen in } C^1$$

Stabilität auch oft Verhalten von  $x$  für  $t \rightarrow \infty$  Stabilität von Fixpunkten,  $x' = f(x), f(\bar{x}) = 0, \bar{x} \in R^n$

$x$  stabil, wenn  $x(t) \rightarrow \bar{x}$  für  $x_0$  hinreichend nahe bei  $\bar{x}$

27.10.3009

---

### 3.2 Einschrittverfahren für gewöhnliche DGL

**Motivation:**  $x' = f(t, x), x(t_0) = x_0$ , AWA

**Problem:** Wir können die Lösung nicht analytisch berechnen

⇒ **Lösung:** approximiere die Lösung  $x \in C^1([0, T], R^d)$  mit Hilfe von numerischen Verfahren.

#### 3.2.1 Eulersches Polygonzugverfahren (1726)

1. Schritt: Lege durch  $P_1 = (t_0, x_0)$  eine Tangente mit der Steigung  $m_0 = f(t_0, x_0)$
2. Schritt: Gehe entlang dieser Tangente  $t_1 = t_0 + \tau_0, x_1 = x_0 + f(t_0, x_0)$
3. Schritt: Berechne Tangente im Punkt  $(t_1, x_1) m_1 := f(t_1, x_1)$

Allgemeine Methode bei **Einschrittverfahren (ESV)**

Zerlege das Zeitintervall  $[z_0, T]$  in  $n + 1$  gewählte Zeitschritte

hier nicht äquidistante Schrittweiten:  $t_{j+1} = t_j + \tau_j, j = 0, \dots, n_\delta - 1$

Gitter  $\delta = t_0, \dots, t_n, t_j = \text{Gitterpunkte}$

**Definiere:** (Feinheit des Gitters)

$$\tau_\delta = \max_{0 \leq j \leq n_\delta - 1} |\tau_j|$$

Es numerisches Verfahren findet so ein  $\Delta$  und ordnet diesem Gitter eine Gitterfunktion  $y_D : \Delta \rightarrow R^d$  zu. Es soll gelten:  $x_\Delta \approx x(t) \forall t \in \Delta$

Wir bestimmen  $x_D$  rekursiv: (2-Term-Rekursion)

- (i)  $x_D(t_0) = x_0$
- (ii)  $x_D(t_{j+1}) = x(t_j) + \tau_j * \psi(t_j, x_\Delta(t_j), \tau_j)$

mit  $\psi(t, x, \tau)$  Verfahrensfunktion/Schrittfunktion

explizites Eulerverfahren:  $\psi(t, x, \tau) = f(t, x)$

- (i)  $x_\Delta(t_0) = x_0$
- (ii)  $x_\Delta(t_{j+1}) = x_\Delta(t_j) - \tau_j f(t_j, x_\Delta(t_j))$

**Beispiel:**  $x' = 1 + x^2 = f(t, x) \Rightarrow x(t) = \tan(x)$

$$x(0) = 0$$

$$I = [0; 0.5]$$

$$S = 0.1$$

$j$	$t_j$	$x_D(t_j)$	$x(t_j)$
0	0	0	0
1	0.1	0.1	0.1003
2	0.2	0.201	0.2027
3	0.3	0.3050	0.3093
4	0.4	0.4193	0.4228
5	0.5	0.5315	0.5463

**Idee von ESV:**  $x(t_{j+1})$  berechnet man mit Hilfe von  $x(t_j)$

**Idee von MSV:**  $x(t_{j+1})$  berechnet man mit Hilfe von  $x(t_j), x(t_{j-1}, \dots)$

### 3.3 Konsistenz

Sei  $x$  die Lösung des AWP und  $\Psi$  Verfahrensfunktion, dann heißt  $\tau(t_i, x(t_i), \tau_i) = \frac{x(t_{i+1}) - x(t_i)}{\tau_i} - \psi(t_i, x(t_i))$ , Verfahrensfehler/Konsistenzfehler an der Stelle  $(t_i, x(t_i))$   
Das KF/VF gib an, wie gut die exakte Lösung das Einschrittverfahren erfüllt.

#### 3.3.1 Definition: Konsistenz

Ein Verfahren heißt konsistent, wenn  $r(t, x(t), \tau) \rightarrow 0$  für  $\tau \rightarrow 0$

Ein Verfahren heißt konsistent von der Ordnung  $p$ , wenn  $\tau(t_i, x(t_i), \tau_i) = O(\tau_i^p)$

Berechnung der Konsistenzordnung mit Hilfe Taylor-Entwicklung.

$$(*) \quad \tau(t_i, x(t_i), \tau_i) = \frac{x(t_{i+1}) - x(t_i)}{\tau_i} - \psi(t_i, x(t_i), \tau_i)$$

$$\text{Taylorentwicklung von } x(t_{i+1}) = x(t_i + \tau_i) = x(t_i) + \frac{x'(t_i)}{1!} \tau_i + \frac{x''(t_i)}{2!} \tau_i^2 + O(\tau^3)$$

$$\text{Einsetzen in } (*) : x(t_i) \tau_i / \tau_i + \frac{x''(t_i) * \tau_i^2}{2\tau_i} + \dots - \psi(t_i, x(t_i), \tau_i)$$

$$x' = f(x, t(x)), x'' = f_t(t, x) + f_x(t, x) * f$$

$$f(t_i, x(t_i)) + 1/2(f_t + f_x f(t_i, x(t_i))) \tau_i + y(t_i, x(t_i), \tau_i)$$

$$\text{für das explizite Eulerverfahren. } \Psi(t_i, x(t_i), \tau_i) = f(t_i, x(t_i), \tau_i)$$

$$\Rightarrow \tau(t_i, x(t_i), \tau_i) = \tau/2(f_t + f_x f(t_i, x(t_i)) + \dots) = \sigma(\tau^1)$$

$\Rightarrow p = 1 \Rightarrow$  das explizite Eulerverfahren ist Konsistent von der Ordnung 1.

Konstruktion von Verfahren hoher Ordnung

**Idee:** verwende die Taylorentwicklung 2. Ordnung:

$$\frac{x(t_{i+1}) - x(t_i)}{\tau_i} = f(t_i, x(t_i)) + 1/2(f_t + f_x f)(t_i, x(t_i), \tau_i) - \tau_i$$

$$\Rightarrow x(t_{i+1}) = x(t_i) + \tau_i f(t_i, x(t_i)) + \tau_i^2 / 2(f_t + f_x f)(t_i, x(t_i))$$

Verfahren der Konsistenzordnung 2

**Beispiel:**  $x' = x = f(t, x)$

$$x(t_{i+1}) = x(t_i) + \tau x(t_i) + \tau^2 / 2(x'(t) + 1x(t)) = x(t_i) + \tau_i(x(t_i)) + \tau_i^2(x(t_i))$$

$$\Psi(t, x(t)) = f(x, t) + \tau/2(f_t + f_x f)(t, x)$$

Ansatz für ein Verfahren 2.Ordnung

Konvergenz

Fehler von  $x$  zu  $x_\Delta$  (Rückwärtsanalyse)

#### 3.3.2 Definition: Gitterfehler

Den Vektor der Approximationsfehler auf dem Gitter  $\Delta$  bezeichnen wir mit  $\epsilon_\Delta : \Delta \rightarrow R^d$  mit  $\epsilon_\Delta^{(t)} = x(t) - x_\Delta(t) \forall t \in \Delta$

Diskretisierungsfehler auf  $\|\epsilon_\Delta\|_\infty = \max_{t \in \Delta} |\epsilon_\Delta(t)|$

#### 3.3.3 Definition(Konvergenz)

Zu jedem Gitter  $\Delta$  auf  $[t_0, T]$  mit  $\tau_\Delta$  hinreichend klein. sei eine Gitterfunktion  $x_\Delta$  gegeben.

Die Familie dieser Gitterfunktionen konvergiert gegen die Lösung  $x \in C^1([t_0, T], R^d)$ , wenn für den Diskretisierungsfehler gilt,  $\|\epsilon_\Delta\|_\infty \rightarrow 0$  für  $\tau_\Delta \rightarrow 0$

Sie konvergiert von der Ordnung  $p$ , wenn gilt  $\|\epsilon_\Delta\|_\infty = O(\tau_\Delta^p)$  für  $\tau_\Delta \rightarrow 0$

**Konvergenzsatz für das ESV** Die Verfahrensfunktion  $\psi(t, x, \tau)$  sei bzgl.  $x$  Lipschitzstetig auf dem Gebiet  $\Omega$ . Dann ist das ESV- konsistent (mit der Ordnung  $P$ ), so ist das ESV konvergent(mit der Ordnung  $P$ )

30.30.2009

$$\begin{aligned} x' &= f(x, t) \\ x(t_0) &= x_0 \end{aligned}$$

$$\begin{aligned}
 x_{\Delta}(t_{j+1}) &= x_{\Delta}(t_j) + \tau_j \psi(\tau_j, x_{\Delta}, \tau_j) \\
 \Delta &= t_0, \dots, t_n = T \\
 \tau_j &= t_{j+1} - t_j \\
 x(t_{j+1}) &= x(t_j) + \tau_j 1/\tau_j (\int_{t_j}^{t_{j+1}} f(x(t), t) dt) \\
 \text{Euler, } \tau_j &\text{ klein} \\
 1/\tau_j \int_{\tau_j}^{\tau_{j+1}} f(x(t), t) dt &\approx f(x_0(t), t)
 \end{aligned}$$

lokaler Konsistenzfehler

$$\begin{aligned}
 r(t, x, \tau) &= \frac{x(t_{i+1}) - x(t_i)}{\tau_i} - \phi(t_i, x(t_i), \tau_i) \\
 x' &= f(x(t), t), x(t_{i+1}))x(t_i) + x'(t_i)\tau_i + 1/2x''(t_i)\tau_i^2 + 1/6x'''(t_i)\tau_i^3 + \dots \\
 &= x(t_i) + f(x(t_i), t_i)\tau_i + 1/2\frac{\partial f}{\partial x}(x(t_i), t_i)f(x(t_i), t_i)\tau_i^2 + 1/2\frac{\partial f}{\partial t}(x(t_i), t_i)\tau_i^2 + \dots \\
 x'' &= d/dt(f(x(t), t)) = \frac{\partial f}{\partial x}x' + \frac{\partial f}{\partial t} = \frac{\partial f}{\partial x}f + \frac{\partial f}{\partial t}
 \end{aligned}$$

Konvergenz:

$$\epsilon_{\Delta} := x(t) - x_{\Delta}(t), t \in \Delta$$

Diskretisierungsfehler:

$$\|\epsilon_{\Delta}\|_{\infty} = \max_{t \in \Delta} |\epsilon_{\Delta}(t)|$$

Konvergenz:  $\|\epsilon_{\Delta}\|_{\infty} \rightarrow 0$  für  $\tau_{\Delta} \rightarrow 0$

Konvergenzordnung:  $\|\epsilon_{\Delta}\|_{\infty} = O(\tau_{\Delta}^P)$

Konvergenz = Konsistenz + Stabilität

Prinzip für lineare Gleichungen:  $A \in R^{n \times n}, A^{-1}$  existiert

$$Ax = b, \tilde{A}\tilde{x} = \tilde{b}$$

$$\text{Konsistenzfehler } r = \tilde{A}x - \tilde{b}$$

$$\text{Konvergenz} = \epsilon = x - \tilde{x}$$

$$\tilde{A}(x - \tilde{x}) = \tilde{A}x - \tilde{b} = r$$

$$\begin{aligned}
 \Rightarrow x - \tilde{x} &= \tilde{A}^{-1} * r \Rightarrow \underbrace{\|x - \tilde{x}\|}_{\text{Konvergenzfehler}} \leq \underbrace{\|\tilde{A}^{-1}\|}_{\text{Stabilitätskonstante}} * \underbrace{\|r\|}_{\text{Konsistenzfehler}} \\
 &= \frac{x(t_{j+1}) - x_{\Delta}(t_{j+1})}{\tau_j} - \frac{x(t_j) - x_{\Delta}(t_j)}{\tau_j} - (\psi(t_j, x(t_j), \tau_j) - \psi(t_j, x_{\Delta}(t_j), \tau_j)) \\
 &= \frac{x(t_j) - x(t_j)}{\tau_j} - \psi(t_j, x(t_j), \tau_j) \\
 &= r(t_j, x, \tau_j)
 \end{aligned}$$

Brauchen noch diskrete Stabilitätsabschätzung

$$(x_{j+1} - \tilde{x}_{j+1}) - (x_j - \tilde{x}_j) - (\psi_j(x_j) - \psi_j(\tilde{x})) = r_j$$

$$\max_j \|x_j - \tilde{x}_j\| \leq C \max_k \|r_k\|$$

Entsprechendes Problem für Dgl:

$$x' - \tilde{x}' - (f(x, t) - f(\tilde{x}, t)) = r$$

$$\text{möchten: } \|x - \tilde{x}\|_{\infty} \leq C \|r\|_{\infty}$$

$f$  Lipschitz nach  $x$

$$\begin{aligned}
 x(t) - \tilde{x}(t) &= \int_{t_0}^t (f(x(s), s) - f(\tilde{x}(s), s)) ds + \int_{t_0}^t r(s) ds \\
 \|x(t) - \tilde{x}(t)\| &= \left\| \int .. + \int .. \right\| \leq \left\| \int .. \right\| + \left\| \int .. \right\| \leq \int_{t_0}^t \|f(x(s), s) - f(\tilde{x}(s), s)\| ds + \int_{t_0}^t \|r(s)\| ds \\
 &\leq \int_{t_0}^t L \|x(s) - \tilde{x}(s)\| ds + \int_{t_0}^t \underbrace{\|r(s)\|}_{\leq B = \max_{s \in [t_0, T]} \|r(s)\|} ds \\
 e(t) &= L \int_{L_0}^t e(s) ds + \int_{t_0}^t B ds
 \end{aligned}$$

Motivation:  $e(t) = L * \int_{t_0}^t e(s) ds + \int_{t_0}^t B ds$

$$e' = Le + B$$

$$(e + B/L)' = L(e + B/L)$$

$$\Rightarrow e(t) + B/L = e^{LT}(e(0) + B/L) = B/Le^{Lt}$$

Substitution:  $f(t) = e(t)e^{-Lt}$   $e(t) = e^{Lt}f(t)$   
 $e^{Lt}f(t) \leq L \int_{t_0}^t e^{Ls}f(s) ds + \int_{t_0}^t B ds$

$$\begin{aligned}
 f(t) &\leq L \int_{t_0}^t e^{L(s-t)} \underbrace{f(s)}_{\leq \max_{s \in [t_0, t]} f(s)} ds + e^{-Lt} \int_{t_0}^t B ds \\
 &\leq (\max_{s \in [t_0, t]} f(s))(1e^{L(t_0, t)}) + e^{-Lt} \int_{t_0}^t B ds \\
 f(t) &\leq (\max_{s \in [t_0, T]} f(s))(1 - e^{-L(T-t_0)}) + e^{-Lt}(T-t_0)B \\
 \Rightarrow \max_{t \in [t_0, T]} f(t) &\leq (\max_{t \in [t_0, T]} f(t))(1 - e^{-L(T-t_0)}) + e^{-Lt}(T-t_0)B \\
 \max_{t \in [t_0, T]} e^{-L(T-t_0)} &\leq e^{-Lt}(T-t_0)B \\
 \max f(t) &\leq e^{Lt-2Lt_0}(T-t_0)B \\
 \max_t e(t) &= \max_t e^{\underbrace{Lt}_{\leq E^{LT}}} f(t) \leq e^{LT} \max_t f(t) \leq e^{2L(T-t_0)}(T-t_0)B \\
 \Rightarrow \max_t \|x(t) - \tilde{x}(t)\| &\leq ((T-t_0)e^{2L(T-t_0)}) \max_t \|r(t)\|
 \end{aligned}$$

### 3.4 Lemma von Gronwall

$$\begin{aligned}
 e(t) &\leq \int_{t_0}^t a(s)e(s)ds + \int_{t_0}^t b(s)ds \\
 \Rightarrow e(s) &\leq e(0)e^{A(t)} + \int_{t_0}^t b(s)e^{A(t)-A(s)}ds
 \end{aligned}$$

$$(A(t) = \int_{t_0}^t a(s)ds)$$

Differenzenverfahren

$$(x_{j+1} - \tilde{x}_{j+1}) - (x_j - \tilde{x}_j) = \psi_j(x_j) - \psi_j(\tilde{x}_j) + r_j$$

Summation über  $j$  von 0 bis  $(k-1)$  Teleskopsumme

$$(x_k - \tilde{x}_k) - (x_0 - \tilde{x}_0) = \sum_{j=0}^{k-1} (\Psi_j(x_j) - \Psi_j(\tilde{x}_j)) + \sum_{j=0}^{k-1} r_j$$

Dreiecksungleichung:

$$\Rightarrow \|x_k - \tilde{x}_k\| \leq \sum_{j=0}^{k-1} \|\Psi_j(x_j) - \Psi_j(\tilde{x}_j)\| + \sum_{j=0}^{k-1} \|r_j\|$$

$$\Psi_j(x_j) = \tau_j \Psi(t_j, x_j, \tau_j)$$

$\Psi$  soll  $f$  approximieren

Annahme:  $\Psi$  Lipschitz-stetig bzgl.  $x$

$$\|\Psi(t, x, \tau) - \Psi(x, y, \tau)\| \leq \tilde{L} \|x - y\| \forall t, \forall \tau > 0$$

$$\Rightarrow \|\Psi_j(x_j) - \Psi_j(\tilde{x}_j)\| = \tau_j \|\Psi(t_j, x_j, \tau_j) - \Psi(t_j, \tilde{x}_j, \tau_j)\| \leq \tau_j \tilde{L} \|x_j - \tilde{x}_j\|$$

$$\underbrace{\|x_k - \tilde{x}_k\|}_{=: e_k} \sum_{j=0}^{k-1} \tau_j \tilde{L} \underbrace{\|x_j - \tilde{x}_j\|}_{=: e_j} + \sum_{j=0}^{k-1} \tau_j \underbrace{\|r(t_j, x, \tau_j)\|}_{=: B}$$

Konsistenzfehler  $B = \max_j \|r(t_j, x, \tau_j)\|$

$$f_k = e_k e^{-L(t_k - t_0)}$$

$$e^{L(t_k - t_0)} f_k \leq \sum_{j=0}^{k-1} L \tau_j e^{L(t_j - t_0)} f_j + B \quad \underbrace{\sum_{j=0}^{k-1} \tau_j}_{=t_{j+1} - t_j = t_k - t_0 \leq T - t_0}$$

$$f_k \leq L \sum_{j=0}^{k-1} \tau_j e^{L(t_j - t_k)} + B(T - t_0) e^{-L(t_k - t_0)}$$

$$f_k \leq L(\sum \tau_j e^{L(t_j - t_k)}) \max_j f_j + B(T - t_0) e^{-L(t_k - t_0)}$$

03.11.2009

$$x' = f(x, t), x(0) = x_0$$

$$\tilde{x}(t_{j+1}) = \tilde{x}(t_j) + \Delta \Psi(t_j, \tilde{x}(t_j), \tau_j)$$

Fehler:  $e_j = x(t_j) - \tilde{x}(t_j)$

$$\frac{e_{j+1} - e_j}{\tau_j} - (\Psi(t_j, x(t_j), \tau_j)) - \Psi(t_j, \tilde{x}(t_j), \tau_j) - r(t_j, x, \tau_j)$$

$$\|e_k\| \leq L \sum_{j=0}^{k-1} \|e_j\| \tau_{j-1} + \underbrace{\sum_{j=0}^{k-1} \tau_j r(t_j, x, \tau_j)}_{\leq (\sum \tau_j) B \leq (T-t_0) B}$$

$$B := \max_j |r(t_j, x, \tau_j)|$$

$$f_k := \|e_k\| e^{-L(t_k - t_0)}, \|e_k\| = f_k e^{L(t_k - t_0)}$$

$$\Rightarrow f_k e^{L(t_k - t_0)} \leq L \sum_{j=0}^{k-1} f_j e^{L(t_j - t_0)} \tau_j + r(T - t_0) B$$

$$\begin{aligned}
 \tau_j e^{L(t_j - t_0)} &\leq \sum_{t_j}^{t_{j+1}} e^{L(t-t_0)} dt \\
 f_k &\leq L \sum_{j=0}^{k-1} f_j \int_{t_j}^{t_{j+1}} e^L(t-t_k) dt + (T-t_0)e^{-L(t_k-t_0)} \\
 &= \sum_{j=0}^{k-1} f_j (e^{L(t_{j+1}-t_k)} - e^{L(t_j-t_k)}) + (T-t_0)e^{-L(t_k-t_0)} B \\
 f_j &\leq \max_{0 \leq i \leq k} f_i \\
 f_k &\leq \max_i f_i \underbrace{\sum_{j=0}^{k-1} (e^{L(t_{j+1}-t_k)} - e^{L(t_j-t_k)})}_{=1} + (T-t_0)e^{-L(t_k-t_0)} \\
 f_k &\leq \max_i f_i \underbrace{(1 - e^{L(t_0-t_k)})}_{\leq 1 - e^{L(t_0-T)}} + (T-t_0) \underbrace{e^{L(t_0-t_k)}}_{\leq 1} B \\
 f_k &\leq (1 - e^{L(t_0-T)}) \max_i f_i + (T-t_0)B \\
 \Rightarrow e^{L(t_0-T)} \max_k f_k &\leq (T-t_0)B \\
 \max_k f_k &\leq (T-t_0) e^{L(T-t_0)} B \\
 \max_k \|e_k\| &= \max_k f_k e^{L(t_k-t_0)} \leq e^{L(T-t_0)} \max_k f_k \\
 &\leq (T-t_0) e^{2L} (T-t_0)B
 \end{aligned}$$

$L$  = Lipschitzkonstante von  $\Psi$

Konvergenzfehler  $\leq C(\tilde{L}, T)$  Konsistenzfehler  $\Rightarrow$  Konvergenzordnung = Konsistenzordnung  
Unterscheidung nach Stabilitätskonstante  $C$

-“Steife Dgl.“ für  $C \gg 1$

-“nicht steift“ wenn  $C \approx e^{L(T-t_0)}$  für  $L$  Lipschitzkonstante von  $f$ )

### 3.5 Steifheit bzgl. Anfangswert

Kondition der Dgl. = kleinste Zahl  $\kappa$  mit  $\|x_1(t) - x_2(t)\| \leq \kappa \|x_1(t_0) - x_2(t_0)\| \forall t \in [t_0, T]$

Diskrete Kondition:  $(x_i^\Delta$  aus Einschrittfunction mit Gitter  $\Delta)$

$$\|x_1^\Delta - x_2^\Delta(t)\| \leq \kappa_\Delta \|x_1^\Delta(t_0) - x_2^\Delta(t_0)\| \forall t \in \Delta$$

„Steifes Problem“  $\kappa^\Delta \gg \kappa$

„nicht steifes Problem“  $K^\Delta \approx \kappa$

(Beachte  $\kappa^\Delta \rightarrow \kappa$  für  $\tau_\Delta \rightarrow 0$ )

Übungsaufgabe 3b:  $y(0) = 1$

Beispiel:  $x' = \lambda x, x(0) = 1$

$\Delta = 0, \tau, 2\tau, \dots$

$$x(t) = e^{\lambda t}$$

Startwert  $x^i(0) \rightarrow x^i(0)e^{\lambda t}$

$$\|x^1(t) - x^0(t)\| = \|e^{\lambda t}(x^1(0) - x^0(0))\| \leq e^{\lambda t} \|x^1(0) - x^0(0)\|$$

$$\kappa = \max : t \in [0, T] e^{\lambda t} = \begin{cases} e^{\lambda T} & \lambda > 0 \\ 1 & \lambda \leq 0 \end{cases} \text{ Eulerverfahren: } \tilde{x}(t - \tau) = \tilde{x}(t) + \tau \lambda \tilde{x}(t) = \tilde{x}(t)(1 + \lambda \tau)$$

$$\tilde{x}(\tau) = \tilde{x}(0)(1 + \tau \lambda)$$

$$\tilde{x}(\lambda \tau) = \tilde{x}(\tau)(1 + \tau \lambda) = \tilde{x}(0)(1 + \lambda \tau)^2$$

$$\tilde{x}(\kappa \tau) = \tilde{x}(1 + \tau \lambda)^2$$

$$\lambda \geq 0 : K_\Delta = (1 + \tau \lambda)^n = (1 + \lambda T/n)^n \rightarrow e^{\lambda T}$$

$\Rightarrow \kappa_\Delta \leq \kappa \rightarrow$  nicht steif

$$\lambda \leq 0 : \tilde{x}(\kappa \tau) = \tilde{x}(1 + \tau \lambda)^k$$

$$|\tilde{x}^1(k\tau) - \tilde{x}^2(\kappa\tau)| = |(1 + \tau \lambda)^k (\tilde{x}^1(0) - \tilde{x}^2(0))| = \begin{cases} (1 + \tau \lambda)^k |\tilde{x}^1(0) - \tilde{x}^2(0)| & \lambda \tau > -1 \\ (\tau |\lambda| - 1)^k |..| & \tau \lambda \leq -1 \end{cases}$$

$$\lambda \tau > -2\tau |\lambda| < 2$$

$$(1 + \tau \lambda)^k \leq 1, \tau \lambda > -1$$

$$(\tau |\lambda| - 1)^k \leq 1, -2 < \tau \lambda \leq -1$$

$$|\tilde{x}^1(k\tau) - \tilde{x}^2(\kappa\tau)| \leq |\tilde{x}^1(0) - \tilde{x}^2(0)| \Rightarrow \kappa_\Delta \leq 1 \leq \kappa$$

$$\tau > 2/(|\lambda|), \kappa_\Delta = \max_k (\tau |\lambda| - 1)^k = (\tau |\lambda| - 1)^n \gg 1 = \kappa$$

### 3.5.1 Beispiel: Implizites Eulerverfahren

$$x(t + \tau) = x(t) + \tau f(x(t + \tau), t + \tau)$$

$$x(t + \tau) = x(t) + \tau \lambda x(t + \tau)$$

$$x(t + \tau) = \underbrace{\frac{1}{1 - \tau \lambda} x(t)}_{\begin{array}{l} > 1 \text{ für } \lambda < 0 \\ < 1 \text{ für jedes } \tau \end{array}}$$

## 3.6 4.2 Explizite Runge-Kutta Verfahren

$$\text{Eulerverfahren: } \frac{x(t+\tau)-x(t)}{\tau} = f(x(t), t)$$

$$\frac{x(t+\tau)-x(t)}{\tau} = \underbrace{x'(t)}_{f(x(t), t)} + 1/2 \underbrace{x''(t)}_{\partial_t f + \partial_x f x' = \partial_t f + f \partial_x f} \tau + 1/6 \underbrace{x''(t)}_{= \partial_{tt} f + 2\partial_{xt} ff + f \partial_{xx} ff + \partial_t f \partial_x f + \partial_x f f \partial_f} \tau^2 + \dots$$

Modifiziertes Eulerverfahren (Runge-Kutta-Verfahren 2. Ordnung)

$$\frac{x(t+\tau)-x(t)}{\tau} = 1/2 f \underbrace{(x(t) + \tau/2 f(x(t), t), t + \tau/2)}_{\approx x(t + \tau/2)}$$

RK mit Stufe s:  $\frac{x(t+\tau)-x(t)}{\tau} = \sum_{i=1}^s b_i k_i$ ,  $b_i$  Koeffizienten,  $k_i$  verschachtelte Funktionswerte

$$K_1 = f(x(t), t)$$

$$K_i = f(x + \tau \sum_{j=1}^{i-1} a_{ij} k_j, t + c_i \tau)$$

$$[K_2 = f(x + \tau a_{21} k_1, t + c_2 \tau)]$$

Freiheitsgrade  $A = (a_{ij})$ ,  $b = (b_i)$ ,  $c = (c_i)$

Butcher-Schema

$$\begin{matrix} c & A \\ & b^T \end{matrix}$$

10.11.2009

## 7 Mehrschrittverfahren für AWP

$$x' = f(x, t), x(t_0) = x_0$$

Runge-Kutta Ordnung p : p-malige Auswertung von  $f$

Verfahren der Ordnung  $p$  mit nur einmaligem Auswerten der Funktion  $f$ . in jedem Zeitschnitt.

### 7.1 Mehrschrittverfahren auf äquidistanten Gittern

$$\Delta = t_j, t_j = t_0 + j\tau, j = 0, \dots, n, \tau = \frac{T-t_0}{n}$$

Wollen nun  $f(x_\Delta(t_j), t_j)$  auswerten

Beispiel 7.1 Motivation 1:

explizites Euler:  $x_\Delta(t_{j+1}) = x_\Delta(t_j) + \tau f(x_\Delta(t_j))$

implizites Euler:  $x_\Delta(t_j) = x_\Delta(t_{j-1}) + \tau f(x_\Delta(t_j), t_j)$

Mittelung:  $1/2x_\Delta(t_{j+1}) + 1/2x_\Delta(t_j) = 1/2 * (x_\Delta(t_j) + \tau f(x_\Delta(t_j))) + 1/2(x_\Delta(t_{j-1}) + \tau f(x_\Delta(t_j), t_j))$   
 $\frac{x_\Delta(t_{j+1}) - x_\Delta(t_{j-1})}{2\tau} = f(x_\Delta(t_j), t_j)$

Zentralen Differenzenquotient

$$\text{Taylor } \frac{x(t_{j+1}) - x(t_{j-1})}{2\tau} - x'(t_j) = \frac{(x(t_{j+1}) - x'(t_j)\tau - x(t_j)) - (x(t_{j+1}) + x'(t_j)\tau - x(t_j))}{2\tau} = \frac{(1/2x''(t_j)\tau^2 + O(\tau^3)) - (1/2x''(t_j)(-\tau)^2 + O(\tau^3))}{2\tau} = O(\tau^2)$$

Konsistenzordnung 2

→ Mehrschrittverfahren  $x_\Delta(t_{j+1})$  berechnet aus  $x_\Delta(t_j), x_\Delta(t_{j-1})$

? Stabilität

Eventueller Nachteil: mehr Speicherbedarf ( bei zentralem Differenzenquotient doppelt)

Motivation 2:

$$x(t_{j+1}) = x(t_{j-1}) + \int_{t_{j-1}}^{t_{j+1}} f(x(t), t) dt$$

Mittelpunktsregel:

$$\int_{t_{j-1}}^{t_{j+1}} f(x(t), t) dt \approx 2\tau f(x(t_j), t_j)$$

Beispiel: 7.2: Milne-Simpson Verfahren ( Simpsom-Regel für die Quadratur)

$$x(t_{j+1}) = x(t_{j-1}) + \int_{t_{j-1}}^{t_{j+1}} f(x(t), t) dt$$

$$\text{Simpson-Regel: } \int_a^b g(t) dt = \frac{b-a}{6} (g(a) + 4g(\frac{a+b}{2}) + g(b))$$

$$x_\Delta(t_{j+1}) = x_\Delta(t_{j-1}) + \tau/3 (f(x_\Delta(t_{j-1}), t_{j-1}) + f(x_\Delta(t_{j+1}), t_{j+1}) + 4(f(x_\Delta(t_j), t_j)))$$

Konsistenzordnung 4, implizit

Allgemeines Lineares Mehrschrittverfahren k-ter Ordnung (k-Schrittverfahren)

$$\alpha_k x_\Delta(t_{j+k}) + \dots + \alpha_0 x_\Delta(t_j) = \tau (\beta_k f(x_\Delta(t_{j+k}), t_{j+k}) + \dots + \beta_0 f(x_\Delta(t_j), t_j))$$

$$f_\Delta(t_j) := f(x_\Delta(t_j), t_j)$$

VS:  $|\alpha_n| + |\beta_k| > 0$

$$|\alpha_0| + |\beta_0| > 0$$

Meist:  $\alpha_k \neq 0$ , Normierung  $\alpha_k = 1$

$\beta_k = 0$ : explizit

$\beta_k \neq 0$ : implizit

Wohldefiniertheit von Mehrschrittverfahren

- Brauchen zusätzliche Anfangswerte  $x_\Delta(t_0), \underbrace{x_\Delta(t_1), \dots, x_\Delta(t_{k-1})}_{\text{Berechnet mit Einschrittverfahren}}$

-  $\beta_k = 0$  induktiv

$$\underbrace{x_\Delta(t_{j+k})}_{\text{wohldefiniert}} = \underbrace{\frac{1}{\alpha_k} \dots}_{\text{abhängig von } x_\Delta(t_{j+k,1}), \dots, x_\Delta(t_j)}$$

$$- \beta_k \neq 0 : x_\Delta(t_{j+k}) = \tau \beta_k / \alpha_k f(x_\delta(t_{j+k}), t_{j+k}) + \underbrace{R_j}_{\text{abhängig von } x_\Delta(t_{j+k,1}), \dots, x_\Delta(t_j)}$$

$$\text{Gleichung: } x = F_j(x) = \tau \beta_k / \alpha_k f(x, t_{j+k}) + R_j$$

Anwendung des Banach'schen Fixpunktsatzes:

Lemma: 7.2 f:  $R^d \times R \rightarrow R$  stetig, Lipschitzbzgl.  $x$   $\|f(x, t) - f(y, t)\| \leq L \|x - y\| \forall t, \forall x, y \in R^n, \tau < 1/L \frac{|\alpha|}{|\beta_n|}$

Dann existiert eine eindeutige Lösung von  $x = F_j(x)$  d.h.  $x_\Delta(t_{j+k})$  ist wohldefiniert.

Beweis: Banach FP:

$$\|F_j(x) - F_j(y)\| = \|\tau \frac{\beta_k}{\alpha_k} (f(x, t_{j+k}) - f(y, t_{j+k}))\| = \tau \frac{|\beta_k|}{|\alpha_k|} \|f(x, t_{j+k}) - f(y, t_{j+k})\| \leq \tau \frac{|\beta_k|}{|\alpha_k|} L \|x - y\|$$

$F_j$  konstruktiv auf  $R^n$ , wenn  $\tau \frac{|\beta_n|}{|\alpha_n|} L < 1$   $\square$

Notation: Shiftoperator  $E$

$$E(x_\Delta(t_j)) = (x_\Delta(t_{j+1}))$$

$$E^2(x_\Delta(t_j)) = (x_\Delta(t_{j+2}))$$

..

$$E^k(x_\Delta(t_j)) = (x_\Delta(t_{j+k}))$$

Mehrschrittverfahren:

$$(\alpha_k E^k + \alpha_{k-1} E^{k-1} + \dots + \alpha_1 E + \alpha_0 I) x_\Delta(t_j) = \tau (\beta_k E^k + \dots + \beta_1 E + \beta_0 I) f_\Delta(t_j)$$

charakteristische Polynome:

$$\rho(z) : 0\alpha_k z^k + \alpha_{k-1} z^{k-1} + \dots + \alpha_0$$

$$\sigma(z) := \beta_k z^k + \beta_{k-1} z^{k-1} + \dots + \beta_0$$

$$\rho(E) x_\Delta = \tau \sigma(E) f_\Delta$$

$$\text{expl. } \rho(z) = z - 1, \sigma(z) = 1$$

$$\text{impl. Euler: } \rho(z) = z - 1, \sigma(z) = z$$

Mittelpunktsregel:  $\rho(z) = z^2 - 1, \sigma(z) = 2z$

Milne-Simpson.  $\rho(z) = z^2 - 1, \sigma(z) = 1/3(z^2 + 4z + 1)$

Extended: Laplace-Transformation:

lineare Dgl mit konstanten Koeffizienten

$$\alpha_k x^{(k)} + \dots + \alpha_1 x' + \alpha_0 x = g(t), t \in R^+$$

$$(Lx)(s) = \int_0^\infty e^{-st} x(t) dt$$

$$(Lx')(s) = \int_0^\infty e^{-st} x'(t) dt = e^{-st} x|_0^\infty - s \underbrace{\int_0^\infty e^{-st} x(t) dt}_{Lx}$$

$L(x') = -x(0) + sL(x)$   
 $x(0) = 0 : Lx' = sL(x)$   
 $x(0) = x'(0) = 0, L(x'') = sL(x') = s^2L(x)$   
 mit einfachsten Anfangswerten  $x(0) = x'(0) = \dots = x^{(k-1)}(0) = 0$   
 $L(x^{(j)}) = s^j L(x)$   
 $\Rightarrow L(\alpha_k x^{(k)} + \dots + \alpha_0 x) = \alpha_k L(x^{(k)}) + \dots + \alpha_0 L(x) = (\alpha_k s^k + \alpha_{k-1} s^{k-1} + \dots + \alpha_0) L(x) = L(g)$   
 $L(x) = \frac{L(g)}{\alpha_k s^k + \dots + \alpha_0}$   
 $x = L^{-1}\left(\frac{L^{-1}(g)}{\alpha_k s^k + \dots + \alpha_0}\right)$   
 $L^{-1}(y) = c \int_0^\infty e^{st} y(s) ds$   
 Integrieren → Partialbruchzerlegung

13.11.2009

---

$$\rho(E)x_\Delta = \tau\sigma(E)f_\Delta \quad \rho(a) = \alpha_k z^k + \alpha_{k-1} z^{k-1} + \dots + \alpha_0$$

$$\text{Laplace-Transformation: } y(s) = \int_s^\infty e^{-st} x(t) dt$$

Analog diskrete Laplace-Transformation

$$y(z) = \sum_{k=0}^{\infty} z^{-k} x_\Delta(t_k)$$

$$y = L(x)$$

$$L(x_\Delta(t_{k+1}) - x_\Delta(t_k)) = L(Ex_\Delta(t_k)) - L(x_\Delta(t_k))$$

$$L(Ex_\Delta) = \sum_{k=0}^{\infty} z^{-k} E(x_\Delta(t_k)) = \sum_{k=0}^{\infty} z^{-k} x_\Delta(t_{k+1})$$

$$= \sum_{l=1}^{\infty} z^{-l} x_\Delta(t_l)$$

$$= z * \sum_{l=1}^{\infty} e^{-l} x_\Delta(t_l)$$

$$= z * \sum_{l=0}^{\infty} z^{-l} x_\Delta(t_l) - tx_\Delta(t_0)$$

$$L(Ex_\Delta) = zL(X_\Delta) - zx_\Delta(t_0)$$

Isbesondere für  $x_\Delta(t_0) = 0$

$$L(Ex_\Delta) = zL(x_\Delta)$$

Alle Anfangswerte = 0 :  $x_\Delta(t_0) = \dots = x_\Delta(t_{k-1}) = 0$

$$L(E^j x_\Delta) = z^j L(x_\Delta)$$

$$L(\rho(E)x_\Delta) - \rho(z)L(x_\Delta)$$

$$L(\sigma(E)f_\Delta) = \sigma(z)L(f_\Delta)$$

$$\Rightarrow \rho(z)L(x_\Delta) = \sigma(z)L(f_\Delta)$$

$$x_\Delta = L^{-1}\left(\frac{\sigma(z)}{\rho(z)}L(f_\Delta)\right)$$

### 3.7 7.1.1 Konsistenz

Definition des Konsistenzfehlers: x Lösung von  $x' = f(x, t)$

$$r(t, x, \tau) = 1/\tau \rho(E)x_\Delta - \sigma(E)f_\Delta$$

### 3.8 Definition 7.6:

Lineare Mehrschrittverfahren ist konsistent von der Ordnung p, wenn für alle Funktionen  $f \in C^p$ ,  $x \in C^{p+1}$ ,  $r(x, t, \tau) = O(\tau^p)$  glm. in t

### 3.9 7.8 Lemma

Konsitenzordnung p, genau dann wenn eine der folgende Bedingungen erfüllt ist.

- (i)  $r(0, Q, \tau) = 0 \forall \tau > 0, \forall$  Polynome Q vom Grad kleiner f.

$$(ii) \quad r(0, exp, \tau) = \rho(e^\tau) - \sigma(e^\tau))O(\sigma^p)$$

$$(iii) \quad \sum_{\alpha_j j^l} j = 0^k \alpha_i = 0 \\ \sum_{j=0} l \sum_{j=0}^k \beta_j^{\beta_i} l = 1, \dots, r$$

Zeigen: Konsistenzordnung  $p \Rightarrow (i) \Rightarrow (ii) \Rightarrow (iii) \Rightarrow$  Konsistenzordnung p

$-KO_p \Rightarrow (i)$

$$r = 1/\tau \rho(E)x - \sigma(E)f = 1/\tau \rho(E) - \sigma(E)x'$$

$$r = (0, q, \tau)$$

$$\text{ko: } p \rightarrow r == 1/2\rho(E)Q - \sigma(E)Q'$$

$$1/\tau \alpha_k Q(k\tau) + \alpha Q(k-1) + \dots + \alpha_0 Q(0) - (\beta_k Q'(k\tau) + \beta_{k-1} W(k-1)\tau + \dots + \beta_0 Q(0)) \\ = c/\tau + \text{Polynom}(\tau), \text{ gradsp .- Polynom(z)} \text{ grad} \leq p-1$$

$$Q(z) = \sum_{i=0}^p l_i z^i$$

$$Q'(z) = \sum_{i=1}^p p^i i z^{i-1} = p_m(i-1)z^i$$

$c/\tau + \text{Polynom in } \tau, \text{ grad} \leq p-1 = r = O(\tau^p)$

$$\Rightarrow r = c/\tau + \sum_{i=0}^{p-1} \delta \tau_i = A\tau^p \forall t > 0$$

$$\dots \Rightarrow \tau = 0$$

$$r(t, x, \tau) = 1/\tau \rho(E)x - \sigma(E)f = 1/\tau \rho(E) - \sigma(E)x'$$

$$(i) \Rightarrow (ii) : r(0, exp, \tau) =$$

$$\exp(\tau) = \underbrace{Q(\tau)}_{\text{Polynom von Grad p}} + R(\tau)\tau^{p+1}$$

$$r = r(0, Q + \tau^{p+1}R, \tau) = 1/\tau \rho(E)(Q + \tau^{p+1}R) - \sigma(E)(Q + \tau^{p+1}R + (p+1)\tau^P R)$$

$$= 1/\tau \rho(E)Q - \sigma(E)Q' + 1/\tau \rho(E)(\tau^{p-1}R) - \sigma(E)(\tau^p(\tau R' + (p+1)R))$$

$$= \tau^p(A(\tau) + B(\tau)) = O(\tau^p)$$

$$r = 1/\tau \rho(e^\tau) - \sigma(e^\tau)$$

$$E^j(\exp) = e^{t+j\tau} = e^{j\tau}e^t$$

$$\rho(E)\exp \sum a_j E^j(\exp) = (\sum a_j e^{j\tau})e^t = \rho(e^\tau)e^t$$

$$t = 0 : r(0, exp, \tau) = 1/\tau \rho(e^\tau) - \sigma(e^\tau)$$

$$-(ii) \Rightarrow (iii) : r = 1/\tau \rho(e^\tau) - \sigma(e^\tau) = O(\tau^p)$$

$$= 1/\tau \left( \sum_{i=0}^k \alpha_i \underbrace{e^{\tau i}}_{= \sum_{l=0}^p 1/l!(\tau i)^l + O(\tau^{p+1}} \right) - \sum_{i=0}^k \beta_i \underbrace{e^{\tau i}}_{= \sum_{l=0}^{p-1} 1/l!(\tau i)^l + O(\tau^p)}$$

$$= 1/\tau \sum_{i=0}^k \sum_{l=0}^p \alpha_i 1/l!(\tau i)^l - \sum_{i=0}^k \sum_{l=0}^p \beta_i 1/l!(\tau i)^l + O(\tau^P)$$

$$= 1/\tau \sum_{i=0}^k \alpha_i + \sum_{l=1}^p \sum_{i=0}^k \alpha_i 1/l! i^l \tau^{l-1} - \sum_{l=0}^{p-1} \sum_{i=0}^k \beta_i 1/l!(\tau i)^l + O(\tau^P)$$

$$= 1/\tau \sum \alpha_i + \sum_{l=1}^p \left( \sum_{i=0}^k \alpha_i 1/l! i^l \tau^{l-1} - \sum_{i=0}^k \beta_i 1/(l-1)! \tau^{l-1} i^{l-1} \right) = 1/\tau \sum \alpha_i + \sum_{l=1}^p \frac{\tau^{l-1}}{l!} \left( \sum_{i=0}^k \alpha_i i^l - l \sum_{i=0}^k \beta_i i^{l-1} \right) + O(\tau^p)$$

Koeffizientenvergleich:  $\sum \alpha_i = 0 \quad \sum \alpha_i i^l = l \sum \beta_i e^{l-1}, l = 1, \dots, L$

(iii)  $\Rightarrow KO_p$

$$r(t, x, \tau) = 1/\tau (\sum \alpha_i c(t + i\tau)) - (\sum \beta_i x'(t + i\tau))$$

$$x(t + i\tau) = x(t) + \sum_{l=1}^p 1/l! x^{(l)}(t) (i\tau)^l + O(\tau^{p+1})$$

$$x'(t + i\tau) = x'(t) + \sum_{l=1}^{p-1} 1/l! x^{(l+1)}(t) (i\tau)^l + O(\tau^P)$$

$$r = 1/\tau ((\sum \alpha_i x(t))) + \sum_{i=0}^k \alpha_i \sum_{l=1}^p 1/l! (i\tau)^l x^{(l)}(t) - (\sum \beta_i \sum_{l=0}^{p-1} 1/l! x^{(l+1)}(i\tau)^l) + O(\tau^P)$$

$$= 1/\tau (\sum \alpha_i) x(t) + \sum_{l=1}^p \frac{\tau^{l-1}}{l!} (\sum \alpha_i i^l - l \sum \beta_i i^{l-1}) x^{(l)}(t)$$

$$= O(\tau^p)$$

$p = 1 : \sum_{i=0}^k \alpha_i = 0, \alpha_k = 1$   
 $\sum_{i=0}^k \alpha_i i = \sum_{i=0}^k \beta_i$   
 $k = 1 : \alpha_0 + 1 = 0 \Rightarrow \alpha_0 = -1$   
 $1 = \alpha_1 = \beta_0 + \beta_1 \Rightarrow \beta_1 = 1 - \beta_0, \beta_1 = \Theta, \beta_0 = 1 - \Theta$   
 $X_\Delta(t + \tau) - x_\Delta(t) = \tau(\Theta f(x_\Delta), r + \tau) + (1 - \Theta)f(x_\Delta(t), t)$   
 $\Theta = 0$ : explizites Euler  
 $\Theta = 1$ : implizites Euler  
 $\Theta/2$ : Crank-Nicholson

---

17.11.2009

---

$\rho(E)x_\Delta = \tau\sigma(E)f_\Delta$   
 $\rho(1) = \sum \alpha_j = 0$   
 $\sum \alpha_j j^l = l \sum \beta_j j^{l-1} l = 1, \dots, p$   
 $\alpha_k = 1$   
 $p+2$  Gleichungen für  $2k+2$  Parameter (Unbekannte)  
 Maximale Konsistenzordnung:  $p = 2k$   
 → implizite Verfahren der Ordnung  $2k$   
 Explizites Verfahren:  $\beta_k = 0 \rightarrow p + 3$  Gleichungen  
 → explizite Verfahren der Ordnung  $2k - 1$   
 ABER: Verfahren maximaler Konsistenzordnung für  $k > 1$  sind immer instabil!

### 3.10 7.1.2 Stabilität

#### 3.10.1 Beispiel: 7.1.1:

$\rho(z) = z^2 + 4z - 5$   
 $\sigma(z) = 4z + 2$   
 $0 = \sum \alpha_k = 1 + 4 - 5 \checkmark$   
 $l = 1 : \sum \alpha_j * j = \sum \beta_j$   
 $-5 * 0 + 4 * 1 + 1 * 2 = 6 = 4 + 2$   
 $l = 2 : \sum \alpha_j j^2 + 2 \sum \beta_j j$   
 $-5 * 0 + 4 * 1 + 1 * 4 = 8 = 2(2 * 0 + 4 * 1)$   
 $l = 3 : \sum \alpha_j j^3 = 3 \sum \beta_j j^2$   
 $(-5 * 0 + 4 * 1 + 1 * 8) = 12 = 3(2 * 0 + 4 * 1) \checkmark$   
 Explizites 2-Schrittverfahren maximaler Konsistenzordnung ( $3 = 2k - 1$ )  
 $x' = 0, x(0) = 1 \rightarrow x(t) = 1$   
 $x_\Delta(0) = 1, x_\Delta(\tau) = 1 + \tau \epsilon$   
 $x_\Delta(t) = 1 + \tau \epsilon / 6 (1 - (-5)^{t/\tau}) \rightarrow \pm \infty$  (für  $t \rightarrow \infty, \tau \rightarrow 0$ )  
 $\rho(z) = z^2 + 4z - 5 = (z - 1)(z - 5)$   
 $\rho$  hat Nullstellen 1 und -5  
 Instabilität in der homogenen Differenzengleichung  $x_{k+2} + 4x_{k+1} - 5x_k = 0$   
 $\rho(E)x_\Delta = 0 \rightarrow$  approximiert  $x' = 0$   
 Forderung:  $x_\Delta$  beschränkt      konstante Lösung und beschränkt

### 3.11 7.12 Definition

Lineares Mehrschrittverfahren heißt stabil, wenn die homogene Differenzengleichung  $\rho(E)x_\Delta = 0$  stabil ist.  
 Die homogene Differenzengleichung heißt stabil, wenn ein  $M > 0$  existiert, sodass  $|x_\Delta(t)| \leq M \forall t = j\tau \forall j > 0$  für alle Anfangswerte  $x_\Delta(0), \dots, x_\Delta((k-1)\tau)$   
 (M darf von Anfangswerten abhängen)  
 Lösung homogenen Differenzengleichungen  
 $\rho(E)x_\Delta = 0$

Ansatz für spezielle Lösungen

$$x_{\Delta}(t) = \lambda^{t/\tau}, t = j\tau \rightarrow \lambda^{t/j} \rightarrow \lambda^j$$

$$E^j(x_{\Delta}(t_k)) = E^j(\lambda^k) = \lambda^{k+j} \lambda^j x_{\Delta}(t_k)$$

$$\rho(E)x_{\Delta}(t_k) = \sum \alpha_j x_{\Delta}(t_k) = \left( \sum_{j=0}^k \alpha_j \lambda^j \right) x_{\Delta}(t_k) = \rho(\lambda)x_{\Delta}(t_k) = 0, \text{ falls } \lambda \text{ Nullstelle von } \rho \text{ ist.}$$

Falls  $\rho$  k Nullstellen mit Einfachheit 1 hat  $\rightarrow$  k linear unabhängige Lösungen

$$x_{\Delta} = \lambda_i^{t/\tau}, i = 1, \dots, k$$

Gleichung linear  $\Rightarrow$  linear Kombinationen von Lösungen sind wieder Lösungen.

$$x_{\Delta} = c_1 \lambda_1^{t/\tau} + c_2 \lambda_2^{t/\tau} + \dots + c_k \lambda_k^{t/\tau}$$

k Koeffizienten  $c_i$  aus den k Anfangswerten  $x_{\Delta}(j\tau), j = 0, \dots, k-1$

lineares Gleichungssystem für  $(c_i)$  mit Vandermonde-Matrix.

Mehrreiche Nullstelle:  $\rho(\lambda) = 0, \rho'(\lambda) = 0$

$$\rho'(\lambda) = \sum \alpha_j j \lambda^{j-1} - 1/\lambda \sum \alpha_j (j \lambda^j)$$

Ansatz:  $x_{\Delta}(t_j) = j \lambda^j$

$$\rho(E)x_{\Delta}(t_j) = \sum \alpha_i E^i x_{\Delta}(t_j)$$

$$= \sum_{j=0}^k \alpha_i (j+i) \lambda^{j+i}$$

$$= j \sum_{i=0}^k \alpha_i \lambda^{j+1} + \sum_{i=0}^k \alpha_i i \lambda^{i+j}$$

$$= j \lambda^j \sum \alpha_i \lambda^i + \lambda^{j+1} \sum \alpha_i i \lambda^{i-1}$$

$$= j \lambda^j \rho(\lambda) + \lambda^{j+1} \rho'(\lambda) > 0$$

$\lambda$  Nullstelle Vielfachheit m:

$$x_{\Delta}(t_j) = j^l \lambda^j \text{ ist Lösung für } l = 0, \dots, m-1$$

$\rho$  hat q verschiedene Nullstellen  $\lambda_i$  mit Vielfachheit  $v_i$

$$\sum_{i=1}^a v_i = k$$

k linear unabhängige Lösungen  $j^l \lambda_i^j, l \in 0, \dots, v_i - 1$

Allgemeine Lösung:

$$X_j = \sum_{i=1}^a \sum_{l=0}^{v_i-1} \lambda_i^j j^l C_{il}$$

k Koeffizienten  $C_{il}$  bestimmt durch die Anfangswerte  $x_{\Delta}(0), \dots, x_{\Delta}((k-1)\tau)$

### 3.12 Satz

Eine lineare Differenzengleichung/lineares Mehrschrittverfahren ist genau dann stabil, wenn für alle Nullstellen  $\lambda$  des char. Polynoms  $\rho$  gilt:  $|\lambda| < 1$  oder  $|\lambda| = 1$  und hat algebraische Vielfachheit 1.

$$\text{Beweis: } x_{\Delta}(t_j) = \sum_{i=1}^a \sum_{l=0}^{v_i-1} c_{il} j^l \lambda_i^j$$

„ $\Rightarrow$ “ Stabilität  $\Rightarrow x_{\Delta}$  beschränkt für alle Anfangswerte  $\Rightarrow x_{\Delta}$  beschränkt für alle  $(c_{il})$

Annahme: es existiert  $\lambda_i$  mit  $|\lambda_i| > 1$  oder  $|\lambda_i| = 1$  und  $v_i \geq 2$

(i)  $|\lambda_i| > 1, c_{i0} = 1$ , alle anderen Koeffizienten = 0.

$$X_{\Delta}(t_j) = \lambda_i^j \Rightarrow \|x_{\Delta}(t_j)\| \rightarrow \infty \text{ für } j \rightarrow \infty$$

„ Widerspruch zur Stabilität

(ii)  $|\lambda_i| = 1$  und  $v_i \geq 2$

$c_{i0} = 0, c_{i1} = 1$ , alle anderen Koeffizienten = 0

$$x_{\Delta}(t_j) = j \lambda_i^j \rightarrow |x_{\Delta}(t_j)| = j |\lambda_i|^j = j \Rightarrow |x_{\Delta}(t_j)| \rightarrow \infty \text{ „ Widerspruch zur Stabilität“}$$

„ $\Leftarrow$ “  $|\lambda_i| < 1, g: s \rightarrow s^l \lambda_i^s, s \in R^+, l \geq 0$

g ist stetige Funktion und  $g(s) \rightarrow 0$  für  $s \rightarrow \infty$

$$\log(g(s)) = l * \log(s) + s \underbrace{\log(\lambda_i)}_{\substack{<0 \\ \rightarrow -\infty}}$$

$$\underbrace{s}_{-\infty} \rightarrow -\infty$$

$s \rightarrow |g(s)|$  stetig

$$(|\lambda_i| = 1, v_i = 1, |\lambda_i^j| = 1 \leq C)$$

$g(0) < \infty, g(\infty) < \infty \Rightarrow g$  beschränkt

$$\Rightarrow EC > 0 : |j^l \lambda_i^j| \leq C$$

$$\begin{aligned}|x_\Delta(t_j)| &> \left| \sum_{i=1}^a \sum_{l=0}^{v_i-1} c_{il} j^l \lambda_i^j \right| \\ &\leq \sum_i \sum_j j |c_{il}| j^2 \lambda_i^j \\ &\leq C \sum_i \sum_j j |c_{il}| =: M\end{aligned}$$

Euler:  $\rho(z) = t - 1\lambda_1 = 1, v_1 = 1 \rightarrow \text{stabil}$

Mittelpunktsregel:  $\rho(z) = z^2 + 1\lambda_1 = 1, v_1 = 1, \text{ oder } \lambda_2 = -1, v_1 = 1 \rightarrow \text{stabil}$

20.11.2009

$$\rho(E)x_\Delta = \sigma(E)f_\Delta$$

Stabilität  $\Leftrightarrow \lambda$  Nullstelle von  $\rho : |\lambda| < 1$

$|\lambda| = 1$  und einfach

HS der Numerik

Konsistenz + Stabilität  $\Leftrightarrow$  Konvergenz

### 3.13 Definition: 7.14

Lineares MS-Verfahren ist konvergent, falls  $\lim_{\tau \rightarrow 0} x_\Delta(t) := x(t)$  für  $t \in \Delta$

Konvergenz von der Ordnung P, wenn  $|x_\Delta(t) - x(t)| = O(\tau p)$

sobald  $x_\Delta(t_0 + j\tau) \rightarrow x_\Delta(t_0) = x_0, j = 0, \dots, k-1$

für  $\tau \rightarrow 0$

### 3.14 Satz 7.15

Ein konvergentes, lineares Mehrschrittverfahren ist stabil und konsistent. Speziell gilt:  $\rho * (1) = \sigma(1) \neq 0$   
 $\rho(1) = 0$

Beweis:

(i) Stabilität:  $\rho(E)x_\Delta = 0$  hat beschränkte Lösungen,  $x' = 0, x(0) = x_0$

Konvergenz  $\Rightarrow x_\Delta(t) \rightarrow x(t) := x_0$

Für  $\tau$  hinreichend klein gilt:

$$|x_\delta(t)| \leq |x_0| + 1$$

$\forall \tau = t_0 : |x_\Delta|$  ist beschränkt durch  $|x_0| + 1$

(ii) Konsistenz  $\rho(1) = \sum \alpha_j = 0$

$$x' = 0, x(t_0) = 1$$

$$x_\Delta(t_0 + j\tau) = 1, j = 0, \dots, k-1$$

$$\rho(E)x_\Delta = \tau\sigma(E)f_\Delta = 0$$

$$\rho(E)x_\Delta = 0$$

$$t = t_0 : \rho(E)x_\Delta(t_0) = \alpha_k x_\Delta(t_k) + \underbrace{\alpha_{k-1} x_\Delta(t_{k-1})}_{=0} + \dots + \alpha_0 x_\Delta(t_0)$$

$$\sum_{\alpha_j} = \alpha_k(1 - x_\Delta(t_k)) \rightarrow \alpha_k(1 - 1) = 0$$

$$\rho(1) = \sum \alpha_j \rightarrow 0 \text{ für } \tau \rightarrow 0$$

$$\Rightarrow \rho(1) = \sum \alpha_j = 0$$

(iii) 1 Nullstelle von  $\rho \Rightarrow 1$  ist einfache Nullstelle  $\Rightarrow \rho'(1) \neq 0$

$$\Theta = \frac{\sigma(1)}{\rho'(1)}$$

$$x_\Delta(t) = \Theta t$$

$$t = 0 : \alpha_k(\Theta k\tau) + \alpha_{k-1}(\Theta(k-1)\tau) + \dots + \alpha_1\Theta\tau = \theta\tau\rho'(1)$$

$$\rho'(1) = \sum j\alpha_j$$

$$\rho(1) = 0 : \alpha_k(\Theta j\tau) + \alpha_{k-1}(\Theta j\tau) + \dots + \alpha_0(\Theta j\tau) = 0$$

$$\alpha_k(\Theta(j+1)\tau) + \underbrace{\alpha_{k-1}(\Theta(j+k-1)\tau)}_{x_\Delta t_{j+k-1}} + \dots + \underbrace{\alpha_0(\Theta j\tau)}_{x_\Delta t_j} = \Theta\tau\rho'(1) = \frac{\sigma(1)}{\rho'(1)}\tau\rho'(1) = \tau\sigma(1)$$

$$= \rho(E)x_\Delta(t_j) = \tau\sigma(1)$$

$$\begin{aligned}
 &= \tau \sum_i \beta_i f_\Delta(r_{j+i}), \text{ wenn } f_\Delta = 1 \\
 &= \tau \sigma(E) f_\Delta \\
 &x_\Delta(t) = \sigma(t) \text{ ist Lösung des linearen Mehrschrittverfahrens für } x' = 1 \\
 &x_\Delta(j\tau) = \Theta j\tau, j = 0, \dots, k-1 \\
 &x_\Delta \text{ ist die Lösung des Mehrschrittverfahrens für } x' = 1, x(0) = 0 \\
 &\text{Lösung: } x(t) = t \\
 &\text{Konvergenz: } x_\Delta(\underbrace{t}_{=\Theta t \rightarrow t \Rightarrow \Theta=1}) \rightarrow x(t) \text{ für } \tau \rightarrow 0
 \end{aligned}$$

### 3.15 Satz 7.23

Ein stabiles und konsistentes, lineares Mehrschrittverfahren ist immer konvergent

Bei Konsistenzordnung  $p$  ist der Diskretisierungsfehler beschränkt durch  $|\epsilon_\Delta(t)| = \Theta(\tau^p + \epsilon_0)$

$$\epsilon_0 = \max_{l=0, \dots, k-1} |x_\Delta(t_l) - x(t_l)|$$

### 3.16 Satz 7.16:

Es existieren stabile Verfahren der Ordnung  $p$ , falls  $p \leq k+2$ ,  $k$  gerade, und  $p \leq k+1$ ,  $k$  ungerade  
 $p \leq k$  für  $\beta_k = 0$

### 3.17 Konstruktion stabiler Mehrschrittverfahren

Standard: Adams-Verfahren

BDF-Verfahren(für steife Probleme)

#### 3.17.1 Adams-Verfahren

Motivation: wollen stabiles Verfahren

wollen maximale Konsistenzordnung (fast)  $\rightarrow p = k+1$  (implizit),  $p = k$  (explizit)

implizit:  $\alpha_k-1, \dots, \alpha_0, \beta_k, \dots, \beta_0 : 2k+1$

$p = k+1$  Gleichungen

Koeffizienten

$k+1$  aus Gleichungen,  $\sum \alpha_j j^l = l \sum \beta_j j^{l-1}$

Idee: Wähle zuerst  $\alpha_{k-1}, \dots, \alpha_0$  so dass das Polynom  $\rho$  einfache Nullstelle 1 hat ( $\sum \alpha_j = 0$ ),  $k-1$  fache Nullstelle 0

$$\rho(z) = z - 1 z^{k-1} = z^k - z^{k-1}$$

Bestimme  $\beta_j$  aus  $\sum \beta_j j^{l-1} = 1/l \sum \alpha_j j^l$

$\rightarrow$  Gleichungssystem mit Vandermonde-Matrix

$\rightarrow$  ex existiert eindeutige Lösung

$\rightarrow$  Lösung darstellbar über dividierte Differenzen wie bei der Interpolation

$k=0 : \sigma(z) = z \rightarrow$  implizite Eulerverfahren

$k=1 : \sigma(z) = \frac{z+1}{2}$  Crichton-Nichelson

$$x_\Delta(t_{j+1}) - x_\Delta(t_j) = \tau/2(f_\Delta(t_{j+1}) + f_\Delta(t_j))$$

$$k=2 : \sigma(z) = \frac{5z^2+8z-1}{12}, \rho(z) = z^2 - z$$

$$k=3 : \sigma(z) = \frac{7z^3+19z^2-5+1}{24}, \rho(z) = z^3 - z^2$$

#### 3.17.2 Adams-Moulton-Verfahren der Ordnung K

explizite Verfahren:  $\beta_k = 0$

$$\rho(z) = z^k - z^{k-1} \beta_{k-1}, \dots, \beta_0 \text{ bestimmt aus } \sum_{j=0}^{k-1} \beta_j j^{l-1} = 1/l \sum \alpha_j j^l, l \leq k-1$$

### 3.17.3 Adams-Bashfort Verfahren

$$k=1 : \sigma(z) = 1: \text{explizites Eulerverfahren}$$

$$k=2 : \sigma(z) = \frac{3z-1}{2} \rho(z) = z^2 - z$$

$$k=3 : \sigma(z) = \frac{23z^2-16z+5}{12} \rho(z) = z^3 - z^2$$

$$k=4 : \sigma(z) = \frac{55z^3-59z^2+37z-9}{24}$$

## 4 §8 Numerische Lösung von Randwertproblemen

Finite-Differenzen-Verfahren

### 4.1 (8.1) Randwertprobleme für lineare gewöhnliche Dgl. zweiter Ordnung

$$\bar{p}(s)x''(s) + \bar{r}(s)x'(s) + \bar{q}(s)x(s) = f(s), s \in (a, b)$$

$$\beta_1 x = \alpha_1, \beta_2 x = \alpha_2, x(a) + \beta_1 x'(a) = \gamma, x(b) + \beta_2 x'(b) = \delta$$

$\beta_1 = \beta_2 = 0 (\alpha_1, \alpha_2 = 1)$ : Dirichlet-Problem

$\alpha_1 = \alpha_2 = 0 (\beta_1 = \beta_2 = 1)$ : Neumann-Problem

Annahme:  $\bar{p}(s) \neq 0, \forall s \in [a, b]$

$$x'' + \frac{\bar{r}}{\bar{p}}x' + \frac{\bar{q}}{\bar{p}}x = \frac{\bar{f}}{\bar{p}}$$

$$\underbrace{px''}_{(px')'} + \underbrace{\frac{\bar{r}}{\bar{p}}px'}_{=:q} + \underbrace{\frac{\bar{q}}{\bar{p}}p}_{=:f}x = \underbrace{\frac{\bar{f}}{\bar{p}}p}_{=:f}$$

$$p' = \frac{\bar{r}}{\bar{p}}p \Rightarrow p(s) = \int_a^s \frac{S(\sigma)}{\bar{p}(\sigma)} d\sigma$$

$$Lx = (px')' + qx = f$$

L ist formal selbstadjungiert in  $L^2$ -Skalarprodukt

$$x, y \text{ mit } x(a) = x(b) = y(a) = y(b) = 0$$

$$\langle Lx, y \rangle = \int_a^b (Lx)(s)y(s)ds = \int_a^b [(px')' + qx]yds$$

$$= \underbrace{px'y|_a^b}_{=0} - \int_a^b px'y'ds + \int_a^b [0, y]xds$$

$$= \underbrace{-pxy'|_a^b}_{=0} + \int_a^b x[(py')' + qy]ds = \langle x, Ly \rangle$$

24.11.2009

---

#### 4.1.1 Randwertprobleme

$$(px')' + \rho x = f, \text{ „Divergenzform“}$$

$$\nabla(-p\nabla x) + \rho x = f$$

$$\nabla(\vec{j})..$$

$$\sum_i \frac{\partial}{\partial s_i} (p \frac{\partial x}{\partial s_i}) + \rho x = f$$

$$(px')' + \rho x = f$$

$$B_1 x = \gamma, B_2 x = \delta$$

$$\alpha_1 x(a) + \beta_1 x'(a), \alpha_2 x(a) + \beta_2 x'(a)$$

### 4.2 8.1: Existenz und Eindeutigkeit von Lösungen

$$px' = \int_a^s (f(\sigma) - \rho(\sigma)x(\sigma))d\sigma + c_1$$

Beispiel:  $\alpha_1 = 0, s = a$

$$x'(a) = \frac{\gamma}{\beta_a}, p(a)x'(a) = c_1 \Rightarrow c_1 = \frac{\gamma}{\beta_1}p(a)$$

Eventuell  $c_1$  zunächst unbestimmt ( $a_1 \neq 0, a_2 \neq 0$ )

Nochmal integrieren:

$$x'(s) = \frac{1}{p(s)}(c_1 + \int_a^s (f - \rho x)d\sigma)$$

$$\Rightarrow x(s) = \int_a^s 1/(p(\sigma))(c_1 + \int_a^\sigma (f - \rho x)d\xi)d\sigma + c_3$$

Bestimme  $c_1, c_2$  aus Nebenbedingungen

Bsp:  $x(a) = \gamma, x(b) = \delta$

$$s = a: \gamma = x(a) = c_2$$

$$s = b: \delta = x(b) = \int_a^b 1/p(\sigma)(c_1 + \int_a^\sigma (f - qx)d\xi)d\sigma + \gamma$$

$$x_1 \int_a^b 1/p(\sigma)d\sigma = \gamma - \delta + \int_a^b 1/p(\sigma) \int_a^\sigma (f - qx)d\xi d\sigma$$

$$x_1 = (\gamma - \delta)\bar{p} + \int_a^b \bar{p}/p(\sigma) \int_a^\sigma (f(\xi) - qx(\xi))d\xi d\sigma$$

$$\bar{p} = 1/(\int_a^b 1/p(\sigma)d\sigma)$$

$$x(s) = (\gamma - \delta) \int_a^s \bar{p}/p(\sigma)d\sigma + \int_a^s 1/p(\sigma)d\sigma \int_a^b \bar{p}/p(\sigma) \int_a^\sigma (f - qx)d\xi d\sigma + \int_a^s 1/p(\sigma) \int_a^\sigma (f - qx)d\xi d\sigma + \gamma$$

$$\int_a^b 1/p(\sigma) \int_a^\sigma (f - qx)d\xi d\sigma = \int_a^b (f(\xi) - q(\xi)x(\xi)) \int_\xi^b 1/p(\sigma)d\sigma d\xi$$

$$= \int_a^b f(\xi) \int_\xi^b 1/p(\sigma)d\sigma + \int_a^b (-q(\xi) \int_x^b 1/p(\sigma)d\sigma)x(\xi)d\xi$$

$$\int_a^s 1/p(\sigma) \int_a^\sigma (f - qx)d\xi d\sigma = \int_a^s f(\xi) \int_\xi^s 1/p(\sigma)d\sigma + \int_a^s (-q(\xi) \int_x^s 1/P(\sigma)d\sigma)x(\xi)d\xi$$

$$\int_a^s Funktion(\xi, s)x(\xi)d\xi = \int_a^b Schlangefunktion(\xi, S)x(\xi)d\xi$$

$$Schlangefunktion = \begin{cases} Funktion(\xi, S)/\zeta \leq s \\ 0 \quad \zeta > s \end{cases}$$

Allgemeine Form:

$$x(s) = g(s) + \int_a^b k(\zeta, s)x(\zeta)d\zeta \text{ mit } k \text{ stetig}$$

Integralgleichung 2. Art (x im Integral und ausserhalb)

$$(\text{integralgleichung 1. Art: } x \text{ nie im Integral, } \int_a^b k(\zeta, s)x(\zeta)d\zeta = g(s))$$

Achtung: i.A. ist  $k(\zeta, s) \neq 0$  für  $\zeta > s$

Vergleichung: AWP:  $x(s) = g(s) + \int_a^s k(\zeta, s)x(\zeta)d\zeta$

mit  $\bar{R}(\zeta, S) = 0$  für  $\zeta > s$ : Volterra-Integralgleichung

RWP:  $k(\zeta, s) \neq 0$  für  $\zeta > s \rightarrow$  Integralgleichung koppelt alle Werte von x im Intervall (a,b), (Fredholm-Integralgleichungen)

$$x(s) = \rho(s) + \underbrace{\int_a^b k(\sigma, s)x(\sigma)d\sigma}_{(Kx)(s)}$$

$$K.L^2([a, b]) \rightarrow L^2([a, b])$$

$$x \rightarrow \int_a^b k(., \sigma)x(\sigma)d\sigma$$

$$x = g + Kx$$

Fixpunktgleichung in Banachraum

Versuch 1: Banach'scher Fixpunktsatz:  $F(x) = g + Kx$

Kontraktivität:  $\|F(x_1) - F(x_2)\| \leq c \|x_1 - x_2\|, c < 1$

$$\|F(x_1) - F(x_2)\|_{L^2} = \|K(x_1 - x_2)\|_{L^2}$$

$$= \sqrt{\int_a^b (\underbrace{\int_a^\sigma k(\sigma, s)x(s)d\sigma}_{\text{Cauchy-Schwarz}})^2 ds}$$

$$\text{aus (1) folgt: } \leq \sqrt{\int_a^b k(\sigma, s)^2 d\sigma} * \sqrt{\int_a^b x(s)^2 d\sigma}$$

$$\|F(x_1) - F(x_2)\|_{L^2} \leq \sqrt{\int_a^b \int_a^b d(\sigma, s)^2 d\sigma \int_a^b x(\sigma)^2 d\sigma ds} \leq \underbrace{\sqrt{\int_a^b \int_a^b k(\sigma, s)^2 d\sigma ds}}_{=c} \|x\|_{L^2}$$

F kontraktiv ( $c < 1$ ), wenn  $k$  klein genug ist

I.A. nicht gegeben, auch für RWP

Versuch 2: Schauderschen Fixpunktsatz (nur Existenz)

VS:

(i) K kompakt ✓

(ii) K bildet  $B_R(x_0)$  in sich selbst ab?? → nur für  $k$  klein

Versuch 3: K kompakt + linear → lineare Operatoren über Eigenwerte charakterisieren Eigenwertproblem für K

$$Kx = \lambda x, (\lambda \in C)$$

$\lambda = 1$  Eigenwert von K,  $x_1$  Eigenvektor

$$Kx_1 = x_1$$

$$x = Kx + g$$

$$(x + \alpha x_1) = K(x + \alpha x_1) + g$$

x Lösung  $\Rightarrow x + \alpha x_1$  Lösung für alle  $\alpha \in R$

Idee: 1 Eigenwert von K + eventuell keine Existenz und Eindeutigkeit

1 kein Eigenwert von K ⇒ Existenz und Eindeutigkeit??

(1 kein Eigenwert  $\Leftrightarrow x = Kx$  hat nur Lösung  $x = 0$ )

Spektrum eines linearen Operators:

$$\sigma(K) := \{\lambda \in C \mid \lambda I - K \text{ stetig invertierbar}\}$$

für uns interessant:  $1 \in \sigma(K)$ ?

$$1 \notin \sigma(K) \quad (I - K)x = f \Rightarrow x = (I - K)^{-1}g$$

Kiesz'sche Sätze für  $\lambda I - K$

1.RS:  $\lambda \neq 0 \Rightarrow \lambda I - K$  surjektiv

2.RS: Es gibt nur abzählbar viele  $\lambda_k \neq 0$  mit  $\lambda_k I - K$ , nicht injektiv ( $\Leftrightarrow N(\lambda_k I - K) \neq 0$ )

$N(\lambda_k I - K)$  ist endlich dimensional

⇒ Spektrum eines kompakten Operators

$$\sigma(K) \subseteq \{0\} \cup \lambda_{k=0}^{\infty}, \lambda_k \text{ ist Eigenwert mit endlicher Vielfachheit}$$

Also  $1 \in \sigma(K) \Leftrightarrow 1$  ist Eigenwert von K  $\Leftrightarrow \exists x \neq 0 : x = Kx$

Integralgleichung 2.Art ist eindeutig lösbar, wenn es kein  $x \neq 0$  mit  $x = Kx$  gibt.

RWP: Analog wie vorher führt  $(px')' + qx = 0$

$$B_1x = 0, B_2x = 0$$

auf  $x = Kx$ .

Satz: Das RWP  $(px')' + qx = f, B_1x = \gamma, B_2x = \delta$

hat genau dann eine (eindeutige) Lösung, wenn das homogene AWP  $(px')' + qx = 0, B_1x = 0, B_2x = 0$  nur die Lösung  $x = 0$  hat.

Überprüfung des homogenen Problems

Beispiel:  $x(a) = 0, x(b) = 0$

$$(px')' + qx = 0$$

$$\Rightarrow (px')'x + qx^2 = 0$$

$$\Rightarrow \int_a^b (px')'x + qx^2 ds = 0$$

$$\Rightarrow \underbrace{px'x|_a^b}_{=0} - \int_a^b px'x' + \int_a^b qx^2 = 0$$

$$\Rightarrow \int_a^b px'^2 - qx^2$$

$$p = \exp(..) > 0, \int x'^2 \geq 0$$

Mögliche Bedingung an q:

$$q(x) < 0 \forall s \in [a, b]$$

$$0 = \int_a^b (px')^2 dx >_0 - \underbrace{qx^2}_{>0} dx < 0 \text{ für } x \neq 0$$

Widerspruch zur Lösung  $x = 0$

---

27.11.2009

$$(px')' + qx = f, s \in [a, b]$$

$$B_1 x = \gamma, B_2 x = \delta$$

$$x = Kx + q$$

$$(Kx)(s) = \int_a^b k(\sigma, s)x(\sigma)d\sigma$$

1 kein Eigenwert von K

⇒ Existenz und Eindeutigkeit

$$1 \text{ kein Eigenwert von } K \Leftrightarrow (px')' + qx = 0, B_1 x + B_2 x = 0$$

1 Eigenwert von K

2. Riesz'schen Satz ⇒ EW 1 hat endliche Vielfachheit m.

Es gibt m linear unabhängige Lösungen von  $x = Kx$  Eigenvektoren (Eigenfunktionen)  $x_1, \dots, x_m$

Falls x Lösung von  $x = Kx + g \Rightarrow \tilde{x} = x + \alpha_1 x_1 + \dots + \alpha_m x_m$ , löse ebenfalls  $\tilde{x} = K\tilde{x} + \rho$

Fragen: Schranke für m?

Überprüfung der Existenz.

Vielfachheit m: Bsp:  $p = 1, q = const.$

$$x'' = qx = 0$$

#### 4.2.1 lineare DGL mit konstanten Koeffizienten

→ Lösungen nur von der Form  $\sum \alpha_i e^{\lambda_i s}$

Was sind mögliche Werte  $\lambda_i \in C$

$$x = e^{\lambda s}$$

$$\Rightarrow x'' + qx = \lambda^2 e^{\lambda s} + qe^{\lambda s} = 0, \forall s \in [a, b]$$

$$\Rightarrow \lambda^2 + q = 0 \Rightarrow \lambda = \pm \sqrt{-q}$$

⇒ maximal zwei linear unabhängige Lösungen

$$x = e^{\sqrt{-q}s}, x = e^{-\sqrt{-q}s} \text{ von } (px') + qx = 0$$

⇒ Vielfachheit  $m \leq 2$

→ Diese Aussage  $m \leq 2$  gilt für allgemeine p,q

Wie überprüfen wir die Existenz?

$$x = Kx + g$$

$$\text{Einfach: } \langle x, x \rangle = \langle Kx, x \rangle + \langle q, x \rangle$$

(wie letztes Mal multiplizieren der Dgl. und integrieren)

Statt mit x nehmen wir Skalarprodukt mit y.

$$\langle x, y \rangle = \langle Kx, y \rangle + \langle q, y \rangle$$

$$\langle g, y \rangle = \langle x, y \rangle - \langle Kx, y \rangle = \langle x, y \rangle - \langle x, K^*y \rangle = \langle x, y - K^*y \rangle$$

$$\Rightarrow \langle g, y \rangle = 0, \text{ falls } y = K^*y$$

Existenz von  $x = Kx + g \Leftrightarrow \langle g, y \rangle = 0 \forall y, y = K^*y$

### 4.3 FREDHOLM-ALTERNATIVE

1 EW von K mit Vielfachheit m  $\Leftrightarrow$  1 EW von  $K^*$  mit der selben Vielfachheit m

#### 4.3.1 Adjungierten Integraloperator

$$\langle Ku, v \rangle = \langle u, K^*v \rangle \quad \forall u, v$$

$$\langle Ku, v \rangle = \int_a^b (Ku)(s)v(s)ds = \int_a^b \int_a^b k(\sigma, s)u(\sigma)v(s)d\sigma ds$$

$$\begin{aligned}
 &= \int_a^b u(\sigma) \int_a^b k(\sigma, s) v(s) ds d\sigma \\
 &= \int_a^b u(\sigma) (K^* v)(\sigma) d\sigma \\
 (Ku)(s) &= \int_a^b K(\sigma, s) u(\sigma) d\sigma \\
 (K^* v)(s) &= \int_a^b K(s, \sigma) v(\sigma) d\sigma
 \end{aligned}$$

$k$  symmetrisch  $\Leftrightarrow k(s, \sigma) = k(\sigma, s) \Rightarrow K = K^*$

Differentialoperator  $L = (px')' + qx$  erfüllt  $L = L^*! \Rightarrow K = K^*$

Also  $y = K^*y \Rightarrow y = Ky \Leftrightarrow (py')' + qy = 0, B_1y = 0, B_2y = 0$

Fredholm - Alternative für RWP

$y$  homogene Lösung  $(py')' + qy = 0$

$(px')' + qx = f < -$  Skalarprodukt

$$\begin{aligned}
 \int_a^b (px')' y + qy ds &= \int_a^b fy ds \\
 px'|_a^b - pxy'|_a^b + \int_a^b x((px')' + qy) ds &= \int_a^b fy ds
 \end{aligned}$$

Bsp.  $x(a) = \gamma, x(b) = \delta, y(a) = y(b) = 0$

$$0 = \int_a^b fy ds + p(b)\delta y'(t) - p(a)\gamma y'(a)$$

#### 4.4 8.2 Numerische Lösung von AWP mit finiten Differenzen

Unterteilung (a,b) in Gitter  $\Delta = \{S_j\}$

$$a = S_0 < S_1 < \dots < S_N = b$$

(am einfachsten: äquidistante Gitter  $S_j = ajh, h = 1/N$ )

Diskretisiere  $(px')'(s_j) + qx(s_j) = f(s_j), j = 1, \dots, N-1$

$$B_1x(S_0) = \gamma, B_2x(S_N) = \delta$$

Approximation erster Ableitungen:

$$\text{z.B. analog Euler, } x'(S_j) \approx \frac{x(S_j) - x(S_{j-1})}{S_j - S_{j-1}} =: D_- x$$

$$x'(S_j) \approx \frac{x(S_{j+1}) - x(S_j)}{S_{j+1} - S_j} =: D_+ x$$

analog Mittelpunktsregel

$$x'(S_j) \approx 1/2D_- x + 1/2D_+ x = \frac{x(S_{j+1}) - x(S_{j-1})}{2h}$$

$$\text{Konsistenz: } x'(S_j) - D_- x(s_j) = x'(S_j) - \frac{x(s_j) - x(s_{j-1})}{h}$$

$$= x'(S_j) - \frac{x(s_j) - (x(s_j) - x'(s_j)h + O(h^2))}{h}$$

$D_-$  und  $D_+$  konsistent von Ordnung 1

$$D_c x := 1/2(D_+ x + D_- x) = x'(s_j) - D_c x(s_j) = x'(s_j) - \frac{1}{2h}(x(S_{j+1}) - x(S_{j-1}))$$

$$= x'(s_j) - \frac{1}{2h}(x(s_j) + x'(s_j)h + 1/2x'(s_j)h^2 + O(h^3) - (x(s_j) - x'(s_j) + 1/2x''(s_j)h^2 + O(h^3))) = x'(s_j) - \frac{1}{2h}(2x'(s_j)h + O(h^3)) = O(h^2)$$

$$(px')' = p'x' + px''$$

$$x'' \approx D_c D_c x \approx D_- D_- x \approx D_+ D_+ x \approx D_+ D_- x \approx D_- D_+ x \approx D_- D_c x \approx D_+ D_c x$$

$$\text{Beispiel: } D_+ D_+ = D_+ \left( \frac{x(s_{j+1}) - x(s_j)}{h} \right) = \frac{x(s_{j+2}) - 2x(s_{j+1}) + x(s_j)}{h^2}$$

$$D_+ D_- x \approx D_+ \left( \frac{x(s_j) - x(s_{j-1})}{h} \right) = \frac{x(s_{j+1}) - 2x(s_j) + x(s_{j-1})}{h^2}$$

$$S_{n-1} = b - h : D_+ D_- x = \frac{x(s_N) - 2x(s_{N-1}) + x(s_{N-2})}{h^2}$$

$$D_+ D_+ x = \frac{x(s_{N+1}) - 2x(s_N) + x(s_{N-1})}{h^2}$$

$$x''(s_j) - D_+ D_- (s_j) = x''(s_j) - 1/h^2(s_{j+1} - 2x(s_j) + x(s_{j-1}))$$

$$= x''(s_j) - 1/h^2(x(s_j) + x'(s_j)h + 1/2x''(s_j)h^2 - 2x(s_j) + x(s_j) - x'(s_j) + 1/2x''(s_j)h^2 + O(h^3)) = O(h^2)1/h^2$$

$$\text{Genauer (1/6}x'''(s_j)h^3 - 1/6x''''(s_j)h^3 + O(h^4)) = O(h^2)$$

$$x''(s_j) - D_+ D_+ x(s_j) = x''(s_j) - 1/h^2(x''(s_{j+1})h^2 + O(h^4)) = x''(s_j) - x''(s_{j+1}) + O(h^2)$$

$$= x''(s_j) - (x''(s_j) + x'''(s_j)h + O(h^2) + O(h^2)) = -x'''(s_j)h + O(h) = O(h)$$

$$\partial_x^t = -\partial_s : D_+^t = -D_-$$

$$\int x'(s)y(s) ds = \int x(s)y'(s) ds$$

$$\langle \partial_s x, y \rangle = \langle x - \partial y \rangle$$

$$D_2 x = D_+ D_- x = D_- D_+$$

04.12.2009

$$(px')' + qx = f, s \in (a, b)$$

Diskretisierung auf äquidistantem Gitter

$$\Delta = a, jh$$

$$x'' \approx D_2x = D_+D_-x = D_-D_+x = \frac{x(s_{j+1}-2x(s_j)+x(s_{j-1})}{h^2}$$

=> Konsistenzordnung 2

$$px'' + p'x' + qx = f$$

Diskretisierung mit Konsistenzordnung 2

$$x_\Delta : \Delta - \rightarrow R^d$$

$$p(s_j)(D_2x_\Delta(s_j)) + p'(S_j)D_cx_\Delta(s_j) + q(s_j)x_\Delta(s_j) = f(s_j), j = 1, \dots, N-1$$

$j = 0$ : Randbedingung  $B_1x = \gamma$

$j = N$ : Randbedingung  $B_2x = \delta$

## 4.5 Diskretisierung der Randbedingung

$$x = a : B_1x = \alpha_1x(a) + \beta_1x'(a) = \gamma$$

$$\text{diskretisiert: } \alpha_1x_\Delta(S_0) + \beta_1D_+x(s_0) = \gamma$$

$$s = b : B_2x = \alpha_2x(b) + \beta_2x'(b) = \delta$$

$$j = N : \alpha_1x_\Delta(S_N) + \beta_1D_-x_\Delta(S_N) = \delta$$

$N+1$  lineare Gleichungen für  $(N+1)$  Unbekannte:  $x_\delta(s_j) =: u_j$

$$j = 0 : \alpha_1u_0 + \beta_1\frac{u_1-u_0}{h} = \gamma$$

$$p_j = p(s_j), p'_j = p'(s_j), f_j = f(s_j), q_j = q(s_j)$$

$$1 \leq j \leq N-1 : P_j \frac{u_{j+1}-2u_j+u_{j-1}}{h^2} + p'_j \frac{u_{j+1}-u_{j-1}}{2h} + q_j u_j$$

$$j = N : \alpha_2u_N + \beta_2\frac{u_N-u_{N-1}}{h} = \delta$$

=> LGS  $Ax = y$

$$y = \begin{pmatrix} \gamma \\ f_1 \\ f_2 \\ .. \\ f_{N-1} \end{pmatrix}$$

$$A = \begin{pmatrix} \alpha_1 - \beta_1/h & \beta_1/h & 0 & & & \\ p_1/h^2 - (p'_1)/(2h) & -\frac{2p}{h^2} + q_1 & \frac{p_1}{h^2} + \frac{p'_1}{2h} & .. & 0 & \\ .. & .. & .. & .. & .. & \\ 0 & .. & .. & 0 & -\frac{\beta_2}{h} & \alpha_2 + \frac{\beta_2}{h} \end{pmatrix}$$

Das ergibt eine tridiagonale Matrix

Max  $\leq 3N-1$  Nichtnullelemente:

$$(N+1)^2 - (3N+1) = N^2 - N \text{ Nulleinträge}$$

Für große N sind die meisten Einträge von A Null.

-> Dünnbesetzte Matrix (sparse)

-> Speicheraufwand Sparse

->  $(i, j, a_{ij})$  für  $a_{ij} \neq 0$

werden  $(3N-1)^2$  Maschinenzahlen und  $2(3N+1)$  Integer

Löse  $Au = f$  mit speziellen Methoden, die Dünnbesetzung von A ausnutzen -> mehr im Januar

Beispiel:

$$\alpha_1 = \alpha_2 = 1$$

$$x(a) = \gamma, x(b) = \delta$$

$$p = 1, x'' + qx = f$$

elimiere:  $u_0 = x_\Delta(a) = \gamma$

$$u_N = x_\Delta(b) = \delta$$

$$\tilde{A} = \begin{pmatrix} -2/h^2 + q_1 & 1/h^2 & .. & \\ 1/h^2 & -2/h^2 + q_2 & 1/h^2 & .. \end{pmatrix}$$

$$\tilde{y} = \begin{pmatrix} f_1 - \gamma/h^2 \\ f_2 \\ .. \\ f_{n-2} \\ f_{N-1} - 2/h^2 \end{pmatrix}$$

$$\tilde{A}\tilde{u} = \tilde{y}$$

$$\tilde{u} = \begin{pmatrix} u_1 \\ .. \\ u_{n-1} \end{pmatrix}$$

$\tilde{A}$  symmetrisch

$B = -\tilde{A}$  hat nur nicht positive Einträge ausserhalb der Diagonalen

$B$  hat für  $h$  klein genug nur positive Einträge auf der Diagonalen

(insbesondere ist  $b_{ii} > 2/h^2$  falls  $q \leq 0$ )

$$b_{ii} \geq 2/h^2 - 1/h^2 + 1/h^2 + 0$$

$$= |a_{i,i-1}| + |a_{i,i+1}| + \sum_{|j-i|>1} |a_{ij}|$$

$$= \sum_{i \neq j} |a_{ij}| = \sum_{j \neq i} |b_{ij}|$$

Definition Eine Matrix  $B \in R^{n \times n}$  heißt diagonaldominant, falls  $b_{ii} \geq \sum_{j \neq i} |b_{ij}|$

(Zeilendiagonaldominant)

Wir wollen garantieren, dass  $A$  invertierbar ist, dann gilt  $\tilde{A}\tilde{u} = \tilde{f} \Rightarrow \tilde{u} = \tilde{A}^{-1}\tilde{f}$  eindeutig + Stabilität

$$\|u\| \leq C\|f\| \text{ mit } C \text{ unabhängig von } h.$$

Wir wissen für  $Au = f$  mit  $A$  regulär gilt:  $\|u\| \leq \underbrace{\|\tilde{A}^{-1}\|}_{\text{zugeordnete Matrixnorm zu } \|f\|} \|f\|$

„Einfachste“ Wahl: euklidische Norm in  $R^n$

zugeordnete Matrixnorm ist Spektralnorm

Alternative: Supremumsnorm:  $\|u\|_\infty = \max_i |u_i|$

zugehörige Matrixnorm ist die Zeilensummennorm

passt zu strukturellen Eigenschaften von  $A$  bzw.  $B$

08.12.2009

$$(px')x' + qx = f, \in (a, b)$$

$$B_1x = \gamma, B_2x = \delta$$

→ Diskretisierung mit Differenzenquotienten auf einem äquisistanten Gleichungssystem für die Werte aus den Gitterpunkten  $\{u_j\}$ .

Stabilität und Konvergenz

Für Konvergenz:

$$Au = y$$

$$(px')' + qx = f$$

$$v = (x(s_j))_{j=0,..,N}$$

$$r = Av - y = O(h)$$

$$A(u - v) = y - (r + y) = -r = O(h^p)$$

$$e = u - v = -A^{-1}r$$

$$\|e\|_\infty = \|u - v\|_\infty \leq \|A^{-1}\|_{ZS} \underbrace{\|r\|_\infty}_{O(h^p)} \leq Ch^p$$

$A^{-1} \in R^{N \times N} \rightarrow \|A^{-1}\|_{ZS}$  abhängig von  $N \rightarrow$  abhängig von  $h$ .

Deshalb benötigen wir Stabilität in folgenden Sinn.

$$\|A^{-1}\|_{ZS} \leq C, \forall N \geq N_0$$

mit  $C$  unabhängig von  $N$  (bzw.  $h = 1/N$ )

Wenn Stabilität erfüllt ist, folgt sofort dass Konvergenzordnung = Konsistenzordnung

Stabilität für das RWP in Supremumsnorm

Vereinfacht:  $(px')' = f \in (a, b)$

$$x(a) = \gamma, x(b) = \delta$$

#### 4.5.1 Maximumsprinzip

Bsp:  $x'' = f$

$f > 0 \in (a, b) \rightarrow x'' > 0$ , x ist konvex

konvexe Funktion nimmt im Maximum, Rand an!

Also:  $x(s) \leq \max \gamma, \delta$

$f < 0 \in (a, b) \rightarrow x'' < 0 \rightarrow x$  konkav

$\rightarrow x$  nimmt Minimum am Rand an.

$x(s) \geq \min \gamma, \delta$

Stabilität aus Maximumsprinzip:

$x'' = f$ , vgl. mit  $y'' = 1$

$x(a) = \gamma, x(b) = \delta$

$(x - cy)'' = f - c, y(a) = \gamma/c; y(b) = \delta/c$

$c \geq \|f\|_\infty, f - c \leq f - \max_\sigma |f(\sigma)| \leq 0$

Mit Randbedingungen  $(x - cy)(a) = 0, (x - cy)(b) = 0$

$(x - cy)'' \leq 0, \rightarrow x - cy$  nimmt Minimum am Rand an

$\Rightarrow (x - cy)(s) \geq 0, \forall s \in (a, b)$

$x \geq cy \geq \|f\|_\infty \min_\sigma y(\sigma)$

$$y'' = 1 \Rightarrow y' = s + a + c_1 \Rightarrow y = \frac{(s-a)^2}{2} + c_1(s-a) + c_2$$

$$\gamma/c = y(a) = c_2$$

$$\delta/c = y(b) = \frac{(b-a)^2}{2} + c_1(b-a) + \gamma/c$$

$$c_1 = \frac{1}{b-a} \frac{\delta-\gamma}{c} - 1/2(b-a)$$

$$y(s) = \frac{(b-a)^2}{2} + \frac{s-a}{b-a} \frac{\delta-\gamma}{c} - 1/2(s-a)(b-a) + \gamma/c$$

$$x(s) \geq c \min_c y(s) = 1/2 \underbrace{(s-a)(s-b)}_{\geq 1/4(b-a)^2} c + \underbrace{\frac{s-a}{b-a} (\delta-\gamma) + \gamma}_{|\cdot| \leq 1}$$

$$\geq 1/8(b-a)^2 c - |\delta| + \gamma$$

$$\geq -\max 1, 1/8(b-a)^2 \max \|f\|_\infty, |\gamma|, |\delta|$$

$$x(s) \geq -\eta \max \|f\|_\infty, \gamma, \delta$$

jetzt:  $c \leq -\|f\|_\infty$

$$(x - cy)'' = f - c \geq -\|f\|_\infty + \|f\|_\infty \geq 0$$

$$(x - cy)(a) = 0, (x - cy)(b) = 0$$

$$\Rightarrow x - cy \leq 0 \in (a, b)$$

$$x \leq |c| \|y\|_\infty \leq \eta \max \|f\|_\infty, |\gamma|, |\delta|$$

$$\Rightarrow \|x\|_\infty = \max_{s \in [a, b]} |x(s)| \leq \eta \max \|f\|_\infty, |\gamma|, |\delta|$$

#### 4.6 Stabilitätsabschätzung für RWP

Supremumsnorm der Lösung  $\leq \text{const} \times$  Supremumsnorm der rechten Seite

Allgemeines Prinzip:

1) Maximumsprinzip für Gleichung  $Lx = f$

2) Konstruieren einfache Lösungen y

mit  $Ly \geq f ((Ly \geq 1) \rightarrow L(cy) \geq f)$

$L(x+y) \geq 0 \Rightarrow \|x\|_\infty \leq \|y\|_\infty \leq \eta \|f\|_\infty$

Allgemeinere Beweise für Maximumsprinzip

(1)  $(px')' = f \rightarrow px' = \int_a^s f d\sigma + c_1$

$x(s) = \int_a^s \frac{1}{p(\sigma)} (\int_a^\sigma f(\eta) d\eta + c_1) d\sigma + c_2; c_1, c_2 \in RB$

Am Ende:

$x(s) = \int_a^b k(s, \sigma) f(\sigma) d\sigma + g(s)$

Vorzeichen von k und g entscheidend!

$\gamma, \delta > 0 \Rightarrow g(s) > 0$

$\gamma, \delta < 0 \Rightarrow g(s) < 0$

$k(s, \sigma) \leq 0$

$\gamma, \delta > 0, f < 0$

$$x(s) = \underbrace{\int_{\leq 0} k(s, \sigma) f(\sigma) d\sigma}_{\geq 0} + \underbrace{f(s)}_{\geq 0} > 0$$

2.) Widerspruchsbeweis:

$(px')' > 0$ , Annahme: Funktion hat ein Maximum in  $(a, b)$

$$x(s_0) \geq x(s) \forall s \in [a, b]$$

$$x'(s_0) = 0, x''(s_0) \leq 0$$

$$0 \leq (px')'(s_0) = p'(s_0) \underbrace{x'(s_0)}_{=0} + p(s_0) \underbrace{x''(s_0)}_{\leq 0} = \leq 0$$

$\Rightarrow$  Widerspruch

$\Rightarrow x$  hat kein Maximum in  $(a, b)$  (starkes Maximumsprinzip)

Schwaches Maximumsprinzip  $(px')' \geq 0$

Dann nimmt  $x$  sein Maximum am Rand an (also  $\{a, b\}$ )

Beweis: über Störungen:  $x_n \rightarrow x$  gleichmäßig mit  $(px'_n)' > 0$

$\Rightarrow X_n$  hat Maximum am Rand

$\Rightarrow x$  hat Maximum am Rand

Wann gilt Max-Prinzip für

$$Au = y, A \in R^n \times R^m m < \in R^n$$

Monotonie:  $y_k \geq 0 \forall j \Rightarrow u_j \geq 0 \forall j$

$$y = e_k, u = A^{-1}y = A^{-1}e_k = k - te \text{ Spalte von } A^{-1}.$$

$k$ -ten Spalte von  $A^{-1}$  hat nur nichtnegative Einträge

Definition:  $A \in R^{n \times n}$  heisst M-Matrix, falls  $A$  regulär und  $A^{-1}$  nur nichtnegative Einträge hat.

---

11.12.2009

$$(px')' = f$$

Diskretisierung FD

$$Au = f, u_j = x_\Delta(s_j)$$

Maximumsprinzip  $\Rightarrow L^\infty$  Stb.

$\Leftrightarrow A$  M-Matrix

Erinnerung:  $A$  ist M-Matrix, wenn  $A^{-1}$  existiert und nur nichtnegative Einträge hat.

Hinreichendes Kriterium für M-Matrix.

Lemma:  $A \in R^{n \times n}$  mit  $a_{ii} > 0 \forall i$

$$a_{ij} \leq 0 \forall i > j$$

$$\text{schwachdiagonaldominant } a_{ii} \geq \sum_{j \neq i} (a_{ij})$$

$A$  regulär

Dann ist  $A$  eine M-Matrix

Beweis:  $u_k = A^{-1}e_k = k$ -te Spalte von  $A^{-1}$ . Müssen also zeigen  $u_k \geq 0 \forall k$ .

Widerspruchsbeweis:  $\exists k : u_k < 0$

D.h. es existiert ein Index  $m$  mit  $(u_k)_m < 0$

o.B.d.A:  $(u_k)_m \leq (U_k)_j \forall j$

1.Fall:

$$(u_k)_m = (u_k)_j \text{ für alle } j$$

$$= c$$

$$u_k = A^{-1}e_k \Leftrightarrow Au_k = e_k$$

$$\Rightarrow \sum_j a_{ij} \underbrace{(u_k)_j}_{=c<0} = (e_k)_i$$

$$i = k : c \sum_{j=1}^n a_{kj} = 1$$

$$c(a_{kk} + \sum_{j \neq k} a_{kj}) = 1$$

$$\Rightarrow c(a_{kk} - \underbrace{\sum_{j \neq k} |a_{kj}|}_{\geq 0}) = 1 \text{ Widerspruch}$$

2.Fall: es existiert:

$$\begin{aligned}
 & j \text{ mit } (a_i)_j > (a_k)_m \\
 & \sum a_{ij} (a_k)_j = (e_k)_j \\
 & i = m : \sum a_{mj} (u_k)_j = (e_k)_m \geq 0 \\
 & = a_{mm} (u_k)_m + \sum_{j \neq m} a_{mj} (a_k)_j \\
 & = a_{mm} (u_k)_m - \sum_{j \neq m} |a_{mj}| (u_k)_j \\
 & < -(u_k)_j \leq -(u_k)_m a_{mm} (u_k)_m - \sum_{j \neq m} |a_{mj}| (u_k)_m \\
 & 0 \leq (u_k)_m (a_{mm} - \sum_{j \neq m} |a_{mj}|) \leq 0
 \end{aligned}$$

Widerspruch

$$A \text{ mit } a_{ii} > 0, a_{ij} \equiv 0, \sum_{i \neq j} |a_{ij}| \leq a_{ij}$$

$$N(A) \subset \text{lin}\{1\}$$

d.h.  $Au = 0 \Rightarrow u_j = c$  für ein  $c \in R$

Korollar: A wie oben ist regulär, wenn  $A1 \neq 0$

$$\Leftrightarrow \exists i : \sum_{j=1}^n a_{ij} \neq 0$$

$$\underbrace{a_{ij} - \sum_{i \neq j} |a_{ij}|}_{}$$

Beweis:

$$Au = 0, U_m = \max_j u_j, \text{ obda: } u_m > 0$$

$$\sum_{i=1}^n a_{mj} u_j = 0$$

$$0 = a_{mm} u_m - \sum_{j \neq m} |a_{mj}| u_j > (a_{mm} - \sum |a_{mj}|) u_m \geq 0$$

Widerspruch, falls  $u_j \neq u_m$  für ein j

$$\Rightarrow u_j = u_m \forall j$$

$$\Rightarrow u \in \lim 1$$

Analoge Eigenschaft für RWP

$$(px')' = f, B_1x = \gamma, B_2x = \delta$$

$$\text{Nullraum: } (px')' = 0, B_1x = 0, B_2x = 0$$

$$\text{Falls } B_1x = x'(a) = 0, B_2(x) = x'(b) = 0$$

$\Rightarrow$  Konstantes x im Nullraum

$$\text{Konvergenz: } Lx = f, Lx = (px')'$$

Diskretisierung auf äquidistanten Gitter  $\rightarrow X_\Delta$

$$u = (x_\Delta(s_j)), v = (x(s_j))$$

wollen Abschätzung für  $|x_\Delta(s_j) - x(s_j)| \rightarrow$  d.h. für  $\|u - v\|_\infty$

Am besten  $\|u - v\|_\infty \leq C \|v\|_\infty$

Konsistenzfehler  $r = Av - y$

Diskretisierung  $Au = y$

Achtung:  $Au = y$  ist diskrete Version von  $-(px')' = -f$

Für Stabilitätskonstante C um  $\tilde{x}$  mit  $L\tilde{x} = -1, B_1\tilde{x} = B_1x, B_2\tilde{x} = B_2x$

$$\tilde{v} = (\tilde{x}(s_j))$$

Konsistenzordnung:  $f = -1 \Rightarrow \tilde{y} = (1, \dots, 1)^t$

$$= A\tilde{v} - 1 = O(H^p)$$

$$Av = y, Au = y + r, A(c\tilde{v}) = c1 + cO(h^p)$$

$$A(u - v - c\tilde{v}) = -r - c1 + cO(h^p)$$

$$\text{c so gewählt, dass } \begin{cases} -r - c1 + cO(h^p) \leq 0 \Rightarrow u - v \leq c\tilde{v} \\ -r - c1 + cO(h^p) \geq 0 \Rightarrow u - v \geq c\tilde{v} \end{cases}$$

$$(i) c = 2\|r\|_\infty$$

$$A(u - v - c\tilde{v}) = \underbrace{-r - \|r\|_\infty 1}_{\leq 0} - \underbrace{\|r\|_\infty 1}_{\geq 0} + \underbrace{2\|r\|_\infty O(h^p)}_{\geq 0} \leq 0$$

$$A(u - v - c\tilde{v}) \leq 0$$

$$\Rightarrow u - v - c\tilde{v} \leq 0$$

$$u - v \leq 2\|r\|_\infty \tilde{v}$$

$$\begin{aligned}
 \text{(ii)} \quad & c = -2\|r\|_\infty \\
 \Rightarrow & A(u - v - c\tilde{v}) \geq 0 \\
 \Rightarrow & u - v \geq c\tilde{v} = -2\|r\|_\infty \tilde{v} \\
 \Rightarrow & -2\|r\|_\infty \tilde{v}_j \leq u_j - v_j \leq 2\|r\|_\infty \tilde{v}_j \\
 \Rightarrow & |u_j - v_j| \leq 2\|r\|_\infty |\tilde{v}_j| \\
 \Rightarrow & \|u - v\|_\infty \leq \underbrace{2\|\tilde{v}\|_\infty}_{=C} \|r\|_\infty
 \end{aligned}$$

$\tilde{v}$  unabhängig von  $h \Rightarrow$  konstante  $C$  ist unabhängig von  $h$ .

Satz:  $A \in R^{n \times n}$  Matrix und  $y \in R^n$  rechte Seite aus der Diskretisierung einer s RWP. Ist  $A$  M-Matrix für alle  $n$ , dann ist dies Verfahren stabil un die Konvergenzordnung entspricht der Konsistenzordnung.

Bsp:  $x'' = f, x(a) = 0, x(b) = 0$

$$-x'' = -f$$

$$D_2 x_\Delta = -f_\Delta = \frac{x_\Delta(s_{j+2}) - 2x_\Delta(s_j) + x_\Delta(s_{j-1})}{h^2}$$

$n = N - 1, N = -1/h$  (elimieren  $x_\Delta(s_0)$  und  $x_\Delta(s_N)$ )

$$\frac{1}{h^2} \begin{pmatrix} 2 & -1 & .. & .. \\ -1 & 2 & -1 & .. \\ .. & -1 & 2 & -1 \\ .. & .. & -1 & 2 \end{pmatrix} \begin{pmatrix} u_1 \\ .. \\ u_n \end{pmatrix} = \begin{pmatrix} -f(s_1) \\ .. \\ -s(S_{N-1}) \end{pmatrix}$$

$A$  ist M-Matrix

15.12.2009

## 4.7 Finite Differenzen

$$Lx = -(px')' + qx = f$$

Diskretisierung

$$Au = y$$

Miximumsprinzip für  $L$  ( $L^\infty$ -Stabilität)  $\rightarrow$  M-Matrix  $A$  ( $l^\infty$ -Stabilität)

Konsistente, monotone ( $A$  M-Matrix)

Diskretisierung  $\Rightarrow$  Konvergenz

Konvergenzordnung = Konsistenzordnung

Wie sieht  $A$  für allgemeines  $p, q$  aus?

$$\underbrace{px''}_{D_2} + \underbrace{p'x'}_{D_c} + qx = f$$

$$s = s_i : p_i 1/h^2 (x_\Delta(s_{i-1}) - 2x_\Delta(s_i) + x_\Delta(s_{i+1})) + p'_i 1/(2h) (x_\Delta(s_{i+1}) - x_\Delta(s_{i-1})) + q_i x_\Delta(s_i) = f_i$$

$$(p_i/h^2 + P'_i/2h) u_{i+1} - (2p_i/h^2 - q_i) u_i + (p_i/h^2 - p'_i/(2h)) u_{i-1} = f_i$$

$$\left( \begin{matrix} 0 & .. & 0 & p_i/h^2 + P'_i/2h & -2p_i/h^2 + q_i & (p_i/h^2 - p'_i/(2h)) & 0 & .. & 0 \end{matrix} \right) \begin{pmatrix} u_1 \\ .. \\ u_n \end{pmatrix} = \begin{pmatrix} -f_i \\ .. \\ -f_i \end{pmatrix}$$

## 4.8 M-Matrix Eigenschaft

$$-q_{ii} > 0 : 2p_i > q_i h^2 \begin{cases} q_i \leq 0 & \text{immer erfüllt} \\ q_i > 0 & \text{für } h \text{ klein genug} \end{cases}$$

$$a_{ij} \leq 0 : 2p_i + hp'_i \leq 0, -2p_i - hp'_i \leq 0 : hp'_i \leq 2p_i \rightarrow p'_i/p_i \leq 2/h; hp'_i \geq -2p_i \rightarrow p'_i/p_i \geq -2/h$$

$$-2/h \leq (\log p)'(s_i) \leq 2/h$$

$$|\log p'| \leq 2/h \rightarrow h \leq 2/(\log p)'$$

Problem, wenn  $p(\log(p))$  sehr steil wird  $\rightarrow$  Extrem kleine Schrittweite nötig, um Stabilität zu bekommen.

$$- \text{Diagonaldominanz } a_{ii} \geq \sum_{j \neq i} |a_{ij}| = |a_{i,i+1}| + |a_{i,i-1}|$$

$$\Rightarrow 2p_i/h^2 - q_i \geq \frac{2p_i - p'_i h}{2h^2} + \frac{2p_i + p'_i h}{2h^2} = 2p_i/h^2$$

$$\rightarrow q_i \leq 0$$

Für  $q \leq 0$  (RB eindeutig lösbar) erhalten wir eine M-Matrix, wenn  $h$  klein genug ist.

Für  $q > 0$  analoge Störungstheorie wie für das RWP.

$$A = \underline{B}_{M-Matrix} + \underline{D}_{Diagonalmatrix}$$

Alternative für  $\|(log p)'\|_\infty$  gross: Verfahren der Konsistenzordnung 1

## 4.9 Upwind-Verfahren:

Statt  $D_c$  für Diskretisierung von  $x'$  wählen wir  $D_+$  oder  $D_-$  abhängig vom Vorzeichen von  $p'$ .

$$px'' + p'x' + qx = f$$

$$D_+ : p_i D_2 x_\Delta(s_i) + p'_i 1/h(x_\Delta(s_{i+1}) - x_\Delta(s_i)) + q_i x_\Delta(s_i) = f_i$$

$$D_- : p_i D_2 x_\Delta(s_i) + p'_i 1/h(x_\Delta(s_i) - x_\Delta(s_{i-1})) + q_i x_\Delta(s_i) = f_i$$

Wenn  $p'_i \geq 0$ , dann Vorwärtsdifferenzenquotienten  $D_+$  verwenden

Wenn  $p'_i \leq 0$ , dann Rückwärtsdifferenzenquotienten  $D_-$  verwenden

$$D_u^{p'}(x) = \max(p'_i, 0) D_+ x + \min(p'_i, 0) D_- x$$

$$A = \begin{pmatrix} 0 & \frac{-p_i + \max(p'_i, 0)}{h^2} & 2p_i/h^2 - q_i + |p'_i|/h & \frac{-p_i - \max(p'_i, 0)}{h^2} & 0 \end{pmatrix} \text{ M-Matrix für } q \leq 0 \text{ für } h \text{ beliebig.}$$

Besonders wichtig:  $p$  nicht bekannt ist, sondern indirekt von  $x$  abhängt.

Beispiel: geladene Teilchen in elektrischen Feld,  $x$  Teilchendichte,  $E$ : elektrisches Feld,  $z$  Ladung:  
 $x'' + z(Ex') = 0$ .

$$E' = zx, E \rightarrow p$$

Ausblick: Lösung partieller DGL

Bsp: 2D,  $u(x,y)$

Poisson-Gleichung:  $-\frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial y^2} = f, \in \Omega \in R^2$  mit RB: z.B.  $u(x,y) = f(x,y)$  für  $(x,y) \in \partial\Omega$

Wärmeleitungsgleichung: in einer Raumdimension und in der Zeit  $u(x,t)$

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}, x \in (0,t), t \in (0,T)$$

RB z.B.:  $u(a,t) = \gamma, u(b,t) = \delta, u(x,0) = u_0(x)$

## 4.10 Anfangswertproblem

### 4.10.1 Numerische Lösung:

Gitter:  $\Delta = (x_i, y_i) \in \Omega$

Approximation von  $\partial\Omega$ :

$\partial\Delta = (x_i, y_j) \in \Omega | (x_{i+1}, y_j) \notin \Omega, \text{ oder } (x_{i-1}, y_j) \notin \Omega, \text{ oder } (x_i, y_{j+1}) \notin \Omega, \text{ oder } (x_i, y_{j-1}) \notin \Omega$

Dort approximieren wir Randwerte, d.h.  $u(x_i, y_j) - f(x_i, y_j)$  für  $(x_i, y_j) \in \partial\Delta$ .

Für  $(x_i, y_i) \notin \partial\Delta$  approximieren wir die partielle DGL, z.B. Poisson-Gleichung auf äquidistantem Gitter  $x_j = jh_1, y_j = jh_2$

$$D_2^x u - D_2^y = f$$

$$-1/h_1^2(u(x_{i+1}, y_j) - 2u(x_i, y_j) + u(x_{i-1}, y_j)) - 1/h_2^2(u(x_i, y_{j+1}) - 2u(x_i, y_j) + u(x_i, y_{j-1})) = f_i$$

Lineares Gleichungssystem für Werte  $u(x_i, y_j)$

Vektor  $U = u(x_i, y_j)$

$$AU = y$$

$A$  ist M-Matrix

Zeile für  $(x_i, y_j)$ : Diagonaleintrag  $2/h_1^2 + 2/h_2^2 > 0$

4 Nebendiagonaleinträge  $\neq 0$ :  $-1/h_1^2(2x), -1/h_2^2(2x)$

Wärmeleitung: Liniennmethode: Diskretisierung zuerst im Ort auf Gitter  $x_j = a + jh$

$$\frac{\partial u}{\partial t}(x_i, t) = 1/h^2(u(x_{i+1}, t) - 2u(x_i, t) + u(x_{i-1}, t))$$

$$\rightarrow u'_i(t) = 1/h^2(u_{i+1}(t) - 2u_i(t) + u_{i-1}(t))$$

System von gewöhnlichen DGL, AWP

$$u'(t) = f(u(t), t)$$

Lipschitzkonstante über  $f \sim 1/h^2$

explizites Eulerverfahren hat Stabilitätskonstante  $\approx e^{L/h^2}$

Zeitdiskretisierung mit implizitem Eulerverfahren oder Crank-Nicholson

18.12.2009

---

## 4.11 Freie Randwertprobleme

Randwertproblem, bei denen ein Teil des Randes unbekannt ist, zusätzliche Gleichungen zur Bestimmung des Randes. Bsp:  $-(px')' = 0 \text{ fin}(a, b)$

$$x(a) = 0, x(b) = 0$$

Eindeutig lösbar für gegebenes a,b

Freier Rand, z.B. b frei

Typsischerweise zwei k Randbedingungen

$$\text{z.B. } x'(b) = 0$$

Beispiel Übergang Eis-Wasser

Unbekannte Funktion  $\mu(x, t)$  Temperatur

Wie stellt sich T ein?

(Wie verändert sich T in der Zeit?)

Wärmeleitungsgleichung:  $\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}$  gilt in Eis und Wasser.

Auf  $\Gamma : u = 0$

zweite Randedingung?

In jeder Phase ist Energie abhängig von der Temperatur, Enthalpie (Energiedichte)

$$e = \begin{cases} a_1 u + b_1 & \text{in Eis} \\ a_2 u + b_2 & \text{in Wasser} \end{cases}$$

Vereinfachung:  $a_1 = a_2 = 1, b_1 - b_2 = L + 0$

Eigentlich: Wärmeleitung durch Energietransport:

$$\frac{\partial e}{\partial t} = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}$$

$$\Xi(x, t) = \begin{pmatrix} 1 & \text{Eis} \\ 0 & \text{Wasser} \end{pmatrix}$$

$$e(x, t) = u(x, t) + b_2 + L\Xi(x, t) = \begin{cases} u + b_1 & \text{Eis} \\ u + b_2 & \text{Wasser} \end{cases}$$

$$\frac{\partial u}{\partial t} + L \frac{\partial \Xi}{\partial t} = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}$$

Im Inneren von Eis/Wasser:  $\frac{\partial \Xi}{\partial t} = 0 \Rightarrow$  Wärmeleitungsgleichung

$$1D: \frac{\partial u}{\partial t} + L \frac{\partial x}{\partial t} = \frac{\partial^2 u}{\partial x^2}$$

$$\text{Betrachten: } \frac{\partial u}{\partial t} + L \frac{\partial \Xi}{\partial t} = \frac{\partial^2 u}{\partial x^2} \text{ in } (\xi(t) + \epsilon, \xi(t) - \epsilon)$$

$$\text{Integration: } \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} \frac{\partial u}{\partial t} dx + L \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} \frac{\partial \Xi}{\partial t} dx = \frac{\partial u}{\partial x}(\xi(t) + \epsilon, t) - \frac{\partial u}{\partial x}(\xi(t) - \epsilon, t)$$

Was ist

$$\frac{d}{dt} \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} f(x, t) dx$$

$$= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} (\int_{\xi(t)+\Delta t-\epsilon}^{\xi(t)+\Delta t+\epsilon} f(x, t + \Delta t) dx - \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} f(x, t) dx)$$

$$= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} (\int_{\xi(t)+\Delta t-\epsilon}^{\xi(t)+\Delta t+\epsilon} f(x, t + \Delta t) - f(x, t) dx + \int_{\xi(t)+\Delta t-\epsilon}^{\xi(t)+\Delta t+\epsilon} f(x, t) dx - \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} f(x, t) dx)$$

$$= \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} \frac{\partial f}{\partial t} f(x, t) dx + \lim_{\Delta t \rightarrow 0} (1/\Delta t \int_{\xi(t)+\Delta t-\epsilon}^{\xi(t)+\Delta t+\epsilon} f(x, t) dx - 1/\Delta t \int_{\xi(t)+\Delta t-\epsilon}^{\xi(t)+\epsilon} f(x, t) dx)$$

$$\frac{\partial F}{\partial x} = f : 1/\Delta t \int_{\xi(t)-\epsilon}^{\xi(t)+\Delta t+\epsilon} f(x, t) dx = \frac{F(\xi(t)+\Delta t+\epsilon, t) - F(\xi(t)+\epsilon, t)}{\delta t}$$

$$\xrightarrow{\Delta t \rightarrow 0} f(\xi(t) + \epsilon, t) \frac{d\xi}{dt}(t)$$

$$\text{Analog } 1/\Delta t \int_{\xi(t)-\epsilon}^{\xi(t)+\Delta t-\epsilon} f(x, t) dx = f(\xi(t) - \epsilon, t) \frac{d\xi}{dt}(t)$$

$$\frac{d}{dt} \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} f(x, t) dx = \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} \frac{\partial f}{\partial t} (x, t) dx + f(\xi(t) + \epsilon) \frac{d\xi}{dt}(t) - f(\xi(t) - \epsilon) \frac{d\xi}{dt}$$

$$f = \xi : \int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} \frac{\partial \Xi}{\partial t} dx = \frac{d}{dx} (\int_{\xi(t)-\epsilon}^{\xi(t)+\epsilon} \Xi(x, t) dx) - \Xi(\xi(t) + \epsilon) \frac{d\xi}{dt} + \Xi(\xi(t) - \epsilon, t) \frac{d\xi}{dt} = d/(dt) \int 1 dx + \frac{d\xi}{dt} = \frac{d\xi}{dt}(t)$$

$$f = u : \int_{\xi-\epsilon}^{\xi+\epsilon} \frac{\partial u}{\partial t} dx = \frac{d}{dt} (\int_{\xi-\epsilon}^{\xi+\epsilon} u(x, t) dx) - u(\xi + \epsilon, t) \frac{d\xi}{dt} + u(\xi - \epsilon) \frac{d\xi}{dt}$$

Achtung:  $u(\xi(t), t) = 0$

$$(t) \text{ für } \epsilon \rightarrow 0 : L \frac{d\xi}{dt}(t) = \frac{\partial u}{\partial x}(\xi(t), t) - \frac{\partial u}{\partial x}(\xi(t), t)$$

Stationäre Lösung:  $\frac{d\xi}{dt} = 0$

$$\frac{\partial u}{\partial t} = 0$$

$$\frac{d\xi}{dt} = 0 \Rightarrow \frac{du}{dx}(\xi_+) = \frac{du}{dx}(\xi_i)$$

$$\frac{d^2 u}{dx^2} = 0 \in (a, \xi) \text{ und } (\xi, b)$$

$$u(x) = c(x - \xi) \text{ in } (a, \xi) \text{ und } (\xi, b)$$

zwei unbekannte, c und  $\xi \Rightarrow$  zu bestimmen aus Randwerten  $u(a) = c(a - \xi), u(b) = c(b - \xi)$

$$c = \frac{u(b) - u(a)}{b - a}$$

$$\xi = a - \frac{u(a)}{c}$$

Energieminimierung unter  $RBu(a) = u(b) = 0$

$$E(u) = 1/2 \int_a^b u'(x)^2 dx$$

Optimierung unter Nebenbedingungen  $u(x) \geq h(x)$

$u$  minimal

$$E(u) \leq E(u + \epsilon\phi), \text{ falls } u + \epsilon\phi \geq h$$

$\phi \geq 0$ , falls  $u = h$

$$\phi(a) = \phi(b) = 0$$

$$0 \leq 1/\epsilon(E(u + \epsilon\phi) - E(u)) = 1/\epsilon \int ((u' + \epsilon\phi)^2 - (u')^2) dx = 2 \int -u'\phi' dx + \epsilon \int (\phi')^2$$

$$0 \leq -2 \int u''\phi dx \rightarrow \int u''\phi dx \leq 0$$

lokale Störung bei  $u > h\epsilon$  klein:  $u + \epsilon\phi \geq h, u - \epsilon\phi \geq h$

$$\int u''\phi \leq 0, -\int u''\phi \leq 0 \Rightarrow \int u''\phi = 0 \rightarrow u'' = 0, \text{ wenn } u(x) > h(x)$$

$u = h \rightarrow$  Störung  $\phi \geq 0$  zulässig ( $u + \epsilon\phi \geq 0$ )

$$\int u''\phi \leq 0 \rightarrow u'' \leq 0 \text{ wenn } u(x) = h(x)$$

Hindernisproblem:  $u \geq h$  in  $(a, b)$

$$u'' \leq 0 \text{ in } (a, b)$$

$$u(a) = u(b) = 0$$

$$u''(u - h) = 0 \text{ in } (a, b)$$

---

05.01.2010

---

## 5 Buch Teil 1, Kapitel 8

Grosse Symmetrische lineare Gleichungssysteme

$$Ax = b, A \in R^{n \times n}, b \in R^n$$

A dünnbesetzt (Sparse), d.h. Anzahl der Nichtnullelementen in  $A \ll n^2$  (Typischerweise  $O(n)$ )

n gross

Zerlegungen (LR, QR, Cholesky) haben hohen Aufwand  $O(n^3)$ , nutzen Dünnbesetztheit nicht aus.

Eventuell hoher Speicherbedarf  $O(n^2)$

z.B.  $A = QR$

A dünnbesetzt, trotzdem Anzahl der Nichtnullelemente in Q von der Ordnung  $n^2$ .

Direkte Verfahren  $x = A^{-1}b$  „exakt“ bis auf Fehlerfortpflanzung

Relative Genauigkeit =  $K(A) * \text{Maschinengenauigkeit}$

Oft ist aber Konditionszahl von A abhängig von A bzw. sehr gross für n gross

Bei Diskretisierung eines RWP  $(px')' + qx = f$  mit n Gitterpunkten entsteht Matrix A mit  $K(A) = O(n^2)$ .

$n \approx 10^4 \rightarrow n^2$  Maschinengenauigkeit  $\approx 10^{-8}$

D.h. wir wollen nur eine Lösung mit Genauigkeit  $\epsilon \gg \text{Maschinengenauigkeit}$

Alternative Iterative Verfahren

In jedem Iterationsschritt: Multiplikation  $Ax$  oder  $A^T x \rightarrow O(n)$  für dünnbesetzte Matrizen

Speicherbedarf  $O(n)$ , nur  $A, b, x^k$

Aufwand  $O(k_{max} * n)$ ,  $K_{max}$  = Anzahl der Iterationen

Umso grösser  $\epsilon$ , umso weniger Iterationen  $k_{max}$  nötig. In der Praxis meist  $k_{max} \ll n^2$

### 5.1 8.1 Klassische Iterationsverfahren

Fixpunktiteration

$$Ax = b \Leftrightarrow 0 = Q^{-1}(b - Ax)$$

$$\Leftrightarrow x = x + Q^{-1}(b - Ax) = \underbrace{(I - Q^{-1}A)}_{=G} x + \underbrace{Q^{-1}b}_{=c}$$

$$x = Gx + c$$

$$x^{k+1} = Gx^k + c - F(x^k)$$

Banach'scher FPS:  $F: R^n \rightarrow R^n$ ,  $F$  kontrahiert  $\Rightarrow x^k \rightarrow \bar{x}$  Fixpunkt + Fehlerabschätzung  
 $\|x^k - \bar{x}\| \leq L^k c$   
 $L < 1$  Lipschitzkonstante  
 $\|F(x) - F(y)\| = \|Gx + c - Gy - c\| = G(x - y)\| \leq \|G\| \|x - y\|$   
 Kontraktivität, falls  $\|G\| < 1$

## 5.2 Satz 8.0:

$G \in R^{n \times n}$ ,  $\|G\| < 1$  in einer Matrixnorm.

Dann konvergiert die Iteration  $x^{k+1} = Gx^k + c$  für jeden Startwert  $x^0 \in R^n$  gegen den eindeutigen Fixpunkt  $\bar{x} = G\bar{x} + c$

## 5.3 Satz 8.1:

$G \in R^{n \times n}$ ,  $\rho(G) := \max_j |\lambda_j(G)| < 1$ .

Dann konvergiert die Iteration  $x^{k+1} = Gx^k + c$  gegen den eindeutigen Fixpunkt  $\bar{x} = G\bar{x} + c$ .

Beweis: (für  $G$  symmetrisch)

$G$  symmetrisch  $\rightarrow \exists Q$  orthogonal

$QGQ^T = \Lambda = \text{diag}(\lambda_i)$

$\lambda_i$  EW von  $G$

$x^{k+1} = Gx^k + c, \bar{x} = G\bar{x} + c$

$x^{k+1} - \bar{x} = G(x^k - \bar{x})$

$$\|x^{k+1} - \bar{x}\| = \|G(x^k - \bar{x})\| = \|Q^T \Lambda Q(x^k - \bar{x})\| = \|Q^T \Lambda Q Q^T \Lambda Q(x^{k-1} - \bar{x})\|$$

$$\|x^{k+1} - \bar{x}\| = \|Q^T \Lambda^2 Q(x^{k-1} - \bar{x})\| = \|Q^T \Lambda^{k+1} Q(X^+ - \bar{x})\| = \|Q^T \Lambda^{k+1} y\|$$

$$= [Q^T \text{orthogonal}] \|\Lambda^{k+1} y\|$$

$$\|(diag(\lambda_j^{k+1}))y\| = \|\lambda_j^{k+1} y_j\| = \underbrace{\|\lambda_j^{k+1}\|}_{\rho(G)^{k+1}} |y_j| \leq \underbrace{\rho(G)^{k+1}}_{\rightarrow 0 \text{ da } \rho(G) < 1} \|y\|$$

$$\Rightarrow \|x^{k+1} - \bar{x}\| \leq \rho(G)^{k+1} \|y\|$$

Es gilt:  $\rho(G) \leq \|G\|$  für jede Matrixnorm

$$\|G\| = \sup_{x \neq 0} \frac{\|Gx\|}{\|x\|} \geq \frac{\|Gx_j\|}{\|x_j\|} = \frac{\|\lambda_j x_j\|}{\|x_j\|} = |\lambda_j| \frac{\|x_j\|}{\|x_j\|}$$

für  $x_j$  EV zum EW  $\lambda_j$ ,  $Gx_j = \lambda_j X_j$

$$\|G\| \geq \max_j |\lambda_j| = \rho(G)$$

## 5.4 Iterative Verfahren:

$$x^{k+1} = Gx^k + c = x^k + Q^{-1}(b - Ax^k)$$

### 5.4.1 RICHARDSON-Verfahren:

$$Q^{-1} = \tau I, \tau \in R$$

$$x^{k+1} = x^k - \tau A x^k + \tau b$$

- Konvergenz:  $G = I - \tau A$

$\rho(G) < 1$ ?

$\mu_j$  EW von  $A$  mit EV  $x_j$

$$\Leftrightarrow 1 - \tau \mu_j \text{ EW von } G, Ax_j = \mu_j x_j$$

$$x_j - \tau Ax_j = x_j - \tau \mu_j x_j$$

$$Gx_j = (I - \tau A)x_j = (1 - \tau \mu_j)x_j$$

Brauchen:  $|1 - \tau \mu_j| < 1$  für alle j

alle  $\mu_j > 0, \tau > 0$

$$-1 < 1 - \tau \mu_j < 1$$

$$\tau \mu_j < 2$$

$$\tau < \frac{2}{\mu_j}, \forall j \Leftrightarrow \tau < \frac{2}{\max_j \mu_j} = \frac{2}{\rho(A)}$$

$$\mu_j < 0, \tau < 0, \tau > \frac{2}{\min_j \mu_j} = -\frac{2}{\rho(A)}$$

$$Re(\mu_j) > 0, Re(\mu_k) < 0$$

$$|1 - \tau\mu_j| < 1, (1 - \tau Re(\mu_j))^2 + (\tau Im(\mu_j))^2 < 1$$

$$|1 - \tau\mu_k| < 1, (1 - \tau Re(\mu_j))^2 + (\tau Im(\mu_k))^2 < 1$$

$$\tau > 0, 1 - \tau Re\mu_k > 1$$

$$\tau < 0, 1 - \tau Re\mu_j > 1$$

$\Rightarrow$  Für jedes  $\tau$  gilt  $\rho(G) > 1 \rightarrow$  nicht konvergent!

Richardson-Verfahren konvergiert, wenn  $Re(\mu_j)$  und  $|\tau|$  klein genug ist und mit richtigem Vorzeichen  $sign(\tau) = sign(Re(\mu_j))$

Insbesondere für A symmetrische-positive-definit, Konvergenz für  $\tau \in (0, \frac{2}{\rho(A)})$

Konvergenzgeschwindigkeit

$$\|x^{k+1} - \bar{x}\| \leq \rho \|x^k - \bar{x}\|$$

$\rightarrow$  lineare Konvergenz

$$\rho(G) = \max_j |1 - \tau\mu_j| = \max |1 - \tau\lambda_{\max}(A)|, |1 - \tau\lambda_{\min}(A)|$$

K(A) gross,  $1 - \tau\lambda_{\min}(A) > 0$

$$\tau < \frac{2}{\rho(A)} = \frac{2}{\lambda_{\max}(A)}$$

$$\lambda_{\max}(A)\tau < 2$$

$$\frac{\lambda_{\max}}{\lambda_{\min}} \lambda_{\min}(A)\tau < 2$$

$$\Rightarrow \lambda_{\min}\lambda < \frac{2}{K_2(A)} < 1$$

$$|1 - \tau\lambda_{\min}| = 1 - \tau\lambda_{\min}$$

$$\tau < \frac{2}{\lambda_{\max}} > 1 - 2\frac{\lambda_{\min}}{\lambda_{\max}}$$

$$= 1 - \frac{2}{K_2(A)} = \frac{K_2(A)-2}{K_2(A)}$$

$$\rho(G) \geq 1 - \tau\lambda_{\min} > \frac{K_2(A)-2}{K_2(A)} \approx 1 \text{ für } K_2(A) \text{ gross}$$

#### 5.4.2 Jacobi-Verfahren:

$$Q = diag(a_{ii})$$

$$A = L + D + R$$

A symmetrisch  $\Rightarrow R = L^t$

D Diagonalmatrix, L = linke untere Dreiecksmatrix

$$Q = D$$

$$x^{k+1} = x^k + D^{-1}(b - Ax^k)$$

$$= D^{-1}(Dx^k - Ax^k + b)$$

$$= D^{-1}(Dx^k - Dx^k - Lx^k - Rx^k + b)$$

$$x^{k+1} = D^{-1}(b - (L + R)x^k)$$

$$x^{k+1} = 1/a_{ii}(b_i - \sum_{j \neq i} a_{ij}x_j^k)$$

09.01.2010

---

$$Ax = b \Leftrightarrow x = x - Q^{-1}(Ax - b)$$

$$x^{k+1} = x^k - Q^{-1}(Ax^k - b)$$

$$= Gx^k + c$$

Konvergenz  $\|G\| < 1$

oder  $\rho(G) < 1$

Richardson-Iteration:  $Q^{-1} = \tau I$

Jacobi-Iteration:  $Q = D, A = L + R + D$

#### 5.5 Satz 8.2:

Das Jacobi-Verfahren konvergiert für jeden Startwert  $x_0 \in R^n$ , wenn A strikt diaognaldominant ist.

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|$$

$$\text{Beweis: } G = I - Q^{-1}A = I - D^{-1}(L + R + D)$$

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

$$\begin{aligned} \text{Sup-Norm in } R^n, \text{ Zeilensummennorm in } R^{n \times n} \|D^{-1}(L+R)\|_\infty &= \left\| \begin{pmatrix} 0 & a_{12}/a_{11} & a_{1n}/a_{11} \\ .. & .. & .. \\ a_{1n}/a_{11} & .. & a_{n-1,n-1} \\ 0 & 0 & 0 \end{pmatrix} \right\| = \\ &\max_j \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|} \\ &= \max_i 1/|a_{ii}| \sum_{j \neq i} |a_{ij}| < 1 \end{aligned}$$

Implementierung:

$$x_i^{k+1} = - \sum_{j \neq i} \frac{a_{ij}}{a_{ii}} x_j^k + \frac{b_i}{a_{ii}}$$

$$x_1^{k+1} = - \sum_{j \neq 1} \frac{a_{1j}}{a_{11}} x_j^k + \frac{b_1}{a_{11}}$$

$$x_2^{k+1} = - \sum_{j \neq 2} \frac{a_{2j}}{a_{22}} x_j^k + \frac{b_2}{a_{22}}$$

Alternative:

$$x_1^{k+1} = - \sum_{j \neq 1} \frac{a_{1j}}{a_{11}} x_j^k + \frac{b_1}{a_{11}}$$

$$x_2^{k+1} = - \frac{a_{21}}{a_{22}} x_1^{k+1} - \sum_{j > 2} \frac{a_{1j}}{a_{11}} x_j^k + \frac{b_2}{a_{11}}$$

..

$$x_i^{k+1} = - \sum_{j < i} \frac{a_{ij}}{a_{ii}} x_j^{k+1} - \sum_{j > i} \frac{a_{ij}}{a_{11}} x_j^k + \frac{b_i}{a_{11}}$$

(Gauss-Seidel-Verfahren)

$$a_{11}x_1^{k+1} = - \sum_{j > 1} a_{1j}x_j^k + b_1$$

$$a_{21}x_2^{k+1} + a_{22}x_2^{k+1} = - \sum_{j > 2} a_{2j}x_j^k + b_2$$

..

$$\sum_{j \leq i} a_{ij}x_j^{k+1} = - \sum_{j > i} a_{ij}x_j^k + b_j$$

Matrix-schreibweise:

$$(D + L)x^{k+1} = -Rx + b$$

$$x^{k+1} = -(D + L)^{-1}Rx^k + (D + L)^{-1}b$$

$$Q = D + L$$

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

Konvergenzbetrachtung in einem geänderten Skalarprodukt.

$$M \text{ spd: } (x, y)_M = y^T M x$$

$$\text{Transponierte Matrix: } y^T A x = (A^T y) x \forall x, y \in R^n$$

$$\text{Adjungierten Matrix: Skalarprodukt induziert durch } M: (Bx, y)_M = (x, B^*y)_M \forall x, y \in R^n$$

$B^*$  adjungierte Matrix

$B$  sei selbstadjungiert:  $B = B^*$

$B$  positiv in  $(., .)_M : (Bx, x)_M > 0 \forall x \in R^n$

Adjungierte Matrix:

$$(Bx, y)_M = y^T MBx - (y^T M B M^{-1}) M x = (B^*y)^T M x, B^* = M^{-1} B^T M$$

## 5.6 Lemma 8.3

$$G \in R^{n \times n}, G^* \text{ adjungierte Matrix in } (., .)_M$$

Falls  $B = I - G^*G$  positiv ist, dann gilt  $\rho(G) < 1$

Beweis:  $B$  positiv

$$0 < (Bx, x) = (x - G^*Gx, x) = (x, x) - \underbrace{(G^*Gx, x)}_{(Gx, Gx)} = \|x\|_M^2 - \|Gx\|_M^2$$

$$\Leftrightarrow \|x\|_M > \|Gx\|_M \forall x \neq 0$$

$$\rho(G) \leq \|G\|_M = \sup_{x \neq 0} \frac{\|Gx\|_M}{\|x\|_M} < 1$$

## 5.7 Satz: 8.4:

Das Gauss-Seidel-Verfahren konvergiert für jeden Startwert  $x^0 \in R^n$ , falls  $A$  spd ist.

Beweis:  $G = -(D + L)^{-1}R = I - (D + L)^{-1}A$

nach Lemma 8.3 betrachten wir  $I - G^*G$

$$M = A$$

$$\begin{aligned} G^* &= A^{-1}G^TA = A^{-1}(I - (D + L)^{-1}A)^TA \\ &= A^{-1}(I - A(\underbrace{D^T}_{=D} + \underbrace{L^T}_{=R})^{-1})A \\ (A^{-1} - (D + R)^{-1})A &= I - (D + R)^{-1}A \\ B &= I - G^*G = I - (I - (D + R)^{-1}A)(I - (D + L)^{-1}A) \\ &= I - (I - (D + R)^{-1}A - (D + L)^{-1}A + (D + R)^{-1}A(D + L)^{-1}A) \\ &= (D + R)^{-1}(I + (D + R)(D + L)^{-1} - A(D + L)^{-1})A \\ &= (D + R)^{-1}(D + \underbrace{L + D + R - A}_{=A})(D + L)^{-1}A \\ &= (D + R)^{-1}D(D + L)^{-1}A \\ (Bx, x)_A &= x^T A(D + R)^{-1}D(D + L)^{-1}Ax \\ &= ((D + L)^{-1}Ax)^T D \underbrace{(D + L)^{-1}Ax}_{=y \neq 0} \\ &= y^T Dy = \sum d_{ii}y_i^2 > 0 \end{aligned}$$

Lemma 8.3  $\Rightarrow \rho(G) < 1 \Rightarrow$  Konvergenz

## 5.8 Relaxierung von FP-Iterationen

Richardson:  $x^{k+1} = x^k - \tau(Ax^k - b) = (1 - \tau)x^k + \tau(x^k - Ax^k + b)$

Allgemein: statt  $x^{k+1} = Gx^k + c$

relaxierte Version:  $x^{k+1}(1 - \omega)x^k + \omega(Gx^k + c)$

$\omega \in (0, 1]$ : Dämpfungsparameter

$\omega > 1$  Überrelaxieren

$$x^{k+1} = \omega(Gx^k + c) + (1 - \omega)x^k$$

$$= \underbrace{(\omega G + (1 - \omega)I)}_{G_\omega} x^k + \omega c$$

Ideal:  $\omega$  so zu wählen, dass das Verfahren möglichst schnell konvergiert.

Fehler geht mit  $\rho(G_\omega)^k \rightarrow 0$

$\omega$  so zu wählen, dass  $\rho(G_\omega)$  minimal

12.01.2010

---

$$x_{k+1} = Gx_k + c \rightarrow x_{k+1} = \omega Gx_k + \omega c + (1 - \omega)x_k$$

$\omega \in (0, 1]$ ,  $\omega = 1$ : FP-Iteration

## 5.9 Def. 8.5

Fixpunktiteration  $x_{k+1} = Gx_k + c$ , wobei  $I = G(A) \sim \text{spd. Matrix für } A \text{ spd}$   
d.h.  $\exists W = W(A), W(I - G(A))W^{-1} \text{ spd}$

## 5.10 Bsp: 8.6:

Richardson  $G(A) = I - A$  symmetrisch,  $I - G(A) = A$ , spd

Jacobi:  $G = I - D^{-1}A, I - G = D^{-1}A$

$W = D^{1/2}, D = \text{diag}(a_{11}, \dots, a_{nn})$

$$\begin{aligned} W(I - G)W^{-1} &= D^{1/2}D^{-1}AD^{-1/2} = D^{-1/2}AD^{-1/2} = (D^{-1/2}AD^{-1/2})^T = (D^{-1/2})^T A (D^{-1/2})^T = D^{-1/2}AD^{-1/2} \\ x^T(D^{-1/2}AD^{-1/2})x &= (\underbrace{D^{-1/2}x}_=)^T A (D^{-1/2}x) = y^T Ay > 0 \text{ für } y \neq 0 \end{aligned}$$

## 5.11 Lemma 8.7:

$$x_{k+1} = Gx_k + c, \text{ sei symmetriesierbar.}$$

Dann sind für  $A$  spd die Eigenwerte von  $G(A)$  reell und  $\lambda < 1$

Konvergenz des Verfahrens für  $(\lambda_i) < 1$

D.h. wir müssen nun erzwingen, dass  $EW > -1$  sind

Relaxierung:  $G_\omega = \omega G + (1 - \omega)I$

$$\lambda_i(G_\omega) = \omega \lambda_i(G) = (1 - \omega) < \omega + (1 - \omega) < 1$$

$$\lambda_i(G_\omega) > \omega \lambda_{min}(G) + (1 - \omega) > -1$$

$$2 > \omega(1 - \lambda_{min}(G))$$

$$\omega < \frac{2}{1 - \lambda_{min}(G)}$$

$$\lambda_{min}(G) > -1 \rightarrow \omega = 1 \text{ erlaubt}$$

$$\lambda_{min}(G) < -1 \rightarrow \omega \text{ nötig}$$

FP-V symmetrisierbar ist → Relaxiertes PF-Verfahren konvergiert für  $0 < \omega < \frac{2}{1 - \lambda_{min}(G)}$   
Was ist die optimale Wahl von  $\omega$ ?

Wehlerabschätzung:  $\|x_{k+1} - \bar{x}\| \leq \rho(G_\omega) \|x_k - \bar{x}\| \leq \rho(G_\omega)^{k+1} \|x_0 - \bar{x}\|$

schnellste Konvergenz, wenn  $\rho(G_\omega)$  minimal

$$\rho(G_\omega) = \max |1 - \omega(1 - \lambda_{min}(G))|, |1 - \omega(1 - \lambda_{max}(G))|$$

$$\bar{\omega} = \min_{0 \leq \omega \leq 1} \rho(G_\omega) \text{ gesucht.}$$

$$0 < 1 - \lambda_{max}(G) \leq 1 - \lambda_{min}(G)$$

$$1 - \omega(1 - \lambda_{max}(G)) \geq 1 - \omega(1 - \lambda_{min}(G))$$

Fallunterscheidung:

1.Fall: beide positiv:  $\rho(G_\omega) = 1 - \omega(1 - \lambda_{max}(G))$

2.Fall: beide negativ:  $\rho(G_\omega) = 1 - \omega(1 + \lambda_{min}(G))$

3.Fall: verschiedene Vorzeichen:  $\rho(G_\omega) = \max 1 - \omega(1 - \lambda_{max}(G)), 1 - \omega(1 - \lambda_{min}(G))$

Optimal:  $-1 + \bar{\omega}(1 - \lambda_{min}(G)) > 0$

$$1 - \bar{\omega}(1 - \lambda_{max}(G)) > 0$$

$$\frac{1}{1 - \lambda_{max}(G)} > \bar{\omega} > \frac{1}{1 - \lambda_{min}(G)}$$

$$\text{Falls } -1 + \bar{\omega}(1 - \lambda_{min}(G)) = \rho(G_\omega) > 1 - \bar{\omega}(1 - \lambda_{max}(G))$$

$\omega = \bar{\omega} - \epsilon, \epsilon \geq \text{KLEIN GENUNG}$

$$-1 + (\bar{\omega} - \epsilon)(1 - \lambda_{min}(G)) < -1 + \bar{\omega}(1 - \lambda_{min}(G))$$

$$\text{aber } -1 + (\bar{\omega} - \epsilon)(1 - \lambda_{min}(G)) < 1 - (\bar{\omega} - \epsilon)(1 - \lambda_{max}(G))$$

$$\text{d.h. } \rho(G_{\bar{\omega}}) = 1 + (\bar{\omega} - \epsilon)(1 - \lambda_{min}(G))$$

Widerspruch zur Optimalität von  $\omega$ .

Analog Widerspruch, falls

$$1 - \bar{\omega}(1 - \lambda_{max}(G)) = \rho(G_{\bar{\omega}}) > -1 + (1 - \lambda_{min}(G))$$

Betrachte  $\bar{\omega} - \epsilon$

D.h.  $\bar{\omega}$  optimal

$$\Rightarrow 1 - \bar{\omega}(1 - \lambda_{max}(G)) = -1 + \bar{\omega}(1 - \lambda_{min}(G))$$

$$\bar{\omega} = \frac{2}{2 - \lambda_{max}(G) - \lambda_{min}(G)}$$

### SOR-Verfahren:

#### Successive Overrelaxation:

Gauss-Seidel:  $(D + L)x^{k+1} = b - Rx^k$

$$x^{k+1} = \omega(D + l)^{-1}b - \omega(D + L)^{-1}Rx^k + (1 - \omega)x^k$$

$$(D + L)x^{k+1} - \omega b - \omega Rx^{k+1} + (1 - \omega)(D + L)x^k$$

SOR

$$(D + \omega L)x^{k+1} = \omega b - \omega Rx^{k+1} + (1 - \omega)Dx^k$$

A spd → Konvergenz für  $\omega \in (0, 2)$

Überrelaxierung,  $\omega \in (1, 5; 2)$ , Beschleunigung des Gauss-Seidel-Verfahrens

## 5.12 8.3 Verfahren der konjugierten Gradienten (CG)

Konstruiere eine Folge von Teilräumen  $U_k \subset \mathbb{R}^n$

$$\dim U_k = k$$

Wählen Startwert  $x_0$

$$V_k = x_0 + U_k$$

$x_k$  Lösung von  $\|x_k - \bar{x}\| = \min_{y \in V_k} \|y - \bar{x}\|$

$\bar{x}$  Lsg von  $Ax = b$

Gute theoretische Eigenschaft:

$$V_n = x_0 + U_n, \dim U_n = n, U_n \subset \mathbb{R}^n \Rightarrow U_n = \mathbb{R}^n$$

$$\Rightarrow V_n = \mathbb{R}^n$$

$$\|x_n - \bar{x}\| = \min_{y \in V_n} \|y - \bar{x}\| = 0$$

Also ist  $x_n = \bar{x} \Rightarrow$  Konvergenz in n Schritten

Was ist passende Norm?

Betrachten  $(\cdot, \cdot)_B, \|x\|_B = \sqrt{(x, x)_B}$

$$\|x_k - \bar{x}\|_B = \min_{y \in V_k} \|y - \bar{x}\|_B$$

$\Rightarrow x_k$  ist orthogonale Projektion von  $\bar{x}$  auf  $V_k$ .

$$(x_k - \bar{x}, u)_B = 0, \forall u \in U_k$$

$$(x_k - \bar{x})^T Bu = 0$$

$$B(x_k - \bar{x})u = 0$$

Für  $B = A : B(x_k - \bar{x}) = Ax_k - A\bar{x} = Ax_k - b = -r_k$

Residuen orthogonal zu  $U_k$ .

Bem.:  $\|x_k - \bar{x}\|_A = \min_{y \in V_k} \|y - \bar{x}\|$  heisst Ritz-Galerkin Approximation

$P_1, \dots, P_k$  sei Orthogonalbasis von  $U_k$

$$x_k = P_k \bar{x} = x_0 + \sum_{j=1}^k \frac{(p_j, \bar{x} - x_0)_A}{(p_j, p_j)_A}$$

$$x_k = x_0 + \sum_{j=1}^k \frac{p_j^T A(\bar{x} - x_0)}{(p_j, p_j)_A}, p_j = x_0 + \sum_{j=1}^k \frac{p_j^T r_0}{(p_j, p_j)_A} p_j$$

$$r_0 = b - Ax_0$$

$$x_k - x_0 + \underbrace{\sum_{j=1}^{k-1} \alpha_j p_j}_{=x_{k-1}} + \alpha_k p_k$$

$$x_k = x_{k-1} + \alpha_k p_k$$

$$r_k = b - Ax_k = b - Ax_{k-1} - \alpha_k Ap_k = r_{k-1} - \alpha_k Ap_k$$

Satz von Cayley-Hamilton (LA)

$\exists P_{n-1}$  Polynom Grad  $n-1$

$$A^{-1} = P_{n-1}(A)$$

$$\bar{x} - x_0 = A^{-1}(b - Ax_0) = A^{-1}r_0 + P_{n-1}(A)r_0$$

Am Ende:  $\bar{x} - x_0 \in U_n - \text{span}\{r_0, \dots, A^{n-1}r_0\}$

Idee  $U_k = \text{span}\{r_0, \dots, A^{k-1}r_0\}$  Krylov-Unterräume

---

15.01.2010

CG-Verfahren  $U_k = \text{span}\{p_1, \dots, p_n\}$

$$x_k = x_{k-1} + \alpha_k p_k$$

$$r_k = r_{k-1} - \alpha_k Ap_k$$

Satz von Cayley-Hamilton  $A^{-1} = P_{n-1}(A), A^{-1}r_0 = x - x_0 = P_{n-1}(A)r_0$

$$U_k = \text{span}\{r_0, Ar_0, \dots, A^{n-1}r_0\} = U_k(Ax_0)$$

$V_k = x_0 + U_k$  Keylov-Unterräume

### 5.12.1 Lemma 8.15:

$r_k \neq 0$ , Dann sind die Residuen  $r_0, \dots, r_k$  paarweise orthogonal, d.h.  $r_i^T r_j = \delta_{ij} r_i^T r_i; i, j = 0, \dots, k$  und spannen den Raum  $U_{k+1}$  auf.

Beweis: Induktion:

$$k = 0 : U_1 = \text{span}\{r_0\}$$

$k > 1$ : Beh.: richtig für  $k = 1$ , d.h.  $r_0, \dots, r_k$  seien paarweise orthogonale Basis von  $U_k$

$$r_k = r_{k-1} - \alpha_k AP_k = \text{span}\{r_0, \dots, A^{k-1}r_0\}$$

$$Ap_k \in AU_k = \text{span}\{Ar_0, \dots, A^k r_0\}$$

$$\Rightarrow r_k = r_{k-1} - \alpha_k AP_k \in \text{span}\{r_0, \dots, A^k p_0\} = U_{k+1}$$

Nach Konstruktion von  $x_k/r_k$  gilt

$$r_k^T u = 0, \forall u \in U_k$$

$$r_k^T r_j = 0 \text{ für } j < k, \text{ da } r_j \in U_k$$

Orthogonalbasis von  $U_{k+1}$

Konstanten der  $p_k$

$$r_0 \neq 0 : P_1 = r_0$$

$$k > 1; r_k = 0 \rightarrow x_k \text{ Lösung von } Ax = b$$

$$r_k \neq 0: r_k \text{ linear unabhängig zu } P_1, \dots, P_k$$

$$p_{k+1} = r_k - \sum_{j=1}^k \frac{(r_k, p_j)_A}{(p_j, p_j)_A} p_j \quad (\text{Projektion von } r_k \text{ auf Orthogonalkomplement von } U_k)$$

$$(p_{k+1}, p_j)_A = 0 \text{ für } j \leq k$$

### 5.12.2 Algorithmus 8.16 (CG-Verfahren)

$$x_0 \in R^n, p_1 = r_0 = b - Ax_0$$

Iteration über  $k$  bis Abbruchbedingung erfüllt: ( $\|r_k\| < \epsilon$ )

$$\alpha_k = \frac{r_{k-1}^T r_{k-1}}{(p_k, p_k)_A} = \frac{r_{k-1}^T r_{k-1}}{p_k^T A p_k}$$

$$x_k = x_{k-1} + \alpha_k p_k$$

$$r_k = r_{k-1} - \alpha_k A p_k$$

$$\beta_{k+1} = \frac{(r_k^T p_k)_A}{(p_k, p_k)_A} = \frac{r_k^T r_k}{r_{k-1}^T r_{k-1}}$$

$$P_{k+1} = r_k + \beta_{k+1} p_k$$

$$\beta_{k+1} = \frac{(r_k, p_k)_A}{(p_k, p_k)_A}$$

$$\alpha_k = \frac{p_k^T r_0}{(p_k, p_k)_A}, p_k^T r_0 = p_k^T (b - Ax_0) = p_k^T A (\bar{x} - x_0) = (p_k, \bar{x} - x_0)_A = (p_k, \bar{X} - X_1 + \alpha_1 P_1)_A = (p_k, x - \bar{x}_1)_A + \alpha_1 (p_k, p_1)_A$$

$$= (p_k, \bar{x} - x_1)_A = \dots = (p_k, \bar{x} - x_{k-1})_A = p_k^T \underbrace{A(\bar{x} - x_{k-1})}_{b - Ax_{k+1}} = p_k^T r_{k-1}$$

$$x_1 = x_0 + \alpha_1 p_1$$

$$p_k^T r_0 = p_k^T r_{k-1} = (r_{k-1} + \beta_k p_{k-1})^T r_{k-1} = r_{k-1}^T r_{k-1} + \beta_k \underbrace{p_{k-1}^T}_{\in U_{k-1}} \underbrace{r_{k-1}}_{\perp U_{k-1}} = r_{k-1}^T r_{k-1}$$

### 5.12.3 Satz: 8.17:

$x_k$  erzeugt durch CG-Verfahren

$$\|\bar{x} - x_k\|_A \leq 2 \left( \frac{\sqrt{K_2(A)} - 1}{\sqrt{K_2(A)} + 1} \right)^k \|\bar{x} - x_0\|_A$$

## 5.13 8.4 Vorkonditionierung

$$Ax = b \rightarrow \tilde{A}\tilde{x} = \tilde{b}$$

$$K_2(\tilde{A}) \ll K_2(A)$$

$$\text{linke Vorkonditionierung } \underbrace{PA}_A x = \underbrace{Pb}_b$$

$$\text{rechte Vorkonditionierung } \underbrace{AB}_A \tilde{x} = b, \tilde{x} = B^{-1}x$$

A spd., P oder B spd

$$K_2(PA), K_2(AB)?$$

$\lambda EW$  von  $PA$

$$PAx = \lambda x \text{ für ein } x \neq 0$$

$Ax = \lambda P^{-1}x$  verallgemeinertes EW-Problem

$$y = p^{-1/2}x, p \text{ spd} \Rightarrow p^{-1} \text{ spd}, \Rightarrow \exists p^{-1/2}, p^{-1} p^{-1/2} p^{-1/2}$$

$$p^{1/2} A p^{1/2} p^{1/2} x = \lambda p^{-1/2} x$$

$$(p^{1/2})^T A p^{1/2} = \lambda y$$

$\lambda$  ist EW von  $(P^{1/2})^T A p^{1/2}$  spd

$$\Rightarrow \lambda \in R_+$$

$$K_2(PA) = \sqrt{\frac{\lambda_{\max}(PA)}{\lambda_{\min}(PA)}}$$

analog ist  $\lambda$  EW von AB auch EW von  $((B^{1/2})^T A B^{1/2})$

$\tilde{A} = AB$  ist selbstadjungiert und positiv in  $(\cdot, \cdot)_B$

$$(x, \tilde{A}y)_B = x^T ABy = x^T B^T A^T By = (ABx)^T By$$

$$= (\tilde{A}x, y)_B$$

$$(x, \tilde{A}x)_B = x^T BABx = (Bx)^T A(Bx) \xrightarrow{y=Bx} y^T Ay > 0, \text{ da } A \text{ spd}$$

Vorkonditionierte Verfahren:

Bsp.: Richardson-Iteration:

$$x_{k+1} = x_k - (ABx_k - b)$$

$$\tilde{x} = Bx_k$$

$$\tilde{x}_{k+1} = \tilde{x}_k - (BA\tilde{x}_k - Bb)$$

Analog Richardson-Iteration mit Vorkonditionierung von links.

$$x_{k+1} = x_k - (PAx_k - Pb)$$

$$P^{-1/2}x_{k+1} = P^{-1/2}x_k - (P^{1/2}Ap^{1/2}P^{-1/2}x_k - P^{1/2}b)$$

$$y_{k+1} = y_k - (P^{1/2}Ap^{1/2}y_k - P^{1/2}b)$$

Für transformierte Variable  $y_k = p^{-1/2}x_k$  standard Richardson-Verfahren mit spd-Matrix

$$\|y_k - \bar{x}\| \leq \left(\frac{K_2(P^{1/2}Ap^{1/2})-1}{K_2(P^{1/2}Ap^{1/2})+1}\right)^k \|y_0 - \bar{y}\|$$

$$\|p^{-1/2}(x_k - \bar{x})\| \leq \underbrace{\left(\frac{K_2(PA)-1}{K_2(PA)+1}\right)^k}_{y_k = P^{-1/2}x_k} \|P^{-1/2}(x_0 - \bar{x})\| \leq \left(\frac{K_2(PA)-1}{K_2(PA)+1}\right)^k \|x_0 - \bar{x}\|_{P^{-1}}$$

$$\|P^{-1/2}z\| = \sqrt{(P^{-1/2}z)^T P^{-1/2}z} = \sqrt{z^T P^{-1/2} P^{-1/2} z} = \sqrt{z^T p^{-1} z} = \sqrt{(z, z)_{P^{-1}}} = \|z\|_{P^{-1}}$$

19.01.2010

---

$Ax = b$ , Konvergenzfaktor wächst mit  $K_2(A)$ .

$$PAx = Pb \rightarrow K_2(PA)$$

$$AB\tilde{x} = b \rightarrow K_2(AB)$$

Konvergenz-Richardson-Verfahren

$$\|x - x_k\| \leq \frac{K_2(AB)-1}{K_2(AB)+1}^K, \text{ oder } \|x - x_k\| \leq \frac{K_2(PA)-1}{K_2(PA)+1}^K$$

c Konstante, die Normäquivalenz zwischen Euklidischer Norm und der von  $B/p^{-1}$  induzierten beschreibt.

## 5.14 Satz 8.22: Für CG-Verfahren mit Vorkonditionierung (PCG)

$$\|x - x_k\|_A \leq 2 \frac{\sqrt{K_2(AB)}-1}{\sqrt{K_2(AB)}+1} \|x - x_0\|_A$$

Voraussetzungen für Vorkonditionierung  $B/P$

- a)  $K_2(AB)$  bzw.  $K_2(PA)$  möglichst klein
- b)  $y \rightarrow By$  bzw.  $y \rightarrow Py$  „effizient“ berechenbar.

Wahl des Vorkonditionierens beruht auf möglichst gutem Kompromiss zwischen a) und b).

## 5.15 Lemma: 8.23:

Falls für  $\mu_0, \mu_1 > 0$  eine der folgenden drei äquivalenten Bedingungen gilt:

$$(i) \quad \mu_0(y, y)_{A^{-1}} = y_0 y^T A^{-1} y \leq y^T B y = (y, y)_B$$

$$(y, y)_B \leq \mu_1(y, y)_{A^{-1}}$$

$$(ii) \quad \mu_0(y, y)_B \leq (y, y)_{BAB}$$

$$(iii) \quad \lambda_{min}(AB) \geq \mu_0, \lambda_{max}(AB) \leq \mu_1$$

Dann gilt:  $K_2(AB) \leq \frac{\mu_1}{\mu_0}$

Beweis: (i)  $\Rightarrow$  (ii)  $y = Au$  in (i)

$$(y, y)_{A^{-1}} = (Au, Au)_{A^{-1}} = (Au)^T A^{-1} Au = (u, u)_A$$

$$y = A^{1/2} B^{1/2} u$$

$$(y, y)_{A^{-1}} = (A^{1/2} B^{1/2} u)^T A^{-1} A^{1/2} B^{1/2} u = u^T B^{1/2} A^{1/2} A^{-1} A^{1/2} B^{1/2} u = u^T Bu = (u, u)_B$$

$$(y, y)_B = y^T By = u^T B^{1/2} A^{1/2} B A^{1/2} B u$$

$$\begin{aligned}
 (i) &\Rightarrow (u, u)_B \leq (u, u)_{B^{1/2} A^{1/2} B A^{1/2} B} \leq \mu_1 (u, u)_B \\
 (ii) &\Rightarrow (iii) \\
 \lambda_{\min}(AB) &= \min_{x \neq 0} \frac{x^T A B x}{x^T x} = \min_y \frac{A B y, y)_B}{(y, y)_B} = \min_y \frac{y^T B A B y}{(y, y)_B} = \min_y \frac{(y, y)_{BAB}}{(y, y)_B} \\
 &\geq \min_y y_0 \frac{(y, y)_B}{(y, y)_B} = \mu_0 \\
 \lambda_{\max}(AB) &= \max_{y \neq 0} \frac{(y, y)_{BAB}}{(y, y)_B} \leq \mu_1 \\
 (ABy, y)_B &= y^T B A B y = y^T B^T A^T B y = (ABy)^T B y = (y, ABy)_B \\
 (iii) &\Rightarrow (ii) \\
 \mu_0 &\leq \lambda_{\min}(AB) = \min_{y \neq 0} \frac{(y, y)_{BAB}}{(y, y)_B} \\
 &\Rightarrow \frac{(y, y)_{BAB}}{(y, y)_B} \geq \mu_0, \forall y \\
 &\Rightarrow (y, y)_{BAB} \geq \mu_0 (y, y) \forall y \\
 (iii) &\Rightarrow K_2(AB) = \frac{\lambda_{\max}(AB)}{\lambda_{\min}(AB)} \leq \frac{\mu_1}{\mu_0} \\
 &\text{(Beweis dieses Lemmas ist falsch)}
 \end{aligned}$$

**Beispiele:**

$$\begin{aligned}
 1: \quad B &= \tau * I \\
 (y, y)_B &= \tau y^T y \\
 (y, y)_{BAB} &= \tau^2 y^T A y \\
 \mu_0 &\leq \lambda_{\min}(AB) = \min_{y \neq 0} \frac{(y, y)_{BAB}}{(y, y)_B} = \min_{y \neq 0} \frac{\tau^2 y^T A y}{\tau y^T y} = \tau \min_{y \neq 0} \frac{y^T A y}{y^T y} = \tau \lambda_{\min}(A) \\
 \lambda_{\max}(AB) &= \tau \lambda_{\max}(A) \\
 &\Rightarrow K_2(A) = K_2(AB)
 \end{aligned}$$

**2. Diagonale Vorkonditionierer**  $B = D^{-1}P = D^{-1}$ ,  $D = \text{diag}(A)$

$$\mu_0 y^T D^{-1} y = \mu(y, y)_{D^{-1}} \leq (y, y)_{D^{-1}AD^{-1}} = y^T D^{-1} A D^{-1} y$$

$$u = D^{-1/2} y, y = D^{1/2} u$$

$$\mu_0 u^T D^{1/2} D^{-1} D^{1/2} u \leq u^T D^{1/2} D^{-1} A D^{-1} D^{1/2} u$$

$$\begin{aligned}
 D &= \begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ 0 & & a_{nn} \end{pmatrix}, D^{-1/2} = \begin{pmatrix} \frac{1}{\sqrt{a_{11}}} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sqrt{a_{nn}}} \end{pmatrix} \\
 D^{-1/2} A D^{-1/2} &= \begin{pmatrix} \frac{a_{11}}{\sqrt{a_{11}}\sqrt{a_{11}}} & \cdots & \cdots & \frac{a_{n1}}{\sqrt{a_{11}}\sqrt{a_{nn}}} \\ \vdots & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \frac{a_{n,n-1}}{\sqrt{a_{n-1,n-1}}\sqrt{a_{nn}}} \\ \frac{a_{1n}}{\sqrt{a_{11}}\sqrt{a_{nn}}} & \cdots & \frac{a_{n-1,n}}{\sqrt{a_{n-1,n-1}}\sqrt{a_{nn}}} & 1 \end{pmatrix}
 \end{aligned}$$

Suchen  $\mu_0$  mit:

$$\mu_0 \sum u_i^2 \leq \sum_{ij} u_i u_j \frac{a_{ij}}{\sqrt{a_{ii}}\sqrt{a_{jj}}} \leq \mu_1 \sum u_i^2$$

Nach oben: Couchy-Schwarz Ungleichung

$$\sum_{ij} \frac{u_i \sqrt{|a_{ij}|}}{\sqrt{a_{ii}}} \frac{u_j \sqrt{|a_{ij}|}}{\sqrt{a_{jj}}} \leq \sqrt{\left(\sum_i u_i^2 \sum_j \frac{a_{ij}}{a_{ii}}\right) \left(\sum_j u_j^2 \sum_i \frac{a_{ij}}{a_{jj}}\right)}$$

$$\Rightarrow \mu_1 = 2$$

A diagonaldominant

$$B = D^{-1}$$

$$\tilde{x}^{k+1} = \tilde{x}^k - (AD^{-1}\tilde{x}^k - b)$$

Konvergenzfaktoren analog Jacobi

Diagonaldominante Matrizen  $\mu_1 \leq 2$

$\mu_0$  abhängig von der Speziellen Matrix A

Diskretisierung:  $-x'' = f, x(a) = a(b) = 0$

$$A = 1/h^2 \begin{pmatrix} 2 & -1 & & \\ -1 & 2 & -1 & \\ & \ddots & \ddots & \ddots \\ & & \ddots & \ddots & -1 \\ & & & -1 & 2 \end{pmatrix}$$

$$D^{-1/2} \frac{\sqrt{2}}{h} I, D^{-1/2} = \frac{h}{\sqrt{2}} I$$

$$D^{-1/2} A D^{-1/2} = \frac{h}{\sqrt{2}} \frac{2}{h^2} \begin{pmatrix} 1 & -\frac{1}{2} & & \\ -\frac{1}{2} & 1 & -\frac{1}{2} & \\ & \ddots & \ddots & \ddots \\ & & \ddots & \ddots & -\frac{1}{2} \\ & & & -\frac{1}{2} & 1 \end{pmatrix} \frac{h}{\sqrt{2}}$$

$$\begin{aligned} u^T D^{-1/2} A D^{-1/2} u &= \sum_{i=1}^n u_i^2 - 1/2 \sum_{i=1}^{n-1} u_i u_{i+1} - 1/2 \sum_{i=2}^n u_{i-1} u_i \\ &= \sum_{i=1}^n u_i^2 - \underbrace{\sum_{i=1}^{n-1} u_i u_{i+1}}_{CS} \leq 2 \sum_{i=1}^n u_i^2 \end{aligned}$$

$$\mu_1 = 2$$

$$\text{nach unten: } \sum_{i=1}^{n-1} (u_i^2 - u_i u_{i+1}) + u_n^2 = \sum_{i=1}^{n-1} u_i (u_i - u_{i+1}) + u_n^2$$

Bsp:  $u = (1, \dots, 1)$

$$u^T u = \sum_{i=1}^n u_i^2 = n$$

$$u^T D^{-1/2} A D^{-1/2} u = 1$$

$\mu_0$  muss  $\mu_0 u \leq 1$  erfüllen

$$\begin{aligned} \mu_0 &\leq \frac{1}{n} \\ \frac{\mu_1}{\mu_0} &\geq \frac{2}{\frac{1}{n}} = 2n \end{aligned}$$

Bessere Vorkonditionierungen:

- unvollständige LR-Zerlegung (ILU)
- unvollständige Cholesky-Zerlegung

$$\tilde{L}^T \tilde{L} \approx A$$

Aufwandsüberlegungen: wie „unvollständig“ durchführen will.

- Vorkonditionierung von CG-durch iterative Verfahren z.B Richardson oder Jacobi

Richardson:  $Ax = b$

$$x^1 = x^0 - (Ax^0 - b)$$

$$x^2 = x^1 - (Ax^1 - b) = (I - A)x^1 + b = (I - A)^2 x^0 + (I - A)b + b$$

Fixpunktform:

$$x = (I - A)^2 x + (2I - A)b$$

$$\underbrace{(I - (I - A)^2) A^{-1}}_P Ax = (2I - A)b$$

Gestaffelte Iteration über  $k = \text{pcg-Verfahren}$

In jedem Iterationsschritt definiere die Anwendung von  $P$  auf  $A$  durch  $m$  Schritte des Richardson-Verfahrens

22.01.2010

## 5.16 Iterative Lösung symmetrisch indefiniter Probleme

Motivation: verallgemeinerte Inverse für unterbestimmte Systeme.

$$Bx = f, B \in R^{m \times n}, m < n, RgB = m$$

Verallgemeinerte Inverse Orthogonale zu  $N(B)$

Für  $m, n$  gross, direkte Berechnung von  $B^+$  unmöglich Berechnung von  $N(B)$  zu aufwendig.

Alternative: Äquivalente Definition „Lösung der minimalen Norm“  $\|x\| \rightarrow \min_{x, Bx=f} \dots$

$$\Leftrightarrow 1/2 \|x\|^2 \rightarrow \min_{x, Bx=f} \dots$$

Allgemeiner:  $1/2 \|x\|_a^2 \rightarrow \min_{x, Bx=f} \dots, A \in R^{n \times n}$  spd

noch allgemeiner:  $1/2\|x\|_A^2 - f^T x \rightarrow \min_{x, Bx=g} \dots$

Einbettung in  $R^{n+m}$ ,  $p \in R^m$

$$L(x, p) = 1/2\|x\|_A^2 - f^T x + (Bx - g)p$$

$L : R^{n+m} \rightarrow R$  Lagrange-Funktional

$p$  Langrange-Parameter, duale Variable

### 5.17 Lemma:

$\bar{x} \in R^n$  löst das beschränkte Optimierungsproblem (A), falls  $\bar{x}$  ein Minimum von  $J : R^n \rightarrow R \cup \infty$

$$J(x) = \sup_{p \in R^m} L(x, p)$$

Beweis:  $Bx = g : J(x) = 1/2\|x\|_A^2 - f^T x$

$Bx \neq g : p_k := k(Bx - g)$

$$\sup_p L(x, p) \geq L(x, p_k) = 1/2\|x\|_A^2 f^T y + k \underbrace{\|Bx - g\|^2}_{\substack{>0 \\ \rightarrow \infty \text{ für } k \rightarrow \infty}}$$

$\rightarrow \infty$  für  $k \rightarrow \infty$

$$\rightarrow \sum_p L(x, p) = \infty$$

$$J(x) = \begin{cases} 1/2\|x\|_A^2 - f^T x & \text{für } Bx = g \\ +\infty & \text{für } Bx \neq g \end{cases}$$

$\bar{x} \in R^n$  minimiert  $J \Rightarrow J(\bar{x}) < \infty \Rightarrow B\bar{x} = f$

und  $J(\bar{x}) = 1/2\|\bar{x}\|_A^2 - f^T \bar{x} \leq J(x) = 1/2\|x\|_A^2 - f^T x$

$$\min_x J(x) = \min_x \sup_p L(x, p)$$

Falls Lösung  $(\bar{x}, \bar{p})$  von  $\min_x \max_p L(x, p)$  existiert, gilt:

$$L(\bar{x}, \bar{p}) = \min_x \max_p L(x, p)$$

$$L(\bar{x}, \bar{p}) \leq L(x, \bar{p}) \forall x \in R^n \rightarrow \frac{\partial}{\partial x} L(\bar{x}, \bar{p}) = 0$$

$$L(\bar{x}, \bar{p}) \geq L(\bar{x}, p) \forall p \in R^m \rightarrow \frac{\partial}{\partial p} L(\bar{x}, \bar{p}) = 0$$

$\Rightarrow m+n$  Gleichungen für  $m+n$  Unbekannte.

$$\begin{aligned} \frac{\partial L}{\partial x}(x, p)y &= \frac{d}{d\epsilon}(L(x + \epsilon y, p)|_{\epsilon=0}) \\ &= \frac{d}{d\epsilon}(1/2(x + \epsilon y)^T A(x + \epsilon y) - f^T(x + \epsilon y) + (Bx + \epsilon By - g)^T p) \\ &= (1/2 y^T A(x + \epsilon y) + 1/2(x + \epsilon y)Ay - f^T y + (By)^T p)|_{\epsilon=0} \\ &= 1/2y^T Ax + 1/2x^T Ay - y^T f + y^T B^T p \\ &= y^T(Ax + B^T p - f) \end{aligned}$$

$$0 = \frac{\partial L}{\partial x} = Ax + B^T p - f$$

$$0 = \frac{\partial L}{\partial p} = \frac{\partial}{\partial p}((Bx - g)^T p) = \frac{\partial}{\partial p}(P^T(Bx - g)^T) \\ = Bx$$

$$Ax + B^T p = f$$

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} x \\ \bar{p} \end{pmatrix} = \begin{pmatrix} f \\ f \end{pmatrix}$$

### 5.18 Satz:

Sei  $(\bar{x}, \bar{p})$  eine Lösung von oben genannten, dann ist  $\bar{x}$  eine Lösung des beschränkten Minimierungsproblems.

Beweis:  $x \in R^n, Bx = g$

$$1/2x^T Ax - f^T x - 1/2\bar{x}^T A\bar{x} - f^T x$$

$$= 1/2(x - \bar{x})^T A(x - \bar{x}) - 1/2\bar{x}^T \bar{x} + 1/2x^T Ax + 1/2x^T A\bar{x} - 1/2\bar{x}^T A\bar{x} - f^T(x - \bar{x})$$

$$= 1/2\|x - \bar{x}\|_A^2 - \bar{x}^T A\bar{x} + \bar{x}^T Ax - f^T(x - \bar{x})$$

$$\geq \bar{x}^T A(x - \bar{x}) - f^T(x - \bar{x}) = (A\bar{x} - f)^T(x - \bar{x})$$

$$= (B^T \bar{p})^T(x - \bar{x}) = \bar{p}^T B(x - \bar{x}) = \bar{p}^T(g - g) = 0$$

$$\Rightarrow 1/2x^T Ax - f^T x \geq 1/2\bar{x}^T A\bar{x} - f^T \bar{x} \quad \forall x, f x = g$$

Beispiel:

$$A = 1, B = 1$$

$$n = m = 1$$

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

$$EW - GG\det \begin{pmatrix} \lambda - 1 & -1 \\ -1 & \lambda \end{pmatrix} = \lambda(\lambda - 1) - 1 = 0$$

$$\lambda^2 - \lambda = 0$$

$$\Rightarrow \lambda = \frac{1 \pm \sqrt{5}}{2}$$

$M = \begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix}$  hat n Eigenwerte  $\geq 0$ , m Eigenwerte  $\leq 0$

$$\text{Aspd: } Ax + B^T = f \rightarrow x = A^{-1}(f - B^T p)$$

$$Bx = g \rightarrow Ba^{-1}(f - B^T p) = g$$

$$BA^{-1}B^T p = BA^{-1}f - g$$

Lemma A spd,  $Rg(B) = m \Rightarrow S^p = BA^{-1}B^T$  ist spd

S = Schur-Komplement

Beweis:  $RgB = m \Rightarrow N(B^T)$

$$p \neq 0 \Rightarrow B^T p = x \neq 0$$

$$p^T B A^{-1} B^T p = x^T A^{-1} x > 0 \text{ wegen } A^{-1} \text{ spd}$$

$$p^T S p > 0 \forall p \neq 0 \Rightarrow S \text{ spd}$$

$S p = h$  wie bisher  $\rightarrow$  alle Verfahren anwendbar

Vorkonditionierte Richardson-Iteration:

$$C p^{k+1} = C p^k - (S p - h), C \in R^{m \times m} \text{ spd}$$

## 5.19 Satz:

$$C \geq BA^{-1}B \quad (\Leftrightarrow C - BA^{-1}B = (C - S) \text{ spd})$$

Dann konvergiert  $p^k \rightarrow \bar{p}$  mit  $x^k = A^{-1}(f - B^T p^k) \rightarrow \bar{x}$

Beweis:  $G = C^{-1}(C - S)$

$$\lambda \text{ EW von } G \text{ mit EW } q : C^{-1}(C - S)q = \lambda q$$

$$(C - S)q = \lambda Cq$$

$$q^T C q - q^T S q = \lambda q^T C q$$

$$\lambda = \frac{q^T(C-S)q}{q^T C q} > 0$$

$$q^T C q (1 - \lambda) = q^T S q$$

$$1 - \lambda = \frac{q^T S q}{q^T C q} > 0 \Rightarrow \lambda < 1$$

$$\Rightarrow \rho(G) < 1$$

26.01.2010

---

## 5.20 Inexaktes Uzawa-Verfahren

$$\hat{A}x^k = (\tilde{A} - A)x^k + f - B^T p^k$$

$$C p^{k+1} = C p^k + B x^{k+1} - g$$

A nicht regulär  $\Rightarrow$  wähle  $\tilde{A}$  regulär

A nicht regulär, wenn ist M regulär

A symmetrisch, positiv semidefinit

Lemma:  $x^T A x > 0 \forall x \neq 0, B x = 0$

Dann ist M regulär und umgekehrt.

Beweis: Ann.: M ist nicht regulär  $\Rightarrow \exists(xp) \neq 0$  mit  $M(xp) = 0$

$$Ax + B^T p = 0 \rightarrow x^T A x = -x^T B^T p = -(Bx)^T p = 0$$

M nicht regulär  $\Rightarrow \exists x \neq 0$  mit  $x^T A x = 0$

( $\Leftrightarrow \forall x \neq 0 : x^T A x > 0 \Rightarrow M$  regulär)

M regulär,  $x \neq 0, B x = 0, M(xp) \neq 0, \forall(xp) \neq 0$

Insbesondere  $M(x_0) \neq 0 \Rightarrow \begin{pmatrix} Ax + B^T p \\ Bx \end{pmatrix} \neq 0 \Rightarrow Ax + B^T p \neq 0$   
 $\Rightarrow x^T Ax + \underbrace{x^T B^T Bx}_C = 0 \neq 0 \Rightarrow x^T Ax \neq 0 \Rightarrow x^T Ax > 0$

Lemma: M regulär  $\Leftrightarrow A + B^T C^{-1} B$  spd für ein C spd

Beweis:

$$A + B^T C^{-1} B \text{ spd}$$

$$\Rightarrow x^T Ax + x^T B^T C^{-1} Bx > 0 \forall x \neq 0$$

$$\Rightarrow x^T Ax > 0 \forall x \neq 0, Bx = 0$$

$\Rightarrow M$  regulär

$(A + B^T C^{-1} B)$  nicht spd

$$\Rightarrow \exists x \neq 0: (A + B^T C^{-1} B)x = 0$$

$$x = y + z, y \in N(B), z^T Ay = y^T Az = 0$$

$$Ay + Az + B^T C^{-1} Bz = 0$$

$$\Rightarrow y^T Ay + \underbrace{y^T Az}_{=0} + \underbrace{y^T B^T C^{-1} Bz}_{=0} = 0$$

$$\Rightarrow y^T Ay = 0, y \neq 0, By = 0$$

$\Rightarrow M$  ist nicht regulär

$$\text{Erinnerung: } L(x, p) = 1/2x^T Ax - f^T x + p^T (Bx - g)$$

Angemerkt Lagrangian:

$$L_c(x, p) = L(x, p) + 1/2||Ax - g||_{C^{-1}}^2$$

$$f_c(x) = \sup_p L_c(x, p) = \begin{cases} 1/2x^T Ax - f^T x & Bx = g \\ \infty & Bx \neq g \end{cases}$$

$$L_c(x, p) = L(x, p) + 1/2(Bx - g)^T C^{-1} (Bx - g) + p^T (Bx - g)$$

$$= 1/2(x^T Ax + x^T B^T C^{-1} Bx) - f^T x - g^T C^{-1} Bx + 1/2g^T C^{-1} g + p^T (Bx - g)$$

$$= x^T (A + B^T C^{-1} B)x$$

$$\begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} x \\ -p \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix}$$

multiplizieren  $Bx = g$  mit  $B^T C^{-1}$  und addieren sie zu  $Ax + B^T y = f$

$$\begin{pmatrix} A + B^T C^{-1} B & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} x \\ -p \end{pmatrix} = \begin{pmatrix} f + B^T C^{-1} C g \\ g \end{pmatrix}$$

$\tilde{A} = A + B^T C^{-1} B$  spd

Falls Lösung von  $\tilde{A}x = b$  mit vernünftigem numerischen Aufwand durchführbar ist  $\rightarrow$  Uzawa-Verfahren anwen-

den auf  $M_c = \begin{pmatrix} \tilde{A} & B^T \\ B & 0 \end{pmatrix} \rightarrow$  Konvergenzbedingung analog,  $A \rightarrow \tilde{A}$

Lösung nicht mit vernünftigem Aufwand durchführbar  $\rightarrow$  Inexakte Uzawa-Verfahren ( $A$  spd)

$$\hat{A}x^{k+1} = (\hat{A} - A)x^k - B^T p^k + f$$

$$\hat{C}p^{k+1} = \hat{C} + Bx^{k+1} - g$$

$$\tilde{C}p^{k+1} = \tilde{C}p^k + B\hat{A}^{-1}((\hat{A} - A)x^k - B^T p^k + f)$$

$$\begin{pmatrix} \hat{A} & 0 \\ 0 & \hat{C} \end{pmatrix} \begin{pmatrix} x^{k+1} \\ p^{k+1} \end{pmatrix} = \begin{pmatrix} \hat{A} - A & -B^T \\ B\hat{A}^{-1}(\hat{A} - A) & \hat{C} - B\hat{A}^{-1}B^T \end{pmatrix} \begin{pmatrix} x^k \\ p^k \end{pmatrix} + \begin{pmatrix} f \\ B\hat{A}^{-1}f - g \end{pmatrix}$$

$$G = \begin{pmatrix} \hat{A}^{-1} & 0 \\ 0 & \hat{C}^{-1} \end{pmatrix} \begin{pmatrix} \hat{A} - A & -B^T \\ .. & .. \end{pmatrix} = \begin{pmatrix} \hat{A}^{-1}(\hat{A} - A) & -\hat{A}^{-1}B^T \\ \hat{C}^{-1}B\hat{A}^{-1}(\hat{A} - A) & \hat{C}^{-1}(\hat{C} - B\hat{A}^{-1}B^T)\hat{C}^{-1}\hat{C} \end{pmatrix}$$

$$G = \begin{pmatrix} -\hat{A}^{-1} & -\hat{A}^{-1}B^T\hat{C}^{-1} \\ -\hat{C}^{-1}B\hat{A}^{-1} & \hat{C}^{-1}(\hat{C} - B\hat{A}^{-1}B^T)C^{-1} \end{pmatrix} \begin{pmatrix} A - \hat{A} & 0 \\ 0 & \hat{C} \end{pmatrix}$$

$$G(xp) = \lambda(xp), (yq) = \begin{pmatrix} A - \hat{A} & 0 \\ 0 & \hat{C} \end{pmatrix} \begin{pmatrix} x \\ p \end{pmatrix}$$

$$\begin{pmatrix} -\hat{A}^{-1} & -\hat{A}^{-1}B^T\hat{C}^{-1} \\ -\hat{C}^{-1}B\hat{A}^{-1} & \hat{C}^{-1}(\hat{C} - B\hat{A}^{-1}B^T)C^{-1} \end{pmatrix} \begin{pmatrix} y \\ q \end{pmatrix} = \lambda \begin{pmatrix} A - \hat{A} & 0 \\ 0 & \hat{C} \end{pmatrix} \begin{pmatrix} y \\ q \end{pmatrix}$$

Verallgemeinertes Symmetrische EWP, falls  $A - \hat{A} > 0 \Rightarrow \lambda < 1$

$$\hat{C} - B^T \hat{A}^{-1} B > 0 \Rightarrow \lambda < 1 \Rightarrow \rho(G) < 1$$

$\Rightarrow$  Konvergenz des exakten Uzawa-Verfahrens

Compressed Sensing/Sparsity

Unterbestimmtes lineares Gleichungssystem  $Bx = g, B \in R^{m \times n}, m < n$

Bisher minimierte quadratisches Funktional  $1/2||x||^2 : 1/2x^T Ax - f^T x$

In der Praxis oft andere Informationen:

Lösung aus wenigen Komponenten einer Basis  $\{b_i\}$

$$x = \sum_{i=1}^n \alpha_i b_i \text{ mit möglichstviele } \alpha_i = 0$$

$(\alpha_1, \dots, \alpha_n)$  sparse

$$\#\{i | x_i \neq 0\} \rightarrow \min_{x, Bx=g}$$

$$S_i \in \{0, 1\}$$

$S_i = 0$  wenn  $x_i = 0$

$S_i = 1$  wenn  $x_i \neq 0$

$$\Rightarrow (x_i(1 - S_i) = 0)$$

$$\text{Äquivalent: } \sum S_i^2 \rightarrow \min_{x, s, Bx=g, x_i(1-S_i)=0, S_i \in \{0, 1\}}$$

Konvexe Relaxierung

$$\sum S_i^2 \rightarrow \min_{x, s, Bx=g, x_i(1-S_i)=0, S_i \in [0, 1]}$$

$$x \text{ gegeben } \sum S_i^2 \rightarrow \min_{x_i(1-S_i)=0} \Leftrightarrow \sum S_i^2 \rightarrow \min_{\sum |x_i|(1-S_i)=0}$$

$$L(x, p) = \sum S_i^2 + p \sum |x_i|(1 - S_i)$$

$$2S_i - p|x_i| = 0$$

$$S_i = p/2|x_i|$$

$$p^2/4 \sum X_i^2 \rightarrow \min_{Bx=g}$$

---

29.01.2010

---

Idee:

$$Bx = g, B \in R^{m \times n}, m < n$$

Suchen Lösung mit  $\#\{i | x_i \neq 0\} \rightarrow \min$

$$l^0 - \text{Minimierung: } \sum |x_i|^0 \rightarrow \min, 0^0 := 0$$

Approximiert durch  $\sum |x_i|^q \rightarrow \min_x, q \rightarrow 0$

$$\sum s_i^2 \rightarrow \min, x_i(1 - s_i) = 0$$

$$\rightarrow \sum |x_i|(1 - s_i) = 0$$

relaxiert ( $s_i = [0, 1]$ )

$$\Rightarrow \sum x_i^2 \rightarrow \min_{Bx=f} \text{ keine gute Relaxierung}$$

Besser:

$$\sum s_i^2 \rightarrow \min_{\sum_i \sqrt{|x_i|}(1-s_i)=0}$$

$$L(s, p) = \sum s_i^2 + p \sum \sqrt{|x_i|}(1 - s_i)$$

$$0 = \frac{\partial L}{\partial s_i} = 2s_i - p\sqrt{|x_i|}, s_i = p/2\sqrt{|x_i|}$$

$$p/2 \sum |x_i| \rightarrow \min_{Bx=g}$$

Unter gewissen Bedingungen an B sind  $l^0$ - und  $l^1$ -Minimierungen äquivalent, (Candès, Tao, Donoho)

$F : R^n \rightarrow R$  konvex

$$\partial F(x) = w |F(x) + w^T(y - x) \leq F(y) \forall y \in R^n$$

[Bei  $x^2$  ist das zum Beispiel die Tangente und liegt immer unterhalb der Funktion ]

Es gilt für F stetig differenzierbar um x

$$\partial F(x) = F'(x)$$

$\partial F(x)$  heisst Subdifferential,  $w \in \partial F(x)$  heisst Subgradient

Beispiel:

$$F(x) = |x|$$

$$x \neq 0, \partial F(x) = \text{sgn}(x)$$

$$x = 0 : \partial F(0) = [-1, 1]$$

Bedingung:  $0 \in \partial F(\bar{x})$  charakterisiert Minimum

## 5.21 Lemma:

$F : R^n \rightarrow R \cup +\infty$  sei konvex.

$\bar{x} \in R^n$  ist ein Minimum von F  $\Leftrightarrow 0 \in \partial F(\bar{x})$

Beweis:  $0 \in \partial F(\bar{x}) \Leftrightarrow F(\bar{x}) + 0^T(y - \bar{x}) \leq F(y), y \in R^n$

$$\Leftrightarrow F(\bar{x}) \leq F(y) \forall y \in R^n$$

$\Leftrightarrow \bar{x}$  ist Minimum von F

$F : R^n \rightarrow R \cup +\infty$  ist nützlich bei Nebenbedingungen.

$G(x) \rightarrow \min_{x \in C}$ ,  $G : R^n \rightarrow R$ , konvex und  $C \subset R^n$  konvex, abgeschlossen

$$F(x) = \begin{cases} G(x) & x \in C \\ +\infty & x \notin C \end{cases}$$

$F(x) \rightarrow \min_{x \in R^n} \Leftrightarrow G(x) \rightarrow \min_{x \in C}$

$\|x\|_1 \rightarrow \min_{Bx=g}$

$$\partial\|x\|_1 = \begin{pmatrix} \partial|x_1| \\ \dots \\ \partial|x_n| \end{pmatrix}$$

$w \in \partial\|x\|_1 \Leftrightarrow w_i \in \partial|x_i| =: sgn(x_i)$

Lösung entspricht Sattelpunkt von  $L(x, p) = \|x\|_1 + p^T(Bx - g)$

$0 \in \partial_x L = \partial\|x\|_1 + B^T p$

$$0 = \frac{\partial L}{\partial p} = Bx - g$$

$$\tau x + \partial\|x\|_1 \ni \tau x - B^T p$$

$$\tau x^{k+1} + \partial\|x^{k+1}\|_1 \ni \tau x^k - B^T p^k$$

Optimalitätsbedingung für

$$\tau/2\|x^{k+1} - x^k + 1/\tau B^T p^k\|^2 + \|x^{k+1}\|_1 \rightarrow \min_{x^{k+1}}$$

$$Cp^{k+1} + Bx^{k+1} - g$$

Für den ersten Schritt müssen wir ein Problem der Form  $\tau/2\|x - f\|^2 + \|x\|_1 \rightarrow \min$

$$\tau/2 \sum_i (x_i - f_i)^2 + \sum_i |x_i| = \sum_i \tau/2(x_i - f_i)^2 + |x_i|$$

$\sum_i$  ist dann minimal, wenn für jedes  $i$ :  $\tau/2(x_i - f_i)^2 + |x_i| \rightarrow \min_{x_i}$

Optimalitätsbedingung:  $0 \in \tau(x_i - f_i) + \partial|x_i|$

1.Fall:  $x_i > 0$ :

$$\partial|x_i| = 1, \tau(x_i - f_i) + 1 = 0 \Leftrightarrow x_i = f_i - 1/\tau > 0$$

2.Fall  $x_i < 0$

$$\partial|x_i| = -1, \tau(x_i - f_i) - 1 = 0 \Leftrightarrow x_i = f_i + 1/\tau > 0$$

3.Fall:  $x_i = 0$

$$0 \in -\tau f_i + [-1; 1] \rightarrow f_i \in [-1/\tau; 1/\tau]$$

$$x = S_{1/\tau}(f_i) = \begin{cases} f_i - 1/\tau & f_i > 1/\tau \\ 0 & |f_i| \leq 1/\tau \\ f_i + 1/\tau & f_i < -1/\tau \end{cases}$$

das S heißt „Shrinkage“