NUMERIK III

Numerische Verfahren für Anfangs- und Anfangsrandwertaufgaben

Ulrich Langer

Institut für Mathematik Johannes Kepler Universität Linz

Vorwort

Das vorliegende Vorlesungsskriptum entstand aus Vorlesungen, die der Autor im Wintersemester 1994/95 und im Wintersemester 1996/97 an der Johannes Kepler Universität Linz gehalten hat. Die Lehrveranstaltung "Numerik III" ist die dritte Vorlesung in einem Zyklus von drei Vorlesungen zur "Höheren Numerischen Mathematik".

Die Vorlesung "Numerik I" stellt das Handwerkszeug zur numerischen Behandlung linearer und nichtlinearer Operatorgleichungen in Banach- bzw. Hilbert-Räumen bereit und gibt eine Einführung in die Theorie moderner Funktionenräume (Sobolev-Räume, Räume von Distributionen) [16]. In der Vorlesung "Numerik II" [17] werden Randwertaufgaben (RWA) für partielle Differentialgleichungen (PDgl.) und die wichtigsten Techniken (FEM, FDM, FVM) zu ihrer Diskretisierung betrachtet, sowie ein Überblick über moderne Verfahren zur Auflösung der bei der Diskretisierung entstehenden Gleichungssysteme gegeben (siehe auch Spezialvorlesung [15]). In der vorliegenden Vorlesung "Numerik III" werden Anfangswertaufgaben (AWA) und Anfangsrandwertaufgaben (ARWA) für gewöhnliche und partielle Differentialgleichungen und die wichtigsten Methoden zu ihrer numerischen Lösung betrachtet. Zur Vorlesung gehört das Praktikum "Zeitintegrationsverfahren" (siehe Anhang A). Im Praktikum werden Übungsaufgaben behandelt und ein praktisches Beispiel auf dem Rechner simuliert. Dieses praktische Beispiel wird durch ein Team von zwei Studenten in seiner Ganzheit (Modellierung, Analysis, numerische Analysis, Implementierung, Simulation, Ergebnisbewertung) bearbeitet. Die einzelnen Teams präsentieren ihre Ergebnisse in seminaristischer Form.

Die vorliegende Vorlesung setzt Kenntnisse aus den Grundvorlesungen zur linearen Algebra, zur Analysis und zur Numerischen Mathematik, sowie die in den Vorlesungen "Numerik I" [16] und "Numerik II" [17] vermittelten Lehrinhalte voraus. Zum anderen liefert die Lehrveranstaltung "Numerik III" Vorkenntnisse für nachfolgende Spezialvorlesungen zur Numerischen und Angewandten Mathematik als auch für Spezialvorlesungen zur Industriemathematik.

Das Skriptum wurde bewußt im Stile eines Vorlesungsmanuskriptes gehalten. Im Gegensatz zu vielen Lehrbüchern wurde auf "belletristische" Ausschmückungen verzichtet. Die Lehrinhalte sollen schnell und kompakt erfaßbar sein. Es wird eine Vielzahl von Abkürzungen eingeführt. Die Abkürzungen Ü x.x und (mms) bedeuten harte Arbeit an der Materie. Das Lösen von Übungsaufgaben und das "Mach-Mal-Selbst" ist angesagt. Das vorliegende Skriptum ist ein Arbeitspapier, es ist kein Ersatz für den Vorlesungsbesuch und auch kein Ersatz für ein Lehrbuch, aber eine gute Vorbereitung auf die allfällige Prüfung.

Der Autor möchte an dieser Stelle Frau Doris Holzer für die Erstellung des IAT_EX-Files und für die umfangreichen technischen Überarbeitungen recht herzlich danken.

Ulrich Langer

Linz, den 1. Februar 1997

Inhaltsverzeichnis

1	Ein	führun	ng	4
2	Par	abolis	che ARWA	6
	2.1	Gewic	chtete zweischichtige Differenzenschemata für die instationäre Wärmeleit-	
		gleich		6
		2.1.1	Konstruktion und Eigenschaften	6
		2.1.2	Lokale Approximationsordnung (Konsistenzordnung)	10
		2.1.3	Zur Stabilitätsproblematik	12
		2.1.4	Zusammenfassung: Approximation + Stabilität ⇒ diskrete Konvergenz	22
	2.2	Eine a	allgemeine Stabilitätstheorie für zweischichtige Schemata in der En-	
		ergien	orm (Energetische Methode)	24
		2.2.1	Allgemeine und kanonische Form	24
		2.2.2	Stabilität in der energetischen Norm $\ \cdot\ _A$	26
	2.3	Galerl	kin–FEM für parabolische ARWA	31
		2.3.1	Verallgemeinerte Formulierungen parabolischer ARWA	31
		2.3.2	Semidiskrete und volldiskrete Ersatzaufgabe	36
		2.3.3	Fehlerabschätzungen für die Lösung der semidiskreten Ersatzaufgabe	
			$(27)_{\mathbf{h}}$ in der $\mathbf{L_2}-$ und in der $\mathbf{W^1_2}-\mathrm{Norm}$	46
		2.3.4	Fehlerabschätzung für volldiskrete Ersatzaufgaben in der ${f L_2}-{ m Norm}$.	53
		2.3.5	Abschließende Bemerkungen	60
3	Ant	fangsw	ertaufgaben für gewöhnliche Differentialgleichungen und Systeme	
	\mathbf{gew}		her Dgl.	63
	3.1		ele	64
	3.2		ılierungen und analytische Resultate	69
	3.3	Einsch	nrittverfahren (= zweischichtig: t_j, t_{j+1})	73
		3.3.1	Das Eulersche Polygonenzugverfahren (EPZV)	73
		3.3.2	Explizite Runge-Kutta-Verfahren (ERKV)	80
		3.3.3	Implizite Runge–Kutta–Verfahren	90
		3.3.4	9	104
		3.3.5		114
	3.4	Mehrs	()	122
		3.4.1	1	123
		3.4.2		130
		3.4.3	Stabilität linearer MSV	134
		3.4.4	Konvergenz linearer MSV	137

	3.5	Steife Differentialgleichung	141		
		3.5.1 A -Stabilität (\longrightarrow ESV)	143		
		3.5.2 L–Stabilität	146		
		3.5.3 B–Stabilität	147		
		3.5.4 A-Stabilität von MSV	149		
4	Nu	merische Behandlung hyperbolischer ARWA	151		
	4.1	Dreischichtige Differenzenschemata für hyperbolische ARWA	151		
		4.1.1 Allgemeine und kanonische Form dreischichtiger Schemata	151		
		4.1.2 Ein allgemeines Stabilitätsresultat für dreischichtige Schemata	152		
		4.1.3 Beispiel:Gewichtete dreischichtige Differenzenschemata für die Saiten-			
		schwingungsgleichung	156		
	4.2	Zurückführung auf AWA mittels Semidiskretisierung	158		
5	Praktikum				
			16 4		
	5.1	Einführung	165		
	5.1 5.2				
		Einführung	165		
		Einführung	$\frac{165}{168}$		
		Einführung Numerische Behandlung parabolischer ARWA 5.2.1 Differenzenverfahren 5.2.2 FEM-Galerkin-Verfahren für parabolische ARWA	168 168 168		
		Einführung	168 168 168 178		
	5.2	Einführung Numerische Behandlung parabolischer ARWA 5.2.1 Differenzenverfahren 5.2.2 FEM-Galerkin-Verfahren für parabolische ARWA	168 168 168 178		
	5.2	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	165 168 168 178 180 181		
	5.2	Einführung Numerische Behandlung parabolischer ARWA 5.2.1 Differenzenverfahren 5.2.2 FEM-Galerkin-Verfahren für parabolische ARWA 5.2.3 Konsultation zur Praktikumsaufgabe Px Numerische Behandlung von AWA für gewöhnliche Dgl. 5.3.1 Einschrittverfahren	165 168 168 178 180 181		
	5.2	Einführung Numerische Behandlung parabolischer ARWA 5.2.1 Differenzenverfahren 5.2.2 FEM-Galerkin-Verfahren für parabolische ARWA 5.2.3 Konsultation zur Praktikumsaufgabe Px Numerische Behandlung von AWA für gewöhnliche Dgl. 5.3.1 Einschrittverfahren 5.3.2 Praktische Durchführung von Einschrittverfahren	165 168 168 178 180 181 183		

Kapitel 1

Einführung

■ Ziel

der Vorlesung NUMERIK III ist das Kennenlernen von Handwerkszeug zur Analysis und zur numerischen Behandlung von

• Anfangswertaufgaben (AWA) für Systeme gewöhnlicher Differentialgleichungen (Dgl.) der Art

(1)
$$\operatorname{Ges.} u(t) = (u_1(t), \dots, u_n(t))^T :$$

$$\operatorname{Dgl.:} \dot{u}(t) = f(t, u(t)), \ t \in I = \bar{T} = [0, T],$$

$$\operatorname{AB:} u(0) = u_0 \text{ (Anfangsbedingung)}.$$

• Anfangsrandwertaufgaben (ARWA) für parabolische <u>partielle Differentialgleichungen</u> (PDgl.) der Art. (→ siehe [17] Numerik II, Pkt. 1.1):

Ges.
$$u(x,t)$$
:
$$\frac{\partial u}{\partial t} + \underbrace{L[u(x,t)]}_{\uparrow} = f(x,t) \quad \forall (x,t) \in Q_T = \Omega \times \mathbf{T},$$
linearer [nichtlinearer]
elliptischer Operator
+ RB: 1. - 4. Art,
+ AB: $u(x,0) = u_0(x), x \in \bar{\Omega}.$

Die Semidiskretisierung von (2) durch die sogenannte (vertikale) Linienmethode

(= Ortsdiskretisierung durch Galerkin-Ansatz $u(x,t) = \sum_{i=1}^{n} u_i(t)p_i(x)$) führt unmittelbar auf (1) zur Bestimmung der zeitabhängigen Koeffizienten $u_i(t), \ldots, u_n(t)$.

• ARWA für hyperbolische PDgl. der Art (→ siehe [17] Numerik II, Pkt. 1.2):

(3)
$$\begin{aligned}
&\operatorname{Ges.} \ u(x,t) : \\
&\frac{\partial^{2} u}{\partial t^{2}} + L \ u(x,t) = f(x,t) \quad \forall (x,t) \in Q_{T} \\
&+ \operatorname{RB:} 1. - 4. \operatorname{Art} \\
&+ \operatorname{AB:} \quad \frac{u(x,0)}{\partial t} = u_{0}(x) \\
&+ \frac{\partial u}{\partial t}(x,0) = u_{1}(x)
\end{aligned} \right\} x \in \overline{\Omega}$$

• Erhaltungsgleichungen (örtlich 1D):

$$\frac{\partial u}{\partial t} + \frac{\partial}{\partial x} [f(u)] = 0 \text{ in } Q_T = \mathbb{R}^1 \times \mathbf{T}$$
$$+ AB: u(x, 0) = u_0(x) \text{ in } \mathbb{R}^1.$$

■ <u>Inhalt:</u>

- ARWA für parabolische PDgl. (Kapitel 2)
- AWA für (Systeme) gewöhnliche Dgl. (Kapitel 3)
- ARWA für hyperbolische PDgl. (Kapitel 4)

■ Wichtigste Literaturquellen:

- zu Kap. 2: [4], [20], [21].
- zu Kap. 3: [1], [5], [6], [19].
- zu Kap. 4: [4], [20].
- zu Erhaltungsgleichungen: [4], [13], [20].
- <u>Praktikum:</u> zur Vorlesung Numerik III zum Thema "Zeitintegrationsverfahren"

 → siehe Anhang A.

Kapitel 2

Numerische Behandlung parabolischer Anfangsrandwertaufgaben

- 2.1 Gewichtete zweischichtige Differenzenschemata für die instationäre Wärmeleitgleichung
- 2.1.1 Konstruktion und Eigenschaften
 - Btr. der Einfachheit halber zunächst die örtlich 1D, instationäre Wärmeleitgleichung in differentieller Form (siehe [17] Numerik II, Pkt. 1.1.2)
 - = Ausgangsaufgabe $\stackrel{\text{hier}}{=}$ 1. ARWA:

Ges.
$$u(x,t) : c\rho \frac{\partial u}{\partial t} - \frac{\partial}{\partial x} \left(\lambda \frac{\partial u}{\partial x}\right) + qu = f(x,t), \quad (x,t) \in Q_T = (a,b) \times (0,T)$$

$$+ \text{RB:} \quad u(a,t) = g_a(t) \\ u(b,t) = g_b(t) \quad t \in \mathbf{T} = (0,T)$$

$$+ \text{AB:} \quad u(x,0) = u_0(x), \quad x \in \bar{\Omega} = [a,b]$$

- Btr. Spezialfall: $c, \rho, \lambda = \text{const.} > 0, q = 0.$
 - Mit $k^2 = \lambda/c\rho$ und $\bar{f} = f/c\rho$ folgt aus (1):

$$\frac{\partial u}{\partial t} - k^2 \frac{\partial^2 u}{\partial x^2} = \bar{f}(x, t), \quad (x, t) \in Q_T$$
+ RB + AB

• O. B. d. Allg.: $k^2 = 1$, a = 0, b = 1. Tatsächlich:

$$\begin{vmatrix}
\tilde{x} &= \frac{x-a}{b-a} = \frac{x-a}{l}, & \tilde{t} &= \frac{k^2}{l^2}t \\
\frac{\partial}{\partial x} &= \frac{\partial}{\partial \tilde{x}} \frac{\partial \tilde{x}}{\partial x} = \frac{1}{l} \frac{\partial}{\partial \tilde{x}}, & \frac{\partial}{\partial t} = \frac{k^2}{l^2} \frac{\partial}{\partial \tilde{t}}
\end{vmatrix} \Rightarrow \frac{\partial u}{\partial \tilde{t}} - \frac{\partial^2 u}{\partial \tilde{x}^2} = \frac{l^2}{k^2} \bar{f} =: \tilde{f}$$

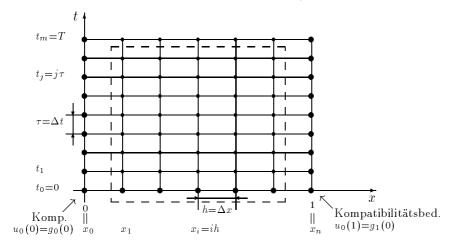
■ Nach eventueller Umbenennung: $\tilde{x} \to x$, $\tilde{t} \to t$, $\tilde{f} \to f$, $\tilde{T} = \frac{k^2}{l^2} T \mapsto T$:

(1) Ges.
$$u(x,t): \frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = f(x,t), \quad x \in (0,1), \quad t \in T,$$

$$+ \text{RB:} \quad \frac{u(0,t) = g_0(t)}{u(1,t) = g_1(t)} \right\} \quad t \in T = (0,T),$$

$$+ \text{AB:} \quad u(x,0) = u_0(x), \quad x \in [0,1].$$

■ <u>Gitterkonstruktion:</u> Raum – Zeit – Gitter (vgl. [17] Numerik II, Pkt. 5.1.1)



$$\begin{array}{ll} \omega_{h\tau} = \omega_h \times \omega_\tau, \\ \omega_h &= \{x_i = ih : i = \overline{1, n-1}\}, \ h = 1/n, \\ \omega_\tau &= \{t_j = j\tau : j = \overline{1, m-1}\}, \ \tau = T/m, \\ \bar{\omega}_\tau &= \{t_j = j\tau : j = \overline{0, m-1}\}. \end{array}$$

■ Bez. von Gitterfunktion:

$$\begin{split} v: \bar{\omega}_{h\tau} &= \bar{\omega}_h \times \bar{\omega}_\tau \longmapsto I\!\!R^1, \quad v_i^j = v\left(x_i, t_j\right); \\ v &= v^j = v\left(\cdot, t_j\right): \bar{\omega}_h \mapsto I\!\!R^1: \text{Gitterfkt. auf } j\text{-ter Zeitschicht}; \\ v: \bar{\omega}_\tau &\longmapsto \left[\bar{\omega}_h \mapsto I\!\!R^1\right] - \text{abstrakte Gitterfkt.}; \\ \hat{v} &= v^{j+1}; \quad \check{v} = v^{j-1}; \end{split}$$

 $\bar{u} = u^{j+0.5} = u\left(\cdot, t_j + \frac{\tau}{2}\right) : \bar{\omega}_h \mapsto I\!\!R^1$, z.B. definiert für $u \in C(\bar{Q}_T)$.

■ explizites Schema

rein implizites Schema

$$v_{t} - v_{\bar{x}x} = \varphi\left(x, t\right)$$

$$x = x_{i} = ih, \ i = \overline{1, n - 1}; \quad t = t_{j} = j\tau, \ j = \overline{0, m - 1}$$

$$\bullet \text{ RB: } v_{0}^{j} = g_{0}(t_{j}), \ t_{j} = j\tau, \ j = \overline{1, m}$$

$$v_{n}^{j} = g_{1}(t_{j}), \ t_{j} = j\tau, \ j = \overline{1, m}$$

$$\bullet \text{ AB: } v_{i}^{0} = u_{0}(x_{i}), \ x_{i} = ih, \ i = \overline{0, n}$$

$$\bullet \text{ Wahl von } \varphi : \varphi = f, \ = \overline{f}, \ = (f + \hat{f})/2 ?$$

$$v_{i}^{j+1} - v_{i}^{j} - v_{i+1}^{j} - 2v_{i}^{j} + v_{i-1}^{j}}{h^{2}} = \varphi_{i}^{j}$$

$$i = 1, 2, \dots, n - 1$$

$$j = 0, 1, \dots, m - 1$$

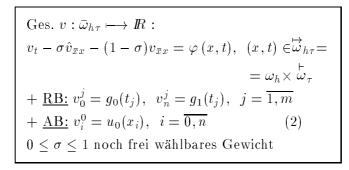
$$+ \text{ RB } (\bullet) + \text{ AB } (\bullet)$$

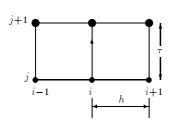
 σ -gewichtetes zweischichtiges DS (=Einschrittverfahren)

■ Kompromiß: Familie von DS = σ -gewichtetes zweischichtiges DS



- Stabilität
- Aufwand





■ Bemerkung:

- * Kompatibilitätsbedingung: $g_0(t_0)=u_0(x_0)=v_0^0$ $g_1(t_0)=u_0(x_n)=v_n^0$
- * Im allgemeinen Fall (veränderliche Koeffizienten, 2D, 3D etc.) wird Ortsdiskretisierung mit Integralbilanzmethode ([17] Numerik II, Pkt. 5.2) bzw. <u>FEM-Linien-</u> methode (siehe Pkt. 2.3) vorgenommen.

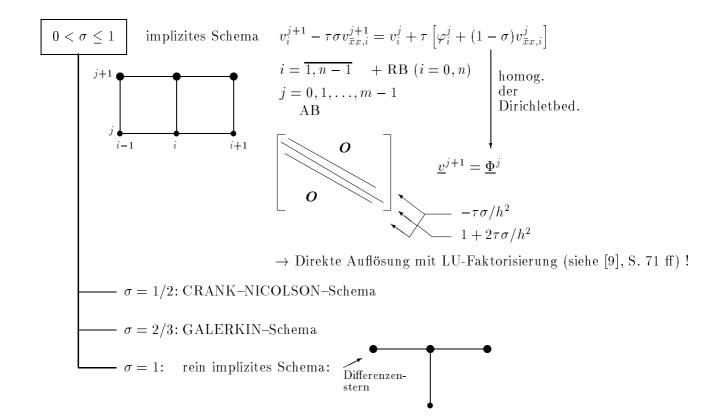
■ Spezialfälle:

explizites Schema
$$v_i^{j+1} = v_i^j + \tau \left[\varphi_i^j + \frac{1}{h^2} \left(v_{i-1}^j - 2v_i^j + v_{i+1}^j \right) \right]$$

$$i = \overline{1, n-1} + \text{RB } (i = 0, n)$$

$$j = 0, 1, \dots, m-1$$
AB

9



- Wahl von σ : in Abhängigkeit von h, τ (\downarrow)
 - Approximation (- ordnung)
 - Stabilität
 - Aufwand, insbes. bei örtl. mehrdim. Probl.

2.1.2 Lokale Approximationsordnung (Konsistenzordnung)

■ Fehlerschema: $z = \stackrel{(1)}{u} - \stackrel{(2)}{v} : \bar{\omega}_{h\tau} \longmapsto \mathbb{R}^1 :$

(3)
$$z_{t} - \sigma \hat{z}_{\overline{x}x} - (1 - \sigma)z_{\overline{x}x} = (u_{t} - \sigma \hat{u}_{\overline{x}x} - (1 - \sigma)u_{\overline{x}x}) - \varphi =: \psi(x, t)$$

$$(x, t) \in \overleftrightarrow{\omega}_{h\tau} = \omega_{h} \times \overleftrightarrow{\omega}_{\tau};$$

$$RB: z_{0}^{j} = 0, \quad z_{n}^{j} = 0, \quad j = \overline{0, m};$$

$$AB: z_{i}^{0} = 0, \quad i = \overline{0, n}.$$

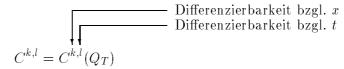
 $=\frac{h^2}{12}\frac{\partial^2 f}{\partial x^2} + O(h^4)$, falls $f \in C^{4,0}$ $(\Rightarrow u \in C^{6,1})$

■ <u>Lemma 2.1.:</u> Zusammenfassung der Approximationsresultate:

DS	Vor.	σ	φ	$\psi(x,t)$
Schema bester Approx.	$u \in C^{6,3}$	$\sigma = \sigma_* = \frac{1}{2} - \frac{h^2}{12\tau}$	$\varphi = \bar{f} + \frac{h^2}{12} \bar{f}_{\bar{x}x} \text{ oder}$ $= \bar{f} + \frac{h^2}{12} \frac{\partial^2 \bar{f}}{\partial x^2}$	$\psi = O(h^4 + \tau^2)$
CRANK- NICOLSON	$u \in C^{4,3}$	$\sigma = \frac{1}{2}$	$\varphi = \bar{f}, = \frac{1}{2}(\hat{f} + f)$	$\psi = O(h^2 + \tau^2)$
	$u \in C^{4,2}$	$0 \le \sigma \le 1$	$\varphi = \bar{f}, = f, = \hat{f}, \dots$	$\psi = O(h^2 + \tau)$
explizites Schema	$u \in C^{4,2}$	$\sigma = 0$	$\varphi = \bar{f}, = f, = \hat{f}, \dots$	$\psi = O(h^2 + \tau)$

Beweis: folgt direkt aus der obigen Darstellung des Approximationsfehlers.





q.e.d.

2.1.3 Zur Stabilitätsproblematik

■ Btr. folgendes Schema mit homogenen RB:

(4)
$$z_{t} - \sigma \hat{z}_{\bar{x}x} - (1 - \sigma)z_{\bar{x}x} = \psi(x, t), \quad (x, t) \in \stackrel{\hookrightarrow}{\omega}_{h\tau} = \omega_{h} \times \stackrel{\vdash}{\omega}_{\tau};$$

$$RB: z_{0}^{j} = 0, \quad z_{n}^{j} = 0, \quad j = 0, 1, \dots, m;$$

$$AB: z_{i}^{0} = w_{i}^{0} \qquad , \quad i = \overline{1, n - 1};$$

Bemerkung 2.2.:

1. Im Fehlerschema (3): \Longrightarrow AB: $z_i^0 = 0$ (homogen), aber dies setzt voraus, daß exakt gilt:

$$v_i^0 = u_0(x_i), \quad i \in \omega_h \qquad (x_i \in \omega_h).$$

<u>praktisch:</u> Rundungsfehler: $v_i^0 = u_0(x_i) - w_i^0$, $i \in \omega_h$. <u>interessant:</u> Untersuchungen der Fortpflanzung solcher Fehler in den AB: \Longrightarrow Stabilität bzgl. AB! 2. RB (1. Art) können o. B. d. Allg. im Fehlerschema homogen angenommen werden, denn durch Homogenisieren

$$u = U + g_0(t) + x(g_1(t) - g_0(t))$$

$$v = V + \tilde{q}_0(t) + x(\tilde{q}_1(t) - \tilde{q}_0(t))$$

kann das Problem der Stabilität bzgl. RB auf die Stabilität bzgl. RS und Stabilität bzgl. AB zurückgeführt werden!

■ <u>Def. 2.3.:</u> (Stabilität bzgl. AB und RS)

Das DS (4) heißt stabil (bzgl. AB z^0 und RS ψ , sowie den gewählten Normen $\|\cdot\|_{(1)}$ und $\|\cdot\|_{(*,j)}!$), falls

(5)
$$||z^{j+1}||_{(1)} \le c_1 ||z^0||_{(1)} + c_2 ||\Psi||_{(*,j)} \quad \forall j = 0, 1, \dots, m-1,$$

wobei $c_1, c_2 = \text{const.} > 0 : c_{\alpha} \neq c_{\alpha}(h, \tau, j, \psi, z^0);$

$$\Psi = (\psi^0, \psi^1, \dots, \psi^j),$$

 $\|\psi\|_{(*,j)}$ – geeignet gewählte Norm:

z.B.:
$$\|\psi\|_{(*,j)} := \max_{k=0,j} \|\psi^k\|_{(2)};$$

 $\|\cdot\|_{(1)}, \|\cdot\|_{(2)}$ – geeignet gewählte Normen auf den

Zeitschichten, d.h. für Gitterfunktionen

$$v:\omega_h\longmapsto I\!\!R^1.$$

Spezialfälle: Stabilität bzgl. AB: $\psi = 0$ in (4) \Rightarrow A-priori Absch. (5)

Stabilität bzgl. RS: $z^0 = 0$ in (4) \Rightarrow A-priori Absch. (5)

- Bemerkung: "geeignet gewählte" Normen
 - soll heißen, die Normen werden eigentlich definiert durch die vom Anwender gewünschten "Konvergenzaussagen"
 - d.h. lok. Approximations ordnung $O(h^p + \tau^q)$

globale Approximationsordnung
$$O(h^p + \tau^q)$$

globale Approximations ordnung
$$O(h^p + \tau^q)$$
, d.h. $\|\psi\|_{(*,j)} \le c_{A,1}(u)h^p + c_{A,2}(u)\tau^q \le c_A(u)(h^p + \tau^q)$

$$\bigoplus_{\text{Stabilität}} \text{ d.n. } \|\psi\|_{(*,j)} \leq c_{A,1}(u)h^p + c_{A,2}(u)\tau^q \leq c_A(u)(h^p + \tau^q)$$
Stabilität

Def 2.2

Def. 2.3. diskrete Konvergenz: $||z^{j+1}||_{(1)} \le c_1 ||w^0||_{(1)} + c_2 c_A(u) (h^p + \tau^q)$.

■ Wir untersuchen Stabilität in folgenden diskreten Normen:

• L₂-Norm: Fourier-Analyse (v. Neumann),

• Energienorm: Energetische Methode (siehe Pkt. 2.2),

• C-Norm: Diskretes Maximumprinzip.

■ Stabilität in diskreten L₂ –Normen:

⇒ Fourier-Analyse nach den Eigenfunktionen des elliptischen Anteils:

• Bezeichnung:

$$\overline{y : \omega_h \longmapsto \mathbb{R}^1} \quad ; \quad y \in L_2(\omega_h) = \overset{\circ}{L_2} (\bar{\omega}_h) :
\|y\|^2 = \|y\|^2_{L_2(\omega_h)} = (y, y)_h := \sum_{i \in \omega_h} h \, y_i^2 = \sum_{i=1}^{n-1} h \, y_i^2
\|\cdot\|_{(1)} = \|\cdot\|_{(2)} := \|\cdot\| = \|\cdot\|_{L_2(\omega_h)} \quad \text{(vgl. Def. 2.3.)}$$

• Btr. EWP (siehe auch PI)

$$-v_{\overline{x}x} = \lambda v(x), \quad x \in \omega_h$$

$$v_0 = v_n = 0$$

$$\Longrightarrow \qquad \underbrace{\frac{\text{EW: } \lambda_k = \frac{4}{h^2} \sin^2 \frac{k\pi h}{2l}, \ k = \overline{1, n-1}}{\text{EFkt.: } \mu_k(x) = \sqrt{\frac{2}{l}} \sin \frac{k\pi x}{l}, \ x \in \overline{\omega}_h, \ k = \overline{1, n-1}}$$

$$\Longrightarrow \qquad \underbrace{\frac{\text{EW: } \lambda_k = \frac{4}{h^2} \sin^2 \frac{k\pi h}{2l}, \ k = \overline{1, n-1}}{\text{EFkt.: } \mu_k(x) = \sqrt{\frac{2}{l}} \sin \frac{k\pi x}{l}, \ x \in \overline{\omega}_h, \ k = \overline{1, n-1}}}_{\text{Es gilt: } \{\mu_k\}_{k = \overline{1, n-1}}} \text{ voll. ONS:}$$

$$\frac{1}{h^{2}}\begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & \mathbf{O} & \\ & \ddots & \ddots & \ddots & \\ & \mathbf{O} & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} \begin{bmatrix} v_{1} \\ v_{2} \\ \vdots \\ v_{n-2} \\ v_{n-1} \end{bmatrix} = \lambda \begin{bmatrix} v_{1} \\ v_{2} \\ \vdots \\ v_{n-2} \\ v_{n-1} \end{bmatrix}$$
• $(\mu_{k}, \mu_{m})_{h} = \delta_{km}$
• Fourier-Zerlegung:
$$y(x) = \sum_{k=1}^{n-1} (y, \mu_{k})_{h} \mu_{k}(x)$$
• PARSEVALsche Gleichung:
$$||y||^{2} = \sum_{n=1}^{n-1} (y, \mu_{k})_{h}^{2}$$

• Fourier-Zerlegung:

$$y(x) = \sum_{n=1}^{n-1} (y, \mu_k)_h \mu_k(x)$$

$$\approx -u''(x) = \lambda u(x), \ x \in (0,1)$$
$$u(0) = u(1) = 0$$

• Entwickeln "Fehler" $z=z(\cdot,t_i)$ auf der j-ten Zeitschicht in eine Fourierreihe nach den Eigenfkt. $\{\mu_k\}_{k=\overline{1,n-1}}$:

$$z = z(x, t_j) = \sum_{k=1}^{n-1} c_k(t_j) \mu_k(x), \quad x = x_i \in \bar{\omega}_h$$
Fourierkoeff.: $c_k = c_k(t_j) = (z(\cdot, t_j), \mu_k(\cdot))_h$

und analog die RS $\psi = \psi(x, t_j)$ des DS (4):

$$\psi = \psi(x, t_j) = \sum_{k=1}^{n-1} d_k(t_j) \mu_k(x) = \sum_{n=1}^{n-1} d_k \mu_k, \quad x = x_i \in \omega_h.$$
Fourierkoeff.: $d_k = (\psi, \mu_k)_h = (\psi(\cdot, t_j), \mu_k(\cdot))_h$

• Btr. nun DS (4) mit diesen Fourierentwicklungen:

$$\underline{\text{Vor.:}} \qquad \boxed{|q_k| \le 1 \quad \forall k = 1, 2, \dots, n-1} \tag{7}$$

$$\|\hat{z}\| = \sqrt{\sum_{k=1}^{n-1} \hat{c}_k^2} \le \sqrt{\sum_{k=1}^{n-1} q_k^2 c_k^2} + \tau \sqrt{\sum_{k=1}^{n-1} \frac{d_k^2}{(1 + \sigma \tau \lambda_k)^2}} \le \uparrow \tag{6}$$

$$\text{PARSEVAL}$$

$$\stackrel{(7)}{\le} \sqrt{\sum_{k=1}^{n-1} c_k^2} + \tau \sqrt{\sum_{n=1}^{n-1} d_k^2} \le \|z\| + \tau \|\psi\|.$$

$$1 + \sigma \tau \lambda_k \ge 1 \qquad \uparrow$$

$$\text{PARSEVAL}$$

• Resultat:
$$||z^{j+1}|| \le ||z^{j}|| + \tau ||\psi^{j}|| \le$$

$$\le ||z^{j-1}|| + \tau \{||\psi^{j-1}|| + ||\psi^{j}||\} \le \dots \le$$

$$\le ||z^{0}|| + \tau \sum_{k=0}^{j} ||\psi^{k}|| \le$$

$$\le ||z^{0}|| + \tau \sum_{k=0,j}^{j} \max_{l=0,j} ||\psi^{l}||$$

$$= (j+1)\tau = t_{j+1} \le T$$

$$\le ||z^{0}|| + T \max_{k=0,j} ||\psi^{k}||$$

$$= : ||\Psi||_{(*,j)}$$

d.h. Stabilität bzgl. AB und RS mit
$$\|\cdot\|_{(1)} = \|\cdot\|_{(2)} = \|\cdot\|$$
, $c_1 = 1$, $c_2 = T$, falls $|q_k| \le 1 \quad \forall k = 1, 2, \dots, n-1$ (vgl. Def. 2.3.).

- Frage: $|q_k| \le 1$? $\forall k = 1, 2, ..., n-1$:
 - (a) $q_k = \frac{1 \tau (1 \sigma) \lambda_k}{1 + \tau \sigma \lambda_k} < 1$ (gilt automatisch),

(b)
$$-1 \le q_k \iff 0 \le 1 + q_k = \frac{1 + \tau \sigma \lambda_k + 1 - \tau (1 - \sigma) \lambda_k}{1 + \tau \sigma \lambda_k}$$

 $\iff 2 - \tau \lambda_k + 2\tau \sigma \lambda_k \ge 0$

$$\iff$$
 $\sigma \ge \frac{\tau \lambda_k - 2}{2\tau \lambda_k} = \frac{1}{2} - \frac{1}{\tau \lambda_k}, \quad \forall k = 1, 2, \dots, n - 1.$

• Resultat:
$$\sigma \ge \frac{1}{2} - \frac{1}{\tau \lambda_k}$$

$$k = 1, 2, \dots, n-1$$
 \iff
$$|q_k| \le 1$$

$$k = 1, 2, \dots, n-1$$

$$\iff$$

$$|q_k| \le 1$$
$$k = 1, 2, \dots, n-1$$

• Aus
$$0 < \frac{8}{l^2} \le \lambda_1 < \lambda_2 < \dots < \lambda_{n-1} = \lambda_{\max} \left(\frac{1}{h^2} \left[\bigotimes \right] \right) < \frac{4}{h^2}$$
 folgt:
$$= \frac{4}{h^2} \sin^2 \frac{(n-1)\pi h}{2}$$

Lemma 2.4.: (L_2 -Stabilität)

• <u>Def. 2.5.</u>: (unbedingte und bedingte Stabilität)

- DS heißt unbedingt stabil, wenn Stabilität vorliegt, unabhängig von den Beziehungen zwischen den Gitterparametern (h, τ) .
- Andernfalls sprechen wir von <u>bedingter Stabilität</u>, d.h. Stabilität unter Zusatzbedingungen an die Gitterparameter.

• Bemerkung 2.6.:

• Bemerkung 2.7.:

Die v. Neumannsche Fourier-Analyse nach den EFkt. des elliptischen Anteils ist immer dann anwendbar, wenn der elliptische Anteil <u>symmetrisch</u> und <u>zeit-unabhängig</u> ist (vgl. auch Pkt. 2.3.).

Aus der oben durchgeführten Analyse ist ersichtlich, daß sich eine Störung in den AB der Form

$$z^0 = w^0 := c\mu_k$$

mit dem Faktor q_k von Zeit- zu Zeitschicht fortpflanzt, d.h. im Falle einer homogenen RS ($\psi \equiv 0$) gilt:

 $z^{j+1} = q_k^{j+1} c \mu_k = q_k^{j+1} z^0.$

Für eine <u>formale</u> Stabilitätsanalyse untersucht man deshalb die Fortpflanzung der harmonischen Störung

$$z_s^0 = e^{i\lambda sh} = \cos(\lambda sh) + i\sin(\lambda sh)$$

mit dem Ansatz

$$z_s^j = (e^{\alpha \tau})^j e^{i\lambda sh}$$

und aus dem betr. DS (ohne RB !) zu bestimmenden α . Das DS ist stabil, falls

$$|e^{\alpha\tau}| \le 1 \quad .$$

Für das explizite Schema mit homogener RS und ohne Berücksichtigung der RB

$$z_s^{j+1} = (1 - 2\gamma)z_s^j + \gamma(z_{s-1}^j + z_{s+1}^j), \quad \gamma = \frac{\tau}{h^2},$$

erhalten wir z.B.

$$(e^{\alpha\tau})^{j+1}e^{i\lambda sh} = (1-2\gamma)(e^{\alpha\tau})^{j}e^{i\lambda sh} + \gamma(e^{\alpha\tau})^{j}\left(e^{i\lambda(s-1)h} + e^{i\lambda(s+1)h}\right)$$

$$e^{\alpha\tau} = (1-2\gamma) + \gamma(e^{-i\lambda h} + e^{i\lambda h}) =$$

$$= 1-2\gamma + \gamma\left(\cos(-\lambda h) + i\sin(-\lambda h) + \cos(\lambda h) + i\sin(\lambda h)\right)$$

$$= 1-2\gamma + 2\gamma\cos(\lambda h)$$

$$= 1-2\gamma\underbrace{(1-\cos(\lambda h))}_{2\sin^{2}\frac{\lambda h}{2}}.$$

Die v. Neumannsche Stabilitätsbedingung $|e^{\alpha\tau}| \leq 1$ ist dann äquivalent zu

$$|e^{\alpha \tau}| = |1 - 2\gamma \underbrace{(1 - \cos(\lambda h))}_{= 2\sin^2 \frac{\lambda h}{2}}| = |1 - 4\gamma \sin^2 \frac{\lambda h}{2}| \le 1, \text{ d.h.}$$

$$4\gamma \underbrace{\sin^2 \frac{\lambda h}{2}}_{\le 1} \le 2.$$

Diese Bedingung ist für $\gamma = \frac{\tau}{h^2} \le \frac{1}{2}$ erfüllt $(\forall \lambda)$!

Das allgemeinere v. Neumannsche Stabilitätskriterium

$$(*) |e^{\alpha \tau}| \le 1 + c\tau$$

mit $c=\mathrm{const.}>0,\ c\neq c(\tau,h)$ läßt exponentielles Wachstum zu:

$$|(e^{\alpha\tau})^m| = |e^{\alpha m\tau}| \le \underbrace{\left(1 + \frac{cT}{m}\right)^{\frac{m}{cT} \cdot cT}}_{\tau = \frac{T}{m}} \le e^{cT}$$

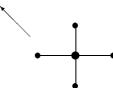
Das Kriterium (*) ist notwendig und hinreichend dafür, daß \exists positive Konstante $K = \text{const.} \geq 1, \ K \neq K(\tau) : |e^{\alpha \tau m}| \leq K. \quad (**)$

- Führen Sie die formale Stabilitätsanalyse nach v. Neumann für die folgenden DS (mit homogener RS) durch:
 - (a) rein implizites Schema,

(b) CRANK-NICOLSON-Schema,



(c) Leapfrog-Schema (=Richardson-Schema) $z_s^{j+1} - z_s^{j-1} - 2\gamma(z_{s-1}^j - 2z_s^j + z_{s+1}^j) = 0; \quad \gamma = \tau/h^2$ (= Dreischichtiges DS = Zweischrittverfahren)



- (d) duFort-Frankel-Schema (= modifiziertes Leapfrog-Schema: $2z_s^j \to z_s^{j+1} + z_s^{j-1}$): $z_s^{j+1} - z_s^{j-1} - 2\gamma \left[z_{s-1}^{j} - (z_s^{j+1} + z_s^{j-1}) + z_{s+1}^{j} \right] = 0$
 - (= Dreischichtiges DS = Zweischrittverfahren).
- $\ddot{\mathrm{U}}$ 2.2 Untersuchen Sie die lokale Approximationsordnung (= Konsistenzordnung) des Leapfrog-Schemas (c) und des duFort-Frankel-Schemas (d) für die homogene Wärmeleitgleichung

$$\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = 0, \ x \in (0, 1), \ t \in \mathbf{T} = (0, T),$$
+ RB: $u(0, t) = g_0(t), \ u(1, t) = g_1(t), \ t \in \mathbf{T},$
+ AB: $u(x, 0) = u_0(x), \ x \in [0, 1].$

Bemerkung: Stabilitäts- und Konsistenzuntersuchungen von (c) und (d) siehe auch $\boxed{P\ II}$: $\ddot{U}\ 5-\ddot{U}\ 7$.

■ Stabilität in diskreten C-Normen:

- \implies Diskretes Maximumprinzip: \implies [17] Numerik II, Pkt. 5.1.4!
 - Btr. DS (4) in der Form

$$\hat{z} - \tau \sigma \hat{z}_{\bar{x}x} = z + \tau (1 - \sigma) z_{\bar{x}x} + \tau \psi =: \eta(x, t)$$

für $x = x_i \in \omega_h$ und $t = t_j \in \overset{\vdash}{\omega}_{\tau}$, d.h.

(8)
$$-\sigma \gamma z_{i-1}^{j+1} + (1 + 2\sigma \gamma) z_i^{j+1} - \sigma \gamma z_{i+1}^{j+1} = \eta_i^j,$$

mit $\gamma = \tau/h^2$ und

(9)
$$\eta_i^j = (1 - \sigma)\gamma z_{i-1}^j + (1 - 2(1 - \sigma)\gamma) z_i^j + (1 - \sigma)\gamma z_{i+1}^j + \tau \psi_i^j.$$

• DS (8) ist offenbar streng monoton mit

$$D(x) := -\sigma \gamma \cdot 1 + (1 + 2\sigma \gamma) \cdot 1 - \sigma \gamma \cdot 1 \equiv 1.$$

Aus Satz II.5.1 aus [17] Numerik II folgt dann sofort:

(10)
$$||z^{j+1}||_{C(\omega_h)} := \max_{i=\overline{1,n-1}} |z_i^{j+1}| \le$$

$$(C(\overline{\omega}_h)) \quad (i = \overline{0,n})$$

$$\le \frac{1}{\min_{x \in \omega_h} |D(x)|} ||\eta^j||_{C(\omega_h)} = \max_{i=\overline{1,n-1}} |\eta_i^j|.$$

Setzt man

$$1 - 2(1 - \sigma)\gamma \ge 0$$

voraus, so erhält man wegen 0 $\leq \sigma \leq$ 1 aus (9) und (10) die Abschätzung

$$\max_i |z_i^{j+1}| \leq \max_i |\eta_i^j| \leq \max_i |z_i^j| + \tau \max_i |\psi_i^j|,$$

deren rekursive Auswertung die gewünschte Stabilitätsabschätzung ergibt:

$$||z^{j+1}||_{C(\omega_h)} \le ||z^0||_{C(\omega_h)} + \tau \sum_{k=0}^{j} ||\psi^k||_{C(\omega_h)}$$

$$\le ||z^0||_{C(\omega_h)} + T \underbrace{\max_{k=0,j} ||\psi^k||_{C(\omega_h)}}_{=: ||\psi||_{(*,j)}}$$

• Lemma 2.8.: (C-Stabilität)

$$\underline{\text{Vor.:}} \ 1 - 2 (1 - \sigma) \gamma \ge 0, \quad \gamma = \tau/h^2, \quad 0 \le \sigma \le 1
\underline{\text{Bh.:}} \ \|z^{j+1}\|_{C(\omega_h)} := \max_{i=\overline{1,n-1}} |z_i^{j+1}| \le \|z^0\|_{C(\omega_h)} + T \max_{k=\overline{0,1}} \|\psi^k\|_{C(\omega_h)}$$

• Bemerkung 2.9.:

 $\sigma=1 \, \Rightarrow$ rein implizites Schema ist unbedingt stabil in der diskreten C–Norm;

 $\sigma = \frac{1}{2} \, \Rightarrow$ C
—Stabilität für $\tau \leq h^2,$ d.h.

 bedingt stabil;

 $\sigma=0 \ \Rightarrow$ C–Stabilität für $\tau \leq \frac{h^2}{2},$ d.h.
 bedingt stabil.

Stabilitätsuntersuchungen in der diskreten C-Norm sind auch für allgemeinere parabolische ARWA mit variablen Koeffizienten vom Typ

$$cg\frac{\partial u}{\partial t} - \left[\frac{\partial}{\partial x}\left(\lambda(x,t)\frac{\partial u}{\partial x}\right) + v(x,t)\frac{\partial u}{\partial t} + q(x,t)u\right] = f(x,t)$$

möglich und relativ einfach.

2.1.4 Zusammenfassung: Approximation + Stabilität \Longrightarrow diskrete Konvergenz

■ Konvergenz in diskreter L_2 -Norm:

Satz 2.10.:
$$||z|| = ||z||_{L_2(\omega_h)} := \left(\sum_{i=1}^{n-1} h z_i^2\right)^{1/2}$$

Vor.	$0 \le \sigma \le 1$	Differenzenstern	Stabilitätsbedingung
		Auflösung	$\sigma \ge \sigma_0 = \frac{1}{2} - \frac{h^2}{4\tau}$
$u \in C^{4,2}$	$\sigma = 1$	•	unbedingt
	rein impl. DS	LU–Zerlegung	stabil
$u \in C^{4,3}$	$\sigma = 1/2$		unbedingt
	CRANK-NICOLSON	LU–Zerlegung	stabil
$u \in C^{6,3}$	$\sigma = \sigma_* = \frac{1}{2} - \frac{h^2}{12\tau}$		unbedingt
	Schema bester Approximation	LU–Zerlegung	stabil
$u \in C^{4,2}$	$1 \stackrel{(>)}{\geq} \sigma \stackrel{(>)}{\geq} \frac{1}{2}$		unbedingt
	fix	LU–Zerlegung	stabil
$u \in C^{4,2}$	$\frac{1}{2} > \sigma > 0$		bedingt stabil, d.h. h, τ :
		LU-Zerlegung	$\sigma \ge \sigma_0 = \frac{1}{2} - \frac{h^2}{4\tau}$
$u \in C^{4,2}$	$\sigma = 0$	T T	bedingt stabil:
	explizites Schema	• • •	$ au \leq \frac{h^2}{2}$

Vor.	rechte Seite	Konsistenz = lokaler	diskrete Konvergenz
		Approximationsfehler	$ \max_{j=0,m} \ u^j - v^j \ \le T \max_{k=0,m-1} \ \psi^k \ \le $
	φ	ψ	$\leq T c_A(u) \left(h^p + \tau^q\right)$
$u \in C^{4,2}$	$\varphi = \bar{f}$	$O(h^2 + \tau)$	$p=2, \ q=1$
	$\varphi = f, \hat{f}$		
$u \in C^{4,3}$	$\varphi = \bar{f}$	$O(h^2 + \tau^2)$	$p=2, \ q=2$
	$\varphi = \frac{1}{2}(\hat{f} + f)$		
$u \in C^{6,3}$	$\varphi = \bar{f} + \frac{h^2}{12} \bar{f}_{\bar{x}x}$	$O(h^4 + \tau^2)$	$p=4, \ q=2$
	$\varphi = \bar{f} + \frac{h^2}{12} \frac{\partial^2 f}{\partial x^2}$		
$u \in C^{4,2}$	$\varphi = \bar{f}$	$O(h^2 + \tau)$	$p=2, \ q=1$
	$\varphi = f, \hat{f}$		
$u \in C^{4,2}$	$\varphi = \bar{f}$	$O(h^2 + \tau)$	$p=2, \ q=1$
	$\varphi = f, \hat{f}$		
$u \in C^{4,2}$	$\varphi = \bar{f}$	$O(h^2 + \tau)$	$p=2, \ q=1$
	$\varphi = f, \hat{f}$		

■ Konvergenz in der diskreten C-Norm:

• Satz 2.11: $||z|| = ||z||_{C(\omega_h)} := \max_{i=1,n-1} |z_i|$

$$\underline{\text{Vor.: }} 1) \ 0 \le \sigma \le 1, \qquad (1 - \sigma) \frac{\tau}{h^2} \le \frac{1}{2} ,
2) \ u \in C^{4,2}(\bar{Q}_T).
\underline{\text{Bh.: }} \max_{j=0,m} \|u(\cdot, t^j) - v^j\|_{C(\bar{\omega}_h)} := \max_{j=0,m} \max_{i=0,n} |u(x_i, t_j) - v_i^j| \le
(1) \quad (2)
\le T c_A(u) (h^2 + \tau).$$

• Satz 2.12:

 \bullet Für das Crank–Nicolson–Schema ($\sigma=1/2$) gilt speziell

$$\max_{j=\overline{0,m}} \|u(\cdot,t^{j}) - v^{j}\|_{C(\bar{\omega}_{h})} \le T c_{A}(u) (h^{2} + \tau^{2}),$$
(1) (2)

falls $\tau/h^2 \leq 1$ (Stabilitätsbed.) und $u \in C^{4,3}(\bar{Q}_T)$.

• Für das Schema bester Approximation ($\sigma=\sigma_*\equiv \frac{1}{2}-\frac{h^2}{12\tau}$ mit $\sigma\in[0,1]$) gilt speziell

$$\max_{j=\overline{0,m}} \|u(\cdot,t^{j}) - v^{j}\|_{C(\bar{\omega}_{h})} \le T c_{A}(u) (h^{4} + \tau^{2}),$$
(1) (2)

falls $\tau/h^2 \le 5/6$ (Stabilitätsbed.) und $u \in C^{6,3}(\bar{Q}_T)$.

2.2 Eine allgemeine Stabilitätstheorie für zweischichtige Schemata in der Energienorm (Energetische Methode)

2.2.1 Allgemeine und kanonische Form

■ Für parabolische ARWA (siehe [17] Numerik II, Pkt. 1.1) der Art

$$\frac{\partial u}{\partial t} + L u(x,t) = f(x,t), \quad \forall (x,t) \in Q_T = \Omega \times \mathbf{T}$$
+ RB + AB, mit $\Omega \subset \mathbb{R}^d$ beschr. Gebiet, $\mathbf{T} = (0,T)$ und mit dem formal s.a. elliptischen Differentialausdruck

$$L u := -\sum_{i,j=1}^{d} \frac{\partial}{\partial x_i} \left(a_j(x) \frac{\partial u}{\partial x_j} \right) + a(x) u(x,t)$$

sei ein zweischichtiges $\left(\frac{\partial u}{\partial t}\right)$ DS aufgeschrieben, dessen <u>allgemeine Form</u> offenbar folgendermaßen aussieht:

(12)
$$C_1 v^{j+1} + C_0 v^j = \tau \varphi^j(x), \quad j = 0, 1, ..., m-1,$$
AB: v^0 geg., $x \in \omega = \omega_h$

wobei
$$v=v^j: \bar{\omega} \longmapsto I\!\!R^1$$
 – Gitterfkt. auf j -ter Zeitschicht
$$(v|_{\gamma=\gamma_1=\bar{\omega}\setminus\omega}=0 \ !);$$

$$H_h^{(1)}, \ \|\cdot\|_{(1)}, \ (\cdot,\cdot)_{(1)} - \text{diskreter } H\text{-Raum};$$

$$\varphi=\varphi^j:\omega\longmapsto I\!\!R^1 - \text{rechte Seite auf } j\text{-ter Zeitschicht};$$

$$H_h^{(2)}, \ \|\cdot\|_{(2)}, \ (\cdot,\cdot)_{(2)} - \text{diskreter } H\text{-Raum};$$

$$C_0,C_1:H_h^{(1)}\longmapsto H_h^{(2)} - \text{lineare Operatoren } (\leftarrow \text{FDM}, \text{FEM},\ldots);$$

$$\text{RB}-\text{sind eingearbeitet:} \bullet \text{RB 1. Art seien homogenisiert} \to \text{RS},$$

$$\bullet \text{Nat. RB: } x \in \gamma_N = \gamma_2 \cup \gamma_3 \Rightarrow \omega = \overset{\circ}{\omega} \cup \gamma_N;$$

$$\tau - \text{Zeitschritt}, \ h - \text{\"{o}rtlicher Diskretisierungsparameter}, \ \omega - \text{Gitter f\"{u}r } \Omega,$$

$$\gamma = \gamma_1 = \gamma_D = \overline{\omega} \setminus \omega - \text{diskreter } \Gamma_1\text{-Rand (Dirichlet-Rand)}.$$

■ Kanonische Form:

(13)
$$B v_t^j + A v^j = \varphi^j, \quad j = \overline{0, m-1}; \quad v^0 \text{ geg.}$$

■ Beziehung zwischen kanonischer und allgemeiner Form:

$$(13) \iff B v^{j+1} - B v^j + \tau A v^j = \tau \varphi^j$$

(14)
$$C_1 = B, C_0 = \tau A - B \text{ bzw. } B = C_1, A = \frac{1}{\tau}(C_0 + C_1)$$

■ Beispiel: σ -gewichtetes Schema aus Pkt. 2.1:

$$\begin{cases} v_t^j - \sigma v_{\bar{x}x}^{j+1} - (1 - \sigma) v_{\bar{x}x}^j = \varphi^j(x), & x \in \omega_h \\ + \text{RB (hom. 1. Art: } g_0 = g_1 = 0) + \text{AB} \end{cases}$$

$$\bar{A} : \bar{A}v := -v_{\bar{x}x}, H_h^{(1)} = \stackrel{\circ}{L_2} (\bar{\omega}_h) = L_2(\omega_h) = H_h^{(2)} = H_h$$

$$\Rightarrow v^{j+1} - v^j + \tau \sigma \bar{A} v^{j+1} + \tau (1 - \sigma) \bar{A} v^j = \tau \varphi^j$$

$$\Rightarrow \underline{\text{allgemeine Form:}} \underbrace{(I + \tau \sigma \bar{A})}_{=: C_1} v^{j+1} + \underbrace{(\tau (1 - \sigma) \bar{A} - I)}_{=: C_0} v^j = \tau \varphi^j$$

$$\Rightarrow \underline{\text{Kanonische Form:}} \ B = C_1 = I + \tau \sigma \bar{A} \\ A = \frac{1}{\tau} (C_0 + C_1) = \frac{1}{\tau} \left(\tau (1 - \sigma) \bar{A} - I + I + \tau \sigma \bar{A} \right) = \bar{A}$$

■ **Def. 2.13:** ($\hat{=}$ Def. 2.3): $(H_h^{(1)}; H_h^{(1)}, H_h^{(2)})$ – Stabilität.

Das zweischichtige DS (12) = (13) heißt stabil (bzgl. AB v^0 und RS φ , sowie den gewählten Normen $\|\cdot\|_{(1)}$ und $\|\cdot\|_{(2)}$), wenn für evtl. hinreichend kleinen h und $\forall j=0,1,\ldots,m-1$ gilt:

(15)
$$||v^{j+1}||_{(1)} \le c_1 ||v^0||_{(1)} + c_2 \max_{0 \le k \le j} ||\varphi^k||_{(2)}$$

$$\begin{split} & \text{mit } c_1, c_2 = \text{const.} > 0, \ \ c_\alpha \neq c_\alpha(h, \tau, j, \varphi, v^0), \ \ \| \cdot \|_{(1)} = \| \cdot \|_{H_h^{(1)}}, \\ & \| \cdot \|_{(2)} = \| \cdot \|_{H_h^{(2)}} - \text{gew. Normen.} \end{split}$$

Spezialfälle: Stabilität bzgl. AB: $\varphi = 0$ in (12): A-priori-Absch. (15), Stabilität bzgl. RS: $v^0 = 0$ in (12): A-priori-Absch. (15).

2.2.2 Stabilität in der energetischen Norm $\|\cdot\|_A$

■ Standardvoraussetzungen an (13):

H – (diskreter) reeller Hilbert–Raum (\rightarrow Grundraum), i. S. Raum von Gitterfkt.: $\|\cdot\|$, (\cdot, \cdot) ;

- (1) $A, B: H \longrightarrow H$ lineare Operatoren ($\leftarrow L$ -linear!).
- (2) $A = A^* > 0$ spd ($\leftarrow L$ formal s.a., ellipt. Operator !); d.h. $(A v, w) = (v, A w), \quad \forall v, w \in H \text{ und } (A v, v) > 0 \quad \forall v \in H, \ v \neq \mathbf{0}$.
- (3) $B > 0 \iff \exists B^{-1} \Rightarrow \text{Schema eindeutig auflösbar !}$.
- (4) A, B seien nicht von t_i abhängig (d.h. Koeff. von L sind zeitunabhängig!).
- Btr. Fehlerschema: $z = u^{(11)} v^{(12)}$
 - (16) $B z_t + A z = \psi, \quad \psi \text{Approximations fehler};$ $AB: z^0 \text{ geg}.$

$$z = y + w$$
: (17) $By_t + Ay = \psi$; $y^0 = 0$: Ziel: Stabilität bzgl. RS \leftarrow RS AB: (18) $Bw_t + Aw = 0$; $w^0 = z^0$: Ziel: Stabilität bzgl. AB \leftarrow $+ \Longrightarrow$ (16): Stabilität bzgl. AB und RS!

■ Die "energetische" Identität: (im parabolischen (zweischichtigen) Falle)

$$\underline{\underline{((16), 2\tau z_t)}} \cdot 2\tau (Bz_t, z_t) + 2\tau (A[z], [z_t]) = 2\tau (\psi, z_t)$$

$$z = \frac{\hat{z}+z}{2} - \frac{\hat{z}-z}{2} = \frac{1}{2}(\hat{z}+z) - \frac{\tau}{2}z_t$$

$$z_t = \frac{\hat{z}-z}{\tau}$$

$$2\tau(Bz_t, z_t) + \underbrace{\tau\left(A(\hat{z}+z), \frac{\hat{z}-z}{\tau}\right) - 2\tau\left(\frac{\tau}{2}Az_t, z_t\right) = 2\tau(\psi, z_t)}_{= (A(\hat{z}+z), \hat{z}-z) = (A\hat{z}, \hat{z}) + \underbrace{(Az, \hat{z}) - (A\hat{z}, z)}_{= 0, \text{ da } A = A^*} - (Az, z)$$

Resultat: Energetische Identität

(19)
$$2\tau \left(\left(B - \frac{\tau}{2} A \right) z_t, z_t \right) + \left(A \hat{z}, \hat{z} \right) - \left(A z, z \right) = 2\tau (\psi, z_t)$$

■ <u>Lemma 2.14.:</u> (Stabilität bzgl. AB)

Bh.: Dann ist DS (16) stabil bzgl. AB (d.h. DS (18)) mit $c_1 = 1$ im Raum H_A (= zu H energetischer Raum):

(20)
$$||w^{j+1}||_A := (Aw^{j+1}, w^{j+1})^{0.5} \le ||w^0||_A.$$

Beweis:
$$\psi = 0$$
, d.h. DS (18)! Für (18) folgt aus (19)

$$\Rightarrow (A\widehat{w}, \widehat{w}) = (Aw, w) - 2\tau \underbrace{\left(\left(B - \frac{\tau}{2}A\right) w_t, w_t\right)} \leq (Aw, w).$$

$$\Rightarrow \|\widehat{w}\|_A \leq \|w\|_A.$$

q.e.d.

■ Beispiel: σ -gewichtetes DS aus Pkt. 2.1.:

$$B = I + \tau \sigma \bar{A}, A = \bar{A}, \bar{A}v := -v_{\bar{x}x};$$

$$H = \mathring{L_2}(\bar{\omega}) = L_2(\omega).$$

- Voraussetzungen (1) (4) offenbar erfüllt, sogar $B = B^*$ (mms)
- Allgemeine Stabilitätsbedingung: $B = I + \tau \sigma \bar{A} \ge \frac{\tau}{2} \bar{A}$!

$$\left\{ \begin{array}{l} \sigma \geq \frac{1}{2} & : \text{offenbar erf\"{u}llt} \\ \sigma \geq \sigma_0 = \frac{1}{2} - \frac{h^2}{4\tau} : \text{ebenfalls erf\"{u}llt} \ \boxed{\ddot{\textbf{U}} \ 2.3} \end{array} \right\} \Longrightarrow \boxed{ \|w^{j+1}\|_A \leq \|w^0\|_A }$$

$$\boxed{\ddot{\textbf{U}} \ 2.3 \ \|>}$$

$$\boxed{ \|w^{j+1}_{\bar{x}}\| \leq \|w^0_{\bar{x}}\|] }$$

• Ü 2.3 Man zeige für das gew. DS mit homog. RS und mit $1 \ge \sigma \ge \sigma_0 = \frac{1}{2} - \frac{h^2}{4\tau}$ die A-priori-Abschätzungen:

(a)
$$\|w_{\bar{x}}^{j+1}\| := \left(\sum_{i=1}^{n} h\left(w_{\bar{x},i}^{j+1}\right)^{2}\right)^{1/2} \le \|w_{\bar{x}}^{0}\|,$$

(b)
$$||w^{j+1}||_{C(\omega_h)} \le \frac{1}{2} ||w^0_{\overline{x}}||$$
.

■ Lemma 2.15.: (Stabilität bzgl. RS)

Aus der Stabilität bzgl. der AB gemäß Lemma 2.14. folgt die Stabilität bzgl. der RS (\rightarrow DS (17)): $(21) \quad \|y^{j+1}\|_A \leq T \max_{0 \leq k \leq j} \|B^{-1}\psi^k\|_A$ $\underline{\qquad \qquad }$ vgl. Def. 2.13.: $\|y^{j+1}\|_{(1)} \quad c_2 \quad \|\psi^k\|_{(2)} = \|\psi^k\|_{\underline{B^{-T}AB^{-1}}}$

Beweis:

- Btr. DS (17): $By_t + Ay = \psi$, $y^0 = 0$ $\tau B^{-1} \cdot (17)$: $\Rightarrow \hat{y} = y - \tau B^{-1} Ay + \tau B^{-1} \psi = (I - \tau B^{-1} A)y + \tau B^{-1} \psi$
- Vor. (2): $A = A^* > 0 \Rightarrow \exists A^{1/2}$; Setzen: $v = A^{1/2}y$, $y = A^{-1/2}v$: $\Rightarrow \hat{v} = \underbrace{(I \tau A^{1/2}B^{-1}A^{1/2})}_{=: S \text{ (Übergangsoperator)}} v + \underbrace{\tau A^{1/2}B^{-1}\psi}_{=: f}$

$$\hat{v} = S v + \tau f$$

• Schätzen nun || S || ab: \uparrow Dazu setzen wir $\psi=0 \Rightarrow f=0$: $\hat{v}=S\,v$

Lemma 2.14. (Beweis):
$$\|\hat{y}\|_A \le \|y\|_A$$

$$\Leftrightarrow y = A^{-1/2}v$$

$$\|\hat{v}\| \le \|v\|$$

$$\Leftrightarrow \|Sv\| \le \|v\| \quad \forall v \in H$$

Resultat:
$$\parallel S \parallel$$

Bem.:
$$||S|| = ||I - \tau B^{-1}A||_A$$
 (mms)

$$\hat{v} = \mathbf{S} \, v + \tau f \qquad \qquad \parallel \mathbf{S} \parallel \le 1$$

• $\|\hat{v}\| \le \|Sv\| + \tau \|f\| \le \|S\| \|v\| + \tau \|f\| \le \|v\| + \tau \|f\|$

$$\implies ||v^{j+1}|| \le ||v^j|| + \tau ||f^j|| \le \ldots \le ||v^0|| + \tau \sum_{k=0}^j ||f^k||$$

$$\iff \|y^{j+1}\|_A \le \underbrace{\|y^0\|_A}_{=0} + \tau \sum_{k=0}^j \|B^{-1}\psi^k\|_A$$
 (22)

$$\leq \max_{0 \leq k \leq j} \|B^{-1}\psi^{k}\|_{A} \cdot \tau \sum_{n=0}^{j} 1 \leq T \max_{0 \leq k \leq j} \|B^{-1}\psi^{k}\|_{A}.$$
$$= (j+1)\tau \leq T$$

q.e.d.

■ Satz 2.16. (Stabilität bzgl. AB und RS)

<u>Vor.:</u> 1. Standardvoraussetzungen (1) - (4).

<u>Bh.:</u> Das DS (16): $Bz_t + Az = \psi$, $z^0 - \text{geg.}$, ist <u>stabil</u> bzgl. AB und RS im energetischen Raum H_A (genauer: H_A ; H_A , $H_{B^{-T}AB^{-1}}$):

(23)
$$||z^{j+1}||_A \le ||z^0||_A + T \max_{0 \le k \le j} ||\psi^k||_{B^{-T}AB^{-1}}.$$

Beweis:

- folgt sofort aus Beweis von Lemma 2.15., Abschätzung (22);

$$\bullet$$
bzw. aus $z=w+y$ – Lemma 2.14. – Lemma 2.15.
$$\|z^{j+1}\|_A \ \leq \ \|w^{j+1}\|_A \ + \ \|y^{j+1}\|_A \ \leq \ \dots$$

q.e.d.

■ Beispiel: σ -gewichtetes DS aus Pkt. 2.1.:

$$B = I + \tau \sigma \bar{A}, \quad A = \bar{A}, \quad \bar{A}v = -v_{\bar{x}x}, \quad H = \stackrel{\circ}{L_2}(\bar{\omega}) = L_2(\omega) = H_h$$
 vertauschbar
$$\|B^{-1}\psi\|_A = \sup_{\substack{v \in H \\ v \neq \boldsymbol{O}}} \frac{((I + \sigma \tau A)^{-1}\psi, v)_A}{\|v\|_A} = \sup_{\substack{v \in H \\ v \in \boldsymbol{H}}} \frac{(A(I + \sigma \tau A)^{-1}\psi, v)}{\|v\|_A} = \sum_{\substack{v \in H \\ v \in \boldsymbol{H}}} \frac{(A(I + \sigma \tau A)^{-1}\psi, v)}{\|v\|_A} = \sum_{\substack{v \in H \\ v \in \boldsymbol{H}}} \frac{((I + \sigma \tau A)^{-1}A^{0.5}\psi, A^{0.5}v)}{\|v\|_A} \leq \underbrace{\|(I + \sigma \tau A)^{-1}\|}_{\text{Satz v. Banach}} \|\psi\|_A \leq 1 \cdot \|\psi\|_A,$$

d.h.
$$(H_A; H_A, H_A)$$
–Stabilität:

Man zeige, daß das σ –gewichtete DS aus Pkt. 2.1 im energetischen Raum $H_{\bar{A}}=A$ stabil ist und daß die a-priori
–Abschätzung

$$||z_{\bar{x}}^{j+1}|| \le ||z_{\bar{x}}^{0}|| + T \max_{0 \le k \le j} ||\psi_{\bar{x}}^{k}||$$

(für das Fehlerschema) gilt, falls $1 \ge \sigma \ge \sigma_0 = \frac{1}{2} - \frac{h^2}{4\tau}$.

Hinweis: 1) Satz 2.16. anwenden.

- 2) $||z||_{\bar{A}} = ||z_{\bar{x}}||$ (Formel der partiellen Summation).
- 3) Satz von Banach $||(I+K)^{-1}|| \le ?$

2.3 Galerkin–FEM für parabolische ARWA

2.3.1 Verallgemeinerte Formulierungen parabolischer ARWA

■ Btr. zunächst parabolische ARWA in klassischer Formulierung (vgl. [17] Nu II, Pkt. 1.1: Instationäre Wärmeleitprobleme bzw. Wärmeleit—Wärmetransport-Probleme):

Ges.
$$u(x,t) \in X = C^{x,t}(Q_T) \cap C^{0,0}(\bar{\Omega} \times [0,T)) \cap C^{1,0}(\Omega \cup \Gamma_2 \cup \Gamma_3 \times (0,T))$$
:
$$\frac{\partial u}{\partial t} - \sum_{i,j=1}^d \frac{\partial}{\partial x_i} \left(a_{ij}(x,t) \frac{\partial u}{\partial x_j} \right) + \sum_{i=1}^d a_i(x,t) \frac{\partial u}{\partial x_i} + a(x,t) u(x,t) = f(x,t)$$

$$\forall (x,t) \in Q_T = \Omega \times T, T = (0,T),$$

$$\Omega \subset \mathbb{R}^d - \text{beschr. Gebiet: } \Gamma = \partial \Omega \in C^2;$$

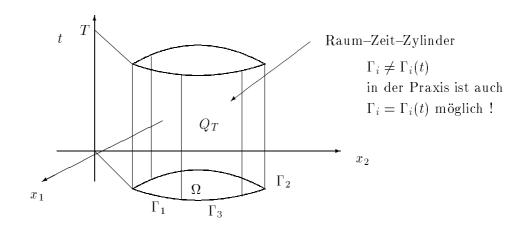
$$0. B. d. A. \underline{\text{(Homogenisieren)}}$$

$$+ \text{RB: } u(x,t) = g_1(x,t) \equiv 0, \forall x \in \Gamma_1$$

$$\frac{\partial u(x,t)}{\partial N} := \sum_{i,j=1}^d a_{ij}(x,t) \frac{\partial u}{\partial x_j} \vec{n}_i = g_2(x,t) , \forall x \in \Gamma_2$$

$$\frac{\partial u(x,t)}{\partial N} + \kappa(x,t) u(x,t) = g_3(x,t) , \forall x \in \Gamma_3$$

$$+ \text{AB: } u(x,0) = u_0(x) \quad \forall x \in \bar{\Omega}.$$



lacksquare Variationsformulierung in Sobolev-Räumen der Art $oldsymbol{W}_{_2}^{^{0,1,k}}(oldsymbol{Q}_T)$ auf $\underline{\mathsf{dem}\,\,\mathsf{Raum}\mathsf{-}\mathsf{Zeit}\mathsf{-}\mathsf{Zylinder}\,\,\boldsymbol{Q}_T:=\boldsymbol{\varOmega}\times\boldsymbol{T}\colon}$

$$\begin{split} & \text{Sei } \overset{\circ}{W}_{2}^{1,k}(Q_{T}) := \left\{ u \in W_{2}^{1,k}(Q_{T}) : u|_{\Gamma_{1} \times [0,T]} = 0 \right\} - \text{Ansatzfkt.} = \text{Lsg.-Raum,} \\ & \dot{W}_{2}^{1,l}(Q_{T}) := \left\{ v \in W_{2}^{1,l}(Q_{T}) \ : \ \begin{aligned} & u|_{\Gamma_{1}} = 0 \\ & u|_{t=T} = 0 \end{aligned} \right. \left. (l \geq 1) \right. \right\} - \text{Testfkt.-Raum,} \\ & \int_{Q_{T}} \overset{(24)}{\text{PDgl.}} \cdot v(x,t) \, dx \, dt \quad \forall v \in \dot{W}_{2}^{1,l}(Q_{T}), \quad k+l = 1 : \end{split}$$

(a) Nur partielle Integration im elliptischen Hauptteil:

Ges.
$$u \in \overset{\circ}{W}_{2}^{1,1}(Q_{T})$$
:
$$\int_{Q_{T}} \left(\dot{u}v + \sum_{i,j=1}^{d} a_{ij} \frac{\partial u}{\partial x_{j}} \frac{\partial v}{\partial x_{i}} + \sum_{i=1}^{d} a_{i} \frac{\partial u}{\partial x_{i}}v + auv \right) dx dt + \int_{0}^{T} \int_{\Gamma_{3}} \kappa uv \, ds \, dt =$$

$$= \int_{Q_{T}} fv \, dx \, dt + \int_{0}^{T} \int_{\Gamma_{2}} g_{2}v \, ds \, dt + \int_{0}^{T} \int_{\Gamma_{3}} g_{3}v \, ds \, dt, \quad \forall v \in \overset{\circ}{W}_{2}^{1,0}(Q_{T})$$

$$+ \text{AB: } u(x,0) = u_{0}(x) \text{ i.S. } L_{2}(\Omega), \text{ d.h.}$$

$$\|u(\cdot,t) - u_{0}(\cdot)\|_{L_{2}(\Omega)} \longrightarrow 0 \text{ für } t \to 0.$$

$$\text{mit } \dot{u} := \frac{\partial u}{\partial t}.$$

mit $\dot{u} := \frac{\partial u}{\partial t}$.

(b) Partielle Integration im elliptischen Hauptteil und im Term mit der Zeitableitung:

(26) Ges.
$$u \in \mathring{W}_{2}^{1,0}(Q_{T})$$
:
$$\int_{Q_{T}} \left(-u\dot{v} + \sum_{i,j=1}^{d} a_{ij} \frac{\partial u}{\partial x_{j}} \frac{\partial v}{\partial x_{i}} + \sum_{i=1}^{d} a_{i} \frac{\partial u}{\partial x_{i}} v + auv \right) dx dt + \int_{0}^{T} \int_{\Gamma_{3}} \kappa uv ds dt =$$

$$= \int_{\Omega} u_{0}(x) v(x,0) dx + \int_{Q_{T}} fv dx dt + \int_{0}^{T} \int_{\Gamma_{2}} g_{2}v ds dt + \int_{0}^{T} \int_{\Gamma_{3}} g_{3}v ds dt$$

$$\forall v \in \mathring{W}_{2}^{1,1}(Q_{T}).$$

Bemerkungen:

- 1. Existenz-, Eindeutigkeits- und Regularitätsaussagen siehe Literatur (z.B. [14] Ladyshenskaja A.: Aufgaben der Mathematischen Physik, Nauka, Moskau 1973).
- 2. VF (25) und (26) sind Ausgangspunkt für die FE-Diskretisierung mit <u>Raum-Zeit-Elementen!</u>

■ Linienvariationsformulierung (LVF):

• $\int_{\Omega} \operatorname{PDgl.} \cdot v(x) dx$, $\forall v \in V_0 = \{v \in V = W_2^1(\Omega) : v|_{\Gamma_1} = 0\}$, \forall f.ü. $t \in T$ $\downarrow \longleftarrow \operatorname{Schritte}(1) - (5)$, [17] Nu II, Pkt. 3.1.2.1! LVF:

$$(27) = (24)_{LVF}$$

$$\begin{aligned} \operatorname{Ges.} u(x,t), \operatorname{soda} & \forall \ \operatorname{f.\ddot{u}.} \ t \in \boldsymbol{T}, \ u \in \boldsymbol{V_0} \ \operatorname{und} \ \dot{u} \in L_2 \left(\Omega\right) : \\ & \underbrace{\int d(x,t)v(x)dx}_{\boldsymbol{\Omega}} + \underbrace{\int \left(\sum_{i,j=1}^d a_{ij} \frac{\partial u}{\partial x_j} \frac{\partial v}{\partial x_i} + \sum_{i=1}^d a_i \frac{\partial u}{\partial x_i} v + auv\right) dx + \int \kappa uv \, dx = \\ & = (\dot{u},v)_0 & = a(u,v) = : < \underbrace{A(t) \ u(t), v >}_{\boldsymbol{\Xi}} \\ & = a(t;u,v) & \in \boldsymbol{V_0^*} \end{aligned} \\ & : = < \dot{u}, v > = \frac{d}{dt}(u,v)_0 \\ & \in \boldsymbol{V_0^*} \\ & = \underbrace{\int fv \, dx + \int g_2 v \, ds + \int g_3 v \, ds}_{\boldsymbol{\Gamma}_2}, \ \forall v \in \boldsymbol{V_0} \ \forall \ \operatorname{f.\ddot{u}.} \ t \in \boldsymbol{T} \end{aligned} \\ & = : < F(t), v > \\ & \in \boldsymbol{V_0^*} \\ & + \operatorname{AB:} u(x,0) = u_0(x) \ \operatorname{in} \ L_2(\Omega). \end{aligned}$$

Nach evtl. Homogenisierung $(u = g_1 + w \uparrow)$ der Dirichletschen RB erhalten wir eine Cauchy-Aufgabe (AWA) für eine Operatordifferentialgleichung (ODgl.)

Ges.
$$u(\cdot): \bar{T} \mapsto V_0 \text{ mit } \dot{u}(\cdot): T \longrightarrow V_0^*:$$

 $\dot{u}(t) + A(t) u(t) = F(t) \text{ in } V_0^*, \quad \forall \text{ f.\"{u}. } t \in T,$
 $+ \text{ AB: } u(0) = u_0 \text{ in } L_2(\Omega).$

• Zur präzisen Formulierung von (27), sowie von Existenz- und Eindeutigkeitsaussagen benötigen wir noch einige funktionalanalytische Hilfsmittel!

■ Funktionalanalytische Formulierung und Untersuchung der LVF:

• L_2 -Räume abstrakter Fkt.:

$$\begin{array}{l} X-\operatorname{Banach-Raum} \text{ (bzw. Hilbert-Raum)}, \quad \boldsymbol{T}=(0,T). \\ L_2\left(\boldsymbol{T},\!X\right) \coloneqq \{u:\boldsymbol{T}\!\rightarrow\!X: \|u\|_{L_2\left(\mathbf{T},X\right)}<\infty\}, \\ \text{wobei } \|u\|_{L_2\left(\mathbf{T},X\right)} \coloneqq \left(\int\limits_0^T \|u(t)\|_X^2 dt\right)^{1/2}. \end{array}$$

Es gilt: (mms)

- 1) $L_2(T,X)$ ist selbst ein B-Raum.
- 2) $L_2(T,X^*) = [L_2(T,X)]^*$, falls X reflexiv, separabel; X^* = Dualraum zu X.

• Evolutionstripel (Gelfand-Dreier):

Das Raumtripel $\{X, H, X^*\}$ heißt Evolutionstripel, falls $X \subseteq H \subseteq X^*$:

- (a) X separabler, reflexibler B-Raum, $\|\cdot\|_X$,
- (b) H separabler H–Raum, $(\cdot, \cdot) = (\cdot, \cdot)_H$, $\|\cdot\|_H$,
- (c) X liegt dicht in $H: ||v||_H \le c||v||_X \quad \forall v \in X$.

• Beispiel: 1)
$$X = \overset{\circ}{H^{1}}(\Omega), H = L_{2}(\Omega), X^{*} = H^{-1}(\Omega)$$
 (1. ARWA),
2) $X = V_{0}$, $H = L_{2}(\Omega), X^{*} = V_{0}^{*}$ (\(\frac{\gamma}{\gamma}\)).

• Man sagt, $u \in L_2(T, X)$ besitzt eine <u>verallgemeinerte Ableitung</u> $w \equiv \dot{u} \in L_2(T, X^*)$, falls

$$\int_{0}^{T} \dot{\varphi}(t) u(t) dt = -\int_{0}^{T} \varphi(t) w(t) dt \text{ in } X^{*} \quad \forall \varphi \in \dot{C}^{\infty}(\mathbf{T}),$$

$$\begin{array}{l} \mathrm{d.h.} \, < \int\limits_0^T \dot{\varphi}(t) \, u(t) \, dt, v > = - < \int\limits_0^T \varphi(t) \, w(t) \, dt, v > \quad \forall v \in X \quad \ \forall \varphi \in \ \dot{C}^\infty(\boldsymbol{T}), \\ \mathrm{wobei} \, < \cdot, \cdot > : X^* \times X \to I\!\!R^1 - \mathrm{Dualit\"{a}tsprodukt}. \end{array}$$

• $W_2^1(T,X) = W_2^1(T;X, H) := \{u \in L_2(T,X) : \exists \dot{u} \in L_2(T,X^*)\}:$ $\|u\|_{W_2^1(T,X)} := \|u\|_{L_2(T,X)} + \|\dot{u}\|_{L_2(T,X^*)}$ (Grapennorm) bzw. $\|u\|_{W_2^1(T,X)}^2 := \|u\|_{L_2(T,X)}^2 + \|\dot{u}\|_{L_2(T,X^*)}^2$, falls X - Hilbert-Raum.

Für $u \in W_2^1(T,X)$ gilt (mms):

- 1) Abb. $u: \bar{T} = [0, T] \mapsto H$ ist stetig (nach Änderung auf einer Menge vom Maße Null), d.h. $W_2^1(T, X) \hookrightarrow C(\bar{T}, H)$ und $u(0) \in H$ ist korrekt definiert.
- 2) Formel der partiellen Integration:

$$\begin{split} &(u(t),v(t))_H - (u(s),v(s))_H = \smallint_s^t [< u'(\tau),v(\tau)> + < v'(\tau),u(\tau)>] d\tau, \\ &\text{wobei} < \cdot,\cdot> : X^* \times X \to I\!\!R^1 - \text{Dualit\"{a}tsprodukt}. \end{split}$$

3)
$$\frac{d}{dt}(u(t), v) = \langle \dot{u}(t), v \rangle \quad \forall v \in X \quad \forall \text{ f.\"{u}. } t \in T.$$

• Damit läßt sich die LVF (27) wie folgt formulieren:

$$\dot{u} + A u = F \qquad \qquad \underline{\text{in } L_2(\mathbf{T}, V_0^*)}$$

 $? \downarrow ?$ in $C(\mathbf{T}, V_0^*)$ mit $u \in C^1(\mathbf{T}, V_0^*) \text{ ges.}$

• Satz 2.17:

<u>Vor.:</u> 1) Die Bilinearform $a(\cdot, \cdot): V_0 \times V_0 \to \mathbb{R}^1$ (der Einfachheit halber sei $a(\cdot, \cdot)$ t-unabhängig, d.h. $a_{ij} = a_{ij}(x), \cdot$) sei V_0 - elliptisch und V_0 - beschr., d.h. $\exists \ \mu_1, \mu_2 = \text{const.} > 0$:

$$\begin{array}{ll} {\rm a)} \ \ \mu_1\|v\|_1^2 \leq a(v,v) \quad \forall v \in V_0, \\ {\rm b)} \ \ |a(u,v)| \leq \mu_2\|u\|_1\,\|v\|_1 \quad \forall u,v \in V_0. \end{array}$$

2)
$$F \in L_2(\mathbf{T}, V_0^*)$$
.

3)
$$u_0 \in L_2(\Omega) = H$$
.

<u>Bh.:</u> 1) $\exists ! u \in W_2^1(T; V_0, H) : (27).$

2)
$$||u(t)||_0 \le ||u_0||_0 + \int_0^t ||F(\tau)||_{V_0^*} d\tau \quad \forall t \in \bar{T}$$
.

<u>Bew.:</u> siehe [23] Zeidler E.: Vorlesungen über nichtlineare Funktionalanalysis II. Teubner-Texte zur Mathematik,

 \rightarrow A-nichtlinear

Leipzig 1977, S. 61 ff. und S. 155 ff.

- Ü 2.5 Man zeige die zweite Aussage (A-priori-Abschätzung) von Satz 2.17. Hinweis: Setzen Sie in (27) v=u ein !
- Regularitätsaussagen siehe Lit., z.B. [21] Thomée, 1984.

Um Diskretisierungsfehlerabschätzungen zu erhalten, benötigen wir in der Regel höhere Glattheit von der Lsg. u als die im \exists ! – Satz 2.17 angegebene. Parabolische ARWA besitzen die sogenannte Glättungseigenschaft!

Z.B. gelten für die 1. ARWA

 $\dot{u} - \Delta u = 0$ in Q_T , u = 0 auf $\partial\Omega \times T$, Ω – glatt, $u = u_0(x)$ für t = 0 und für $u_0 \in L_2(\Omega)$ und $t \geq \delta > 0$ die Aussagen:

- 1) $u \in H^k(\Omega) \quad \forall k \ge 1$,
- 2) $||u(t)||_k \le ct^{-\frac{1}{2}k}||u_0||_0 \quad \forall t > 0.$

2.3.2 Semidiskrete und volldiskrete Ersatzaufgabe

■ Ausgangspunkt: = LVF (27):

(27) Ges.
$$u \in W_2^1(T, V_0)$$
 mit $u(0) = u_0 \in L_2(\Omega)$ geg.:
$$\frac{d}{dt}(u(t), v)_0 + a(u(t), v) = \langle F(t), v \rangle \quad \forall v \in V_0, \ t \in T$$

$$\dot{u} + A u = F \text{ in } L_2(T, V_0^*)$$

- Zwei Diskretisierungsstrategien:
 - 1) Die (vertikale) <u>Linienmethode</u>: Erst $\,x$, dann $\,t\,$!

 \uparrow \uparrow FEM ODE–Solver

Diese Methode wird in dieser Vorlesung ausführlich behandelt!

2) Die <u>Rothe-Methode</u> (horizontale Linienmethode): Erst t, dann x! (siehe auch Pkt. 2.3.3 "Zusammenfassung") \uparrow \uparrow

Impl. Euler FEM
Operator - für ellip.
ODE-Solver Problem

• Zerlegen $\bar{T} = [0,T] = \bigcup_{\delta=0}^{m-1} [t_j,t_{j+1}], \quad t_j = j\tau$ (ohne Beschränkung der Allgemeinheit: gleichmäßig!), $\tau = T/m$:

$$\bar{T} = [0, T] \quad t_0 \quad t_1 \quad t_2 \qquad \qquad t_j = j\tau \qquad \tau \qquad t_m$$

• Def.
$$\varphi_j(t) = \begin{cases} (t - t_{j-1})/\tau, & t \in [t_{j-1}, t_j] \\ (t_{j+1} - z)/\tau, & t \in [t_j, t_{j+1}], & j = \overline{1, m-1} \\ 0, & \text{sonst} \end{cases}$$

mit offensichtlichen Modifikationen für j = 0 und j = m.

• Die Näherung an die Lösung u(x,t) wird durch die Rothe-Funktion

$$u_{\tau}(x,t) = \sum_{j=0}^{m} u_{j}(x) \varphi_{j}(t),$$

gegeben, wobei $u_{j+1} \in V_0 \ (\approx u(\cdot, t_{j+1}))$ aus der Variationsgleichung

$$\begin{cases} \left(\frac{u_{j+1} - u_j}{\tau}, v\right)_0 + a(t_{j+1}; u_{j+1}, v) = \langle F(t_{j+1}), v \rangle & \forall v \in V_0, \\ (?) & (? \text{ Stetigkeit}) \\ j = 0, 1, \dots, m-1; \ u_0 \text{ geg}. \end{cases}$$

bestimmt werden.

- <u>Literatur</u>: [18] Rektorys K.: The Method of Discret. in Time and PDEs. Dordrecht, Boston, 1982.
 - [10] Kačur J.: Method of Rothe in Evolution Equations. Teubner, Leipzig, 1985.

■ Die semidiskrete Ersatzaufgabe mittels (vertikaler) Linienmethode:

= FEM-Galerkin-Approximation mit zeitabhängigen Koeffizienten! (vgl. [17] Nu II)

Ansatz:
$$u_h = u_h(x, t) := \sum_{i \in \omega_h} \underbrace{u^{(i)}(t)} p^{(i)}(x) \in W_2^1 (\mathbf{T}, V_{0h}) \subset W_2^1 (\mathbf{T}, V_0),$$

$$(28) \qquad \qquad \in W_2^1 (\mathbf{T}) \nearrow$$
wobei $V_{0h} = \operatorname{span} \{ p^{(i)} : i \in \omega_h \} \subset V_0 - \operatorname{FE-UR}.$

 $(27)_h$

$$\begin{split} a(t;u_h,v_h) \\ \text{Ges. } u_h(x,t): &\frac{d}{dt}(u_h,v_h)_0 + a(u_h,v_h) = < F(t), v_h > \quad \forall v_h \in V_{0h} \quad \forall \text{ f.\"{u}. } t \in \textbf{\textit{T}} \\ &\underline{\text{AB: z.B. } L_2\text{-Projektion:}} \\ &(u_h(\cdot,0),v_h)_0 = (u_0(\cdot),v_h)_0 \quad \forall v_h \in V_{0h}. \end{split}$$

Ges.
$$\underline{u}_h = [u^{(i)}(t)]_{i \in \omega_h} \in [W_2^1(0,T)]^N$$
; $N = N_h = |\omega_h|$:

$$\sum_{i \in \omega_h} (p^{(i)}, p^{(k)})_0 \dot{u}^{(i)}(t) + \sum_{i \in \omega_h} a(p^{(i)}, p^{(k)}) u^{(i)}(t) = \langle F(t), p^{(k)} \rangle, \quad k \in \omega_h,$$

$$\underline{AB}: \sum_{i \in \omega_h} (p^{(i)}, p^{(k)})_0 u^{(i)}(0) = (u_0, p^{(k)})_0, \quad k \in \omega_h.$$

 $(27)_h$

Ges.
$$\underline{u}_h(t) \in [W_2^1(0,t)]^N$$
: $M_h \dot{u}_h(t) + K_h(t)\underline{u}_h(t) = \underline{f}_h(t)$, \forall f.ü. $t \in \mathbf{T}$,
$$M_h \underline{u}_h(0) = \underline{g}_h$$
,

mit
$$M_h = \left[\int_{\Omega} p^{(i)}(x) \ p^{(k)}(x) \ dx\right]_{k,i\in\omega_h}$$
 – Massenmatrix,
$$K_h(t) = [a(t;p^{(i)},p^{(k)})]_{k,i\in\omega_h}$$
 – Steifigkeitsmatrix,
$$\uparrow \qquad \uparrow$$
 falls Koeff. der PDgl. bzw. der RB 3. Art zeitabhängig sind!
$$\underline{f}_h(t) = [\langle F(t),p^{(k)}\rangle]_{k\in\omega_h} \in [L_2(0,T)]^N - \text{Lastvektor},$$

$$\underline{g}_h = \left[\int_{\Omega} u_0(x)p^{(k)}(x) \ dx\right]_{k\in\omega_h} \in \mathbb{R}^N - \text{Vektor der "Momente" der AB.}$$

Resultat:

Zur Bestimmung der Vektorfkt. (Koeffizienten) $\underline{u}_h(t) = [u^{(i)}(t)]_{i \in \omega_h} \in [W_2^1(0,T)]^N$ erhalten wir das System $\underline{(27)}_h$ gewöhnlicher Dgl. mit den AB $\underline{u}_h(0) = M_h^{-1}\underline{g}_h$, d.h. eine AWA (Cauchy–Aufgabe).

■ Probleme:

- 1.) $\exists ! \underline{u}_h(t) \in [W_2^1(0,T)]^N : (27)_h.$
- 2.) Weitere Eigenschaften der AWA $(27)_h$: ? Eventuell <u>steifes</u> Dgl.-System ?
- 3.) Numerische Verfahren zur Lösung von (27)_h:
 ⇒ volldiskrete Ersatzaufgabe ?
 ? Stabilität, Approximation, Konvergenz, Aufwand ?
- 4.) Gesamtfehlerabschätzung.
- Zu Problem 1.): $\exists ! \underline{u}_h(t) \in [W_2^1(0,T)]^N$: $\underline{(27)}_h$

Satz 2.18: (Satz von Picard und Lindelöf)

Vor.: Btr. AWA für folgendes System gew. Dgl.

(29) $\begin{cases} \frac{d\underline{y}(t)}{dt} + A(t)\underline{y}(t) = \underline{f}(t) & \forall \text{ f.ü. } t \in \mathbf{T} = (0, T) \\ \underline{y}(0) = \underline{y}_0 & \text{mit } \underline{y} = \underline{y}(t) = (y_1(t), \dots, y_N(t))^T \text{ ges.,} \\ \underline{y}_0 = (y_{01}, \dots, y_{0N})^T \in \mathbb{R}^N \text{ geg. AW,} \\ A = A(t) = [a_{ki}(t)]_{k,i=\overline{1,N}} : a_{ki}(\cdot) \in L_{\infty}(0, T), \\ \underline{f}(t) \in [L_2(0, T)]^N. \end{cases}$

 $\underline{\mathrm{Bh}.:} \ \exists \ ! \ y(t) \in [W_2^1(0,T)]^N \! \colon (29).$

 $\underline{\text{Beweis:}} \to \boxed{\ddot{\text{U}} \text{ 2.4}}$ (siehe auch [16] Nu I, Bsp. I.2.4)

 $\ddot{\mathbb{U}}$ 2.6 Man beweise den Satz 2.18 von Picard und Lindelöf (für nichtstetige Daten) ! \overrightarrow{PA} \overrightarrow{V}

Hinweise:

1) Übergang zur äquivalenten Igl. \int_{0}^{t} (29) $d\xi$

$$\underline{y}(t) = -\int_{0}^{t} A(\xi) \ \underline{y}(\xi) \ d\xi + \int_{0}^{t} \underline{f}(\xi) \ d\xi + \underline{y}_{0}$$

$$\begin{array}{l} \underline{y} = B\underline{y} \quad (= \text{Fixpunktgleichung}) \\ \text{mit } B: X = [C(\bar{T})]^N \longmapsto [W_2^1(0,T)]^N \subset X \\ & \uparrow \\ B\text{-Raum } ! \end{array}$$

2) Wenden Sie auf Fixpunktgleichung

Ges.
$$y \in X : y = By$$
 in X

den verallgemeinerten Banachschen Fixpunktsatz I.2.5 an (vgl. auch Bsp. I.2.4!)!

Anwendung von Satz 2.18 auf $(27)_h$:

• $M_h: M_h = M_h^T$ p.d. (da GRAMER–Matrix bzgl. $(\cdot, \cdot)_0$)

$$M_h \neq M_h(t) \longrightarrow \exists M_h^{-0.5}$$

Setzen: $\underline{u}_h = M_h^{-0.5} \underline{y}$ $\Rightarrow \underline{\dot{u}}_h = M_h^{-0.5} \underline{\dot{y}}$

$$\Rightarrow \underline{(27)_h} \iff \boxed{\frac{\underline{y}(t) + A(t)\underline{y}(t) = \underline{f}(t)}{\underline{y}(0) = \underline{y}_0}}$$
(30)

$$\begin{split} & \text{mit } \underline{y} = M_h^{0.5} \underline{u}_h \\ & A(t) = M_h^{-0.5} K_h(t) \, M_h^{-0.5} = [a_{ki}(t)]_{k,i \in \omega_h \; (k,i = \overline{1,N_h})} \\ & \underline{f}(t) = M_h^{-0.5} \underline{f}_h(t) = [f_k(t)]_{k \in \omega_h \; (k = \overline{1,N_h})} \\ & \underline{y}_0 = M_h^{-0.5} \underline{g}_h \end{split}$$

Vor.: 1) Koeff. der Bilinearform $a(\cdot,\cdot) \in L_{\infty}(Q_T)$.

 $2) f \in L_2(Q_T), \quad g_i \in L_2(\Gamma_i \times T), \quad i = 2, 3.$ $\underline{Bh} : 1) \ a_{ki}(\cdot) \in L_\infty(0, T).$ $2) \ f_k(\cdot) \in L_2(0, T).$

• Offenbar:

Bew.: (mms)

• Nach Satz 2.18 gilt dann:

$$\exists \, ! \, \underline{y}(t) \in [W_2^1(0,T)]^N : \quad (30)$$

$$\exists \, ! \, \underline{u}_h(t) = M_h^{0.5} \underline{y}(t) \in [W_2^1(0,T)]^N : \quad \underline{(27)}_h$$

- Zu Problem 2.): Weitere Eigenschaften der AWA $(27)_h = (30)$:
 - Seien die Voraussetzungen von Satz II.4.4 (gleichmäßig in t) erfüllt:

$$(31) \left\{ \begin{array}{c} 1) \ a(\cdot,\cdot): V_0 \times V_0 \to I\!\!R^1: V_0 - \text{elliptisch} \quad \text{gleichmäßig in} \quad t \ ! \\ V_0 - \text{beschränkt} \quad \text{gleichmäßig in} \quad t \ ! \\ V_0 - \text{symmetrisch}. \end{array} \right.$$

- 2) Triangularisierung sei regulär i. S. der Def. II.4.3.
- Nach Satz II.4.4 gilt dann:

(32)
$$\begin{cases} \exists \underline{c}_{E}, \overline{c}_{E} = \text{const.} > 0 : c_{E} \neq c_{E}(h) : \\ \underline{\gamma} = \underline{c}_{E}h^{d} \leq \lambda(K_{h}) \leq \overline{c}_{E}h^{d-2} = \overline{\gamma} \Rightarrow \kappa(K_{h}) \leq (\overline{c}_{E}/\underline{c}_{E}) h^{-2}. \\ \text{EW} \end{cases}$$

Diese Abschätzungen sind ordnungsgemäß scharf (vgl. Ü II.4.6), d.h. $\kappa(K_h) = O(h^{-2}).$

Man zeige, daß die Massenmatrix unter der Voraussetzung (31) 2) gut konditioniert ist, d.h. $\exists \nu_1, \nu_2 = \text{const.} > 0 : \nu_i \neq \nu_i(h)$:

(33)
$$\begin{cases} \nu_1 h^d \le \lambda(M_h) \le \nu_2 h^d, \ \nu_2^{-1} h^{-d} \le \lambda(M_h^{-1}) \le \nu_1^{-1} h^{-d}, \\ \kappa(M_h) \le \nu_2 / \nu_1 = O(1). \end{cases}$$

ullet Ü 2.8 Man zeige für $A_h=M_h^{-0.5}K_hM_h^{-0.5}$ die Abschätzungen

(34)
$$\begin{cases} \delta_1 = \underline{c}_E \nu_2^{-1} \le \lambda(A_h) \le \bar{c}_E \nu_1^{-1} h^{-2} = \delta_2 h^{-2}, \\ \text{d.h. } \kappa(A_h) \le (\delta_2/\delta_1) h^{-2}. \end{cases}$$

Darüberhinaus zeige man, daß $\kappa(A_h) = O(h^{-2})$, d.h. die Abschätzungen sind ordnungsgemäß scharf!

Hinweis: (32) und natürlich auch (33) sind ordnungsgemäß scharf!

- Folgerung: Für kleine (?) h ist Dgl.-System $(27)_h = (30)$ steif (?)! (⇒ unbedingt stabile Schemata zur Zeitintegration!)
- Zu Problem 3.): Numerische Verfahren zur Lösung des evtl. steifen Dgl.-Systems (AWA) $(27)_h = (30): \rightarrow$ siehe Kapitel 3!

• z.B.
$$\sigma$$
-gewichtete Differenzenschemata aus Pkt. 2.1:
$$\bar{T} = [0,T] \longmapsto \bar{\omega}_{\tau} = \{t_j = j\tau : j = \overline{0,m}, \ \tau = T/m\} :$$

$$0 = t_0 < t_1 < \ldots < t_j = j\tau < \ldots < t_m = T$$

$$\underline{\text{Zeitschritt:}} \ \tau = t_{j+1} - t_j = T/m,$$
 auch variabler Zeitschritt möglich!
$$\rightarrow \text{Zeitschrittsteuerung!}$$

• <u>Bez.:</u>

$$v^{j} = [v^{(i)}(t_{j})]_{i \in \omega_{h}} := \underbrace{U_{h}^{j}}_{h} = \underline{U_{h}}(t_{j}) = [U^{(i)}(t_{j})]_{i \in \omega_{h}} \in \mathbb{R}^{N_{h}} \text{ bzw. } v^{j} : \omega_{h} \to \mathbb{R}^{1}$$
Gitterfkt. auf j-ter Zeitschicht
$$U_{h}^{j} = U_{h}(t_{j}) = U_{h}(x, t_{j}) = \sum_{i \in \omega_{h}} U^{(i)}(t_{j}) p^{(i)}(x) \in V_{0h}$$

• Resultat: Volldiskrete Ersatzaufgabe = σ -DS für AWA (27) $_h$ = (30):

$$(27)_{h\tau} \quad \text{Ges. } v^j = \underline{U}_h^j \in I\!\!R^{N_h} : M_h v_t^j + \sigma K_h(t_{j+1}) \ v^{j+1} + (1-\sigma) \ K_h(t_j) v^j = \varphi^j$$

$$j = 0, 1, \dots, m-1,$$

$$\underline{AB:} \ v^0 = \underline{u}_h(0) := M_h^{-1} \underline{g}_h,$$

$$\text{wobei z.B. } \varphi^j := \sigma \underline{f}_h(t_{j+1}) + (1-\sigma) \underline{f}_h(t_j),$$

$$\text{falls } \underline{f}_h(\cdot) \text{ stetig ist (sonst Mittlungen)}.$$

$$(27)_{h\tau} \quad \text{Ges.} \qquad U_h^j \in V_{0h} : (U_{h,t}^j, v_h)_0 + \sigma \, a(t_{j+1}; U_h^{j+1}, v_h) + (1 - \sigma) \, a(t_j, U_h^j, v_h) =$$

$$= \langle \sigma F(t_{j+1}) + (1 - \sigma) \, F(t_j), v_h \rangle \quad \forall v_h \in V_{0h},$$

$$j = 0, 1, \dots, m - 1,$$

$$\underline{AB:} \, (U_h^0, v_h)_0 = (u_0, v_h)_0 \quad \forall v_h \in V_{0h}.$$

• Beispiel: $\sigma = 0$ Explizites Schema (= Euler vorwärts):

$$M_h \frac{v^{j+1} - v^j}{\tau} = \varphi^j - K_h(t_j) v^j \| \text{GS !}$$

Mass-Lumping, d.h.

Übergang zu einer \longrightarrow \longleftrightarrow (Vorsicht für Elemente höherer Ordnung, d.h. $k \ge 2$!)

Diagonalmatrix

$$D_h = \operatorname{diag}\left[\sum_{i \in \omega_h} m_{ki}\right]$$
bzw.
$$M_h = [m_{ki}]_{k,i \in \omega_h}$$



lineare Ansätze (k=1):
Integrationspkt.
= Knotenpunkt

 $\begin{array}{c|c} \sigma = 1/2 & \text{CRANK-NICOLSON (= Trapezregel)} \searrow \\ & \text{unbed. stabil bzgl.} \\ & L_2 \ \& \ \text{Energie} \\ & (\downarrow) \\ \end{array}$

 $\sigma = 1$ Rein implizites Schema (= Euler rückwärts) \nearrow

• Bemerkung:

1. Mass-Lumping wird auch für $\sigma \neq 0$ durchgeführt, um für System-Matrix M-Matrix-Eigenschaft zu erhalten:

2.
$$M_h \underline{\dot{u}}_h(t) + K_h(t) \underline{u}_h(t) = \underline{f}_h(t)$$

$$\downarrow \frac{1}{\tau} \int_{t_j}^{t_{j+1}} \dots dt$$

$$M_h \frac{\underline{u}_h(t_{j+1}) - \underline{u}_h(t_j)}{\tau} = \frac{1}{\tau} \int_{t_j}^{t_{j+1}} (\underline{f}_h(t) - K_h(t) \underline{u}_h(t)) dt \approx$$

$$\downarrow \qquad \qquad \sigma = 0: \text{Rechteckregel: } t = t_j$$

$$D_h \qquad \approx \xrightarrow{\nearrow} \sigma = 1/2: \text{Trapezregel}$$

$$\sigma = 1: \text{Rechteckregel: } t = t_{j+1}$$

$$\underline{u}_h(t_j) \longmapsto \underline{U}_h^j$$

• Anwendung der allgemeinen Stabilitätstheorie (Energienormstabilität) aus Pkt. 2.2:

<u>Vor.:</u> $K_h \neq K_h(t)$, d.h. Koeff. der Bilinearform sind t-unabhängig! $K_h = K_h^T$ p.d., d.h. Bilinearform ist sym., V_0 -ellipt., V_0 -beschränkt!

Allgem. Form:
$$\begin{array}{|c|c|c|c|c|c|}\hline C_1 \, v^{j+1} + C_0 \, v^j = \tau \, \varphi^j & & \updownarrow \\ v^0 \, \text{geg.} & & & C_1 = M_h + \tau \, \sigma \, K_h \\ & & & C_0 = -M_h + \tau \, (1-\sigma) K_h \\ \hline \end{array}$$

<u>Räume:</u> $H = I\!\!R^n$ mit Euklidischem Skalarprodukt $(\cdot, \cdot) := (\cdot, \cdot)_{I\!\!R^{N_h}}$ bzw. $L_2(\omega_h)$ mit FE-diskretem L_2 -Skalarprodukt $(\cdot, \cdot) := (\cdot, \cdot)_{M_h} = (M_h \cdot, \cdot)_{I\!\!R^{N_h}}$ oder $(\cdot, \cdot) := (\cdot, \cdot)_{D_h}$ und der entsprechenden Norm $\|\cdot\| = (\cdot, \cdot)^{0.5}$ (im weiteren sei $(\cdot, \cdot) = (\cdot, \cdot)_{I\!\!R^{N_h}}$).

Stabilitätssatz 2.16:

$$\underline{\mathbf{Bh}} : \|v^{j+1}\|_{A} \le \|v^{0}\|_{A} + T \max_{k=0,j} \|B^{-1}\varphi^{k}\|_{A}$$

Überprüfen Stabilitätsbedingung:

$$B \ge \frac{\tau}{2}A$$

wobei $\lambda_{\max} = \text{Maximaler EW des EWP: } K_h \underline{u}_h = \lambda M_h \underline{u}_h.$

Resultat: 1) $\sigma \geq 1/2 \implies$ Stabilitätsbed. trivialerweise erfüllt ! \Rightarrow unbed. stabil !

2)
$$\sigma = 0 \implies \text{bedingt stabil: } \tau \leq \frac{h^2}{2\delta_2} \leq \frac{1}{2\lambda_{\text{max}}} \left(\to \boxed{\ddot{U} \ 2.6} \right).$$

Bem.: $\|\underline{u}_h^j\|_A = \|\underline{u}_h^j\|_{K_h} \cong \|u_h^j\|_1$.

■ Zu Problem 4.): Gesamtfehlerabschätzung:

- 1) $\sigma = 1$: Pkt. 2.3.4.1:
- 2) $\sigma = \frac{1}{2}$: Pkt. 2.3.4.2:

■ Zusammenfassung:

(27)_{LVF} Ges.
$$u \in W_2^1(T, V_0)$$
 mit AB: $u(0) = u_0 \in L_2(\Omega)$ geg.:
$$\frac{d}{dt}(u(t), v)_0 + a(t; u, v) = \langle F(t), v \rangle \quad \forall v \in V_0 \quad \forall \text{ f.\"{ii.}} \ t \in T$$
Ausgangsaufgabe

(vertikale) Linienmethode / Ortsdiskr.

Zeitdiskr. \(\square\) Methode von Rothe

AWA für System gew. Dgl.

$$(27)_{h} \text{ Ges. } u_{h} = \sum_{i \in \omega_{h}} u^{(i)}(t) p^{(i)}(x) \in W_{2}^{1}(\boldsymbol{T}, V_{0}):$$

$$(\dot{u}_{h}, v_{h})_{0} + a(t; u_{h}, v_{h}) = \langle F, v_{h} \rangle \quad \forall v_{h} \in V_{h}$$

$$\forall \text{ f.ü. } t \in \boldsymbol{T}$$

$$+ \underline{AB:} (u_{h}(\cdot, 0), v_{h})_{0} = (u_{0}(\cdot), v_{h})_{0} \quad \forall v_{h} \in V_{h}$$

$$\underline{(27)_{h}} \text{ Ges. } \underline{u_{h}} = \underline{u_{h}}(t) \in [W_{2}^{1}(0, T)]^{N}:$$

$$M_{h}\underline{\dot{u}}(t) + K_{h}(t)\underline{u_{h}}(t) = \underline{f_{h}}(t) \quad \forall \text{ f.ü. } t \in \boldsymbol{T}$$

$$+ \underline{AB:} M_{h}\underline{u_{h}}(0) = \underline{g_{h}}$$

Methode von Rothe Semidiskrete Ersatzaufgabe

 $\begin{array}{c}
(27)_{h\tau} \\
\updownarrow \\
\underline{(27)}_{h\tau}
\end{array}$

Ges.
$$U_h^{j+1} \in V_{0h}: (U_{h,t}^j, v_h)_0 + \sigma \, a(t_{j+1}; U_h^{j+1}, v_h) +$$

$$+ (1 - \sigma) \, a(t_j; U_h^j, v_h) =$$

$$= \langle \sigma F(t_{j+1}) + (1 - \sigma) F(t_j), v_h \rangle \quad \forall v_h \in V_0$$

$$\underline{j} = 0, 1, \dots, m-1; \, \underline{AB}: (U_h^0, v_h)_0 = (u_0, v_h)_0 \quad \forall v_h \in V_0$$
Ges. $\underline{U}_h^{j+1} \in \mathbb{R}^{N_h}: M_h \underline{U}_{h,t}^j + \sigma K_h(t_{j+1}) \underline{U}_h^{j+1} +$

$$+ (1 - \sigma) K_h(t_j) \underline{U}_h^j = \underline{\varphi}_h^j,$$

$$\underline{j} = 0, 1, \dots, m-1$$

$$+ \underline{AB}: M_h \underline{u}_h^0 = \underline{g}_h, \, \text{mit } \underline{\varphi}_h^j = \sigma \underline{f}_h(t_{j+1}) + (1 - \sigma) \underline{f}_h(t_j)$$
Folge von $(m+1)$ linearen Gleichungssystemen

Volldiskrete Ersatzaufgabe

Methode von Rothe

$$(27)_{\tau} \text{ Ges. } U^{j+1} \in V_0 : (U^{j+1} \equiv u_{j+1}(\uparrow))$$

$$(U^j_t, v)_0 + \sigma a(t_{j+1}; U^{j+1}, v) + (1 - \sigma)a(t_j; U^j, v) = \langle \sigma F(t_{j+1}) + (1 - \sigma)F(t_j), v \rangle$$

$$\forall v \in V_0$$

$$j = 0, 1, \dots, m - 1; \quad 0 \le \sigma \le 1$$

$$+ \underline{AB:} U^0 = u_0 \in L_2(\Omega) \text{ für } \sigma = 1 \quad (\uparrow) \text{ bzw. } U^0 = u_0 \in V_0 \quad (!) \text{ sonst.}$$
Folge von m elliptischen RWA

2.3.3 Fehlerabschätzungen für die Lösung der semidiskreten Ersatzaufgabe $(27)_h$ in der L_2- und in der W_2^1- Norm

■ Btr. ARWA (27)_{LVF}:

(27)_{LVF} Ges.
$$u \in W_2^1(\mathbf{T}, V_0)$$
 mit AB: $u(0) = u_0 \in L_2(\Omega)$ geg.:
$$\frac{d}{dt}(u(t), v)_0 + a(u(t), v) = \langle F(t), v \rangle \quad \forall v \in V_0 \quad \forall \text{ f.\"{u}. } t \in \mathbf{T}$$

unter den Standardvoraussetzungen:

$$(35) \begin{array}{l} \left\{ \begin{array}{l} 1) \ a(\cdot,\cdot) : V_0 \times V_0 \mapsto I\!\!R^1 \ \ \text{sei} \ \underline{t\text{-unabh\"{a}ngige}} \ \ \text{Bilinearform:} \\ 2) \ V_0 - \text{elliptisch:} \ a(v,v) \geq \mu_1 \|v\|_1^2 \ \ \forall v \in V_0 \subset V = W_2^1(\Omega) \,, \\ 3) \ V_0 - \text{beschr\"{a}nkt:} \ |a(u,v)| \leq \mu_2 \|u\|_1 \ \|v\|_1 \ \ \forall u,v \in V_0 \,, \\ 4) \ V_0 - \text{symmetrisch:} \ a(u,v) = a(v,u) \ \ \forall u,v \in V_0 \ \text{(der Einfachheit halber !)} \end{array} \right.$$

und die semidiskrete Ersatzaufgabe $(27)_h$:

Ges.
$$u_h(x,t) = \sum_{i \in \omega_h} u^{(i)}(t) p^{(i)}(x) \in W_2^1(T, V_{0h}) \subset W_2^1(T, V_0)$$
:
$$\frac{d}{dt}(u_h, v_h)_0 + a(u_h, v_h) = \langle F(t), v_h \rangle \quad \forall v_h \in V_{0h} \quad \forall \text{ f.ü. } t \in T$$

$$\underline{AB:} (u_h(\cdot, 0), v_h)_0 = (u_0(\cdot), v_h(\cdot))_0 \quad \forall v_h \in V_{0h}$$

unter den Voraussetzungen:

$$\begin{cases} 1) \text{ Triangularisierung sei regulär i.S. der Def. II.4.3,} \\ 2) \tau(\Delta) = \operatorname{span}\{p^{(\alpha)}(\xi): \alpha \in A\} \supset P_k(\Delta), \\ 3) \text{ Regularität von } u \text{ und } \dot{u} \text{ } (\downarrow) \end{cases}$$

des Approximationssatzes (Satz II.4.5) aus [17] Nu II, Pkt. 4.4.2.

■ Ziel: 1)
$$||u(\cdot,t) - u_h(\cdot,t)||_0 \le .?$$
. Satz 2.19/2.19*, 2) $||u(\cdot,t) - u_h(\cdot,t)||_1 \le .?$. Satz 2.21.

■ Dazu benötigen wir den Ritz-Galerkin-Projektor (siehe [16] Pkt. I.4.1.3)

$$(37)_0 \quad R_h: V_0 \longmapsto V_{0h} \subset V_0: a(R_h, u, v_h) = a(u, v_h) \quad \forall v_h \in V_{0h}.$$

$$=:$$

$$< F_u, v_h >, \quad F_u \in V_0^* \text{ (mms)}$$

siehe \square : $R_h(=P_R$ in Nu I) ist Orthoprojektor bzgl. energetischen Skalarproduktes $[\cdot,\cdot]:=a(\cdot,\cdot)$, und es gilt in der Energienorm $|||\cdot|||^2=[\cdot,\cdot]=a(\cdot,\cdot)$:

(37)₁
$$|||u - R_h u||| = \inf_{v_h \in V_{0h}} |||u - v_h|||$$

$$\text{bzw. (Cea)}$$

$$||u - R_h u||_1 \le \sqrt{\frac{\mu_2}{\mu_1}} \inf_{v_h \in V_{0h}} ||u - v_h||_1 \le \bar{c}_{1,k+1} h^k |u|_{k+1},$$

falls $a(\cdot, \cdot)$ – symmetrisch.

Approximationssatz II.4.5

■ Zunächst L₂-Abschätzung (im Unterschied zu ellipt. RWA):

$$\parallel \underbrace{\underbrace{\underbrace{u(\cdot,t) - u_h(\cdot,t)}_{U_h(\cdot,t)}}_{\text{Eehler}} \parallel_0 \leq .?.$$

$$z(x,t) = u - u_h = \underbrace{u - R_h u}_{=: \rho = \rho(t) = \rho(x,t)} - \underbrace{\underbrace{(u_h - R_h u)}_{=: \theta_h = \theta_h(t) = \theta_h(x,t)}_{=: \theta_h = \theta_h(t) = \theta_h(x,t) \in V_{0h} !}_{(b)}$$

(a) Unter den Voraussetzungen von Satz II.4.9 (L_2 -Abschätzung für H^2 -koerziative elliptische RWA: \Rightarrow Nitsche-Trick !):

1) Reguläre Triangularisierung
2)
$$\tau(\Delta) \supset P_k(\Delta)$$
3) $u, \dot{u} \in W_2^{k+1}(\Omega) \quad \forall \quad (\text{f.\"{u}.}) \quad t \in \mathbf{T}$
4) $W_2^2 (= H^2)$ -Koerzitivität $(\Rightarrow (37)_0)$

gilt offenbar (siehe Satz II.4.9)

(b)
$$\|\theta_{h}(t)\|_{0} = \|u_{h}(\cdot,t) - R_{h}u(\cdot,t)\|_{0} \le .?.$$

Btr. $(\dot{\theta}_{h},v_{h})_{0} + a(\theta_{h},v_{h}) =$

$$= \underbrace{(\dot{u}_{h},v_{h})_{0} + a(u_{h},v_{h})}_{(27)_{h}} - ((R_{h}\dot{u}),v_{h})_{0} - a(R_{h}u,v_{h})$$

$$= \underbrace{(\dot{u}_{h},v_{h})_{0} + a(u_{h},v_{h})}_{(27)_{h}} - ((R_{h}\dot{u}),v_{h})_{0} - a(R_{h}u,v_{h})$$

$$= \underbrace{(\dot{u}_{h},v_{h})_{0} + a(u_{h},v_{h})}_{a(u,v_{h})} - ((R_{h}\dot{u}),v_{h})_{0} - a(R_{h}u,v_{h})$$

$$=\underbrace{\langle F(t), v_h \rangle - a(u, v_h)}_{\text{LVF}} - ((R_h \dot{u}), v_h)_0 \qquad a(\cdot, \cdot) \quad t \text{ -unabhängig !}$$

$$(27)_{\text{LVF}}, \quad V_{0h} \subset V_0 \rightarrow \underbrace{}_{\text{LVF}} = \frac{d}{dt}(u, v_h)_0 = \langle \dot{u}, v_h \rangle = (\dot{u}, v_h)_0$$

$$\uparrow \qquad \qquad \qquad \qquad \qquad \qquad \qquad \uparrow$$

$$\text{da } \dot{u} \in V_0^* \cap H^{k+1} \subset L_2(\Omega)$$

$$= \left(\frac{d}{dt}(u - R_h u), v_h\right)_0 = (\dot{\rho}(t), v_h)_0$$

Setzen hier
$$v_{h} = \theta_{h} := u_{h} - R_{h}u \in V_{0h}$$

$$\Rightarrow (\dot{\theta}_{h}, \theta_{h})_{0} + (\dot{\theta}_{h}, \theta_{h})_{0} = (\dot{\rho}, \theta_{h})_{0}$$

$$\frac{1}{2} \frac{d}{dt} (\theta_{h}, \theta_{h})_{0} \geq \mu_{1} \|\theta_{h}\|_{1}^{2} \geq \mu_{1} \|\theta_{h}\|_{0}^{2} \geq 0$$

$$\Rightarrow \frac{1}{2} \frac{d}{dt} \|\theta_{h}\|_{0}^{2} + \mu_{1} \|\theta_{h}\|_{0}^{2} \leq \|\dot{\rho}\|_{0} \|\theta_{h}\|_{0}$$

(40)
$$\frac{d}{dt} \|\theta_h\|_0 + \mu_1 \|\theta_h\|_0 \le \|\dot{\rho}\|_0 \quad \forall \text{ f.ü. } t \in \mathbf{T}$$

(40),
$$\frac{d}{dt} \|\theta_h\|_0 \le \|\dot{\rho}\|_0 \quad \forall \text{ f.ü. } t \in \mathbf{T}$$

$$\int_{0}^{t} (40)' ds \Rightarrow \|\theta_{h}(t)\|_{0} \leq \|\theta_{h}(0)\|_{0} + \int_{0}^{t} \|\dot{\rho}(s)\|_{0} ds =$$

$$= \|\underbrace{u_{h}(\cdot,0) - R_{h}u(\cdot,0)}_{0}\|_{0} + \int_{0}^{t} \|\dot{u}(\cdot,s) - R_{h}\dot{u}(\cdot,s)\|_{0} ds$$

$$= 0, \text{ falls } u_{h}(\cdot,0) = R_{h}u_{0}(\cdot), \text{ d.h. } a(u_{h}(\cdot,0),v_{h}) = a(u_{0},u_{h}) \quad \forall v_{h} \in V_{0h}$$

$$\leq \|u_{h}(\cdot,0) - u_{0}(\cdot)\|_{0} + \|u_{0} - R_{h}u_{0}\|_{0} + \int_{0}^{t} \|\dot{u} - R_{h}\dot{u}\|_{0} ds$$

$$\leq \|u_{h}(\cdot,0) - u_{0}(\cdot)\|_{0} + c_{0,k+1}h^{k+1}|u_{0}|_{k+1} + c_{0,k+1}h^{k+1} \int_{0}^{t} |\dot{u}(s)|_{k+1} ds$$

$$\leq a_{0,k+1}h^{k+1}|u_0|_{k+1} + 2c_{0,k+1}h^{k+1}\left\{|u_0|_{k+1} + \int\limits_0^t |\dot{u}(s)|_{k+1} \, ds\right\}.$$

Damit haben wir den folgenden Satz bewiesen:

■ Satz 2.19: $(L_2$ -Konvergenz: $O(h^{k+1})$)

 $\forall t \in \bar{T} = [0, T].$

\blacksquare Verbesserung der L_2 -Abschätzung aus Satz 2.19:

→ exponentielles Abfallen der Fehler in den AB! Dazu benötigen wir das Gronwall'sche Lemma:

Lemma 2.20: (Gronwall'sches Lemma)

 $u_h(\cdot,0) = P_{0h}u_0(\cdot)$

$$\begin{array}{c} \underline{\mathrm{Vor}.:}\ 1)\ f\in W^1_1(0,T)\ (\mathrm{d.h.\ absolut\ stetig})\colon f(t)\geq 0\quad\forall t\in\bar{T};\\ 2)\ \frac{d\ f(t)}{dt}\leq c_1(t)\ f(t)+c_2(t)\quad\forall\ \mathrm{f.\ddot{u}.}\ t\in T=(0,T);\\ 3)\ c_i\in L_1(0,T)\ \mathrm{f\ddot{u}r}\ i=1,2.\\ \\ \underline{\mathrm{Bh.:}}\ \mathrm{Dann\ gilt}\ \forall t\in\bar{T}=[0,T]\ \mathrm{die\ Absch\"{a}tzung};\\ 0\leq f(t)\leq e^\circ\qquad f(0)+e^\circ\qquad \int\limits_{-\infty}^t c_1(s)ds \ t\qquad -\int\limits_{-\infty}^\xi c_1(s)ds \ d\xi\,. \end{array}$$

Beweis:
$$\frac{d f(t)}{dt} e^{-\int_{0}^{t} c_{1}(s)ds} \leq c_{1}(t)f(t) e^{-\int_{0}^{t} c_{1}(s)ds} + c_{2}(t) e^{-\int_{0}^{t} c_{1}(s)ds}$$

$$= \frac{d}{dt} \left(f(t) e^{-\int_{0}^{t} c_{1}(s)ds} \right)$$

$$\Rightarrow c_{2}(t) e^{-\int_{0}^{t} c_{1}(s)ds}$$

$$\Rightarrow c_{3}(t) e^{-\int_{0}^{t} c_{1}(s)ds}$$

$$\Rightarrow c_{4}(t) e^{-\int_{0}^{t} c_{1}(s)ds}$$

$$\Rightarrow c_{5}(t) e^{-\int_{0}^{t} c_{1}(s)ds}$$

$$\Rightarrow c_{6}(t) e^{-\int_{0}^{t} c_{1}(s)ds}$$

$$\Rightarrow c_{7}(t) e^{-\int_{0}^{t} c_{1}(s)ds}$$

$$\Rightarrow c_{8}(t) e^{-\int_{0}^{t} c_{1}(s)ds}$$

q.e.d.

Satz 2.19*: (Verbesserte L_2 -Abschätzung)

Beweis: Analog zu Satz 2.19 mit folgenden Modifikationen:

• (38)
$$\|\rho(t)\|_0 = \|u - R_h u\|_0 \le c_{0,k+1} h^{k+1} |u(\cdot,t)|_{k+1}.$$

• Verwenden anstelle von

(40),
$$\frac{d}{dt} \|\theta_h\|_0 \le \|\dot{\rho}\|_0 \quad \forall \text{ f.ü. } t \in T$$

die verschärfte Abschätzung

(40)
$$\frac{d}{dt} \|\theta_h\|_0 \le -\mu_1 \|\theta_h\| + \|\dot{\rho}\|_0 \quad \forall \text{ f.\"{u}. } t \in T.$$

• Anwendung des Gronwall'schen Lemmas mit

$$f(t) = \|\theta_h(t)\|_0, \quad c_1 = -\mu_1, \quad c_2(t) = \|\dot{\rho}(t)\|_0$$

ergibt:

$$\|\theta_h(t)\|_0 \le e^{-\mu_1 t} \|\theta_h(0)\|_0 + \int_0^t e^{-\mu_1 (t-\xi)} \|\dot{\rho}(\xi)\|_0 d\xi.$$

q.e.d.

■ W₂-Abschätzungen:

Satz 2.21: $(W_2^1$ -Konvergenz: $O(h^k)$)

Beweis:
$$z = u - u_h = \underbrace{u - R_h u}_{\text{(a)} =: \rho} - \underbrace{(u_h - R_h u)}_{\text{(b)} =: \theta_h \in V_{0h}}$$

(a)
$$\|\rho(t)\|_1 = \|u - R_h u\|_1 \le c_{1,k+1} h^k |u(\cdot,t)|_{k+1}$$

(b) Setzen in (39)
$$(\dot{\theta}_{h}, v_{h})_{0} + a(\theta_{h}, v_{h}) = (\dot{\rho}, v_{h})_{0} \quad \forall v_{h} \in V_{0h}$$

$$v_{h} = \dot{\theta}_{h} \in V_{0h}:$$

$$\Longrightarrow (\dot{\theta}_{h}, \dot{\theta}_{h})_{0} + a(\theta_{h}, \dot{\theta}_{h}) = (\dot{\rho}, \dot{\theta}_{h})_{0} \leq \frac{\varepsilon}{2} \parallel \dot{\rho} \parallel_{0}^{2} + \frac{1}{2\varepsilon} \parallel \dot{\theta}_{h} \parallel_{0}^{2}$$

$$= \parallel \dot{\theta}_{h} \parallel_{0}^{2} = \frac{1}{2} \frac{d}{dt} a(\theta_{h}, \dot{\theta}_{h}) = \frac{1}{2} \frac{d}{dt} \parallel |\theta_{h}| \parallel^{2} \qquad \uparrow$$

$$\text{Cauchy}$$

$$\frac{\varepsilon - \text{Ungleichung:}}{ab \leq \frac{\varepsilon}{2} a^{2} + \frac{1}{2\varepsilon} b^{2}}$$

$$\varepsilon = \frac{1}{2}$$

$$\Longrightarrow (\frac{d}{dt} \parallel |\theta_{h}| \parallel^{2} \leq \frac{1}{2} \parallel \dot{\rho} \parallel_{0}^{2}) \text{ mit Bez. } \parallel |\theta_{h}| \parallel^{2} = a(\cdot, \cdot).$$

$$\int_{0}^{t} \implies \| |\theta_{h}(t)| \|^{2} \leq \| |\theta_{h}(0)| \|^{2} + \frac{1}{2} \int_{0}^{t} \| \dot{\rho}(s) \|_{0}^{2} ds$$

$$\Rightarrow \mu_{1} \|\theta_{h}(t)\|_{1}^{2} \leq \mu_{2} \|\theta_{h}(0)\|_{1}^{2} + \frac{1}{2} \int_{0}^{t} \|\dot{\rho}(s)\|_{0}^{2} ds$$

$$\Rightarrow \|\theta_{h}(t)\|_{1} \leq \sqrt{\frac{\mu_{2}}{\mu_{1}}} \|\theta_{h}(0)\|_{1} + \frac{1}{\sqrt{2\mu_{1}}} \left(\int_{0}^{t} \|\dot{\rho}(s)\|_{0}^{2} ds\right)^{1/2}$$

$$\uparrow \qquad \qquad \parallel$$

$$u_{h}(\cdot,0) - R_{h}u(\cdot,0) = u_{h}(\cdot,0) - R_{h}u_{0}(\cdot) = u_{h} - u_{0} + u_{0} - R_{h}u_{0}$$

$$\leq \sqrt{\frac{\mu_{2}}{\mu_{1}}} \|u_{h}(\cdot,0) - u_{0}(\cdot)\|_{1} + \sqrt{\frac{\mu_{2}}{\mu_{1}}} \bar{c}_{1,k+1}h^{k} \|u_{0}\|_{k+1} + \frac{1}{\sqrt{2\mu_{1}}} c_{0,k}h^{k} \left(\int_{0}^{t} |\dot{u}(\cdot,s)|_{k}^{2} ds\right)^{1/2}$$

q.e.d.

Bemerkung 2.22:

$$L_2$$
-Abschätzung — Inverse — W_2^1 -Abschätzung!

Unter Verwendung der inversen Ungleichung (siehe Lemma II.4.10)

$$(45) ||v_h||_{1,\Omega} \le c h^{-1} ||v_h||_{0,\Omega} \forall v_h \in V_{0h}$$

kann aus einer L_2 -Abschätzung immer auch eine W_2^1 -Abschätzung abgeleitet werden:

$$\begin{array}{c}
\operatorname{Satz} 2.19/2.19^* \\
\downarrow \\
\|u(\cdot,t) - u_h(\cdot,t)\|_{1} \leq (\bar{a}_{1,k+1} + c \, a_{0,k+1}) h^k |u(t)|_{k+1} + c \, h^{-1} \|u(\cdot,t) - u_h(\cdot,t)\|_{0}
\end{array}$$

Approximationssatz II.4.5 (Beweis)

Nun kann $||u - u_h||_0$ sowohl mit Satz 2.19 als auch mit Satz 2.19* weiter abgeschätzt werden. Falls $u_h(\cdot, 0) = P_{0h}u_0(\cdot)$, dann gilt:

$$(47) ||u_h(\cdot,0) - u_0(\cdot)||_0 \le a_{0,k+1}h^{k+1}|u_0|_{k+1},$$

und damit erhalten wir in beiden Fällen insgesamt eine $O(h^k)$ -Abschätzung. Im Falle der Anwendung von Satz 2.19* mit Abklingverhalten $e^{-\mu_1 t}$ des Fehlers in den AB.

Ü 2.7 | Man zeige, daß für die L_2 -Projektion der AB,

$$u_h(\cdot,0) = P_{0h} u_0(\cdot), \text{ d.h.}$$

$$(u_h(\cdot,0),v_h(\cdot))_0 = (u_0(\cdot),v_h(\cdot))_0 \quad \forall v_h \in V_{0h},$$

die Abschätzung

$$||u_h(\cdot,0) - u_0(\cdot)||_1 \le (2c \, a_{0,k+1} + \bar{a}_{1,k+1})h^k |u_0|_{k+1}$$

gilt!

Hinweis: Verwenden Sie wieder die inverse Ungleichung (45)!

Beweis:
$$u_h = P_{0h}u_0$$
 Satz II.4.5
a) $\|u_h(\cdot,0) - u_0(\cdot)\|_0 = \inf_{v_h \in V_{0h}} \|u_0 - v_h\|_0 \le a_{0,k+1}h^{k+1}|u_0|_{k+1}$

b)
$$||u_h(\cdot,0) - u_0(\cdot)||_1 \le ||u_h - I_h u_0||_1 + ||I_h u_0 - u_0||_1$$

 $\in V_{0h}$
 $\le c h^{-1} ||u_h - I_h u_0||_0 + ||I_h u_0 - u_0||_1 \le$
 $\le c h^{-1} (||u_h - u_0||_0 + ||u_0 - I_h u_0||_0) + ||I_h u_0 - u_0||_1 \le$
 $\le c h^{-1} (a_{0,k+1} h^{k+1} ||u_0|_{k+1} + a_{0,k+1} h^{k+1} ||u_0|_{k+1}) + \bar{a}_{1,k+1} h^k ||u_0|_{k+1} =$
 $= (2c a_{0,k+1} + \bar{a}_{1,k+1}) h^k ||u_0|_{k+1} = O(h^k).$

2.3.4 Fehlerabschätzung für volldiskrete Ersatzaufgaben in der L₂-Norm

2.3.4.1 Rein implizites Schema (Euler rückwärts)

■ Btr. rein implizites Schema $(27)_{h\tau} \leftrightarrow (\underline{27})_{h\tau} \ \ (\boldsymbol{\sigma}=\mathbf{1})$:

$$(27)_{h\tau} \text{ Ges. } U_h^{j+1} \in V_{0h} : (U_{h,t}^j, v_h)_0 + a(U_h^{j+1}, v_h) = \langle F(t_{j+1}), v_h \rangle \quad \forall v_h \in V_{0h},$$

$$\downarrow \qquad \qquad j = 0, 1, \dots, m - 1,$$

$$\underline{AB} : (U_h^0, v_h)_0 = (u_0, v_h)_0 \quad \forall v_h \in V_{0h}.$$

$$\underline{(27)}_{h\tau} \text{ Ges. } \underline{U}_h^{j+1} \in \mathbb{R}^{N_h} : M_h \underline{U}_{h,t}^j + K_h \underline{U}_h^{j+1} = \underline{\varphi}_h^j := \underline{f}_h(t_{j+1}), \quad j = \overline{0, m-1},$$

$$M_h \underline{U}_h^0 = \underline{g}_h$$

$$V_{0h} \ni U_h^j = \sum_{i \in \omega} U^{(i)}(t_j) \, p^{(i)}(x) \longleftrightarrow \underline{U}_h^j = [U^{(i)}(t_j)]_{i \in \omega_h} \in \mathbb{R}^{N_h}.$$

■ Satz 2.23: $(L_2$ -Konvergenz: $O(h^{k+1} + \tau)$ für impl. Schema)

$$\begin{array}{c} \underline{\text{Vor.:}} \ 1. \ \text{Vor. 1}) - 4) \ \text{aus Satz 2.19}. \\ 2. \ u, \dot{u} \in W_2^{k+1}(\Omega), \quad \ddot{u} \in L_2(\Omega) \quad \forall \ (\text{f.\"{u}.}) \ t \in \textbf{\textit{T}} \ (\downarrow). \end{array}$$

 $\underline{\operatorname{Bh}.:}$ Dann gilt die Diskretisierungsfehlerabschätzung

$$(27)_{\text{LVF}} \quad (27)_{h\tau} \leq a_{0,k+1}h^{k+1}|u_{0}|_{k+1}$$

$$(48) \quad ||u(\cdot,t_{j}) - U_{h}^{j}(\cdot)||_{0} \leq ||U_{h}^{0} - u_{0}||_{0} +$$

$$+2c_{0,k+1}h^{k+1}\{|u_{0}|_{k+1} + \int_{0}^{t_{j}}|\dot{u}(s)|_{k+1}\,ds\} + \tau \int_{0}^{t_{j}}||\ddot{u}(s)||_{0}\,ds,$$

$$j = 0, 1, \dots, m.$$

Beweis:

• Bez.
$$u(t_j) = u(\cdot, t_j) \in V_0 : (27)_{LVF}$$

$$U_h^j = U_h^j(\cdot) \in V_{0h} : (27)_{h\tau}$$

• Btr. Fehler

$$u(t_j) - U_h^j = \underbrace{u(t_j) - R_h \, u(t_j)}_{\mathbf{a}) =: \rho^j \in V_0} - \underbrace{(U_h^j - R_h \, u(t_j))}_{\mathbf{b}) =: \theta^j = \theta_h^j \in V_{0h}$$

• a)
$$\|\rho^j\|_0 = \|u(t_j) - R_h u(t_j)\|_0 \le c_{0,k+1} h^{k+1} |u(t_j)|_{k+1} =$$

$$= c_{0,k+1} h^{k+1} |u(0) + \int_0^{t_j} \dot{u}(s) \, ds|_{k+1} \le c_{0,k+1} h^{k+1} \{|u(0)|_{k+1} + \int_0^{t_j} |\dot{u}(s)|_{k+1} \, ds\}.$$

• b)
$$(\theta_t^j, v_h)_0 + a(\theta^{j+1}, v_h) =$$

$$= \underbrace{(U_{h,t}^j, v_h)_0 + a(U_h^{j+1}, v_h)}_{\parallel (U_h^j, v_h)_0 - a(R_h u(t_{j+1}), v_h)} - \underbrace{(R_h u)_t^j, v_h)_0 - a(R_h u(t_{j+1}), v_h)}_{\parallel (L_h^j, v_h)_0 - a(R_h u(t_{j+1}), v_h)} =$$

$$< F(t_{j+1}), v_h > = (\dot{u}(t_{j+1}), v_h)_0 + a(u(t_{j+1}), v_h)$$

$$\uparrow (27)_{\text{LVF}}, V_{0h} \subset V_0$$

$$= (\dot{u}(t_{j+1}) - R_h u_t(t_j), v_h)_0 = (\tilde{\rho}^j, v_h)_0 \quad \forall v_h \in V_{0h},$$

$$=: \tilde{\rho}^j$$

Resultat:

(50)

$$(\theta_t^j, v_h)_0 + a(\theta^{j+1}, v_h) = (\tilde{\rho}^j, v_h)_0 \quad \forall v_h \in V_{0h}, \text{ mit}$$

$$\tilde{\rho}^j = \dot{u}(t_{j+1}) - R_h \, u_t(t_j) = \underbrace{\dot{u}(t_{j+1}) - u_t(t_j)}_{=:: \tilde{\rho}_1^j} - \underbrace{(R_h - I) \, u_t(t_j)}_{=:: \tilde{\rho}_2^j} = \tilde{\rho}_1^j - \tilde{\rho}_2^j$$

Setzen in (50)
$$v_h = \theta^{j+1} \equiv \theta_h^{j+1} \in V_{0h}$$
:

(51)
$$\boxed{ (\theta_t^j, \theta^{j+1})_0 + \underbrace{a(\theta^{j+1}, \theta^{j+1})}_{\downarrow} = (\tilde{\rho}^j, \theta^{j+1})_0 } \le \|\tilde{\rho}^j\|_0 \|\theta^{j+1}\|_0$$

1. Variante: ≥ 0

$$\begin{array}{l} \underline{\text{2. Variante:}} \geq \mu_1 \|\theta^{j+1}\|_1^2 \geq \mu_1 \|\theta^{j+1}\|_0^2 \\ \Rightarrow \text{Satz 2.23* (mms)} \\ \rightarrow \text{Abklingverhalten !} \end{array}$$

$$\frac{1. \text{ Variante: } (\theta_t^j, \theta^{j+1})_0 \leq \|\tilde{\rho}^j\|_0 \|\theta^{j+1}\|_0}{\|\frac{\theta^{j+1} - \theta^j}{\tau}}$$

$$\implies \|\theta^{j+1}\|_{0}^{2} - \underbrace{(\theta^{j}, \theta^{j+1})_{0}}_{\leq \|\theta^{j}\|_{0} \|\theta^{j+1}\|_{0}} \leq \tau \|\tilde{\rho}^{j}\|_{0} \|\theta^{j+1}\|_{0}$$

$$\implies \|\theta^{j+1}\|_0^2 \leq \tau \|\tilde{\rho}^j\|_0 \ \|\theta^{j+1}\|_0 + \|\theta^j\|_0 \ \|\theta^{j+1}\|_0$$

$$\implies \|\theta^{j+1}\|_0 \le \|\theta^j\|_0 + \tau \|\tilde{\rho}^j\|_0$$

⇒ Rekursiv anwenden:

(52)
$$\|\theta^{j}\|_{0} \leq \|\theta^{0}\|_{0} + \tau \sum_{k=0}^{j-1} \|\tilde{\rho}^{k}\|_{0} \leq$$

$$\leq \|\theta^{0}\|_{0} + \tau \sum_{k=0}^{j-1} \|\tilde{\rho}_{1}^{k}\|_{0} + \tau \sum_{k=0}^{j-1} \|\tilde{\rho}_{2}^{k}\|_{0}$$

$$1) \qquad 2) \qquad 3)$$

2.3.4.2 CRANK-NICOLSON-Schema

■ Btr. CRANK-NICOLSON-Schema $(27)_{h\tau} \leftrightarrow (\underline{27})_{h\tau} \ (\sigma = \mathbf{0.5})$:

$$(27)_{h\tau} \operatorname{Ges.} U_{h}^{j+1} \in V_{0h} : (U_{h,t}^{j}, v_{h})_{0} + a(\frac{1}{2}(U_{h}^{j+1} + U_{h}^{j}), v_{h}) = \langle F(t_{j+\frac{1}{2}}), v_{h} \rangle \quad \forall v_{h} \in V_{0h}$$

$$\downarrow \qquad \qquad j = 0, 1, \dots, m - 1,$$

$$\underline{\operatorname{AB:}} (U_{h}^{0}, v_{h})_{0} = (u_{0}, v_{h})_{0} \quad \forall v_{h} \in V_{0h}.$$

$$\underline{(27)}_{h\tau} \operatorname{Ges.} \underline{U}_{h}^{j+1} \in \mathbb{R}^{N_{h}} : M_{h}U_{h,t}^{j} + K_{h}(\frac{1}{2}(\underline{U}_{h}^{j+1} + \underline{U}_{h}^{j})) = \underline{f}_{h}(t_{j+\frac{1}{2}}), \quad j = \overline{0, m - 1}$$

$$M_{h}\underline{U}_{h}^{0} = \underline{g}_{h}$$

■ Satz 2.24: $(L_2$ -Konvergenz: $O(h^{k+1} + \tau^2)$ für CN-Schema)

Vor.: 1. Vor. 1) – 4) aus Satz 2.19.
$$2.\ u,\dot{u}\in W_2^{k+1}(\Omega),\quad \ddot{u}\in W_2^2(\Omega),\quad u'''\in L_2(\Omega)\quad \forall \text{ f.\"{u}. }t\in \boldsymbol{T}\ (\downarrow),$$
wobei $u'''=d^3u|dt^3.$

Bh.: Dann gilt die Diskretisierungsfehlerabschätzung

$$(27)_{\text{LVF}} \qquad (27)_{h\tau} \qquad \leq a_{0,k+1} h^{k+1} |u_{0}|_{k+1}$$

$$(53) \qquad \| u(\cdot,t_{j}) - U_{h}^{j}(\cdot) \|_{0} \leq \underbrace{\| U_{h}^{0} - u_{0} \|_{0}}_{+} + 2c_{0,k+1} h^{k+1} \{ |u_{0}|_{k+1} + \int_{0}^{t_{j}} |\dot{u}(s)|_{k+1} ds \} +$$

$$+ \left(\frac{1 + 2c_{2}}{8} \right) \tau^{2} \left[\int_{0}^{t_{j}} (\|\ddot{u}(s)\|_{0} + \|\ddot{u}(s)\|_{2}) ds \right],$$

wobei $c_2 = W_2^2$ -Beschränktheitskonstante (\Rightarrow partielle Integration):

(54)
$$|a(w,v)| \leq c_2 ||w||_2 ||v||_0 \quad \forall w \in W_2^2(\Omega) \cap V_0, \quad \forall v \in V_0,$$
$$c_{0,k+1} = \text{Konstante aus Satz II.4.9 } (L_2\text{-Abschätzung}),$$
$$a_{0,k+1} = \text{Konstante aus Approximationssatz II.4.5}.$$

Beweis:

• Btr. Fehler

$$z^{j} = u(t_{j}) - U_{h}^{j} = \underbrace{u(t_{j}) - R_{h} u(t_{j})}_{\mathbf{a}) =: \rho^{j} \in V_{0}} - \underbrace{(U_{h}^{j} - R_{h} u(t_{j}))}_{\mathbf{b}) := \theta^{j} = \theta_{h}^{j} \in V_{0h}$$

• a)
$$\|\rho^j\|_0 \le c_{0,k+1} h^{k+1} |u(t_j)|_{k+1} \le c_{0,k+1} h^{k+1} \{|u_0|_{k+1} + \int_0^{t_j} |\dot{u}(s)|_{k+1} ds\}.$$

• b)
$$(\theta_t^j, v_h)_0 + a\left(\frac{1}{2}(\theta^{j+1} + \theta^j), v_h\right) = a(\cdot, \cdot) = t$$
-unabhängig!

unter Benutzung der Beziehungen

(56)
$$\dot{u}(t_{j+1/2}) - u_t(t_j) = \frac{1}{2\tau} \left(\int_{t_j}^{t_{j+1/2}} (s - t_j)^2 u'''(s) \, ds + \int_{t_{j+1/2}}^{t_{j+1}} (s - t_{j+1})^2 u'''(s) \, ds \right)$$

und

$$(57)$$

$$u(t_{j+1/2}) - \frac{1}{2}(u(t_{j+1}) - u(t_{j})) = \frac{1}{2} \left[\int_{t_{j+1/2}}^{t_{j+1}} (s - t_{j+1}) \ddot{u}(s) ds - \int_{t_{j}}^{t_{j+1/2}} (s - t_{j}) \ddot{u}(s) ds \right]$$

$$1 \times \text{ partiell integrieren !}$$

erhalten wir die Abschätzungen

$$= \frac{1}{\tau} \int_{t_{j}}^{t_{j+1}} \dot{u}(s) \, ds \, \uparrow$$

$$a(\cdot, \cdot) \, t\text{-unabhängig}$$

$$\leq \|\frac{1}{\tau} \int_{t_{j}}^{t_{j+1}} (I - R_{h}) \, \dot{u}(s) \, ds\|_{0} \, \|v_{h}\|_{0} \leq$$

$$\leq \frac{1}{\tau} \int_{t_{j}}^{t_{j+1}} \|(I - R_{h}) \, \dot{u}(s) \, ds\|_{0} \, \|v_{h}\|_{0} \leq$$

$$\leq \frac{1}{\tau} \, c_{0,k+1} \, h^{k+1} \int_{t_{j}}^{t_{j+1}} |\dot{u}(s)|_{k+1} \, ds \cdot \|v_{h}\|_{0}.$$

(iii)
$$\left| a \left(u(t_{j+1/2}) - \frac{1}{2} (u(t_{j+1}) + u(t_j)), v_h \right) \right| \stackrel{(57)}{=}$$

$$= \left| a \left(\frac{1}{2} \left[\int_{t_{j+1/2}}^{t_{j+1}} (s - t_{j+1}) \ddot{u}(s) ds - \int_{t_{j}}^{t_{j+1/2}} (s - t_{j}) \ddot{u}(s) ds \right], v_{h} \right) \right| (54)$$

$$\leq c_{2} \left\| \frac{1}{2} \left[\int_{t_{j+1/2}}^{t_{j+1}} (s - t_{j+1}) \ddot{u}(s) ds - \int_{t_{j}}^{t_{j+1/2}} (s - t_{j}) \ddot{u}(s) ds \right] \right\|_{2} \|v_{h}\|_{0} \leq$$

$$\leq \frac{c_{2}}{4} \tau \int_{t_{j}}^{t_{j+1}} \|\ddot{u}(s)\|_{2} ds \cdot \|v_{h}\|_{0}.$$

Aus (i) - (iii) folgt sofort die Abschätzung

$$(58) \quad |\langle \hat{\rho}^{j}, v_{h} \rangle| \leq \left[\frac{\tau}{8} \int_{t_{j}}^{t_{j+1}} ||u'''(s)||_{0} ds + \frac{1}{\tau} c_{0,k+1} h^{k+1} \int_{t_{j}}^{t_{j+1}} |\dot{u}(s)|_{k+1} ds + \frac{c_{2}}{4} \tau \int_{t_{j}}^{t_{j+1}} ||\ddot{u}(s)||_{2} ds\right] ||v_{h}||_{0} = [(58)^{j}] ||v_{h}||_{0}$$

Setzen in (55) $v_h = \frac{1}{2}(\theta^j + \theta^{j+1}) \in V_{0h}: \Rightarrow$

$$\left(\theta_t^j, \frac{1}{2}(\theta^j + \theta^{j+1})\right)_0 + \underbrace{a\left(\frac{\theta^j + \theta^{j+1}}{2}, \frac{\theta^j + \theta^{j+1}}{2}\right)}_{\geq \mu_1 \|\frac{\theta^j + \theta^{j+1}}{2}\|_1^2 \geq \frac{\mu_1}{4} \|\theta^j + \theta^{j+1}\|_0^2 \geq 0 !$$

$$\underbrace{\|\theta^{j+1}\|_{0}^{2} - \|\theta^{j}\|_{0}^{2}}^{(58)} \stackrel{\downarrow}{\stackrel{\checkmark}{\leq}} \tau[(58)^{j}] (\|\theta^{j}\|_{0} + \|\theta^{j+1}\|_{0})$$

$$= (\|\theta^{j+1}\|_{0} - \|\theta^{j}\|_{0}) (\|\theta^{j}\|_{0} + \|\theta^{j+1}\|_{0})$$

$$\Rightarrow \boxed{ \|\theta^{j+1}\|_{0} \leq \|\theta^{j}\|_{0} + \tau[(58)^{j}] \leq \|\theta^{0}\|_{0} + \tau \sum_{l=0}^{j} [(58)^{l}] }$$

Wegen $\|\theta^0\|_0 \le \|U_h^0 - u_0\|_0 + \|u_0 - R_h u_0\|_0 \le \|U_h^0 - u_0\|_0 + c_{0,k+1} h^{k+1} \|u_0\|_{k+1}$ erhalten wir

$$(59) \|\theta^{j}\|_{0} \leq \|\theta^{0}\|_{0} + \tau \sum_{l=0}^{j-1} [(58)^{l}] \leq$$

$$\leq \|U_{h}^{0} - u_{0}\|_{0} + c_{0,k+1} h^{k+1} |u_{0}|_{k+1} + \frac{\tau^{2}}{8} \int_{0}^{t_{j}} \|u'''(s)\|_{0} ds +$$

$$+ c_{0,k+1} h^{k+1} \int_{0}^{t_{j}} |\dot{u}(s)|_{k+1} ds + \frac{c_{2}}{4} \tau^{2} \int_{0}^{t_{j}} \|\ddot{u}(s)\|_{2} ds.$$

Aus a) und (59) folgt (53).

q.e.d.

2.3.5 Abschließende Bemerkungen

lacksquare Fehlerabschätzungen in anderen Normen (z.B. $H^{\scriptscriptstyle 1}, L_{\scriptscriptstyle \infty}$):

1. Über <u>inverse Ungleichung</u> unter Benutzung der L_2 -Abschätzungen. <u>Vorsicht:</u> z.B. liefert diese Technik für CN-Schema offenbar die folgende H^1 -Abschätzung:

$$||u(\cdot,t_j) - u_h^j(\cdot)||_1 = O(h^k + h^{-1}\tau^2) = O(h)!$$

? $t = 0$

- 2. Direkte Abschätzung ($e^{-\mu_1 t}$ bzw. $e^{-\mu_1 (t-s)}$ Abklingen ?)
 - a) H^1 –Norm: analog zu Satz 2.21 !
 - b) L_{∞} -Norm: siehe [4] S. 311, [21] S. 62 ff.
- 3. Mass-Lumping ist nur für k=1 (lineare Dreieckselemente) begründet: siehe [4] S. 312 ff., [21] S. 166 ff.

■ Auflösung:

- 1. Explizites Schema $(\sigma = 0)$: $M_h \underline{U}_{h,t}^j + K_h(t_j) \underline{U}_h^j = \underline{f}_h(t_j)$
 - (60) $M_h \underline{U}_h^0 = \underline{g}_h =: \underline{d}_h^0$ $j = 0, 1, \dots, m-1$
 - (60) $M_{h} \underline{U}_{h}^{j+1} = \underline{d}_{h}^{j} := M_{h} \underline{U}_{h}^{j} + \tau (\underline{f}_{h}(t_{j}) K_{h}(t_{j}) \underline{U}_{h}^{j})$ $\downarrow \qquad \qquad \longleftarrow \text{Mass-Lumping für } k = 1$ $D_{h} \qquad \qquad D_{h} \Rightarrow \text{Lösung ist explizit hinschreibbar } !$

Bemerkung: $\kappa(M_h) \leq \nu_2/\nu_1 = O(1)$, d.h. M_h ist gut konditioniert. Folglich sind obige GS mit Systemmatrix offenbar selbst mit klassischen Iterationsverfahren (z.B. Jacobi-Verfahren) schnell auflösbar. Startnäherung für (60) ist Näherung \underline{U}_h^{j,h_j} von letzter Zeitschicht.

2. Implizite Schemata (0 < $\sigma \le 1$):

(60)
$$M_h \underline{U}_h^0 = \underline{d}_h^0 := \underline{g}_h^0$$
$$j = 0, 1, \dots, m - 1$$

Offenbar gelten für die Systemmatrizen $M_h + \tau \sigma K_h(t_{j+1})$ die folgenden Eigenwert- und Konditionsabschätzungen:

$$\begin{array}{lll}
\nu_1 \, h^d + \tau \, \sigma \, \underline{c}_E \, h^d & \leq \lambda (M_h + \tau \, \sigma \, K_h) & \leq \nu_2 \, h^d + \tau \, \sigma \, h^{d-2} \\
& \geq & EW & \leq \\
(\nu_1 + \tau \, \sigma) \, h^d & (\nu_2 + \sigma \, \tau \, h^{-2}) \, h^d \\
& \geq & \\
\nu_1 \, h^d
\end{array}$$

Also:
$$\kappa(M_h + \tau \sigma K_h) \le \frac{\nu_2 + \sigma \tau h^{-2}}{\nu_1 + \sigma \tau} \le \frac{\nu_2}{\nu_1} + \frac{\sigma}{\nu_1} \tau h^{-2}$$

Damit gilt z.B. für k=1 und $K_h \neq K_h(t_j)$ o. B. d. Allg.:

a)
$$\underline{\sigma=1}$$
: L_2 -Konvergenz: $O(h^2+\tau)$, d.h. $\tau=O(h^2)$
$$\Longrightarrow \qquad \kappa(M_h+\tau\,K_h) \leq \frac{\nu_2}{\nu_1} + \frac{1}{\nu_1} \cdot \mathbf{c} = O(1)$$

$$\uparrow \qquad \qquad \qquad \qquad \uparrow$$

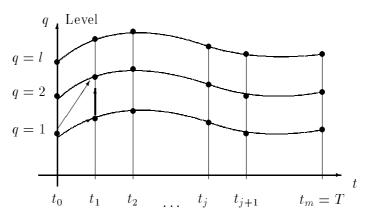
<u>Praktisch:</u> $\tau \le c \, h^2$, wobei c = const. > 0 wesentlich größer sein kann als Konstante in der <u>Stabilitätsbedingung</u> für explizites Schema!

b)
$$\underline{\sigma = 1/2}$$
: L_2 -Konvergenz: $O(h^2 + \tau^2)$, d.h. $\tau = O(h)$

$$\implies \kappa \left(M_h + \frac{\tau}{2} K_h \right) \le \frac{\nu_2}{\nu_1} + \frac{1}{2\nu_1} \cdot \mathbf{c} \ h^{-1} = O(h^{-1})$$

$$\uparrow \qquad \qquad \qquad \uparrow$$

⇒ Lösung von (61) durch Multigrid-Verfahren bzw. MG-Präkonditionierte CG-Verfahren + Nested Iteration (siehe auch [17]):



- Nichtlineare parabolische ARWA: (⇒ Kap. 3)
 - Nichtlinearitäten: a) Koeffizienten hängen von der Temperatur ab
 b) Strahlung
 - $\{V_0, H, V_0^*\}$ sei Evolutionstripel, $T = (0, T), \quad \frac{1}{p} + \frac{1}{q} = 1$

(62) Ges.
$$u \in W_p^1(T; V_0, H) := \{u \in L_p(T, V_0) : \exists \ \dot{u} \in L_q(T, V_0^*)\} :$$

$$\dot{u}(\cdot) + A(\cdot) \ u(\cdot) = F(\cdot) \text{ in } L_q(T, V_0^*)$$

$$\underline{AB} : \ u(0) = u_0 \in H$$

$$\text{mit geg. } F \in L_q(T, V_0^*), \quad u_0 \in H,$$

$$A(\cdot) : V_0 \to V_0^*$$

$$\text{Genauer: } A : W_p^1(T; V_0, H) \longrightarrow L_q(T, V_0^*)$$

FE-Galerkin-Semidiskretisierung

$$(62)_h \qquad \begin{array}{l} \operatorname{Ges.} \ \underline{u}_h(t) \in [W^1_p(0,T)]^N: \\ M_h \underline{\dot{u}}_h(t) + K_h(t,\underline{u}_h(t)) = \underline{f}_h(t) \quad \forall \text{ f.\"{u}. } t \in \mathbf{T} \\ M_h \underline{u}_h(0) = \underline{g}_h \quad \parallel \text{ oft} \\ \hat{K}_h(t;\underline{u}_h(t)) \ \underline{u}_h(t) - \text{quasilinear} \end{array}$$

AWA für System nichtlinearer gew. Dgl. 1. Ordnung Ges. $u(\cdot): \bar{T} \to \mathbb{R}^N: \dot{u}(t) = f(t, u(t)), \quad t \in T$ AB: $u(0) = u_0$ geg.

• <u>Literatur:</u> [2], [10], [23].

Kapitel 3

Anfangswertaufgaben für gewöhnliche Differentialgleichungen und Systeme gewöhnlicher Dgl.

■ <u>Literatur:</u>

- [1] Deuflhard P., Bornemann F.: Numerische Mathematik II: Integration gewöhnlicher Differentialgleichungen. de Gruyter Lehrbuch; Berlin · N.Y. 1994.
- [5] Hairer E., Nørsett S.P., Wanner G.: Solving Ordinary Differential Equation I: Nonstiff Problems.
 Springer-Verlag, Berlin · Heidelberg 1987.
- [6] Hairer E., Wanner G.: Solving Ordinary Differential Equation II: Stiff and Differential-Algebraic Problems. Springer-Verlag, Berlin · Heidelberg 1991.
- [19] Stetter H.J.: Analysis of Discretization Methods for Ordinary Differential Equations. Springer-Verlag, Berlin · Heidelberg · N.Y. 1973.

■ Bezeichnung:

$$\dot{u}(t) = u'(t) := \frac{d}{dt}u(t)$$

3.1 Beispiele

■ Bsp. 3.1: Chemische Reaktionen (Brusselator)

Brusselator (von R. Lefever und G. Nicolis, 1971) ist ein Modell einer chemischen Reaktion, die aus den folgenden Einzelreaktionen der Substanzen A, B, D, E, X, Y besteht:

wobei k_i – Reaktionsgeschwindigkeitskoeffizient (reaction rate coefficient).

Aufgrund des <u>Massenwirkungsgesetzes</u> ergibt sich die folgende Reaktionskinetik (siehe [1], S. 8 ff.), beschrieben durch ein System von Dgl. für die Konzentrationen $c_A = c_A(t)$, $c_B = c_B(t)$, $c_D = c_D(t)$, $c_E = c_E(t)$, $c_X = c_X(t)$ und $c_Y = c_Y(t)$ der einzelnen Substanzen als Fkt. der Zeit t:

+ AB:
$$c_A(0)$$
, $c_B(0)$, $c_D(0)$, $c_E(0)$, $c_X(0)$, $c_Y(0)$ geg. $(c_A + c_B + c_D + c_E + c_X + c_Y = 1)$.
Literatur: [1], S. 8 – 14; [5], S. 111 f.

<u>Ü</u> 3.1 Der <u>Oregonator</u> (Zhabotinski-Belousov-Reaktion einer chemischen Oszillation) wird durch das Reaktionsschema

$$BrO_{3}^{-} + Br^{-} \xrightarrow{k_{1}} HBrO_{2}$$

$$HBrO_{2} + Br^{-} \xrightarrow{k_{2}} P$$

$$BrO_{3}^{-} + HBrO_{2} \xrightarrow{k_{3}} 2HBrO_{2} + Ce(IV)$$

$$2HBrO_{2} \xrightarrow{k_{4}} P$$

$$Ce(IV) \xrightarrow{k_{5}} Br^{-}$$

beschrieben, mit $k_1=1.34, k_2=1.6\cdot 10^9, k_3=8.0\cdot 10^3, k_4=4.0\cdot 10^7, k_5=1.0$ und geg. AB für die ges. Konzentrationen $c_1=c_{BrO_3^-}(t), c_2=c_{Br^-}(t), c_3=c_{HBrO_2}(t), c_4=c_P(t), c_5=c_{Ce(IV)}(t).$

Stellen Sie das Differentialgleichungssystem für die ges. Konzentrationen $c_1(t),\ c_2(t),\ c_3(t),\ c_4(t),\ c_5(t)$ auf, und lösen Sie es (später) numerisch mit einem geeigneten (?) Integrationsverfahren für die Anfangskonzentrationsverteilung: $c_1(0)=0.25,\ c_2(0)=0.25,\ c_3(0)=0,\ c_4(0)=0,\ c_5(0)=0.5$.

Ü 3.2 Die Robertson-Reaktion (1966) wird durch das folgende Reaktionsschema beschrieben:

$$A \xrightarrow{0.04} B \qquad \text{(langsame Reaktion)}$$

$$B + B \xrightarrow{3 \cdot 10^7} C + B \quad \text{(sehr schnell)}$$

$$B + C \xrightarrow{10^4} A + C \quad \text{(schnell)}$$

Stellen Sie das Dgl.-System zur Bestimmung der Konzentrationen $c_A = c_A(t)$, $c_B = c_B(t)$, $c_C = c_C(t)$ auf, und lösen Sie es (später) numerisch mit einem geeigneten (?) Verfahren für die AB: $c_A(0) = 1$, $c_B(0) = 0$, $c_C(0) = 0$ (vgl. Willoughby, 1974).

■ Bsp. 3.2: Semidiskrete (\Rightarrow Vertikale Linienmethode: siehe Pkt. 2.3.) lineare parabolische A(R)WA:

Ges.
$$\underline{u}_h(t) \in [W_2^1(0,T)]^N$$
: $M_h \, \underline{\dot{u}}_h(t) + K_h(t) \, \underline{u}_h(t) = \underline{f}_h(t)$, $\forall \text{ f.\"{u}. } t \in \mathbf{T} = (0,T)$
 $+ \text{AB: } M_h \, \underline{u}_h(0) = \underline{g}_h$

Besonderheiten: • $\underline{u}_h(t) = [u^{(i)}(t)]_{i \in \omega_h} : u^{(i)}(\cdot) \in W_2^1(T)$ (!) mit geg. Daten $\underline{f}_h(\cdot) \in [L_2(T)]^N$ und L_{∞} -Koeffizienten der Matrix $K_h(t)$.

• $h \to 0 \Rightarrow N = N_h = |\omega_h| \to \infty$ (keine fixierte Dimension !!).

Bemerkung: \exists ! (Satz 2.18), Integrationsverfahren (σ -gew. DS), Stabilität, Approximation und Fehlerabschätzungen $O(\tau^q + h^p)$ siehe Pkt. 2.3.

■ Bsp. 3.3: Semidiskrete (⇒ Vertikale Linienmethode) <u>nichtlineare</u> parabolische A(R)WA (z.B. nichtlineare Wärmeleitprobleme):

$$\begin{split} &\text{Ges. } u \in W^1_p(\boldsymbol{T}; V_0, H) := \{u \in L_p(\boldsymbol{T}, V_0) : \exists \ \dot{u} \in L_q(\boldsymbol{T}, V_0^*), p^{-1} + q^{-1} = 1\} : \\ &\dot{u}(\cdot) + A(\cdot) \ u(\cdot) = F(\cdot) \text{ in } L_q(\boldsymbol{T}, V_0^*) \ + \text{AB} : u(0) = u_0 \text{ in } H \\ &\text{mit geg. } F \in L_q(\boldsymbol{T}, V_0^*), \ u_0 \in H, \ A : W^1_p(\boldsymbol{T}, V_0, H) \longrightarrow L_q(\boldsymbol{T}, V_0^*) \text{ nichtlinear}. \end{split}$$

 $\leftarrow \text{FE-Galerkin-Semidiskretisierung (vertikale Linienmethode)}$

Ges.
$$\underline{u}_h(t) \in [W^1_p(0,T)]^{N_h}: M_h \, \underline{\dot{u}}_h(t) + K_h(t,\underline{u}_h(t)) = \underline{f}_h(t), \quad \forall \text{ f.\"{u}. } t \in \mathbf{T}$$

$$+ \text{ AB: } M_h \underline{u}_h(0) = \underline{g}_h$$
oft liegt der quasilineare Fall vor: $K_h(t,\underline{u}_h(t)) = \hat{K}_h(t,\underline{u}_h(t)) \, \underline{u}_h(t).$

AWA für System nichtlinearer gew. Dgl. 1. Ordnung: Ges. $u(\cdot): \overline{T} \to \mathbb{R}^N: u'(t) = f(t, u(t)), \ t \in T$ mit AB: $u(0) = u_0$ geg., die in diesem Kap. behandelt werden (allerdings in Räumen stetig diff. Fkt.)

■ Bsp. 3.4: van der Polsche Differentialgleichung ([5], S. 107 – 111) Die van der Polsche Dgl.

$$y''(t) - \varepsilon(1 - y^2(t))y'(t) + y(t) = 0, t > 0$$

+ AB: y(0), y'(0) geg.

beschreibt Kippschwingungen in einem elektrischen Schaltkreis. Es handelt sich hier um eine <u>nichtlineare</u> Dgl. 2. Ordnung.

 $\ddot{\mathbf{U}}$ 3.3 In [5] S. 110, Figure 16.3, wird die Lösung der van der Polschen Dgl. für $\varepsilon=10$ und für die konkreten Anfangsbedingungen y(0)=0 und y'(0)=0 gezeigt. Man bestimme die Lösung mit einem geeigneten (?) Integrationsverfahren numerisch.

■ Bsp. 3.5: Newtonsche Himmelsmechanik $(mx'' = F(x, x'), F = -\nabla U, U - \text{Gravitationspotential})$: Das restringierte Dreikörperproblem (siehe [1], S. 3 – 8 bzw. [5], S. 127 – 129):

Die Bahn eines Satelliten in der Ebene des Erde-Mond-Systems läßt sich durch das folgende System von Dgl. 2. Ordnung beschreiben:

$$y_1'' = y_1 + 2y_2' - (1 - \mu) \frac{y_1 + \mu}{D_1} - \mu \frac{y_1 - (1 - \mu)}{D_2}, \quad t > 0,$$

$$y_2'' = y_2 - 2y_1' - (1 - \mu) \frac{y_2}{D_1} - \mu \frac{y_2}{D_2}, \quad t > 0$$

$$+ \text{AB: } y_1(0), y_1'(0), y_2(0), y_2'(0) \text{ geg.},$$

$$\text{mit } D_1 = [(y_1 + \mu)^2 + y_2^2]^{3/2}, \quad D_2 = [(y_1 - (1 - \mu))^2 + y_2^2]^{3/2},$$

$$\mu = 0.012277471.$$

Bei geeigneter Setzung der AB ergeben sich periodische Lösungen!

 $\ddot{\mathbf{U}}$ 3.4 Man bestimme die Lösung $(y_1(t), y_2(t))$ mit einem geeigneten (?) Integrationsverfahren numerisch für die Anfangswerte:

 $y_1(0) = 0.994$

 $y_1'(0) = 0$

 $y_2(0) = 0$

 $y_2'(0) = -2.00158510637908252240537862224$

<u>Hinweis:</u> $t_{per} = 17.0652165601579625588917206249$, $\tau = \frac{t_{per}}{m}$ mit m > 6000 (RK).

■ Bsp. 3.6: Semidiskrete (⇒ Vertikale Linienmethode) <u>lineare</u> bzw. <u>nichtlineare</u> hyperbolische A(R)WA (Schwingungsgleichung, Wellengleichung):

Ges.
$$u(x,t)$$
: $(\ddot{u},v)_0 + a(t;u,v) = \langle F(t),v \rangle \quad \forall v \in V_0 \quad \forall$ f.ü. $t \in T$ + AB: $u(0) = u_0, \dot{u}(0) = u_1$ (vgl. auch [17] Numerik II, Pkt. 1.2.3)

 $\Big\| \longleftarrow \text{FE-Galerkin-Semidiskretisierung (vertikale Linienmethode)}$

$$\begin{split} \text{Ges.}\,\,\underline{u}_h(t): & \ M_h\underline{\ddot{u}}_h(t) + K_h(t;\underline{u}_h(t)) = \underline{f}_h(t), \quad \forall \text{ f.\"{u}.} \, t \in \mathbf{T} \\ & + \text{AB:} \quad M_h\underline{u}_h(0) = \underline{u}_0, \\ & \ M_h\underline{\dot{u}}_h(0) = \underline{u}_1, \\ \text{wobei} & \ K_h(t;\underline{u}_h(t)) := \hat{K}_h(t,\underline{u}_h(t))\underline{u}_h(t) \text{ im quasilinearen Fall und} \\ & \ K_h(t;\underline{u}_h(t)) := \hat{K}_h(t)\underline{u}_h(t) \text{ im linearen Fall.} \end{split}$$

AWA für Systeme nichtlinearer gew. Dgl. 2. Ordnung: Ges. $u(\cdot): \bar{T} \to I\!\!R^N: \quad u''(t) = f(t,u(t)), \ t \in T = (0,T),$ mit AB: $u(0) = u_0, \ u'(0) = u_1.$

- Bsp. 3.7: Evolutionsgleichungen für interne Parameter in Prozessen, die von der Prozeßgeschichte abhängen:
 - z.B. Elastisch-plastische Fließprobleme in der Festkörpermechanik (vgl. auch [17] Numerik II, Pkt. 3.1.2.2: lineares Elastizitätsproblem \Rightarrow quasistatisch) [12]:

Ges. Verschiebungen
$$u \in C^1(T, V_0)$$
,
Spannungen $\sigma \in C^1(T, S)$,
Verfestigungsparameter $\kappa \in C^1(T, H)$:

- (i) das verallgemeinerte Kräftegleichgewicht $\smallint_{\Omega} \dot{\sigma}^T \varepsilon(v) \; dx = \, < F, v > := \smallint_{\Omega} \dot{f}^T v \; dx + \smallint_{\Gamma_2} \dot{g}^T v \; ds, \quad \forall v \in V_0, \;\; t \in \boldsymbol{T},$
- (ii) die differentiellen Spannungs-Verzerrungsbez. (Evolutionsgl.) $\dot{\sigma} = D\dot{\varepsilon} D\alpha(\sigma,\kappa,\dot{\varepsilon}),$
- (iii) die Evolutionsgleichung $\dot{\kappa} = \beta(\sigma, \kappa, \dot{\varepsilon})$ für die Verfestigungsparameter,
- (iv) die geometrischen Verzerrungs-Verschiebungsbeziehungen $\dot{\varepsilon} = \varepsilon(\dot{u}) \equiv [\varepsilon_{ij}(\dot{u})]_{ij=\overline{1,3}}, \ \varepsilon_{ij}(\dot{u}) := \frac{1}{2} \left(\frac{\partial \dot{u}_i}{\partial x_j} + \frac{\partial \dot{u}_j}{\partial x_i} \right) = 0.5 \ (\dot{u}_{i,j} + \dot{u}_{j,i}),$
- (v) die AB: $u(x,0)\stackrel{\text{z.B.}}{=} 0$, $\sigma(x,0)=0$, $\kappa(x,0)=0$, $x\in\Omega$,

erfüllt werden.

mit
$$V_0 = \{v \in V = [H^1(\Omega)]^3 : v = \mathbf{O} \text{ auf } \Gamma_1\}, (u, v)_V = \int_{\Omega} \varepsilon^T(u) D \varepsilon(v) dx,$$

 $S = \{\sigma \in [L_2(\Omega)]^9 : \sigma_{ij} = \sigma_{ji}\}, (\sigma, \tau)_S = \int_{\Omega} \sigma^T D^{-1} \tau dx,$
 $H = [L_2(\Omega)]^l, (\kappa, \lambda)_H = \int_{\Omega} \kappa \lambda dx, \|\cdot\|_X = (\cdot, \cdot)_X^{0.5}, D - \text{Matrix der elast. Konst.}$

FE-Diskretisierung von (ii) bzw. (iii) ⇒ AWA für System gew. Dgl.!

3.2 Formulierungen und analytische Resultate

- Im weiteren wird von folgender klassischen Formulierung einer AWA für ein System gew. Dgl. 1. Ordnung ausgegangen:
 - (1) Sei $I = \bar{T} = [0, T] \subset \mathbb{R} \equiv \mathbb{R}^1$ kompaktes Intervall, $0 < T < \infty$; $u_0 \in \mathbb{R}^N$ geg. Vektor der Startwerte, $f: D \subset I \times \mathbb{R}^N \longmapsto \mathbb{R}^N$ geg. stetige Fkt.

$$f: D \subset I \times \mathbb{R}^{n} \longmapsto \mathbb{R}^{n} - \text{geg. stetige FKI.}$$

$$(=)$$

Ges. stetig differenzierbare Fkt. $u: I \longrightarrow \mathbb{R}^N$, d.h. $u \in C^1(I)$:

$$u'(t) = f(t, u(t)), \quad t \in I = [0, T],$$

AB: $u(0) = u_0,$

wobei N als fixiert (?) angenommen wird!

■ Bez.:

- $u \in C^1(I) \cong C^1(I, \mathbb{R}^N)$ bedeutet für Vektorfkt. $u = (u_1(\cdot), u_2(\cdot), ..., u_N(\cdot))^T \in [C^1(I)]^N$, d.h. $u_i(\cdot) \in C^1(I) \ \forall i = \overline{1, N}$.
- Analog: $f \in C(D)$ bedeutet $f = (f_1, \ldots, f_N)^T \in [C(D)]^N$.
- Bemerkung 3.8: zu verallgemeinerten Formulierungen!
 - 1. In den Bsp. 3.2 und 3.3 können die Daten, d.h. f, oft nicht mehr als stetig vorausgesetzt werden. Folglich kann <u>nicht</u> $u \in C^1(I)$ erwartet werden. Typischerweise gilt: <u>Bsp. 3.2</u>: $u \in [W_2^1(T)]^N$ <u>Bsp. 3.3</u>: $u \in [W_p^1(T)]^N$. Analoge Aussagen gelten für <u>Bsp. 3.6</u>.
 - 2. AWA für Operatordgl. (vgl. Bsp. 3.7): Ges. $u \in C^1(I,Y) : \dot{u}(t) = F(t,u(t))$ in $Y, \forall t \in I$ mit AB: $u(0) = u_0$ in Y, wobei Y ein Banach-Raum ist: \Rightarrow Galerkin-Diskr. von $Y \longrightarrow (1)$ (in allen Gauß-Pkt.).

■ Spezialfälle:

1. Die Einschränkung auf Systeme gew. Dgl. 1. Ordnung ist <u>nicht</u> wesentlich, denn jedes System gew. Dgl. höherer Ordnung läßt sich auf ein System 1. Ordnung transformieren:

z.B. für allgem. Dgl. (bzw. System von Dgl.) 2. Ordnung

$$\begin{cases} u''(t) = f(t, u(t), u'(t)), \\ \text{mit AB: } u(0) = u_0, u'(0) = u_1 \end{cases}$$

erhält man mit der Setzung

$$u_1(t) = u(t), \ u_2(t) = u'(t)$$

ein äquivalentes System 1. Ordnung

$$\begin{cases} u'_1(t) = u_2(t), \\ u'_2(t) = f(t, u_1(t), u_2(t)), \\ \text{mit AB: } u_1(0) = u_0, \ u_2(t) = u_1. \end{cases}$$

2. Das System (1) nennt man (affine) linear, falls

$$(1)_{\text{lin.}} \quad f(t, u(t)) = A(t) \ u(t) + f(t)$$

$$A(\cdot) : I \longmapsto_{\text{stetig}} \mathbb{R}^N \times \mathbb{R}^N - N \times N - \text{Matrix.}$$

Bsp. 3.2 (nach M_h^{-1} -Multiplikation) ist ein affine lineares System gew. Dgl. 1. Ordnung.

- 3. Hängt die rechte Seite der Dgl. nicht von u ab, d.h.
- (1)_{integr.} $u'(t) = f(t), t \in I \text{ mit AB: } u(0) = u_0,$

so läßt sich die exakte Lsg. mit Hilfe des Integrals

$$u(t) = u_0 + \int_0^t f(s) \, ds$$

darstellen. Die numerische Lsg. des AWP führt dann auf die Fragestellung der numerischen Integration !

- 4. Hängt die rechte Seite der Dgl. nicht explizit von t ab, d. h.
- (1)_{autonom} $u'(t) = f(u(t)), t \in I \text{ mit AB: } u(0) = u_0,$

dann nennt man die Dgl. (System v. Dgl.) autonom.

Die Dgl. (1) u'(t) = f(t, u(t)) kann durch die Setzung

$$u_1(t) = t, \ u_2(t) = u(t)$$

immer in das äquivalente autonome System

$$\begin{cases} u_1'(t) &= 1 \\ u_2'(t) &= f(u_1(t), u_2(t)) \\ + AB: & u_1(0) = 0 \\ & u_2(0) = u(0) \equiv u_0 \end{cases}$$

transformiert werden.

lacksquare lacksquare $\ddot{ ext{U}}$ 3.5 lacksquare Transformieren Sie

Bsp.
$$3.1 - 3.3 \longrightarrow (1)_{autonom}$$
,

Bsp.
$$3.4 - 3.6 \longrightarrow (1) \longrightarrow (1)_{\text{autonom}}$$
.

■ Formulierung als Operatorgleichung:

Sei $X = C^{1}(I, \mathbb{R}^{N}) \cong [C^{1}(I)]^{N} \quad \text{mit geeignet definierter Norm } \|\cdot\|_{X},$ $\text{z.B.} \quad \|u\|_{X} := \max_{t \in I} |u(t)| + \max_{t \in I} |u'(t)|,$ $\text{mit} \quad |v| := \max_{i = 1, N} |v_{i}| \text{ oder}$ $|v|^{2} := \sum_{i=1}^{N} |v_{i}|^{2},$

 $Y = C(I, \mathbb{R}^N) \times \mathbb{R}^N$ mit geeignet definierter Norm $\| \cdot \|_Y$,

$$F: X \mapsto Y$$
 : $F(u) := \begin{bmatrix} u'(\cdot) - f(\cdot, u(\cdot)) \\ u(0) - u_0 \end{bmatrix}$.

Das AWP (1) wird dann äquivalent zur Operatorgleichung

(2) Ges.
$$u \in X : F(u) = 0 \text{ in } Y$$

und es können die Techniken aus [16] Numerik I (Kap. 6) angewendet werden.

• Eine andere Möglichkeit zur Formulierung des AWP (1) als <u>Operatorgleichung</u> in einem <u>B-Raum</u> basiert auf der Integralbeziehung für Lsg. von (1):

(3)
$$u(t_2) = u(t_1) + \int_{t_1}^{t_2} u'(t) dt = u(t_1) + \int_{t_1}^{t_2} f(t, u(t)) dt$$

für beliebige $t_1, t_2 \in I = \bar{T} = [0, T].$

Mit der speziellen Wahl $t_1 = 0$ und $t_2 = t$ erhält man aus AB und (3):

(4)
$$u(t) = u_0 + \int_0^t f(s, u(s)) ds, \quad \forall t \in I.$$

Also ist die Lösung von (1) offenbar äquivalent zur Lösung der <u>Fixpunktgleichung</u> (= Volterrasche Integralgleichung)

(4) Ges.
$$u \in X : u = G(u)$$
 in X ,

$$\mathrm{mit}\ G: X:=C(I,I\!\!R^N) \longrightarrow X:$$

(5)
$$G(u) \equiv G(u)(t) := u_0 + \int_0^t f(s, u(s)) ds.$$

Existenz— und Eindeutigkeitsaussagen liefern eine Verallgemeinerung des Satzes von Picard–Lindelöf auf der Basis des Banachschen Fixpunktsatzes I.6.1 [16] (\rightarrow Satz 3.1: \exists !) und des Satzes von Peano auf der Basis des Fixpunktsatzes von Schauder (\rightarrow \exists : siehe [22], S. 26 – 27, 35 – 36).

71

Satz 3.1: (Verallgemeinerter Satz v. Picard-Lindelöf)

 $\begin{array}{ll} \underline{\text{Vor.:}} & 1) & f: \tilde{D} \longrightarrow Y := I\!\!R^N \ (B\text{-Raum}) - \text{stetig}, \\ & \text{wobei } \tilde{D} := \{(t,v) \in I\!\!R \times Y : |t| \le a \equiv T, \|v-u_0\|_Y \le b\}, \\ & a,b \in (0,\infty), \end{array}$

- 2) $||f(t,v) f(t,\overline{v})||_Y \le L ||v \overline{v}||_Y$, $||f(t,v)||_Y \le K$ auf \tilde{D} mit festen Konstanten $L, K \in [0,\infty]$,
- 3) c = const. > 0: 0 < c < a, cK < b, cL < 1.

Bh.: 1) Dann besitzt

$$\begin{split} &(\tilde{1}) \ u'(t) = f(t,u(t)), \ t \in [-c,c], \ \underline{AB} \colon u(0) = u_0 \\ &\text{bei festem } u_0 \in Y \text{ genau eine Lösung } u(\cdot) \text{ in der Kugel} \\ &M := \{v \in X := C([-c,+c],Y) : \|v-u_0\|_X \leq b\} \\ &\text{mit } \|v\|_X := \max_{t \in [-c,+c]} \|v(t)\|_Y. \\ &u(\cdot) \text{ ist auf } [-c,+c] \text{ stetig differenzierbar}. \end{split}$$

- 2) Die Fixpunktiteration $u_{n+1}(t) = u_0 + \int\limits_0^t f(s,u_n(s))\,ds,\; u_0(t) := u_0,$ konvergiert auf [-c,+c] gleichmäßig gegen $u(\cdot)$.
- 3) Die Lsg. $u(\cdot)$ hängt in der Norm auf X stetig vom Anfangswert $u_0 \in Y$ ab.
- 4) Ist f stetig auf \tilde{D} für alle b>0 und gilt Vor. 2) für alle b>0 mit einem einheitlichen L, dann besitzt $(\tilde{1})$ genau eine Lsg. auf C([-a,+a],Y), folglich auch (1) auf $C(I,\mathbb{R}^N)$. Vor. 3) entfällt hier!

Beweis:

- 1. Zunächst für N=1, d.h. $Y=I\!\!R^1$ (mms) <u>Hinweis:</u> Bh. 4) $M=X=C[-a,a], \ \|v\|_X:=\max_{|t|< a}|v(t)|e^{-L|t|}.$
- 2. Dann übertragen auf $Y = I\!\!R^N$ mit N>1 und $B ext{-Raum }Y.$
- Literatur: [22] Zeidler E.: Vorlesungen über nichtlineare Funktionalanalysis I, Fixpunktsätze.

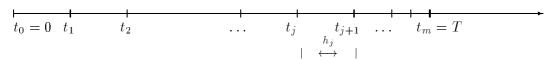
 Teubner-Texte zur Mathematik, Teubner-Verlag, Leipzig, 1976: $Y = I\!\!R^1$: S. 19 21 und Y B-Raum: S. 34 35.

Einschrittverfahren (= zweischichtig: t_j, t_{j+1}) 3.3

Das Eulersche Polygonenzugverfahren (EPZV) 3.3.1

- Eulersches Polygonenzugverfahren:
 - = explizites Eulerverfahren ($\hat{=} \sigma = 0$ in Kap. 2)
 - = Euler vorwärts (vorwärtige Differenzen)
 - = linksseitige Rechteckregel
- Unterteilen Zeitintervall $I = \overline{T} = [0, T]$ i. allg. ungleichmäßig (z.B. adaptiv durch Schrittweitensteuerung) in m Teilintervalle:

— Schrittweitensteuerung (?) \longrightarrow



$$0 = t_0 < t_1 < \ldots < t_j < t_{j+1} < \ldots < t_{m-1} < t_m = T$$

Gitterpkt.

$$I_h := \{t_0, t_1, \dots, t_m\} \equiv \bar{\omega}_h$$
 – (Zeit-)Gitter.

$$\begin{split} I_h := \{t_0, t_1, \dots, t_m\} &\equiv \bar{\omega}_h & - & (\text{Zeit-}) \text{Gitter}, \\ h_j &= t_{j+1} - t_j & - & \text{Schrittweite im Gitterpkt.} \ t_j \ \ (h_j = \tau_j \ (\uparrow)) \end{split}$$

$$h = \vec{h} = (h_0, h_1, \dots, h_{m-1})$$
 – Schrittweitenvektor

$$h=\vec{h}=(h_0,h_1,\ldots,h_{m-1})$$
 – Schrittweitenvektor,
 $h=T/m=h_j \ \ \forall j=\overline{0,m-1}$ – Konstante Schrittweite bei äquidistanter Unterteilung $(h=\tau\ (\uparrow))$

$$h = |\vec{h}| := \max_{j=0, m-1} h_j$$

der Einfachheit halber

Frühere Bezeichnung: $\tau_j = h_j, \ \tau = h, \ \bar{\omega}_{\tau} = I_h, \dots (\uparrow)$

■ Methoden zur Herleitung des Eulerschen Polygonenzugverfahrens:

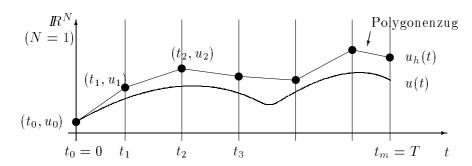
73

1. Taylorentwicklung (\Rightarrow Polygonenzug):

$$u(t) \approx u(t_0) + (t - t_0) u'(t_0) \stackrel{\text{(1)}}{=} u(t_0) + (t - t_0) f(t_0, u(t_0))$$

Lsg. (1)

(6)
$$u_h(t) = u_j + (t - t_j) f(t_j, u_j), \quad t \in [t_j, t_{j+1}]$$
$$u_{j+1} \equiv u_h(t_{j+1}) = u_j + h_j f(t_j, u_j)$$
$$j = 0, 1, \dots, m-1; \text{ AB: } u_0 \text{ geg.}$$



Btr. $u_h: I_h \longrightarrow \mathbb{R}^N$ – als Gitterfkt. $u_h \in X_h = \{v_h: I_h \longrightarrow \mathbb{R}^N\}.$

Mit dieser Schreibweise läßt sich das EPZV (6) als diskretes Ersatzproblem (Näherungsgleichung)

$$(2)_h \qquad \boxed{\text{Ges. } u_h \in X_h : F_h(u_h) = 0 \text{ in } Y_h}$$

für (2) F(u) = 0 interpretieren, wobei

 $Y_h = X_h$ zumindest mengenmäßig, aber $\|\cdot\|_{X_h}, \ \|\cdot\|_{Y_h}$ (?),

$$F_h(v_h)(t_{j+1}) := \begin{cases} D_h v_h(t_j) - f(t_j, v_h(t_j)), & j = 0, 1, \dots, m-1, \\ v_h(t_0) - u_0, & j = -1, \end{cases}$$

mit

$$D_h v_h(t_j) = D_h v_j = v_{t,j} := \frac{v_{j+1} - v_j}{h_j}.$$

2. Vorwärtiges (Explizites) Differenzenverfahren:

(1)
$$u'(t) = f(t, u(t)) + \underline{AB}: u(0) = u_0$$

ersetzen durch
vorwärtigen Differenzquotienten

$$\downarrow$$
(1) $h = \frac{u(t+h) - u(t)}{h} \approx f(t, u(t)), \qquad t = t_j, \quad j = 0, 1, ..., m-1$
(6) $\Rightarrow t = t_j: h \longmapsto h_j, \quad u \longmapsto u_h, \quad u_h(t_j) = u_j, \quad \approx \longmapsto = !$

3. Linksseitige Rechteckregel:

 $\int \mapsto \text{linksseitige Rechteckregel}$

■ Die Sätze von Cauchy und Peano zur Konvergenz des EPZV:

Satz 3.2: (Cauchy)

1) $f: D \mapsto Y = \mathbb{R}^N$ – stetig, wobei $|\cdot| := ||\cdot||_{Y=\mathbb{R}^N}$, Vor.: $D = \{(t, v) : t \in [0, T], |v - u_0| \le b\}, b \in (0, \infty) \text{ fix. } (\Rightarrow \forall b).$ 2) $|f(t,v)-f(t,\bar{v})| \leq L|v-\bar{v}|, f(t,v) \leq K$ auf D mit festen Konstanten $L, K \in [0, \infty)$. 3) $K \cdot T \leq b$ (fällt weg, falls 1) $\forall b > 0$ gilt, d.h. $b = \infty$, $,, \leq " \mapsto ,, < "$).

Bh .: Dann konvergiert $u_h(\cdot)$ aus (6) gleichmäßig gegen eine stetig differenzierbare Fkt. $u(\cdot)$ für $h \to 0$. Die Fkt. $u(\cdot)$ ist eindeutige Lsg. des AWP (1) in D.

Beweis: siehe Literatur [5].

Das EPZV liefert ebenfalls konstruktive Möglichkeit zum Beweis des Bemerkung: Satzes 3.1 von Picard-Lindelöf über $\exists + ! \#$

75

Satz 3.3: (Peano)

 $\underline{\text{Vor.:}} \quad 1) \ f:D:=\{(t,v): t\in [0,T], \ |v-u_0|\leq b\} \to I\!\!R^N \text{ -- stetig, } b\in (0,\infty) \text{ fix.}$

2) $|f(t,v)| \leq K$ auf D mit festen Konstanten $K \in [0,\infty)$.

3) $K \cdot T \leq b$ (fällt weg, falls 1) $\forall b > 0$ gilt, d.h. $b = \infty$, $", \leq " \mapsto ", < ")$.

Bh.: Dann konvergiert $u_h(\cdot)$ aus (6) gleichmäßig gegen eine Lsg. $u(\cdot)$ des AWP (1) für $h \to 0$.

Beweis: siehe Literatur [5].

Bemerkung: Das EPZV liefert ebenfalls konstruktive Möglichkeit zum Beweis der \exists der Lsg. des AWP (1), falls f nicht Lipschitz-stetig ist. #

■ Fehlerbegriffe:

• globaler Diskretisierungsfehler (globaler Fehler):

(7)
$$e_h(t) := u(t) - u_h(t) \lesssim \frac{\text{SL: } t \in I_h = \{t_0, t_1, \dots, t_m\}}{\text{NL: } t \in I = [0, T]}$$

$$\begin{array}{ll} \underline{\text{Satz 3.2 und 3.3:}} & e_h \rightarrow 0 \text{ in } C(I, I\!\!R^N) \text{ für } h \equiv |h| \rightarrow 0, \\ & \text{d.h. } \|e_h\|_{C(I, I\!\!R^N)} := \max_{t \in I} \|e_h\|_{I\!\!R^N} \xrightarrow[h \rightarrow 0]{} 0. \end{array}$$

Btr. im weiteren

$$e_h(\cdot): I_h \longrightarrow I\!\!R^N, \qquad e_h \in X_h$$

als Gitterfunktion, d.h. als Element $X_h := \{v_h : I_h \to \mathbb{R}^N\}.$

Wir nennen das EPZV konvergent (\hat{=} diskrete Konvergenz), falls

$$||e_h||_{X_h} \xrightarrow{h \to 0} 0 \text{ für } h \to 0,$$

wobei $\|\cdot\|_{X_h}$ geeignet definierte Norm in X_h , z.B.

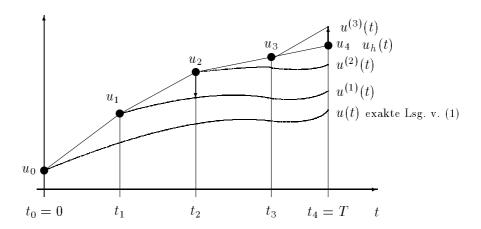
$$||v_h||_{X_h} := \max_{j=0,1,\dots,m} |v_h(t_j)| + \max_{j=0,1,\dots,m-1} |D_h v_h(t_j)|,$$

d.h.
$$X_h = C^1(I_h, \mathbb{R}^N) = X = C^1(I, \mathbb{R}^N).$$

Nach Sätze 3.2 u. 3.3 konvergiert das EPZV in $X_h = C(I_h, I\!\!R^N)$.

• Lokale Diskretisierungsfehler (lokale Fehler):

= Unterschied zwischen exakter Lsg. der Dgl. und Näherungslösung (bzw. SL) nach einem Schritt des Verfahrens bei gleichen AB im Pkt. $t = t_j$, j = 0, 1, ..., m-1:



- o EPZV $t = t_1$: Lokaler Fehler = globaler Fehler = $u(t_1) u_h(t_1)$

$$u^{(0)}(t) - u^{(1)}(t)$$
 mit $u^{(0)}(t) = u(t)$

die Fortpflanzung des lokalen Fehlers $u(t_1)-u_h(t_1)$ nach dem ersten Schritt des EPZV durch die Dgl. Auch bei exakter Integration der Dgl. bleibt ab t_1 dieser Fehler erhalten !

o Dem nächsten Schritt des EPZV läßt sich der lokale Fehler $u^{(1)}(t_2) - u_h(t_2)$

zuordnen, dessen Fortpflanzung analog durch die Lsg.

$$u^{(2)}(t)$$
 des AWP:
$$\begin{cases} \dot{u}^{(2)}(t) = f(t, u^{(2)}(t)), & t \in [t_2, T] \\ \text{mit AB: } u^{(2)}(t_2) = u_2 \end{cases}$$

beschrieben wird, u.s.w. (siehe Abb.).

<u>Fazit:</u> Der <u>globale Fehler</u> setzt sich aus der Fortpflanzung der oben eingeführten <u>lokalen Fehler</u> durch die Dgl. zusammen (siehe Abb.).

Abschätzung des lokalen Fehlers für EPZV:

$$v(t) = u^{(j)}(t) : v'(t) = f(t, v(t)), \ t \in [t_j, T]$$

$$\underline{AB}: \ v(t_j) = u_j$$

 \circ Lokaler Fehler: $h = h_i$

(8)
$$u^{(j)}(t_{j+1}) - u_h(t_{j+1}) = v(t_j + h) - u_h(t_j + h) =$$

$$= v(t_j) + hv'(t_j) + \frac{h^2}{2}v''(t_j) + \dots - u_h(t_j + h)$$

$$= v(t_j) + hf(t_j, v(t_j)) + \frac{h^2}{2} \left(\frac{\partial f}{\partial t} + \frac{\partial f}{\partial v}f\right) (t_j, v(t_j)) + \dots - (u_j + hf(t_j, u_j))$$

$$= \frac{h^2}{2} (f_t + f_v f)(t_j, u_j) + o(h^2) = O(h^2).$$

Man spricht in diesem Fall von einer Methode der Konsistenzordnung 1 (d.h. $\hat{=}$ lokaler Fehler $= O(h^2)$!).

Falls die Fortpflanzung der lokalen Fehler durch die Dgl. "stabil" erfolgt, darf erwartet werden, daß der globale Fehler höchstens von der Größenordnung

$$O(h_0^2) + O(h_1^2) + \ldots + O(h_{m-1}^2) = O(h)$$
 ist!

Um in diesem Sinne zumindest mit Konvergenz rechnen zu dürfen, muß für die <u>lokalen Fehler</u> gelten

$$u^{(j)}(t_i + h) - u_h(t_i + h) = o(h).$$

Man spricht dann von einer konsistenten Methode.

• Lokale Abschneidefehler (= lokaler Approximationsfehler ψ (↑)): engl.: = local truncation error:

(9)
$$\tau_{h}(t_{j+1}) = F_{h}(u)(t_{j+1}) \equiv \begin{cases} \frac{u(t_{j+1}) - u(t_{j})}{h} - f(t_{j}, u(t_{j})) & , j = \overline{0, m-1}, \\ 0 & , j = -1, \end{cases}$$

$$(2)_{h} \text{ Lsg. v. (1)}$$

d.h. der <u>Abschneidefehler</u> mißt, inwieweit die exakte Lsg. $u(\cdot)$ des AWP (1) die Näherungsgl. (2)_h erfüllt.

Bemerkung:
$$\tau_h(t_{j+1}) = \psi_h(u)(t_{j+1}) = F_h(u)(t_{j+1}) - F(u)(t_{j+1})$$
 (\(\frac{1}{2}\)

ullet Zusammenhang: lokaler Abschneidefehler \leftrightarrow lokaler Fehler

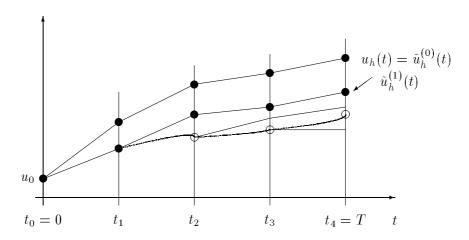
(10)
$$\tau_h(t_{j+1}) = \frac{1}{h} \underbrace{\left\{ u(t_{j+1}) - \left[u(t_j) + hf(t_j, u(t_j)) \right] \right\}}_{\text{elokaler Fehler, wenn man EPZV in } t = t_j$$
 mit exakter Lsg. der Dgl. $u(t_j)$ starten würde !

Wegen (8), $u_j \mapsto u(t_j)$, gilt: $\tau_h = O(h)$.

Im speziellen folgt

(11)
$$\|\tau_h\|_{Y_h} := \max_{j=0,1,\dots,m} |\tau_h(t_j)| = O(h) \xrightarrow{h \to 0} 0.$$

Dies ist eine weitere Möglichkeit, die Konsistenz bzw. die Konsistenzordnung 1 zu definieren:



Der globale Fehler kann auch aus der Fortpflanzung von $u(t_{j+1}) - \tilde{u}_h^{(j)}(t_{j+1})$ durch EPZV bestimmt werden.

■ Ü 3.6 | Man beweise den folgenden Konvergenzsatz!

Satz 3.4:

Vor.: Es gelten die Voraussetzungen des Satzes 3.2 mit $b = \infty$.

Bh.: Dann gilt die folgende Fehlerabschätzung für das EPZV auf äquidistantem Gitter $I_h = \{t_j = jh, j = \overline{0, m}, h = T/m\}$:

$$\begin{split} |u(t_j)-u_h(t_j)| &\leq e^{Lt_j} \left\{ |e_0| + \frac{\tau}{L} \right\}, \ j=1,2,\ldots,m, \\ \text{wobei $L-$Lipschitz-Konstante}, \\ e_0 &= u(t_0) - u_h(t_0) = u_0 - u_0 = 0, \text{ d.h. im betrachteten Fall ist der Anfangsfehler } 0, \ \tau &= \max_{j=1,m} |\tau_j|, \\ \tau_j &= \tau_h(t_j) - \text{lokale Abschneidefehler}. \end{split}$$

Hinweise zum Beweis: (siehe P VII):

- 1) Schreiben Sie unter Benutzung der Darstellung $(\tau_{j+1} \text{ nach } u(t_{j+1}) \text{ auflösen})$
- $(*) \ u(t_{j+1}) = u(t_j) + h f(t_j, u(t_j)) + h \underbrace{\left[\frac{u(t_{j+1}) u(t_j)}{h} f(t_j, u(t_j))\right]}_{=: \tau_h(t_{j+1}) = \tau_{j+1}}$

und des EPZV

- (**) $u_{j+1} = u_j + hf(t_j, u_j)$ eine Rekursionsbeziehung ((*) (**)!) für den Fehler $e_{j+1} = u(t_{j+1}) u_h(t_{j+1}) = u(t_{j+1}) u_{j+1}$ auf, und schätzen Sie diese ab.
 - 2) Verwenden Sie dabei die elementare Beziehung (mms) $(1+hL)^{j+1} \leq e^{(j+1)hL} = e^{Lt_{j+1}} \leq e^{LT}$.

q.e.d.

■ Folgerung 3.5:

2)
$$f$$
, $\frac{\partial f}{\partial t}$, $\frac{\partial f}{\partial v}$ seien auf D beschränkt.

Vor.: 1) Es gelten die Voraussetzungen von Satz 3.4, 2) f, $\frac{\partial f}{\partial t}$, $\frac{\partial f}{\partial v}$ seien auf D beschränkt.

Bh.: Dann gilt die folgende Fehlerabschätzung für das EPZV auf äquidistantem Gitter $I_h = \{t_j = jh\}$:

$$|u(t_j) - u_h(t_j)| \le e^{Lt_j} \frac{c}{L} h.$$

Beweis: folgt sofort aus Satz 3.4 und den Fakten

•
$$e_0 = 0$$

• $(11) \|\tau_h\|_{Y_h} := \max_{j=0,\dots,m} |\tau_h(t_j)| \le c h$
 \uparrow
 $(8), (10), \text{Vor. 2})$

q.e.d.

■ Siehe auch Pkt. 3.3.4: "Allgemeine Konvergenztheorie für Einschrittverfahren"!

Explizite Runge-Kutta-Verfahren (ERKV) 3.3.2

■ Motivation:

- EPZV hat Konsistenzordnung $1 \Rightarrow \|u(\cdot) u_h(\cdot)\|_{C(I_h,\mathbb{R}^N)} = O(h)$!
- Es sind sehr kleine Schrittweiten $\{h_j\}$ notwendig, um genaue Resultate zu erzielen!
- Deshalb: Ges. sind Konstruktionsprinzipien für Verfahren höherer Konsistenzordnung!

3.3.2.1Konstruktionsprinzip

■ Ausgangspunkt: = dabei jene Interpretation des EPZV, die auf der linksseitigen Rechteckregel für

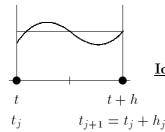
$$\int_{t}^{t+h} f(s, u(s)) ds \approx h f(t, u(t))$$

 $\int\limits_t^{t+h} f(s,u(s))\,ds \approx h\,f(t,u(t))$ beruht, d.h. $\int\limits_t^{t+h} f(s,u(s))\,ds\,\mathrm{durch}\,\,\underline{\mathrm{bessere}}\,\mathrm{Quadratformel}\,(\mathrm{QF})$

■ Eine genauere QF ^{z.B.} Mittelpunktsregel (Gauß 1):

$$\int_{t}^{t+h} f(s, u(s)) ds \approx h f\left(t + \frac{h}{2}, u\left(t + \frac{h}{2}\right)\right)$$

zunächst unbekannt!



<u>Idee:</u> Approximieren u(t + 0.5 h) durch EPZV:

$$u\left(t + \frac{h}{2}\right) \approx u(t) + \frac{h}{2}f(t, u(t))$$

Als Resultat erhält man das folgende Verfahren (Runge, 1895):

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

$$K_2 = f\left(t + \frac{h}{2}, u + \frac{h}{2}K_1\right)$$

$$u_h(t+h) = u + h K_2$$

bzw. vollständig aufgeschrieben:

AB:
$$u_h(t_0) \equiv u_h(0) = u_0$$

 $j = 0, 1, \dots, m - 1$
 $K_{1,j} = f(t_j, u_j)$
 $K_{2,j} = f(t_j + 0.5 h_j, u_j + 0.5 h_j K_{1,j})$
 $u_{j+1} \equiv u_h(t_{j+1}) = u_j + h_j K_{2,j}$

Dieses Verfahren wird

verbesserte Euler-Methode bzw. Euler-Cauchy-Methode

genannt.

Durch <u>Taylorentwicklung</u> des <u>lokalen Fehlers</u> läßt sich nun leicht die <u>Ordnung</u> des Verfahrens bestimmen:

•
$$u_h(t+h) = u(t) + h f \left(t + \frac{h}{2}, u(t) + \frac{h}{2} f(t, u(t))\right)$$

$$= u + h \left[f(t, u) + f_t \frac{h}{2} + f_u \cdot \frac{h}{2} f + \frac{1}{2!} \left(f_{tt} \left(\frac{h}{2}\right)^2 + 2 f_{tu} \cdot \frac{h}{2} \cdot \frac{h}{2} f + f_{uu} \left(\frac{h}{2} f\right)^2\right) + \dots\right]$$

$$= u + h f(t, u) + \frac{h^2}{2} (f_t + f_u f)(t, u) + \frac{h^3}{8} (f_{tt} + 2 f_{tu} f + f_{uu} f^2)(u, t) + \dots$$

$$\underline{\text{Bez.:}} \ f_t = \frac{\partial f}{\partial t}, \ f_u = \frac{\partial f}{\partial u} \ \text{usw.}$$

•
$$u(t+h) = u(t) + h \cdot u'(t) + \frac{h^2}{2} u''(t) + \frac{h^3}{6} u'''(t) + \dots$$

$$Dgl.u'(t) = f(t,u(t))$$

$$\downarrow$$

$$(12) = u(t) + h f(t,u) + \frac{h^2}{2} (f_t + f_u f)(t,u) + \dots$$

$$+ \frac{h^3}{6} (f_{tt} + f_{tu} f + f_{ut} f + f_{ut} f^2 + f_u f_t + f_u f_u f)(t,u) + \dots$$

•
$$u(t+h) - u_h(t+h) =$$

$$= \frac{h^3}{24} [4f_{tt} + 8f_{tu}f + 4f_{uu}f^2 + 4f_uf_t + 4f_u^2f - 3f_{tt} - 6f_{tu}f - 3f_{uu}f^2](u,t) + \dots$$

$$= \frac{h^3}{24} [f_{tt} + 2f_{tu}f + f_{uu}f^2 + 4(f_uf_t + f_u^2f)](u,t) + \dots$$

$$= O(h^3)$$

d.h. Verfahren hat die Ordnung 2!

■ Verallgemeinerung dieses Konstruktionsprinzips:

•
$$\int_{t}^{t+h} f(s, u(s)) ds \approx h \sum_{i=1}^{l} b_{i} f(t + c_{i}h, u(t + c_{i}h))$$

$$\uparrow \qquad \qquad \uparrow$$

$$l-\text{stufige QF} \qquad \text{unbekannt !}$$
wobei $\{b_{i}\}, \{c_{i}\} \uparrow - \text{noch frei w\"{a}hlbar mit } c_{1} = 0.$

• Anstelle der unbekannten Fkt-Werte $u(t+c_ih)$ werden Näherungen

$$g_i \approx u(t + c_i h)$$

rekursiv durch (i-1)-stufige QF berechnet:

$$g_{1} = u,$$

$$g_{2} = u + h a_{21} f(t, g_{1}),$$

$$g_{3} = u + h [a_{31} f(t, g_{1}) + a_{32} f(t + c_{2}h, g_{2})],$$

$$\vdots$$

$$g_{l} = u + h [a_{l1} f(t, g_{1}) + a_{l2} f(t + c_{2}h, g_{2}) + \dots + a_{l,l-1} f(t + c_{l-1}h, g_{l-1})].$$

• Setzt man

$$K_i = f(t + c_i h, g_i),$$

so erhält man die Darstellung

$$K_{1} = f(t, u),$$

$$K_{2} = f(t + c_{2}h, u + ha_{21}K_{1}),$$

$$K_{3} = f(t + c_{3}h, u + h[a_{31}K_{1} + a_{32}K_{2}]),$$

$$\vdots$$

$$K_{l} = f(t + c_{l}h, u + h[a_{l1}K_{1} + a_{l2}K_{2} + \ldots + a_{l,l-1}K_{l-1}]).$$

• Die nächste Näherung hat dann die Form

$$u_h(t+h) = u(t) + h(b_1K_1 + b_2K_2 + \dots + b_lK_l).$$

= $u(t) + h \sum_{i=1}^{l} b_i f(t+c_ih, g_i).$

Man nennt diese Methode ein

$l{\rm -stufiges}$ explizites Runge-Kutta-Verfahren/Formel

und ordnet dieser Methode das folgende Tableau zu, durch das die Methode eindeutig beschrieben wird:

■ Beispiele:

(a) EPZV: $u_h(t+h) = u + hb_1K_1$, $K_1 = f(t, u)$ d.h. EPZV = 1-stufiges RK-Verfahren mit dem Tableau

(b) Verbesserte Euler-Methode:

$$\begin{array}{ll} u_h(t+h) = & u + h \big(b_1 K_1 + b_2 K_2 \big) \\ & K_1 = f(t,u), \ K_2 = f \big(t + \frac{1}{2} h, u + \frac{1}{2} h \ K_1 \big), \ \ b_1 = 0, \ \ b_2 = 1, \end{array}$$

d.h. VEM = 2-stufiges RK-Verfahren mit dem Tableau

$$\begin{array}{c|cc}
0 & & \\
1/2 & 1/2 & \\
\hline
& 0 & 1
\end{array}$$

■ Koeffizienten werden aus gewünschter Konsistenzordnung bestimmt!

83

3.3.2.2 Konsistenzordnung der expliziten Runge-Kutta-Formeln

■ Definition 3.6:

Eine Runge-Kutta-Formel hat die Konsistenzordnung p, falls $u(t+h) - u_h(t+h) = O(h^{p+1}),$ wobei $u(t+h): u'(s) = f(s, u(s)), s \in [t, t+h]$ u(t) geg. $(!u_h(t)$ geg. anstelle von u(t)) $\updownarrow \leftarrow \text{ Start mit gleichen Werten zum Zeitpunkt } t.$ $u_h(t+h) = u(t) + h \sum_{i=1}^l b_i K_i \text{ mit } K_i = f(t+c_i h, g_i) \; (\uparrow)$

Beispiele:

- (a) EPZV = ein 1-stufiges ERKV der Ordnung 1.
- (b) Die verbesserte Euler-Methode = ein 2-stufiges ERKV der Ordnung 2.
- Ü 3.7 Verwendet man anstelle der Mittelpunktsregel die Trapezregel zur Berechnung des Integrals $\int\limits_{\cdot}^{t+h}f(s,u(s))\,ds$, so erhält man das <u>Verfahren von Heun:</u>

$$\smallint_{t}^{t+h} f(s,u(s)) \; ds \overset{TR}{\approx} \tfrac{h}{2} [f(t,u(t)) + f(t+h,u(t+h))]$$

$$K_{1} = f(t, u(t))$$

$$K_{2} = f(t + h, u + h K_{1})$$

$$u_{h}(t + h) = u + \frac{h}{2}K_{1} + \frac{h}{2}K_{2}$$

$$EPZV$$

$$u(t + h) \approx u(t) + hf(t, u(t))$$

$$\begin{array}{c}
\uparrow \\
\text{EPZV} \\
\hline
 u(t+h) \approx u(t) + h f(t, u(t))
\end{array}$$

Damit ergibt sich für das Verfahren von Heun das folgende Tableau:

$$\begin{array}{c|cccc}
0 & & & \\
1 & 1 & & \\
\hline
& 1/2 & 1/2 & \\
\end{array}$$

d.h. das Verfahren von Heun ist eine 2-stufige Runge-Kutta-Formel, und es gilt:

$$u(t+h)=u(t)+\tfrac{h}{2}[f(t,u)+f(t+h,u+hf(t,u))].$$

Man zeige, daß dieses Verfahren die Konsistenzordnung 2 hat.

■ Frage: Existiert eine 2-stufige Runge-Kutta-Formel

$$\begin{array}{c|cccc}
0 & & & \\
c_2 & a_{21} & & \\
\hline
& b_1 & b_2 & \\
\end{array}$$

der Ordnung 3 oder höher?

Zur Klärung dieser Frage entwickeln wir den lokalen Fehler $u(t+h) - u_h(t+h)$ wieder in eine Taylor-Reihe:

•
$$u(t+h) \stackrel{(12)}{\stackrel{\downarrow}{=}} u(t) + hu'(t) + \frac{h^2}{2}u''(t) + \frac{h^3}{6}u'''(t) + o(h^3) =$$

$$\stackrel{(12)}{\stackrel{\downarrow}{=}} u(t) + hf(t,u) + \frac{h^2}{2}(f_t + f_u f)(t,u) +$$

$$+ \frac{h^3}{6}(f_{tt} + 2f_{tu}f + f_{uu}f^2 + f_u f_t + f_u^2 f)(t,u) + o(h^3)$$

$$\begin{aligned} \bullet \ \, u_h(t+h) &= u + h[b_1f(t,u) + b_2f(t+c_2h,u+ha_{21}f(t,u))] = \\ &= u + h\,b_1f + hb_2[f + f_tc_2h + f_uha_{12}f + \\ &\quad + \frac{1}{2!}(f_{tt}(c_2h)^2 + 2f_{tu}(c_2h)(ha_{21}f) + f_{uu}(ha_{12}f)^2) + o(h^2)] \\ &= u + h(b_1 + b_2)f + h^2b_2(c_2f_t + a_{12}f_uf) + \\ &\quad + h^3b_2\left(\frac{c_2^2}{2}f_{tt} + f_{tu}c_2a_{21}f + \frac{a_{12}^2}{2}f_{uu}f^2\right) + o(h^3) \end{aligned}$$

• Koeffizientenvergleich: ⇒ (Notwendige) Bedingungen für Ordnung 2:

$$\begin{bmatrix} b_1 + b_2 = 1 \\ b_2 c_2 = 1/2 \\ b_2 a_{12} = 1/2 \end{bmatrix} \xrightarrow{\text{allgem. Lsg.}} \begin{bmatrix} b_1 = 1 - b_2 \\ c_2 = 1/2 b_2 \\ a_{12} = 1/2 b_2 \end{bmatrix}$$

mit beliebiger Wahl von b_2 !

Damit ergibt sich der <u>lokale Fehler</u> zu:

$$u(t+h) - u_h(t+h) = \frac{h^3}{6} \left[\left(1 - \frac{3}{4b_2} \right) (f_{tt} + 2f_{tu}f + f_{uu}f^2) + \underline{f_u f_t + f_u^2 f} \right] (t, u) + \dots,$$

woraus sofort erkennbar ist, daß die Ordnung 3 i. a. nicht erreichbar ist.

Für

$$b_2 = 3/4$$

ist die Summe der Beträge der einzelnen Beiträge zum lokalen Fehler minimal. Man spricht von einer (in diesem Sinne) optimalen Formel. Falls $f_u=0$ (\Rightarrow reine Integration), dann folgt für $b_2=3/4$ Ordnung 3!

■ Runge-Kutta-Verfahren der Ordnung 3:

• RK-Verfahren der Ordnung 3 muß mindestens 3-stufig sein (†):

$$\begin{array}{c|cccc}
0 & & & & & \\
c_2 & a_{21} & & & & \\
c_3 & a_{31} & a_{32} & & & & \\
\hline
& b_1 & b_2 & b_3 & & & \\
\end{array}$$

• Formales Herleiten von Bedingungsgleichungen für die Koeffizienten $\{c_i\}$, $\{a_{ij}\}$, $\{b_j\}$ ohne weiteres durch <u>Taylor-Entwicklung</u> und <u>Koeffizientenvergleich</u> möglich:

$$u(t+h) \stackrel{(12)}{\stackrel{\downarrow}{=}} u + hf + \frac{h^2}{2}(f_t + f_u f) + \frac{h^3}{6}(f_{tt} + 2f_{tu} f + f_{uu} f^2 + f_u f_t + f_u^2 f_t) + \dots \text{ usw.}$$

$$\downarrow u_h(t+h) = u + h(b_1 f(t,u) + b_2 f(t+c_2 h, u+ha_{21} f(t,u)) + b_3 f(t+c_3 h, u+\dots))$$

$$= u + h(b_1 + b_2 + b_3) f + \frac{h^2}{2}(\dots) + \frac{h^3}{6}(\dots) + \dots$$

→ Vorgehen ist technisch zu kompliziert!

■ Systematisches Vorgehen zur Ermittlung der Bedingungsgleichungen für die Ordnung von RK-Formeln:

• Die folgende Ü 3.8 besagt, daß man sich bei der Analyse der Ordnung von RK-Formeln auf autonome Dgl.

$$(14) u'(t) = f(u(t))$$

beschränken kann.

Ü 3.8 Falls die Bedingungen

(15)
$$c_i = \sum_{j=1}^{i-1} a_{ij}, \quad i = \overline{2, l}$$

erfüllt sind, erhält man für den Fall allgem. Dgl. der Form u'(t)=f(t,u(t)) die gleiche Ordnung wie für den Fall autonomer Dgl. (14).

- Für die weitere Diskussion können wir annehmen, daß (15) gilt und daß nur autonome Dgl. (14) betrachtet werden !
 - Btr. dazu RK-Formeln in der Form (für autonome Dgl. (14))

(16)
$$g_{i} \equiv g_{i}(h) = u + \sum_{j=1}^{i-1} a_{ij} h f(g_{j}(h)), \quad i = \overline{1,l} \quad \left(\sum_{j=1}^{0} = 0\right)$$

Taylor-Entwicklung an der Stelle $\underline{h=0}$

$$u_{h}(t+h) = u + \sum_{j=1}^{l} b_{j} \underbrace{h f(g_{j}(h))}_{h \varphi(h)} = \dots$$

$$h \varphi(h)$$

• Lemma 3.7: (Leibnitz)

Falls
$$\varphi \in C^q(I\!\! R)$$
, dann gilt: $[h\varphi(h)]^{(q)}(0) = q \, \varphi^{(q-1)}(0)$.

Beweis: durch Induktion nach q (für q = 1 o.k.).

$$[h\varphi(h)]^{(q+1)}(0) = [(h\varphi(h))']^{(q)}(0) = [\varphi(h) + h\varphi'(h)]^{(q)}(0) =$$

$$= \varphi^{(q)}(0) + (h\varphi'(h))^{(q)}(0) = \varphi^{(q)}(0) + q\varphi'^{(q-1)}(0) = (1+q)\varphi^{(q)}(0)$$

q.e.d.

• Mit Lemma 3.7 folgt aus (16) für
$$g_i^{(q)}(0) := \frac{d^q g_i}{dh^q}\Big|_{h=0}$$
:
$$\underline{q=0:} \quad g_i(0) = u \equiv u(t),$$

$$\underline{q=1:} \quad g_i'(0) = \sum_{j=1}^{i-1} a_{ij} \cdot 1 \cdot f(g_j(0)) = \sum_{j=1}^{i-1} a_{ij} f(u),$$

$$\underline{q=2:} \quad g_i''(0) = 2 \sum_{j=1}^{i-1} a_{ij} (f \circ g_j)'(0) =$$

$$\begin{bmatrix} (f \circ g_j)'(h) = f'(g_j(h))g_j'(h) \\ \vdots \\ i^{-1} \\ i^{-1} \\ i^{-1} \end{bmatrix} = 2 \sum_{j=1}^{i-1} a_{ij} f'(u) g_j'(0) =$$

$$= 2 \sum_{j=1}^{i-1} a_{ij} f'(u) \sum_{k=1}^{j-1} a_{jk} f(u) =$$

$$= 2 \sum_{j=1}^{i-1} \sum_{k=1}^{j-1} a_{ij} a_{jk} f'(u) f(u)$$

$$\underline{q=3:} \quad g_i'''(0) = 3 \sum_{j=1}^{i-1} a_{ij} (f \circ g_j)''(0) =$$

$$\begin{bmatrix} f(g_j(h))]'' = [f'(g_j(h))g_j'(h)]' = \\ & = f''(g_j(h))g_j'(h)g_j'(h) + f'(g_j(h))g_j''(h) \end{bmatrix} =$$

$$= 3 \sum_{j=1}^{i-1} \sum_{k=1}^{j-1} \sum_{n=1}^{j-1} a_{ij} a_{jk} a_{jn} f''(u) f(u) f(u) +$$

$$+ 6 \sum_{j=1}^{i-1} \sum_{j=1}^{j-1} \sum_{k=1}^{k-1} a_{ij} a_{jk} a_{kn} \cdot f'(u) \cdot f'(u) \cdot f(u)$$

usw.

• $u_h(t+h)$ unterscheidet sich von $g_i(h)$ nur durch eine andere Gewichtung:

$$\Rightarrow \frac{d^q u_h(t+h)}{dh^q}\bigg|_{h=0} = u_h^{(q)}(t)$$

 \Rightarrow in obigen Formeln a_{ij} durch b_j ersetzen:

$$u'_{h}(t) = \sum_{j=1}^{l} b_{j} f(u)$$

$$u''_{h}(t) = 2 \sum_{j=1}^{l} \sum_{k=1}^{j-1} b_{j} a_{jk} f'(u) f(u)$$

$$u'''_{h}(t) = 3 \sum_{j=1}^{l} \sum_{k=1}^{j-1} \sum_{n=1}^{j-1} b_{j} a_{jk} a_{jn} f''(u) f(u) f(u) + 6 \sum_{j=1}^{l} \sum_{k=1}^{j-1} \sum_{n=1}^{k-1} b_{j} a_{jk} a_{kn} f'(u) f'(u) f(u)$$
usw.

• Aus der Dgl. u'(t) = f(u(t)) folgen die entsprehenden Abl. für u:

$$u'(t) = f(u(t))$$

$$u''(t) = f'(u(t))u'(t) = f'(u(t))f(u(t))$$

$$u'''(t) = f''(u(t))u'(t)f(u(t)) + f'(u(t))f'(u(t))u'(t) =$$

$$= f''(u(t))f(u(t))f(u(t)) + f'(u(t))f'(u(t))f(u(t))$$
usw.

• Daraus erhält man die Taylor-Entwicklung des lokalen Fehlers:

$$u(t+h) - u_h(t+h) = u(t) + hu'(t) + \frac{h^2}{2}u''(t) + \frac{h^3}{6}u'''(t) + \dots$$

$$-(u''_h(t) + hu'_h(t) + \frac{h^2}{2}u''_h(t) + \frac{h^3}{6}u'''(t) + \dots)$$
(18)

$$\begin{array}{l} \stackrel{(17)}{=} & h (1 - \sum\limits_{j=1}^{l} b_j) f(u) + \\ & + \frac{h^2}{2} (1 - 2 \sum\limits_{j=1}^{l} \sum\limits_{k=1}^{j-1} b_j a_{jk}) f'(u) f(u) + \\ & + \frac{h^3}{6} (1 - 3 \sum\limits_{j=1}^{l} \sum\limits_{k=1}^{j-1} \sum\limits_{n=1}^{j-1} b_j a_{jk} a_{jn}) f''(u) f(u) f(u) + \\ & + \frac{h^3}{6} (1 - 6 \sum\limits_{j=1}^{l} \sum\limits_{k=1}^{j-1} \sum\limits_{n=1}^{k-1} b_j a_{jk} a_{kn}) f'(u) f'(u) f(u) + \\ & + \dots \text{ usw.} \end{array}$$

• Damit ergeben sich die folgenden 4 Bedingungen für RK-Verfahren der Ordnung 3:

(19)
$$\begin{bmatrix}
\sum_{j=1}^{l} b_{j} & = 1 & h \\
2 \sum_{j=1}^{l} \sum_{k=1}^{j-1} b_{j} a_{jk} & = 1 & \frac{h^{2}}{2!} \\
3 \sum_{j=1}^{l} \sum_{k=1}^{j-1} \sum_{n=1}^{j-1} b_{j} a_{jk} a_{jn} & = 1 & \frac{h^{3}}{3!} \\
6 \sum_{j=1}^{l} \sum_{k=1}^{j-1} \sum_{n=1}^{k-1} b_{j} a_{jk} a_{kn} & = 1 & \frac{Bez.:}{2} \\
usw. für höhere Ordnung, z.B. + 4 Bd. für $\frac{h^{4}}{4!}$ usw. $\frac{Bez.:}{2} = 0$$$

• Ein 3-stufiges RK-Verfahren (l=3) wird durch das Tableau

beschrieben und hat somit <u>6 frei wählbare Parameter</u>, da die Koeffizienten $\{c_i\}$ bereits durch (15) festgelegt sind.

Diese 6 Parameter $\{a_{21},a_{31},a_{32},b_1,b_2,b_3\}$ lassen sich so wählen, daß die Bedingungen (19) für l=3 erfüllt werden, d.h. es gibt 3-stufige RK-Formeln der Ordnung 3!

■ Runge-Kutta-Verfahren der Ordnung 4:

- Um Ordnung 4 zu erreichen, müssen weitere 4 Bedingungen erfüllt werden.
- Es läßt sich zeigen, daß dies mit 3-stufigen Formeln nicht möglich ist!
- 4-stufige RK-Formeln werden durch das Tableau

beschrieben und haben somit 10 freie Parameter, die tatsächlich so gewählt werden können, daß die 8 Bedingungen für Ordnung 4 erfüllt sind (Allgem. Lösung des polynominalen GS: eine 2-parametrige Lösungsschar und drei 1-parametrige Lösungsscharen).

89

- Beispiele 4-stufiger RK-Formeln der Ordnung 4:
 - (1) Das sogenannte "klassische" Runge-Kutta-Verfahren:

Für ein Quadraturproblem u'(t) = f(t) wird aus dem klassischen RK-Verfahren die Simpson-Regel:

$$u_h(t+h) = u(t) + h\left[\frac{1}{6}f(t) + \frac{1}{3}f\left(t + \frac{h}{2}\right) + \frac{1}{3}f\left(t + \frac{h}{2}\right) + \frac{1}{6}f(t+h)\right]$$

(2) Die 3/8-Regel:

Für ein Quadraturproblem u'(t) = f(t) wird aus der 3/8-Regel die sogenannte Newton-3/8-Regel:

$$u_h(t+h) = u(t) + h \left[\frac{1}{8} f(t) + \frac{3}{8} f\left(t + \frac{h}{3}\right) + \frac{3}{8} f\left(t + \frac{2}{3}h\right) + \frac{1}{8} f(t+h) \right]$$

- Mit 4-stufigen RK-Formeln läßt sich die Ordnung 5 nicht erreichen!
- Ergebnisse zur Analyse von <u>s -stufigen RK-Formeln</u> ($s \ge 5$) siehe [1], Punkt 4.2.3, S. 124-130.

3.3.3 Implizite Runge-Kutta-Verfahren

3.3.3.1 Beispiele

■ Implizites Euler-Verfahren ($\hat{\sigma} = 1$ in Kap.2):

$$J = \int\limits_t^{t+h} f(s,u(s)) \, ds \qquad \circlearrowleft$$

$$\lim_{t \to \infty} \frac{\text{linksseitige Rechteckregel: } J \approx hf(t,u(t))$$

$$\Rightarrow \text{ Expliziter Euler: } u_{j+1} = u_j + h_j f(t_j,u_j)$$

$$\underbrace{\text{rechtsseitige Rechteckregel: } J \approx hf(t+h,u(t+h))}_{\Rightarrow \text{ Impliziter Euler: } u_{j+1} = u_j + h_j f(t_j + h_j,u_{j+1})}$$

Resultat: Implizites Euler-Verfahren ($\sigma = 1$)

= Euler rückwärts

= rechtsseitige Rechteckregel

(20)
$$\mathbf{u_{j+1}} = u_j + h_j f(t_j + h_j, \mathbf{u_{j+1}}), j = 0, 1, \dots, m-1; \ u_0 \text{ geg.}$$

Zur Bestimmung von $\mathbf{u_{j+1}}$ muß i. allgem. ein nichtlineares Gleichungssystem gelöst werden, z.B. durch Fixpunktiteration oder durch Newton–Iteration mit der Startnäherung $u_{j+1}^0 = u_j$!

Bemerkung:

- Nachteile impliziter Verfahren: Nichtlineare GS zu lösen (†)
- Vorteile impliziter Verfahren:
 - Bessere Stabilität (↓)
 - Genauigkeit (↓)
- Implizite Mittelpunktsregel:

Resultat: Implizite Mittelpunktsregel

$$\mathbf{g_1} = u + \frac{h}{2} f\left(t + \frac{h}{2}, \mathbf{g_1}\right)$$

$$u_h(t+h) = u + h f\left(t + \frac{h}{2}, g_1\right)$$

$$\mathbf{K_1} \equiv f\left(t + \frac{h}{2}, g_1\right) = f\left(t + \frac{h}{2}, u + \frac{h}{2}\mathbf{K_1}\right)$$

$$u_h(t+h) = u + h K_1$$

$$g - \text{Form}$$

$$K - \text{Form}$$

■ Beide Verfahren sind <u>Beispiele</u> von sogenannten <u>impliziten Runge-Kutta-Verfahren</u> (Formeln).

3.3.3.2 Allgemeines Konstruktionsprinzip

- Die allgemeine Form einer impliziten <u>l</u> -stufigen Runge-Kutta-Formel läßt sich entweder
 - in der g-Form

$$\mathbf{g_1} = u + h[a_{11}f(t + c_1h, \mathbf{g_1}) + \dots + a_{1l}f(t + c_lh, \mathbf{g_l})]$$

$$\mathbf{g_2} = u + h[a_{21}f(t + c_1h, \mathbf{g_1}) + \dots + a_{2l}f(t + c_lh, \mathbf{g_l})]$$

$$\vdots$$

$$\mathbf{g_l} = u + h[a_{l1}f(t + c_1h, \mathbf{g_1}) + \dots + a_{ll}f(t + c_lh, \mathbf{g_l})]$$

$$u_h(t + h) = u + h[b_1f(t + c_1h, g_1) + \dots + b_lf(t + c_lh, g_l)]$$

bzw.

 \bullet in der K-Form

$$\mathbf{K_{1}} = f(t + c_{1}h, u + h[a_{11}\mathbf{K_{1}} + a_{12}\mathbf{K_{2}} + \dots + a_{1l}\mathbf{K_{l}}])$$

$$\mathbf{K_{2}} = f(t + c_{2}h, u + h[a_{21}\mathbf{K_{1}} + a_{22}\mathbf{K_{2}} + \dots + a_{2l}\mathbf{K_{l}}])$$

$$\vdots$$

$$\mathbf{K_{l}} = f(t + c_{l}h, u + h[a_{l1}\mathbf{K_{1}} + a_{l2}\mathbf{K_{2}} + \dots + a_{ll}\mathbf{K_{l}}])$$

$$u_{h}(t + h) = u + h[b_{1}K_{1} + b_{2}K_{2} + \dots + b_{l}K_{l}]$$

schreiben.

■ Das Verfahren wird durch das folgende <u>Tableau</u> eindeutig definiert:

■ Definition 3.8:

Eine durch das Tableau (21) beschriebene <u>RK-Formel</u> heißt

- ullet explizit, falls A eine echte linke untere Dreiecksmatrix ist,
- implizit, falls A nicht explizit ist.

<u>Bem.:</u> Diese Def. einer expliziten RK-Formel ist etwas allgemeiner, da nicht mehr $c_1=0$ a-priori angenommen wird.

■ Beispiele:

1. Implizites Euler–Verfahren = 1-stufige implizite RK–Formel:

2. Implizite Mittelpunktsregel = 1-stufige implizite RK–Formel:

$$\begin{array}{c|c}
1/2 & 1/2 \\
\hline
& 1
\end{array}$$

3. $\underline{\text{Implizite Trapezregel}} = 2\text{-stufige implizite RK-Formel:}$

$$u_{j+1} = u_j + \frac{h_j}{2} [f(t_j, u_j) + f(t_j + h_j, u_{j+1})]$$

 $(\sigma = 1/2 \text{ in Kap. 2: CRANK-NICOLSON-Verfahren})$

$$\begin{array}{c|cc}
0 & 0 & 0 \\
1 & 1/2 & 1/2 \\
\hline
& 1/2 & 1/2
\end{array}$$

3.3.3.3 Durchführbarkeit

- Bei allen impliziten RK-Verfahren stellt sich die Frage nach der Durchführbarkeit, d.h. nach der Lösbarkeit des i. allgem. nichtlinearen Gleichungssystems zur Bestimmung der nächsten Näherung:
 - *g*–Form:

(22)
$$\mathbf{g} = \Phi(\mathbf{g}, t, u, h)$$

• Ges.
$$g = (g_1, \dots, g_l)^T$$
 als Lösung der Fixpunktgleichung
$$(22) \qquad \boxed{\mathbf{g} = \Phi(\mathbf{g}, t, u, h)}$$
mit $\Phi = (\Phi_1, \dots, \Phi_l)^T$, $\Phi_i = u + h \sum_{j=1}^l a_{ij} f(t + c_j h, g_j)$
• $u_h(t+h) = u + h \sum_{j=1}^l b_j f(t + c_j h, g_j)$

- K-Form:
 - Geg. (t, u, h)
 - Ges. $K = (K_1, \dots, K_l)^T$ als Lösung der Fixpunktgleichung

(23)
$$\mathbf{K} = \Psi(\mathbf{K}, t, u, h)$$

(23)
$$\mathbf{K} = \Psi(\mathbf{K}, t, u, h)$$

$$\text{mit } \Psi = (\Psi_1, \dots, \Psi_l)^T, \quad \Psi_i = f(t + c_i h, u + h \sum_{j=1}^l a_{ij} K_j)$$

$$\bullet \quad u_h(t+h) = u + h \sum_{j=1}^l b_j K_j$$

•
$$u_h(t+h) = u + h \sum_{j=1}^l b_j K_j$$

■ Der folgende Satz zeigt unter geeigneten Voraussetzungen die Existenz und Eindeutigkeit der Lösung der Fixpunktgleichungen (22) bzw. (23), die z.B. durch die Banachsche Fixpunktiteration iterativ bestimmt werden kann:

Satz 3.9: (Existenz, Eindeutigkeit, Konstruktion)

- 1) $f: D \longrightarrow \mathbb{R}^N \text{stetig},$
- 2) f sei auf D beschränkt, d.h. $|f(t,v)| \leq M$, $\forall (t,v) \in D$, und f sei auf D im zweiten Argument Lipschitz-stetig, d.h. $|f(t,v)-f(t,u)| \leq L|v-u| \ \forall (t,v), (t,u) \in D$ mit fixierten Konstanten $M,L \in [0,\infty)$,
- 3) $h||A||_{\infty} M < b$, $h||A||_{\infty} L < 1$.

Bh.: 1) Dann besitzt die Fixpunktgleichung

(22)
$$\mathbf{g} = \Phi(\mathbf{g}, t, u, h)$$

für bel. (aber fixierte) (t, u, h) mit $t, t + c_i h \in [0, T]$ $(i = \overline{1, l})$ eine eindeutig bestimmte Lösung

$$g \in \mathcal{K} := \{g \in (\mathbb{R}^N)^l := \mathbb{R}^N \times \ldots \times \mathbb{R}^N : |g_i - u| \le b \ \forall i = \overline{1, l}\}.$$

2) Die Fixpunktiteration

(24)
$$\mathbf{g}^{(n+1)} = \Phi(\mathbf{g}^{(n)}, t, u, h), n = 0, 1, \dots$$

mit dem Startwert $g^{(0)} = (u, u, ..., u)^T \in \mathcal{K}$ konvergiert linear gegen diese Lösung.

Beweis: beruht natürlich auf Banachschem Fixpunktsatz I.6.1 [16]:

$$\bullet \ (I\!\!R^N)^l - B \text{--Raum}, \, \|\cdot\| := \|\cdot\|_\infty := \max_{i=\overline{1,l}} \, |\cdot| \ .$$

- $\mathcal{K} = \bar{\mathcal{K}} \subset (I\!\!R^N)^l$ abgeschlossen, nicht leer.
- $\Phi \mathcal{K} \subset \mathcal{K}$, tatsächlich

$$\begin{split} |\Phi_i(g)-u| &= \left| \left[u + h \sum_{j=1}^l a_{ij} f(t+c_jh,g_j) \right] - u \right| = \\ &= h \sum_{j=1}^l |a_{ij}| \underbrace{|f(t+c_jh,g_j)|}_{\leq M} \leq h \|A\|_{\infty} M \leq b \ \forall g \in \mathcal{K}. \end{split}$$

• Kontraktivität von Φ:

$$\begin{split} \|\Phi(\bar{g}) - \Phi(g)\|_{\infty} &= \max_{i=1,\bar{l}} |u - h \sum_{j=1}^{l} a_{ij} f(t + c_{j}h, \bar{g}_{j}) - [u - h \sum_{j=1}^{l} a_{ij} f(t + c_{j}h, g_{j})]| \\ &= \max_{i=1,\bar{l}} |h \sum_{j=1}^{l} a_{ij} (f(t + c_{j}h, g_{j}) - f(t + c_{j}h, \bar{g}_{j}))| \\ &\leq h \max_{i=1,\bar{l}} \sum_{j=1}^{l} |a_{ij}| |f(t + c_{j}h, g_{j}) - f(t + c_{j}h, \bar{g}_{j})| \\ &\leq h \max_{i=1,\bar{l}} \sum_{j=1}^{l} |a_{ij}| L |\bar{g}_{j} - g_{j}| \\ &\leq h \|A\|_{\infty} L \|\bar{g} - g\|_{\infty} = q \|\bar{g} - g\|_{\infty}. \\ &=: q < 1. \end{split}$$

• Aussagen von Satz 3.9 folgen nun unmittelbar aus Satz I.6.1.

q.e.d.

■ Bemerkungen:

- 1. Aussagen von Satz 3.9. gelten auch für Fixpunktgl. (23), (mms).
- 2. Durchführbarkeit impl. RK-Formeln ist somit für hinreichend klein h garantiert!
- 3. Fixpkt.-Iteration: $g^{(n+1)} = \Psi(g^{(n)}, t, u, h) : \|g g^{(n)}\|_{\infty} \le q^n \|g g^{(0)}\|_{\infty}$ usw. Möglich: Newton-Iteration, falls $f \in C^1$, sowie andere Iterationsverfahren.

3.3.3.4 Konsistenzordnung impliziter Runge-Kutta-Formeln

■ Eine RK-Formel läßt sich auch folgendermaßen aufschreiben:

(25)
$$u_{h}(t+h) = u + h \varphi(t, u, h),$$

$$\varphi(t, u, h) = \sum_{j=1}^{l} b_{j} K_{j}(t, u, h) = \sum_{j=1}^{l} b_{j} f(t + c_{j}h, g_{j})$$

$$K_{j} = K_{j}(t, u, h) : K = \Psi(K, t, u, h) \text{ bzw. } g_{j} = g_{j}(t, u, h) : g = \Phi(g, t, u, h)$$

$$(23) \qquad (22)$$

- Fehlerbegriffe werden analog zu den expliziten RK-Formeln definiert:
 - globaler Diskretisierungsfehler: $e_h(t) = u(t) u_h(t)$,

- lokale Diskretisierungsfehler: $u(t+h) u_h(t+h)$ (vgl. Def. 3.6) Dgl. RK-Formel mit gleichen Startwerten $u(t) = u_h(t)$
 - im Pkt. t, wobei $t \in \{t_0 = 0, t_1, \dots, t_{m-1}\}.$
- lokaler Abschneidefehler (vgl. (9)):

$$\tau_h(t+h) = \frac{u(t+h) - u(t)}{h} - \varphi(t, u, h), \quad t \in \{t_0 = 0, t_1, \dots, t_{m-1}\},$$

$$\tau_h(t_0 \equiv 0) = 0.$$

- Es gilt: $\tau_h(t+h) = \frac{1}{h}[u(t+h) (u(t) + h\varphi(t, u(t), h))].$
- Konsistenzordnung p (vgl. auch Def. 3.6)
 - entspricht der Bedingung $u(t+h) - u_h(t+h) = O(h^{p+1})$
 - an den lokalen Fehler
 - bzw. der Bedingung $\tau_h = O(h^p)$
 - an den lokalen Abschneidefehler.
- Die Konsistenzordnung impliziter RK-Formeln läßt sich analog zum expliziten Fall für autonome Dgl. durch Taylor-Entwicklung des lokalen Fehlers bestimmen (siehe Pkt. 3.3.2.2):
 - 1 Bedingung für Ordnung 1:

$$\sum_{i=1}^{l} b_i = 1$$

1 zusätzliche Bed. für Ordnung 2:

$$2\sum_{j,k}b_ja_{jk}=1$$

•
$$\frac{2}{2}$$
 zusätzliche Bed. für Ordnung 3: $\frac{3}{3} \sum_{j,k,n} b_j a_{jk} a_{jn} = 1$ $\frac{6}{3} \sum_{j,k,n} b_j a_{jk} a_{kn} = 1$

4 zusätzliche Bed. für Ordnung 4:

$$\frac{4 \sum_{j,k,n,m} b_{j} a_{jk} a_{jn} a_{jm} = 1}{8 \sum_{j,k,n,m} b_{j} a_{jk} a_{jn} a_{nm} = 1}$$

$$12 \sum_{j,k,n,m} b_{j} a_{jk} a_{kn} a_{km} = 1$$

$$24 \sum_{j,k,n,m} b_{j} a_{jk} a_{kn} a_{nm} = 1$$

• usw., mit $\sum_{\dots} = \sum_{\dots=1}^{i}$

Die für die allgem. Dgl. u'(t) = f(t, u(t)) notwendigen Koeffizienten $\{c_i\}_{i=\overline{1,l}}$:

$$(15) c_i = \sum_{j=1}^l a_{ij}.$$

■ Unter Verwendung von (15) $c_i = \sum\limits_{j=1}^l a_{ij}$ kann man die obigen Ordnungsbedingungen für den allgemeinen Fall u'(t) = f(t, u(t)) umschreiben: $(\sum\limits_{\dots} = \sum\limits_{\dots=1}^l)$

$$\sum_{j} b_{j} c_{j} = 1$$
 Ordn. 1
$$\sum_{j} b_{j} c_{j} = 1/2$$
 Ordn. 2
$$\sum_{j} b_{j} c_{j}^{2} = 1/3$$
 Ordn. 3
$$\sum_{j} b_{j} a_{jk} c_{k} = 1/6$$
 Ordn. 4
$$\sum_{j} b_{j} c_{j}^{3} = 1/4$$
 Ordn. 4
$$\sum_{j} b_{j} c_{j} a_{jn} c_{n} = 1/8$$
 Ordn. 4
$$\sum_{j,n} b_{j} a_{jk} c_{k}^{2} = 1/12$$
 Ordn. 4

■ Beispiele:

1. Impliziter Euler
$$(l=1)$$
: $\begin{array}{c|c} \hline 1 & 1 \\ \hline & 1 \\ \end{array}$ \Rightarrow $\begin{array}{c|c} \hline \text{Ordnung} = 1 \\ \hline \end{array}$

2. Implizite MP-Regel
$$(l=1)$$
: $1/2$ $1/2$ \Rightarrow Ordnung = 2

- Frage: ∃ 2-stufige, implizite RK-Formeln, die eine höhere Ordnung als 2 besitzen?
 - Für $l=2\,$ haben die 4 Bedingungen für Ordnung 3 eine 2-parametrige Lösungsschar!
 - Die zusätzlichen 4 Bedingungen für Ordnung 4 können durch eine geeignete Wahl der 2 freien Parameter ebenfalls noch erfüllt werden.

Die einzige Lösung hat das folgende Tableau:

| Stützstellen der

Gauß-QF der Stufe 2

für das Intervall [0,1]

- $\Rightarrow \exists$! 2-stufige (l=2) implizite RK–Formel der Ordnung 4 ! \longrightarrow RK–Formel vom Gauß–Typ !
- Resümee: Mit l-stufigen impliziten RK-Formeln sind offenbar höhere Konsistenzordnungen erreichbar als mit den entsprechenden expliziten RK-Formeln!
- Die folgenden beiden Sätze geben Aussagen über die mit *l* -stufigen, impliziten RK-Formeln erreichbare Konsistenzordnung:
 - Satz 3.10:

Vor.: 1. Es seien die Voraussetzungen von Satz 3.2 erfüllt:

- $f: D = \{(t, v): t \in [0, T], |v u| \le b\} \to \mathbb{R}^N$ stetig, $b \in (0, \infty)$ fix,
- $\begin{array}{l} \bullet \; |f(t,v)-f(t,\bar{v})| \leq L|v-\bar{v}|, |f(t,v)| \leq M \; \text{auf} \; D; \\ L,M \in [0,\infty), \end{array}$
- $M \cdot T < b$ (entfällt bei " $b = \infty$ ").

2.
$$r, l \in \mathbb{N} : r \leq l; f \in C^r(D, \mathbb{R}^N); [h||A||_{\infty} L < 1].$$

Bh.: Falls für eine l-stufige RK-Formel die Bedingungen

erfüllt sind, besitzt die Formel die Konsistenzordnung p = r.

<u>Beweis:</u> \Rightarrow direkte Abschätzung des <u>lokalen Fehlers</u> $u(t+h) - u_h(t+h) = ?$ für <u>autonome Dgl.</u> (\Rightarrow (15) $c_i = \sum_{j=1}^{l} a_{ij}$)!:

•
$$u(t+h) = u + \underbrace{\int_{t}^{t+h} f(u(s)) ds}_{\text{Integral}} = u + \underbrace{h \sum_{j=1}^{l} b_{j} f(u(t+c_{j}h))}_{\text{QF}} + \underbrace{Restglied}_{\text{Restglied}}$$

Aus der Taylorentwicklung

$$f(u(s)) = u'(s) = \sum_{k=0}^{r-1} \frac{1}{k!} u^{(k+1)}(t) (s-t)^k + O(h^r), \quad s \in [t, t+h],$$

folgt für das Restglied

$$R_{0} = \int_{t}^{t+h} f(u(s)) ds - h \sum_{j=1}^{l} b_{j} f(u(t+c_{j}h)) =$$

$$= \sum_{k=0}^{r-1} \frac{1}{k!} u^{(k+1)}(t) \int_{t}^{t+h} (s-t)^{k} ds - h \sum_{j=1}^{l} b_{j} \left(\sum_{k=0}^{r-1} \frac{1}{k!} u^{(k+1)}(t) (c_{j}h)^{k} \right) +$$

$$+ O(h^{r+1})$$

$$= \sum_{k=0}^{r-1} \frac{1}{k!} u^{(k+1)}(t) \frac{h^{k+1}}{k+1} - \sum_{k=0}^{r-1} \frac{1}{k!} u^{(k+1)}(t) h^{k+1} \sum_{j=1}^{l} b_{j} c_{j}^{k} + O(h^{r+1}) =$$

$$= \sum_{k=0}^{r-1} \frac{1}{k!} u^{(k+1)}(t) h^{k+1} \underbrace{\left[\frac{1}{k+1} - \sum_{j=1}^{l} b_{j} c_{j}^{k} \right]}_{0} + O(h^{r+1}) = O(h^{r+1})$$

• Mit

$$u_h(t+h) = u + h \sum_{j=1}^{l} b_j f(g_j)$$

erhalten wir für den lokalen Fehler

$$u(t+h) - u_h(t+h) = h \sum_{j=1}^{l} b_j [f(u(t+c_jh)) - f(g_j)] + O(h^{r+1})$$

die Abschätzung

(28)
$$|u(t+h) - u_h(t+h)| \le h \sum_{j=1}^{l} |b_j| |f(u(t+c_jh)) - f(g_j)| + ch^{r+1}$$

$$\le h \sum_{j=1}^{l} |b_j| L |u(t+c_jh) - g_j| + ch^{r+1}.$$

• Schätzen nun völlig analog wie oben

$$|u(t+c_ih)-g_i|\leq\ldots$$

für i = 1, 2, ..., l ab:

$$u(t+c_ih) = u + \int_t^{t+c_jh} f(u(s)) ds = u + h \sum_{j=1}^l a_{ij} f(u(t+c_jh)) + R_i.$$

Für das Restglied R_i erhalten wir (\uparrow) :

$$R_{i} = \sum_{k=0}^{r-2} \frac{1}{k!} u^{(k+1)}(t) h^{k+1} \underbrace{\left[\frac{c_{i}^{k+1}}{k+1} - \sum_{j=1}^{l} a_{ij} c_{j}^{k} \right]}_{= \frac{r}{(26)} = 0} + O(h^{r}) = O(h^{r}).$$

Mit

$$g_i = u + h \sum_{j=1}^{l} a_{ij} f(g_j)$$

erhalten wir

$$u(t+c_ih) - g_i = h \sum_{j=1}^{l} a_{ij} [f(u(t+c_jh)) - f(g_j)] + O(h^r).$$

Wegen der Lipschitz–Stetigkeit von f folgt:

$$\max_{i=1,l} |u(t+c_ih) - g_i| \le h \|A\|_{\infty} L \max_{j=1,l} |u(t+c_jh) - g_j| + ch^r,$$

und daher gilt:

(29)
$$\max_{i=\overline{1,l}} |u(t+c_i h) - g_i| \le \frac{1}{1-h||A||_{\infty} L} ch^r = O(h^r).$$

• Aus (28) und (29) folgt nun unmittelbar:

$$u(t+h) - u_h(t+h) = O(h^{r+1}), h \to 0.$$

q.e.d.

• Bemerkung 3.11:

1. Für r=l folgt aus Satz 3.10, daß eine l-stufige, implizite RK-Formel die Ordnung p=l hat, falls

(26)
$$\sum_{j=1}^{l} a_{ij} c_j^k = \frac{c_i^{k+1}}{k+1}, \quad k = \overline{0, l-2}$$

für alle $i = 1, 2, \ldots, l$ und

gilt. Diese Bedingungen bedeuten für die dazugehörigen QF, daß sie den algebraischen Genauigkeitsgrad l-1 besitzt, d.h. z.B. für

(30)
$$\int_{t}^{t+h} g(s) ds = h \sum_{j=1}^{l} b_{j} g(t + c_{j}h),$$
$$\forall g(s) = (s - t)^{k}, \quad k = 0, 1, \dots, l - 1.$$

Tatsächlich,

$$\int_{t}^{t+h} (s-t)^{k} ds = \frac{(s-t)^{k+1}}{k+1} \Big|_{t}^{t+h} = \frac{h^{k+1}}{k+1},$$

$$h \sum_{j=1}^{l} b_{j} (t+c_{j}h-t)^{k} = h^{k+1} \sum_{j=1}^{l} b_{j} c_{j}^{k},$$

d.h. (30) ist äquivalent zu (27).

2. Falls $\{c_j\}_{j=\overline{1,l}}$ paarweise verschieden, d.h. $c_i \neq c_j \ \forall \ i \neq j$, vorgegeben werden, dann lassen sich die Vektoren $\{a_{ij}\}_{j=\overline{1,l}}$ aus den l GS (26) mit $k=\overline{0,l-1}$ (!) und der Vektor $\{b_j\}_{j=\overline{1,l}}$ aus den GS (27) eindeutig bestimmen, da die Systemmatrix

$$\begin{bmatrix} 1 & \dots & 1 \\ c_1 & \dots & c_l \\ c_1^{l-1} & \dots & c_l^{l-1} \end{bmatrix}, \det [\dots] = \bigcap_{\substack{i,j=1 \\ i>j}}^{l} (c_i - c_j) \neq 0,$$

(Van der Monde'sche Determinante)

dieser GS offenbar regular ist.

Die damit erzeugten QF beruhen auf Polynominterpolation an den Stützstellen $t + c_i h$, i = 1, 2, ..., l.

3. Man kann nun durch geeignete <u>Wahl</u> der Stützstellen $\{t + c_i h\}_{i=\overline{1,l}}$ versuchen, die Ordnung weiter zu erhöhen. Das erinnert, so wie auch schon die ersten 3 Bedingungen für Ordnung 3, an die Gauß-QF (siehe Satz 3.12).

• Satz 3.12:

Vor.: 1. Es seien die Voraussetzungen von Satz 3.2 erfüllt (siehe auch Satz 3.10).

$$2. \quad l \in I\!\!N: f \in C^l(D, I\!\!R^N).$$

Bh.: Falls für eine l-stufige RK-Formel die Bedingungen

(26)
$$\sum_{j=1}^{l} a_{ij} c_j^k = \frac{c_i^{k+1}}{k+1}, \quad k = 0, 1, \dots, l-1,$$
 für alle $i = 1, 2, \dots, l$ und

(31)
$$\sum_{j=1}^l b_j c_j^k = \frac{1}{k+1}, \ k=0,1,\ldots,2l-1$$
 erfüllt sind, besitzt die RK-Formel die Konsistenzordnung $p=2l$.

<u>Beweis:</u> siehe z.B. [3] Grigorieff R.D.: Numerik gewöhnlicher Differentialgleichungen I. Teubner-Verlag, Stuttgart 1972.

• Bemerkung 3.13:

Bedingung (31) bestimmt genau die Gauß-QF der Ordnung l, die bekanntlich den algebraischen Genauigkeitsgrad 2l-1 besitzt (vgl. auch Bem. 3.11.1). Bedingungen an die $\{a_{ij}\}$ mit $k=\overline{0,l-1}$ besitzen eine eindeutige Lsg. bei geg. $\{c_i\}$:

⇒ Runge-Kutta-Formeln vom Gauß-Typ

• Bemerkung 3.14:

- 1. QF von Radau := QF maximaler algebr. Genauigkeit (2l-1): $c_1 = 0 \text{ oder } c_l = 1$
 - \Rightarrow Runge-Kutta-Formeln vom Radau-Typ :

(32) •
$$\sum_{j=1}^{l} b_j c_j^k = \frac{1}{k+1}, \ k = 0, 1, \dots, p-1 \text{ mit } p = 2l-1.$$

(33) • a)
$$c_1 = 0$$
: $a_{11} = a_{12} = \ldots = a_{1l} = 0$ (Bd. !)
 $(\Rightarrow K_1 = f(t, u) \text{ bzw. } g_1 = u)$

(34) • b)
$$c_l = 1$$
: $a_{1l} = a_{2l} = \ldots = a_{ll} = 0$ (Bd. !)
 ($\Rightarrow K_l$ bzw. g_l explizit berechenbar !)

- $\Rightarrow \underline{\text{Konsistenzordnung}} = \underline{2l-1}$
- 2. QF von Lobatto := QF maximaler algebr. Genauigkeit (2l-2): $c_1 = 0$ und $c_l = 1$.

- $\bullet \ (32) \ \text{für} \ p = 2l-2 \qquad \bullet \ (33) \ + \ (34) \qquad \bullet \ (35) \ \text{für} \ r = l-1$
- \Rightarrow Konsistenzordnung = $\underline{2l-2}$.

3.3.4Konvergenztheorie für Einschrittverfahren

■ **Def. 3.15:** (ESV)

Ein Verfahren zur Lösung des AWP

(1)
$$\begin{cases} u'(t) = f(t, u(t)), & t \in I = [0, T] \\ \text{AB: } u(0) = u_0 \end{cases}$$
heißt Einschrittverfahren (ESV), falls es von der Form

(36)
$$u_{j+1} = u_j + h_j \varphi(t_j, u_j, h_j), \quad j = 0, 1, \dots, m-1,$$

ist, mit geg. AB u_0 .

Bsp.: Die bisher behandelten RK-Formeln (expliziten wie auch impliziten) sind offenbar ESV. Tatsächlich,

EPZV:
$$\varphi = f$$
, RK: $\varphi(t_j, u_j, h_j) = \sum_{i=1}^l b_j f(t_j + c_i h_i, g_i)$, $g = \Phi(g, t_j, u_j, h_j)$.

■ Analog zum Pkt. 3.3.1 (\Rightarrow EPZV) können wir auch (36) als Operatorgleichung in der Form

(36) Ges.
$$u_h \in X_h : F_h(u_h) = 0 \text{ in } Y_h$$

schreiben mit

$$F_h(v_h)(t_{j+1}) := \begin{cases} \underbrace{\frac{=:D_h v_h(t_j)}{1}}_{=:D_h v_h(t_j)} - \varphi(t_j, v_j, h_j), j = 0, 1, \dots, m-1, \\ v_0 - u_0, & j = -1, \end{cases}$$

und (z.B.) $X_h = C^1(I_h)$ und $Y_h = C(I_h)$.

Der lokale Abschneidefehler (Approximationsfehler) wird dann durch (vgl. (9))

(37)
$$\tau_h(t_{j+1}) = F_h(u)(t_{j+1}) \equiv \begin{cases} \frac{u(t_{j+1}) - u(t_j)}{h_j} - \varphi(t_j, u(t_j), h_j), j = \overline{0, m-1}, \\ 0, j = -1, \end{cases}$$

eingeführt oder kurz: $\tau_h = F_h(u)$, wobei u Lsg. des AWP (1) ist.

■ Konsistenz (lokale und globale):

• **Definition 3.16:** (globale Konsistenz(-ordnung))

Ein ESV (36) heißt mit dem AWP (1) konsistent, falls

$$\|\tau_h\|_{Y_h} \longrightarrow 0$$
 für $h \equiv |h| \to 0$.

Ein ESV (36) hat die Konsistenzordnung p, falls

$$\|\tau_h\|_{Y_h} = O(h^p).$$

• Bemerkungen:

- In den Pkt. 3.3.2.2 (explizite) und 3.3.3.4 (implizite) wurde die (lokale) Konsistenzordnung von RK-Formeln studiert.
- RK-Formeln: $\sum_{j=1}^{l} b_j = 1 \Rightarrow$ Konsistenz (Ordnung = 1, falls $f \in C^1$).
- Für (36) bedeutet Konsistenz (vgl. auch Bew. Satz 3.17)

$$\varphi(t,u,h) \longrightarrow f(t,u)$$
 für $h \to 0$

bzw. präziser ($\|\tau_h\|_{Y_h} := \max_{j=0,m} |\tau_j| \equiv \max_{j=1,m-1} |\tau_{j+1}|, \ m = m(h)$!)

(38)
$$\max_{j=0,m-1} |\varphi(t_j, u(t_j), h_j) - f(t_j, u(t_j))| \xrightarrow{h \to 0} 0.$$

• Satz 3.17:

 $\underline{\text{Vor.:}} \quad 1) \quad f \in C(I \times \mathbb{R}^N) \equiv [C(I \times \mathbb{R}^N)]^N - \text{stetig};$

2) Es gelte (38).

Bh.: Dann ist das ESV konsistent.

Beweis: folgt sofort aus der Darstellung

$$\tau_h(\underbrace{t+h}_{t_{j+1}}) = \left[\frac{1}{h}(u(t+h) - u(t)) - u'(t)\right] + [f(t, u(t)) - \varphi(t, u(t), h)].$$

$$t_{j+1}$$

$$t = t_j$$

q.e.d.

■ Stabilität:

 Btr. zunächst die Auswirkung eines <u>Fehlers</u> im Startwert (⇒ Stabilität bzgl. der AB):

(39)
$$\begin{cases} v_0 = u_0 + \delta_0, \\ v_{j+1} = v_j + h_j \varphi(t_j, v_j, h_j), & j = 0, 1, \dots, m-1. \end{cases}$$

• Frage:
$$|u_j - v_j| \le ?$$
 $j = 1, 2, ..., m$.
 $\uparrow \qquad \uparrow$
 $(36) \quad (39)$

Lemma 3.18:

Vor.: φ sei Lipschitz-stetig im zweiten Argument, d.h. $\exists L = \text{const.} \in [0, \infty)$:

$$\begin{aligned} (40) & |\varphi\left(t,v,h\right)-\varphi\left(t,w,h\right)| \leq L \,|v-w| \\ & \forall t \in I, \ \forall v,w \in I\!\!R^N, \ \forall \ \text{zul\"{assigen}} \ h \in \{h_0,h_1,\ldots,h_{m-1}\}. \end{aligned}$$

Bh.: Dann gilt für (39) die Abschätzung

(41)
$$|u_j - v_j| \le e^{Lt_j} |\delta_0|, |D_h u_j - D_h v_j| \le L |u_j - v_j|,$$

mit $D_h v_j = v_{t,j} = h_j^{-1} (v_{j+1} - v_j).$

Beweis: Btr. Fehler $z_i = u_i - v_i$:

q.e.d.

• Btr. nun die Fortpflanzung der einzelnen lokalen Fehler $\delta_{j+1} = h_j y_{j+1}$ durch das ESV:

(42)
$$\begin{cases} v_0 = u_0 + \mathbf{y_0} \\ v_{j+1} = v_j + h_j \varphi(t_j, v_j, h_j) + \underbrace{\mathbf{h_j y_{j+1}}}_{\delta_{j+1}}, \quad j = \overline{0, m-1} \end{cases}$$

$$\frac{\text{Frage:}}{\uparrow} \quad |u_j - v_j| \le ?$$

$$\uparrow \quad \uparrow$$

$$(36) \quad (42)$$

Lemma 3.19:

$$\underline{\text{Vor} ::} \quad |\varphi(t, v, h) - \varphi(t, w, h)| \le L|v - w| \quad \forall \ t; v, w; h.$$

Bh:. Dann gelten für (42) die Abschätzungen:
$$(43) \quad |u_j-v_j| \leq e^{Lt_j}|y_0| + \frac{e^{Lt_j}-1}{L} \max_{i=\overline{1,j}}|y_i|,$$

$$(44) |D_h u_j - D_h v_j| \le L |u_j - v_j| + |y_{j+1}|.$$

Beweis:

o Def. für fixiertes $i \in \{0, 1, ..., m-1\}$ Gitterfkt. $v^{(i)}$:

(45)
$$\begin{cases} v_{j+1}^{(i)} = v_j^{(i)} + h_j \varphi(t_j, v_j^{(i)}, h_j), & j = i, i+1, \dots, m-1, \\ v_i^{(i)} = v_i := \underbrace{v_{i-1} + h_{i-1} \varphi(t_{i-1}, v_{i-1}, h_{i-1})}_{=: v_i^{(i-1)}} + h_{i-1} y_i, \end{cases}$$

wobei $h_{-1} = 1$ und $v_0^{(-1)} = u_0$ gesetzt wird.

 \circ Btr. $v_i^{(i-1)}$:

(46)
$$\begin{cases} v_{j+1}^{(i-1)} = v_j^{(i-1)} + h_j \varphi(f_j, v_j^{(i-1)}, h_j), & j = i-1, i, i+1, \dots, m-1, \\ v_{i-1}^{(i-1)} = v_{i-1}. \end{cases}$$

• Aus <u>Lemma 3.18:</u> folgt:

(47)
$$|v_j^{(i-1)} - v_j^{(i)}| \le e^{L(t_j - t_i)} h_{i-1} |y_i|, \quad j = \overline{i, m} \text{ und } i = \overline{1, m},$$

wobei $v_m^{(m)} = v_m := v_{m-1} + h_{m-1} \varphi(t_{m-1}, v_{m-1}, h_{m-1}) + h_{m-1} y_m.$

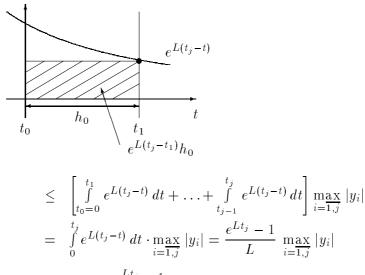
o Für den Gesamtfehler folgt daraus:

$$|u_j - v_j| \le \underbrace{|u_j - v_j^{(0)}|}_{\le e^{Lt_j}|y_0|} + |v_j^{(0)} - \underbrace{v_j^{(j)}|}_{=v_j}$$

<u>Lemma 3.18</u>

Btr.

$$\begin{aligned} |v_{j}^{(0)} - v_{j}^{(j)}| & \leq |v_{j}^{(0)} - v_{j}^{(1)}| + |v_{j}^{(1)} - v_{j}^{(2)}| + \ldots + |v_{j}^{(j-1)} - v_{j}^{(j)}| \\ & \leq e^{L(t_{j} - t_{1})} h_{0} |y_{1}| + e^{L(t_{j} - t_{2})} h_{1} |y_{2}| + \ldots + e^{L(t_{j} - t_{j})} h_{j-1} |y_{j}| \\ & \leq [e^{L(t_{j} - t_{1})} h_{0} + e^{L(t_{j} - t_{2})} h_{1} + \ldots + e^{L(t_{j} - t_{j})} h_{j-1}] \max_{i=1,j} |y_{i}| \end{aligned}$$



$$\Longrightarrow |u_j - v_j| \le e^{Lt_j} |y_0| + \frac{e^{Lt_j} - 1}{L} \max_{i=1,j} |y_i|$$

$$0 D_h u_j - D_h v_j = \varphi(t_j, u_j, h_j) - \varphi(t_j, v_j, h_j) - y_{j+1}$$

$$\implies |D_h u_j - D_h v_j| \le L |u_j - v_j| + |y_{j+1}|.$$

q.e.d.

• Das gestörte ESV (42) läßt sich kompakt in der Form

$$(48) F_h(v_h) = y_h$$

schreiben. Aus Lemma 3.19 folgt sofort:

Satz 3.20:

$$\underline{\text{Vor} ::} \quad |\varphi(t, v, h) - \varphi(t, w, h)| \le L|v - w| \quad \forall t; v, w; h$$

Bh.: Dann gibt es eine positive Konstante C:

(49)
$$\|v_h - w_h\|_{X_h} \le C \|F_h(v_h) - F_h(w_h)\|_{Y_h} \quad \forall v_h, w_h \in X_h,$$

mit $X_h = C^1(I_h) = C^1(I_h, \mathbb{R}^N)$ und $Y_h = C(I_h)$.

Beweis:

o Sei zunächst:
$$w_h = u_h$$
 : $F_h(u_h) = 0$ (36)
 v_h : $F_h(v_h) = y_h$

Dann folgt aus Lemma 3.19 sofort:

$$\max_{j=\overline{0,m}} |v_j - u_j| \overset{(43)}{\leq} C_1 \max_{j=\overline{0,m}} |y_j|.$$
mit $C_1 = e^{LT} + \left(e^{LT} - 1\right) / L$. Aus (44) und (43) folgt
$$\max_{j=\overline{0,m-1}} |D_h u_j - D_h v_j| \leq \underbrace{\left[Le^{LT} + e^{LT} - 1 + 1\right]}_{j=\overline{0,m}} \max_{j=\overline{0,m}} |y_j|.$$

Daraus ergibt sich dann:

$$||v_h - u_h||_{X_h} := \max_{j = \overline{0, m}} |v_j - u_j| + \max_{j = \overline{0, m - 1}} |D_h v_j - D_h u_j| \le \underbrace{(C_1 + C_2)}_{C} ||y_h||_{Y_h}.$$

o Für fix. bel. $w_h \in X_h$ folgt (49) aus dem eben bew. Spezialfall angewandt auf

$$\widetilde{F}_h(\cdot) := F_h(\cdot) - z_h \text{ mit } z_h = F_h(w_h) \Rightarrow \widetilde{F}_h(v_h) = y_h - z_h \equiv F(v_h) - F(w_h)$$
$$\widetilde{F}_h(w_h) = 0$$

q.e.d.

• **Definition 3.21:** (Stabilität)

Ein Einschrittverfahren heißt stabil, falls $\exists C = \text{const.} > 0, C \neq C(h)$:

$$(50) ||v_h - w_h||_{X_h} \le C||F_h(v_h) - F_h(w_h)||_{Y_h} \forall v_h, w_h \in X_h.$$

• Bemerkung:

Seien $y_h, z_h \in Y_h$ beliebig, und erfüllen $v_h, w_h \in X_h$ die Gleichungen

$$F_h(v_h) = y_h$$
 und $F_h(w_h) = z_h$.

Dann gilt für ein stabiles ESV

$$||v_h - w_h||_{X_h} \le C||y_h - z_h||_{Y_h}$$
.

Die Stabilität des ESV bedeutet also, daß sich die Lösung der Näherungsgleichung (Lipschitz-) stetig und gleichmäßig in h mit den Daten (den rechten Seiten) ändert!

Satz 3.20 zeigt also, daß die

Lipschitz-Stetigkeit von $\varphi(\cdot, \bullet, \cdot)$

ein hinreichendes Kriterium für die Stabilität des ESV ist.

■ Konvergenz: $e_h = u - u_h \longrightarrow 0$ in X_h für $h = |h| \to 0$.

• **Definition 3.22:** (diskrete Konvergenz)

Ein ESV heißt konvergent, falls

$$||e_h||_{X_h} \equiv ||u - u_h||_{X_h} \stackrel{h \to 0}{\longrightarrow} 0.$$

Falls

$$||e_h||_{X_h} = O(h^p), \quad h = |h|,$$

nennt man das ESV ein Verfahren der Konvergenzordnung p. Hierbei ist $e_h=u\ -\ u_h$ der Fehler.

$$(1)$$
 (36)

 $\bullet \ \ \frac{\text{Approximation}}{\text{Satz 3.23:}} \left(\begin{array}{c} \text{Approximation} \\ \text{Konsistenz}_{Y_h} \end{array} + \text{Stabilit\"at}_{(Y_h,X_h)} \Rightarrow \text{diskr. Konvergenz}_{X_h} \right)$

<u>Vor.:</u> 1) ESV (36) sei stabil i. S. Def. 3.21.

2) ESV(36) sei konsistent (von der Ordnung p) i. S. Def. 3.16.

<u>Bh.:</u> Dann ist das ESV (36) konvergent (von der Ordnung p).

Beweis: trivial!

$$\circ \ u \in X|_{I_h} \subset X_h : F_h(u) = \tau_h$$

$$(1)$$

$$\circ u_h \in X_h : F_h(u_h) = 0$$

$$\begin{array}{c} \circ \ \underline{\text{Stabilit"at:}} \Rightarrow \| \underbrace{u - u_h} \|_{X_h} \leq C \, \| \tau_h \|_{Y_h} \\ + \\ \circ \ \underline{\text{Konsistenz:}} \Rightarrow \text{a)} \, \| \tau_h \|_{Y_h} \to 0, \ h \to 0 \\ & \downarrow \qquad \text{b)} \, \| \tau_h \|_{Y_h} = O(h^p) \\ & \Rightarrow \| e_h \|_{X_h} \to 0 \\ \Rightarrow \| e_h \|_{X_h} = O(h^p) \\ \end{array}$$

diskrete

Konvergenz

q.e.d.

Anwendung der Theorie auf (implizite) RK-Formeln:

• RK-Formeln sind ESV mit

(*)
$$\varphi(t, u, h) = \sum_{j=1}^{l} b_j f(t + c_j h, g_j(t, u, h)),$$

wobei $g = (g_1, \ldots, g_l)^T$ Lösung der Fixpunktgleichung

$$(**) g = \Phi(g, t, u, h)$$

ist, mit $\Phi = (\Phi_1, \dots, \Phi_l)^T$, $\Phi_i = \Phi_i(g, t, u, h) = u + h \sum_{i=1}^l a_{ij} f(t + c_j h, g_j)$.

• Konsistenz (\uparrow) :

Satz 3.24: (Konsistenz)

f sei stetig, beschränkt und erfülle die Lipschitz-Bedingung

$$|f(t,v) - f(t,w)| \le L|v - w| \quad \forall \ t; v, w.$$

Für die Gewichte b_j gelte:

$$\sum_{j=1}^{l} b_j = 1.$$

Dann ist die RK-Formel konsistent mit dem AWP (1). Bh .:

Beweis: \longrightarrow Satz 3.17 !

o Sei $g_i = g_i(t, u(t), h)$. Dann folgt wie im Beweis von Satz 3.9 aus

$$g_{i} - u(t) \stackrel{\stackrel{(**)}{\downarrow}}{=} h \sum_{j=1}^{l} a_{ij} f(t + c_{j}h, g_{j}) =$$

$$= h \sum_{j=1}^{l} a_{ij} [f(t + c_{j}h, g_{j}) - f(t + c_{j}h, u(t))] + h \sum_{j=1}^{l} a_{ij} f(t + c_{j}h, u(t))$$

die Abschätzung

$$\max_{i=1,l} |g_i - u(t)| \le h ||A||_{\infty} L \max_{i=1,l} |g_i - u(t)| + h ||A||_{\infty} M$$

$$\mathrm{mit}\ M = \sup_{\tau,t\in I}\ |f(\tau,u(t))|.$$

Daher gilt für hinreichend kleine h

$$\max_{i=\overline{1,l}} |g_i - u(t)| \le \frac{h ||A||_{\infty} M}{1 - h ||A||_{\infty} L}.$$

o Aus

Aus
$$(*), \sum b_{j}=1$$

$$\varphi(t, u(t), h) - f(t, u(t)) =$$

$$= \sum_{j=1}^{l} b_{j} [f(t + c_{j}h, g_{j}) - f(t + c_{j}h, u(t))] + \sum_{j=1}^{l} b_{j} [f(t + c_{j}h, u(t)) - f(t, u(t))]$$

folgt dann

$$\begin{aligned} |\varphi\left(t,u(t),h\right) - f(t,u(t))| &\leq \frac{\sum\limits_{\|b\|_{1}} |b\|_{1} |Lh\|A\|_{\infty} M}{1 - h\|A\|_{\infty} L} + \\ &+ \|b\|_{1} \max\limits_{j=\overline{1,l}} |f(t + c_{j}h,u(t)) - f(t,u(t))| \end{aligned}$$

und daher (beachte: $f(\tau, u(t))$ ist gleichmäßig stetig!)

$$\max_{k=0,m-1} |\varphi(t_k,u(t_k),h_k) - f(t_k,u(t_k))| \longrightarrow 0 \text{ für } |h| \to 0.$$

o Der Rest folgt sofort aus Satz 3.17!

q.e.d.

Ü 3.9 Man zeige, daß im Falle expliziter RK-Formeln die Lipschitz-Bd. für den Konsistenzbew. <u>nicht</u> benötigt wird!

<u>Bem.:</u> Falls f entsprechend oft differenzierbar ist, dann kann man Aussagen über die Konsistenzordnung von RK-Formeln gewinnen:

siehe Pkt. 3.3.2.2: explizite RK-Formeln,

Pkt. 3.3.3.4: implizite RK-Formeln.

• Stabilität:

Satz 3.25: (Stabilität)

Vor.: f sei stetig und erfülle die Lipschitz-Bedingung.

$$|f(t,v) - f(t,w)| \le L|v - w| \quad \forall \ t; v, w.$$

Bh.: Dann ist die RK-Formel stabil.

 $\underline{\text{Beweis:}} \longrightarrow \text{Satz } 3.20 \Rightarrow \text{z.Z.:} \left| \varphi \left(t, v, h \right) - \varphi \left(t, w, h \right) \right| \leq \widetilde{L} |v - w| \ !$

 \circ Seien t, v, w, h bel.: \Rightarrow

$$\begin{aligned} |\varphi\left(t,v,h\right) - \varphi\left(t,w,h\right)| &\overset{(*)}{\leq} \sum_{j=1}^{l} |b_{j}| \left| f(t + c_{j}h, g_{j}(t,v,h)) - f(t + c_{j}h, g_{j}(t,w,h)) \right| \leq \\ &\leq \sum_{j=1}^{l} |b_{j}| L \left| g_{j}(t,v,h) - g_{j}(t,w,h) \right|. \end{aligned}$$

• Wie im Beweis von Satz 3.9 folgt:

$$\begin{split} \max_{j=\overline{1,l}} |g_j(t,v,h) - g_j(t,w,h)| &\overset{(**)}{\leq} |v-w| + h \, \|A\|_{\infty} \, L \max_{j=\overline{1,l}} |g_j(t,v,h) - g_j(t,w,h)| \\ & \downarrow \downarrow \\ \max_{j=\overline{1,l}} |g_j(t,v,h) - g_j(t,w,h)| &\leq \frac{1}{1 - h \|A\|_{\infty} \, L} \, |v-w|. \end{split}$$

 \circ Insgesamt folgt damit die Lipschitz–Bed. für $\varphi\colon$

$$|\varphi\left(t,v,h\right)-\varphi\left(t,w,h\right)|\leq \underbrace{\frac{L\|b\|_{1}}{1-h\|A\|_{\infty}L}}_{=:\widetilde{L}}|v-w|.$$

o Der Rest folgt aus Satz 3.20!

q.e.d.

Konvergenz:

Aus den Sätzen Konsistenz+Stabilität \Rightarrow Konvergenz Aus den Sätzen 3.24, 3.25 und $\underline{3.23}$ folgt sofort:

Satz 3.26: (\$\hat{2}\$ Satz 3.23 !)

 $\underline{\text{Vor.:}}$ 1) f sei stetig und erfülle die Lipschitz-Bedingung:

$$|f(t,v) - f(t,w)| \le L|v - w| \quad \forall \ t; v, w.$$

2) Für die Gewichte b_j gelte:

$$\sum_{j=1}^{l} b_j = 1.$$

Bh.: Dann ist die RK-Formel für das AWP (1) konvergent.

Bemerkung:

Falls die RK-Formel die Konsistenzordnung p hat, und f entsprechend oft stetig differenzierbar ist, dann ist auch die Konvergenzordnung p.

Bemerkung:

o Die obigen Aussagen gelten auch unter der lokalen Lipschitz-Bd.

$$|f(t,v) - f(t,w)| < L|v-w| \quad \forall (t,v), (t,w) \in U,$$

wobei $U \subset I \times I\!\!R^N$ eine Umgebung des Graphen von f

$$\{(t,f(t,u(t))):t\in I\}$$

ist und u(t) die exakte Lsg. des AWP (1) bezeichnet.

 Für eine lokale Variante des <u>Stabilitätssatzes 3.25</u> genügt ebenfalls eine lokale Lipschitz-Bed. der Art

$$|\varphi(t,v,h) - \varphi(t,w,h)| < L|v-w| \quad \forall (t,v), (t,w) \in U, \quad h < H,$$

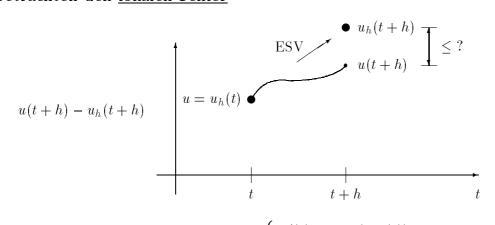
falls zusätzlich die Konsistenz des Verfahrens vorausgesetzt wird. Man kann die Aussagen $\forall v_h, w_h \in X_h : \|F_h(v_h)\|_{Y_h}, \|F_h(w_h)\|_{Y_h} \leq \eta$ mit η – hinr. klein, zeigen.

3.3.5 Praktische Durchführung von Einschrittverfahren

- Die Wahl der **richtigen Schrittweiten** $\{h_j\}$ ist offenbar von entscheidender Bedeutung für eine effiziente Durchführung eines ESV. Die Aufgabe einer Schrittweitensteuerung ist es, den jeweils neu entstehenden <u>lokalen Fehler</u> in einer gewünschten Größenordnung zu halten.
- Im folgenden werden zwei Möglichkeiten zur Schätzung der lokalen Fehler angegeben (Pkt. 3.3.5.1), darauf aufbauend eine Schrittweitensteuerung vorgeschlagen (Pkt. 3.3.5.2) und schließlich eine Technik zur Berechnung von Näherungen an vorgegebenen Punkten $t + \theta h$, die nicht von der Schrittweitensteuerung herrühren (Pkt. 3.3.5.3).

3.3.5.1 Schätzung der lokalen Fehler

■ Betrachten den <u>lokalen Fehler</u>



wobei
$$u(t+h)$$
 Lösung d. AWP:
$$\begin{cases} u'(s) &= f(s,u(s)), s \geq t, \\ u(t) &= u \end{cases}$$
 im Pkt. $t+h$, und
$$u_h(t+h) = u + h\varphi(t,u,h) - \text{ESV}.$$

Es wird nun vorausgesetzt, daß das ESV die Konsistenzordnung p hat, d.h.

$$u(t+h) - u_h(t+h) = O(h^{p+1}),$$

und der Fehlerterm die Gestalt

(51)
$$u(t+h) - u_h(t+h) = \underbrace{c(t,u) \, h^{p+1}}_{\text{eighrender Fehlerterm}} + O(h^{p+2})$$

$$= \text{führender Fehlerterm}$$

$$(\text{engl.: principal error term})$$

hat, wobei c unabhängig von der Schrittweite h sein soll.

Für Runge-Kutta-Formeln läßt sich diese Darstellung des lokalen Fehlers durch Taylor-Entwicklung leicht nachweisen, z.B.:

• explizites Euler-Verfahren (vgl. Pkt. 3.3.1):

$$u(t+h) - u_h(t+h) = \underbrace{\frac{1}{2} (f_t + f_u f) (t, u)}_{=: c (t, u)} h^2 + O(h^3),$$

• explizite Mittelpunktsregel (vgl. Pkt. 3.3.2):

$$u(t+h) - u_h(t+h) = \underbrace{\frac{1}{24} \left(f_{tt} + 2f_{tu}f + f_{uu}f^2 + 4\left(f_uf_t + f_u^2f \right) \right)(t,u)}_{=: c(t,u)} h^3 + O(h^4).$$

Im folgenden wird nur von der Existenz eines von h unabhängigen, in t und u <u>hin</u>reichend glatten (!) Koeffizienten c (t,u), nicht aber von seiner speziellen Gestalt Gebrauch gemacht !!

■ Schätzung durch Extrapolation:

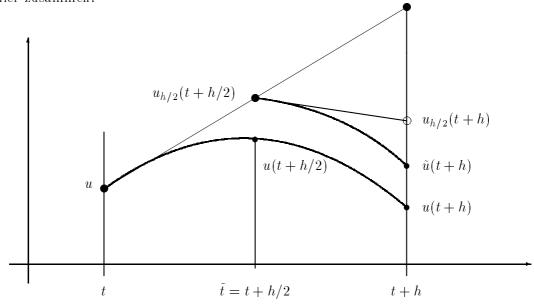
• Zuerst wird ein Schritt mit der Schrittweite h durchgeführt. Man erhält eine Näherung u_h im nächsten Gitterpunkt t + h und wegen (51) gilt:

$$u(t+h) - u_h(t+h) = c(t,u)h^{p+1} + O(h^{p+2}).$$

• Zwei Schritte mit halber Schrittweite h/2 führen ebenfalls zum Gitterpunkt t+h, erzeugen aber eine andere Näherung $u_{h/2}(t+h)$ im Punkt t+h. Für den ersten Schritt gilt wegen (51):

$$u(t+h/2) - u_{h/2}(t+h/2) = c(t,u)\left(\frac{h}{2}\right)^{p+1} + O(h^{p+2}).$$

Der Fehler nach dem zweiten Schritt setzt sich aus der Fortpflanzung des ersten (↑) lokalen Fehlers (= Fehler in den AB: siehe Lemma 3.18) und dem neuen lokalen Fehler zusammen:



$$(52) \qquad u(t+h) - u_{h/2}(t+h) = \\ = [u(t+h) - \tilde{u}(\tilde{t}+h/2)] + [\tilde{u}(t+h/2) - u_{h/2}(t+h)] = \\ \text{(mms)} \\ \downarrow \\ = (1+O(h)) c(t,u) \left(\frac{h}{2}\right)^{p+1} + (c(t,u) + O(h)) \left(\frac{h}{2}\right)^{p+1} + O(h^{p+1}) = \\ = 2c(t,u) \left(\frac{h}{2}\right)^{p+1} + O(h^{p+1})$$

• Aus (51) und (52) kann nun $c \equiv c(u,t)$ offenbar eliminiert werden:

Daraus folgt sofort, daß

$$\hat{u}_h(t+h) := u_{h/2}(t+h) + \frac{u_{h/2}(t+h) - u_h(t+h)}{2^p - 1}$$

eine bessere Näherung für die Lösung

$$u(t+h) = \hat{u}_h(t+h) + O(h^{p+2})$$

ist, d.h. mit einem um eine Ordnung verbesserten Fehler.

• Eine Schätzung err (\rightarrow führender Fehlerterm) des lokalen Fehlers erhält man nun, indem man u(t+h) durch den extrapolierten Wert $\hat{u}_h(t+h)$ ersetzt:

$$err := \frac{|\hat{u}_h(t+h) - u_h(t+h)|}{d} ,$$

wobei

$$d = \left\{ \begin{array}{ll} 1 & \text{f\"{u}r den absoluten Fehler,} \\ |u_{h/2}(t+h)| & \text{f\"{u}r den relativen Fehler,} \\ \max\{|u_{h/2}(t+h)|,1\} & \text{f\"{u}r ein gemischtes Fehlermaß} \end{array} \right.$$

ein Skalierungsfaktor ist. Gebräuchlich ist auch eine komponentenweise Anwendung der obigen Skalierungsfaktoren.

• Um den Mehraufwand zur Schätzung des lokalen Fehlers auch zumindestens zum Teil für das Verfahren zu nutzen, betrachtet man entweder $u_{h/2}(t+h)$ oder sogar den extrapolierten Wert $\hat{u}_h(t+h)$ als die Näherung im nächsten Gitterpunkt t+h (\rightarrow lokale Extrapolation).

Der berechnete Fehler err kann jedoch nur als Schätzung des auf die ursprünglichen Werte definierten Verfahrens interpretiert werden. Es darf allerdings erwartet werden, daß die Verwendung von $\hat{u}_h(t+h)$ anstelle von $u_h(t+h)$ die Genauigkeit des Verfahrens insgesamt eher erhöht. Die konsequente Verfolgung dieser Idee führt zu den Extrapolationsverfahren (siehe z.B. [1], [5]).

■ Schätzung mittels eingebetteter Runge-Kutta-Formeln:

• Grundidee:

Neben der das ESV bestimmende RK–Formel der Ordnung p wird eine weitere $\widehat{\text{RK}}$ –Formel mindestens der Ordnung p+1 konstruiert. Die erste Formel liefert einen lokalen Fehler

$$u(t+h) - u_h(t+h) = O(h^{p+1}),$$

für die zweite Formel gilt

$$u(t+h) - \hat{u}_h(t+h) = O(h^{p+2}),$$

und somit folgt

$$u(t+h) - u_h(t+h) = \hat{u}_h(t+h) - u_h(t+h) + O(h^{p+2}).$$

Daher bietet sich für den lokalen Fehler folgende Schätzung an:

$$err = \frac{|\hat{u}_h(t+h) - u_h(t+h)|}{d}$$

mit d = 1 – absoluter Fehler, $d = |\hat{u}_h(t+h)|$ – relativer Fehler, $d = \max\{|\hat{u}_h(t+h)|, 1\}$ – gemischtes Fehlermaß.

• Eingebettete RK-Formeln:

Um den Aufwand zur Berechnung von $\hat{u}_h(t+h)$ möglichst gering zu halten, wählt man die gleichen Werte für c und A in den beiden Tableaus. Man spricht dann von einer eingebetteten Runge-Kutta-Formel (engl.: embedded RK-formula) und stellt die Koeffizienten in einem gemeinsamen Tableau dar, z.B. im Fall expliziter Formeln:

Die beiden Näherungen sind durch

$$u_h(t+h) = u + h(b_1 K_1 + b_2 K_2 + \ldots + b_l K_l)$$

und

$$\hat{u}_h(t+h) = u + h(\hat{b}_1 K_1 + \hat{b}_2 K_2 + \ldots + \hat{b}_l K_l)$$

gegeben.

• <u>Beispiel:</u> Runge-Kutta-Fehlberg Formel [RKF 2(3)]: Ausgangspunkt ist das allgemeine (expl.) Tableau für l = 3:

Die Bedingung, daß die erste (eigentliche) Formel die Ordnung 2 besitzt, lautet (siehe Pkt. 3.3.2.2: $c_i = \sum a_{ij}$):

$$b_1 + b_2 + b_3 = 1,$$

 $b_2 c_2 + b_3 c_3 = 1/2.$

Damit die zweite (zur Schätzung benutzte) Formel die Ordnung 3 besitzt, muß gelten (siehe Pkt. 3.3.2.2: $c_i = \sum a_{ij}$):

$$\hat{b}_1 + \hat{b}_2 + \hat{b}_3 = 1,
\hat{b}_2 c_2 + \hat{b}_3 c_3 = 1/2,
\hat{b}_2 c_2^2 + \hat{b}_3 c_3^2 = 1/3,
\hat{b}_3 a_{32} c_2 = 1/6.$$

Für die Wahl

$$c_2 = 1$$
, $c_3 = 1/2$, $b_3 = 0$

erhält man eine sogenannte Runge-Kutta-Fehlberg-Formel mit der Kurzbezeichnung **RKF 2(3)**. Das Symbol **2(3)** bringt zum Ausdruck, daß das eigentliche Verfahren von der Ordnung 2 ist und daß für die Schätzung des lokalen Fehlers ein Verfahren der Ordnung 3 verwendet wird. Entsprechend ist ein Symbol p(q) zu interpretieren. Das Tableau für RKF 2(3) lautet:

• Bemerkungen:

1. Vom Blickwinkel der Effizienz sind eingebettete Runge-Kutta-Formeln mit $a_{li}=b_i,\ i=1,2,\ldots,l$ ($\begin{cases} \begin{cases} \begin{case$

$$K_{1}(t+h) \equiv f(t+h, u_{h}(t+h))$$

$$= f(t+c_{l}h, u_{h}(t) + h \sum a_{li}(t) K_{i}(t)) \equiv K_{l}(t)$$

$$\downarrow c_{l} = \sum a_{li} = \sum b_{i} = 1 \quad \text{(Konsistenzbed.)}$$

$$u_{h}(t) + h \sum a_{li}(t) K_{i}(t) = u_{h} + h \sum b_{i} K_{i} = u_{h}(t+h)$$

- 2. Am Beispiel der 3-stufigen, expliziten, eingebetteten RK-Formeln sieht man, daß die Lösung keineswegs eindeutig ist. Es wurde daher versucht, die frei wählbaren Parameter so zu wählen, daß der führende Fehlerterm von $u_h(t+h)$ möglichst klein wird. Das führt allerdings dazu daß err den lokalen Fehler eher unterschätzt.
- 3. Es liegt nahe, nicht $u_h(t+h)$, sondern den genaueren Wert $\hat{u}_h(t+h)$ für das (eigentliche) Verfahren und $u_h(t+h)$ nur für die Schätzung des lokalen Fehlers zu verwenden.

Man versucht dann natürlich, die frei wählbaren Parameter so zu wählen, daß der führende Fehlerterm von $\hat{u}_h(t+h)$ möglichst klein wird. Dann darf damit gerechnet werden, daß err den lokalen Fehler (für die erste Formel) eher überschätzt.

4. Ein sehr wichtiges, 7-stufiges, eingebettetes Runge-Kutta-Verfahren 4(5) mit den Eigenschaften aus Bemerkung 1 und 3 wurde von J. R. Dormand und P. J. Prince konstruiert [7], [8]:

0							
$\frac{1}{5}$	$\frac{1}{5}$						
$\frac{3}{10}$	$\frac{3}{40}$	$\frac{9}{40}$					
$\frac{4}{5}$	$\frac{44}{45}$	$-\frac{56}{15}$	$\frac{32}{9}$				
$\frac{8}{9}$	$\frac{19372}{6561}$	$-\frac{25360}{2187}$	$\frac{64448}{6561}$	$-\frac{212}{729}$			
1	$\frac{9017}{3168}$	$-\frac{355}{33}$	$\frac{46732}{5247}$	$\frac{49}{176}$	$-\frac{5103}{18656}$		
1	$\frac{35}{384}$	0	$\frac{500}{1113}$	$\frac{125}{192}$	$-\frac{2187}{6784}$	$\frac{11}{84}$	
	35 384	0	500 1113	125 192	$-\frac{2187}{6784}$	11 84	0
	5179 57600	0	7571 16695	393 640	<u>92097</u> 339200	$\frac{187}{2100}$	$\frac{1}{40}$

In der Literatur hat dieses eingebettete Runge-Kutta-Verfahren die Kurzbezeichnung **DOPRI 4(5)**.

3.3.5.2 Schrittweitensteuerung

■ Idee:

Es wird nun davon ausgegangen, daß eine Schätzung err des lokalen Fehlers vorliegt, die auf einer gewählten Schrittweite h basiert und für die gilt (vgl. Pkt. 3.3.5.1):

$$err = ch^{p+1}$$
.

Eigentliches Ziel ist es, daß der lokale Fehler einen vorgegebenen Wert tol nicht überschreitet. Also wäre die "optimale" Schrittweite h_{neu} durch die Bedingung

$$tol = ch_{neu}^{p+1}$$

definiert. Aus diesen beiden Bedingungen läßt sich die Unbekannte c eliminieren, und man erhält:

$$(*) h_{\text{neu}} := h \cdot \sqrt[p+1]{\frac{tol}{err}} .$$

Daraus leitet sich sofort der folgende Algorithmus zur Schrittweitensteuerung ab.

■ Algorithmus zur Schrittweitensteuerung:

- 1. Es wird ein Schritt mit einer vorgegebenen Schrittweite h durchgeführt und gleichzeitig eine Schätzung err für den lokalen Fehler berechnet.
- 2. Falls $err \leq tol$, wird dieser Schritt akzeptiert und mit der neuen Schrittweite h_{neu} das Verfahren fortgesetzt.
- 3. Andernfalls wird der Schritt verworfen und noch einmal mit der neuen Schrittweite h_{neu} gestartet.

■ Praktische Hinweise:

- 1. Um sicher zu sein, daß die neue Schrittweite auch tatsächlich den lokalen Fehler unterhalb des Wertes tol hält, wird die optimale Schrittweite um einen Sicherheitsfaktor fac reduziert (z.B. fac = 0.8).
- 2. Weiters ist es zweckmäßig, zwischen h und $h_{\rm neu}$ einen maximalen Vergrößerungsfaktor facmax und einen minimalen Verkleinerungsfaktor facmin einzuführen, um allzu dramatische Änderungen der Schrittweite zu verhindern. Mit diesen Modifikationen wird aus (*) die Formel

$$h_{\text{neu}} = h \cdot \min \left\{ facmax, \max \left\{ facmin, fac \cdot \left(\frac{tol}{err} \right)^{1/(p+1)} \right\} \right\}.$$

3. Außerdem ist es ratsam, nach einem verworfenen Schritt zwei erfolgreiche Schritte abzuwarten, bevor wieder eine Vergrößerung der neuen Schrittweite zugelassen wird.

3.3.5.3 Berechnung von Näherungen an vorgegebenen Punkten

■ Problematik:

Häufig ist es erwünscht, Näherungen an vorgegebenen Punkten zu berechnen (z.B. für eine graphische Ausgabe mit gegebener Feinheit). Berücksichtigt man solche Punkte bereits bie der Schrittweitensteuerung (z.B. indem man grundsätzlich höchstens bis zum nächsten gegebenen Punkt fortschreitet), so verliert man unter Umständen an Effizienz.

■ Kontinuierliche Runge-Kutta-Formeln:

- Eine Möglichkeit, Näherungen an vorgegebenen Punkten effizient zu berechnen, bieten die sogenannten kontinuierlichen (engl.: continuous) RKF. Diese Formeln beinhalten einen Parameter $\theta \in (0,1]$ und liefern Näherungen für $u(t+\theta h)$. Für $\theta=1$ erhält man also die eigentliche RKF. Durch spezielle Setzungen für θ kann man Näherungen an vorgegebenen Punkten mitrechnen, ohne die Schrittweitensteuerung zu beeinflussen.
- Aus Effizienzgründen fordert man, daß die Werte von c und A im Tableau unabhängig von θ sind:

$$\begin{array}{c|c} c & A \\ \hline & b(\theta) \end{array}$$

• Beispiel: Damit eine explizite, 3-stufige RKF für alle $\theta \in (0,1]$ von der Ordnung 3 ist, müssen folgende Bedingungen gelten (vgl. Pkt. 3.3.2.2: $c_i = \sum a_{ij}$):

$$b_1 + b_2 + b_3 = \theta,$$

$$b_2 c_2 + b_3 c_3 = \frac{\theta^2}{2},$$

$$b_2 c_2^2 + b_3 c_3^2 = \frac{\theta^3}{3},$$

$$b_3 a_{32} c_2 = \frac{\theta^3}{6}.$$

Man sieht leicht, das es <u>nicht</u> möglich ist, c_2, c_3 und a_{32} unabhängig von θ zu wählen. Man fordert daher Ordnung 3 nur für $\theta = 1$ und begnügt sich sonst mit Ordnung 2. Für $c_2 = 1/2$ und $c_3 = 1$ erhält man dann folgendes Tableau einer kontinuierlichen RKF:

3.4 Mehrschrittverfahren (MSV)

ESV z.B. RK-Verfahren:
$$u_{n+1} = u_n + h_n \sum_{j=1}^l b_j f(\underbrace{t_n + c_j h_n}, g_j), \quad n = \overline{0, m-1}$$

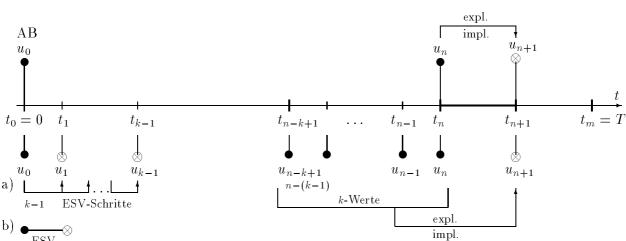
$$g = \Phi(g, t_n, u_n, h_n)$$

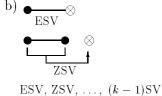
$$g = (g_1, \dots, g_l)^T$$

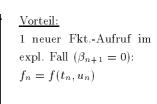
Nachteil:

Funktionsauswertung in den Zwischenpkt. $t_n + c_j h_n$ auch









MSV z.B. lineare MSV
$$(h_n = h)$$
:
$$\sum_{j=n-k+1}^{n+1} \alpha_j u_j = h \sum_{j=n-k+1}^{n+1} \beta_j f_j$$
$$f_j = f_j(t_j, u_j)$$

1. Bez. der Näherungen im MSV:
$$\underbrace{u_{n-k+1}, \dots, u_{n-1}, u_n}_{k-(k-1)} \longrightarrow u_{n+1}$$
$$\underbrace{u_{n-k+1}, \dots, u_{n-1}, u_n}_{u_{n+k-2}, u_{n+k-1}} \longrightarrow u_{n+k}$$

2. Im Startschritt:
$$u_0, \underbrace{u_1, \dots, u_{k-1}}_{\text{unbekannt}} u_k$$
 MSV:

- a) k-1 ESV-Schritte entsprechender (\triangleq MSV) Genauigkeit,
- b) ESV, ZSV, ..., (k-1)SV entsprechender (\triangleq MSV) Genauigkeit.

Spezielle Mehrschrittverfahren

■ Ausgangspkt.:

1) Dgl.:
$$\underline{u'(t)} = f(t, u(t))$$

Numerische Differentiation eines Interpolations polynoms $\Rightarrow t = t_{n+1}$

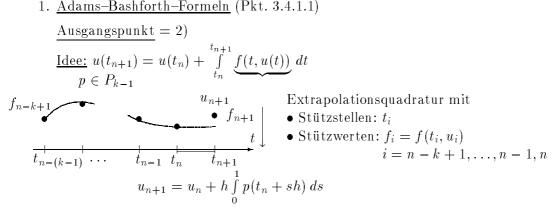
2) Igl.:
$$u(t) = u(t-s) + \int_{\frac{t-s}{2}}^{t} f(\tau, u(\tau)) d\tau$$

Numerische Integration beruhend auf Interpolation bzw. Extrapolation von f $\Rightarrow t = t_{n+1}$

■ Btr. <u>folgende</u> 5 Klassen von MSV für $h = h_n, n = \overline{0, m-1}$:

1. Adams-Bashforth-Formeln (Pkt. 3.4.1.1)

$$\underline{\underline{\text{Idee:}}} \ u(t_{n+1}) = u(t_n) + \int_{t_n}^{t_{n+1}} \underbrace{f(t, u(t))} dt$$



$$u_{n+1} = u_n + h \int_{0}^{1} p(t_n + sh) ds$$

2. Adams-Moulton-Formeln (Pkt. 3.4.1.2)

$$Ausgangspunkt = 2$$

$$\underline{\underline{\text{Idee: }} u(t_{n+1})} = u(t_n) + \int_{t_n}^{t_{n+1}} \underline{f(t, u(t))} dt$$

$$\underbrace{\frac{\operatorname{Idee:} u(t_{n+1}) = u(t_n) + \int\limits_{t_n}^{t_{n+1}} f(t,u(t))}_{p^* \in P_k} dt}_{f_{n-k+1}} \underbrace{f_{n-1} \int\limits_{t_n}^{t_{n+1}} f(t,u(t))}_{t_{n-k+1}} dt$$

$$\underbrace{f_{n-k+1} \int\limits_{t_{n-1}}^{t_{n-1}} f_n \int\limits_{t_{n+1}}^{t_{n+1}} f_{n+1}}_{t_{n+1}} \underbrace{f_{n-k+1} \int\limits_{t_{n-1}}^{t_{n+1}} f(t,u(t))}_{t_{n-1} \int\limits_{t_{n-1}}^{t_{n+1}} f(t,u(t))} dt$$

$$\underbrace{f_{n-k+1} \int\limits_{t_{n-1}}^{t_{n-1}} f_n \int\limits_{t_{n-1}}^{t_{n+1}} f(t,u(t))}_{t_{n-1} \int\limits_{t_{n-1}}^{t_{n-1}} f(t,u(t))} dt$$

$$\underbrace{f_{n-k+1} \int\limits_{t_{n-1}}^{t_{n-1}} f_n \int\limits_{t_{n-1}}^{t_{n-1}} f(t,u(t))}_{t_{n-1} \int\limits_{t_{n-1}}^{t_{n-1}} f(t,u(t))} dt$$

$$\underbrace{f_{n-k+1} \int\limits_{t_{n-1}}^{t_{n-1}} f(t,u(t))}_{t_{n-1} \int\limits_{t_{n-1}}^{t_{n-1}} f(t,u(t))}_{t_{n-1} \int\limits_{t_{n-1}}^{t_{n-1}} f(t,u(t))} dt$$

$$\underbrace{f_{n-k+1} \int\limits_{t_{n-1}}^{t_{n-1}} f(t,u(t))}_{t_{n-1} \int\limits_{t_{n-1}}^{t_{n-1}} f(t,u(t))}_{t_{n-1} \int\limits_{t_$$

$$i = n - k + 1, \dots, n, n + 1$$

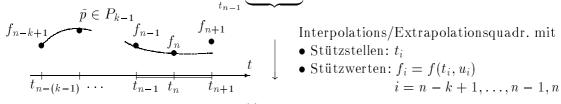
$$u_{n+1} = u_n + h \int_0^1 p^*(t_n + sh) ds$$

3. Nyström-Formeln (Pkt. 3.4.1.3)

$$Ausgangspunkt = 2$$

$$\underline{\underline{\text{Idee:}}} \ u(t_{n+1}) = u(t_{n-1}) + \int_{t_{n-1}}^{t_{n+1}} \underbrace{f(t, u(t))}_{dt} dt$$

$$\tilde{p} \in P_{k-1}$$

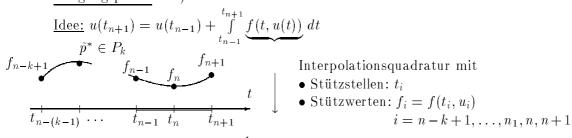


$$u_{n+1} = u_{n-1} + h \int_{-1}^{+1} \tilde{p}(t_n + sh) ds$$

4. Milne-Simpson-Formeln (Pkt. 3.4.1.4)

$Ausgan\underline{gspunkt} = 2)$

$$\underline{\underline{\text{Idee: }} u(t_{n+1}) = u(t_{n-1}) + \int_{t_{n-1}}^{t_{n+1}} \underbrace{f(t, u(t))}_{t} dt$$



$$i = n - k + 1, \dots, n_1, n,$$

$$u_{n+1} = u_{n-1} + h \int_{-1}^{1} \tilde{p}^*(t_n + sh) ds$$

5. BDF-Verfahren (Pkt. 3.4.1.5)

$$Ausgangspunkt = 1$$

$$\underline{\text{Idee:}}\ \underline{u'(t_{n+1})} = f(t_{n+1}, u(t_{n+1}))$$

Interpolation von $u(\cdot)$ mit

- Stützstellen: t_i Stützwerten: u_i

$$u_{n-1}$$
 u_{n-1}
 u_{n-1}
 u_{n-1}
 u_{n-1}
 u_{n-1}

$$i = n - k + 1, \dots, n - 1, n, n + 1$$

$$q'(t_{n+1}) = f_{n+1} := f(t_{n+1}, u_{n+1})$$

3.4.1.1 Die Adams-Bashforth-Formeln

■ Idee: Extrapolationsquadratur mit

• Stützstellen
$$t_i$$

• Stützwerten $f_i = f(t_i, u_i)$ $i = n - k + 1, \dots, n - 1, n$

aus Newton-Darstellung des
Interpolationspolynoms mit
Hilfe dividierter Differenzen =

$$u(t_{n+1}) = u(t_n) + \int_{t_n}^{t_{n+1}} \underbrace{f(t, u(t))}_{t} dt$$

$$p \in P_{k-1}$$

$$u_{n+1}$$

$$\downarrow f_{n-k+1}$$

$$\downarrow f_$$

■ Resultat: Adams-Bashforth-Formeln (k = 1, 2, ...)

$$u_{n+1} = u_n + h \sum_{j=0}^{k-1} \gamma_j \nabla^j f_n \text{ mit } \gamma_j = (-1)^j \int_0^1 \begin{pmatrix} -s \\ j \end{pmatrix} ds$$

$$j \quad 0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8$$

$$\gamma_j \quad 1 \quad \frac{1}{2} \quad \frac{5}{12} \quad \frac{3}{8} \quad \frac{251}{720} \quad \frac{95}{288} \quad \frac{19087}{60480} \quad \frac{5257}{17280} \quad \frac{1070017}{3628800}$$

■ Beispiele:

$$k = 1 : u_{n+1} = u_n + hf_n \Rightarrow \text{Explizite Euler-Methode},$$

$$k = 2 : u_{n+1} = u_n + h \left[\frac{3}{2} f_n - \frac{1}{2} f_{n-1} \right],$$

$$k = 3 : u_{n+1} = u_n + h \left[\frac{23}{12} f_n - \frac{16}{12} f_{n-1} + \frac{5}{12} f_{n-2} \right].$$

- 1. Adams-Bashforth-Formeln sind explizite MSV.
- 2. Konsistenzordnung = k (siehe Pkt. 3.4.2).
- 3. Stabilität: 0-stabil (siehe Pkt. 3.4.3).

3.4.1.2 Die Adams-Moulton-Formeln

■ Idee: Interpolationsquadratur mit

• Stützstellen
$$t_i$$

• Stützwerten $f_i = f(t_i, u_i)$ $i = n-k+1, \dots, n-1, n, n+1$

aus Newton-Darstellung des Interpolationspolynoms mit Hilfe dividierter Differenzen =

Hilfe dividierter Differenzen
$$\Rightarrow$$

$$u(t_{n+1}) = u(t_n) + \int_{t_n}^{t_{n+1}} \underbrace{f(t, u(t))}_{t_n} dt$$

$$p^* \in P_k$$

$$f_{n-k+1} \qquad f_{n-1} \qquad f_n \qquad f_{n+1}$$

$$t \qquad t_{n-(k-1)} \qquad t_{n-1} \qquad t_n \qquad t_{n+1}$$

$$u_{n+1} = u_n + h \int_{1}^{1} p^* (t_n + sh) ds$$

$$p^*(t_n + sh) = \sum_{j=0}^{k} (-1)^j \binom{-s+1}{j} \nabla^j f_{n+1}$$

$$mit \nabla f_{n+1} = f_{n+1} - f_n, \nabla^{j+1} = \nabla^j \nabla, \nabla^0 = I,$$

$$\binom{x}{j} = \frac{x(x-1) \cdot \ldots \cdot (x-j+1)}{1 \cdot 2 \cdot \ldots \cdot j} = \frac{x!}{(x-j)!j!}$$

■ Resultat: Adams-Moulton-Formeln (k = 0, 1, ...)

	u	$u_{n+1} = $	$= u_n +$	$h \sum_{j=0}^{k} \gamma_{j}$	$_{j}^{*}\nabla^{j}f_{n+1}$	$_{\scriptscriptstyle \parallel}$ mit γ_j^*	$\dot{z} = (-1)^j \int_0^t$	$\int_{0}^{1} \left(-s+1\right)^{s}$	$\int ds$
j	0	1	2	3	4	5	6	7	8
γ_j^*	1	$-\frac{1}{2}$	$-\frac{1}{12}$	$-\frac{1}{24}$	$-\frac{19}{720}$	$-\frac{3}{160}$	$-\frac{863}{60480}$	$-\frac{275}{24192}$	$\frac{33953}{3628800}$

■ Beispiele:

$$\begin{split} k &= 0: u_{n+1} = u_n + h f_{n+1} & \Rightarrow \text{ Implizite Euler-Methode}, \\ k &= 1: u_{n+1} = u_n + h \left[\frac{1}{2} f_{n+1} + \frac{1}{2} f_n \right] & \Rightarrow \text{ CRANK-NICOLSON}, \\ k &= 2: u_{n+1} = u_n + h \left[\frac{5}{12} f_{n+1} + \frac{8}{12} f_n - \frac{1}{12} f_{n-1} \right]. \end{split}$$

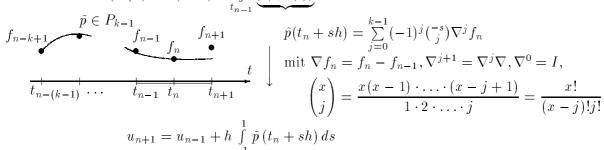
- 1. Adams-Moulton-Formeln sind implizite MSV: Für hinr. kleines h ist Auflösbarkeit nach u_{n+1} durch Fixpunktiteration garantiert (f = Lipschitz-stetig)!
- 2. Konsistenzordnung = k + 1 (siehe Pkt. 3.4.2).
- 3. Stabilität: 0-stabil (siehe Pkt. 3.4.3).

3.4.1.3 Die Nyström-Formeln

- <u>Idee:</u> Interpolations/Extrapolationsquadratur mit
 - Stützstellen t_i • Stützwerten $f_i = f(t_i, u_i)$ $\}$ $i = n - k + 1, \dots, n - 1, n$

aus Newton-Darstellung des Interpolationspolynoms mit Hilfe dividierter Differenzen ⇒

$$u(t_{n+1}) = u(t_{n-1}) + \int_{t_{n-1}}^{t_{n+1}} \underbrace{f(t, u(t))}_{t} dt$$



■ Resultat: Nyström-Formeln (k = 1, 2, ...)

u_{n+}	₋₁ =	$= u_r$	ı−1 ·	+h	$\sum_{j=0}^{k-1} \kappa$	$_{j} abla^{j}f$	$_n$ mit $_{ extit{ heta}}$	$\kappa_j = (-1)$	$(1)^{j} \int_{-1}^{+1} \begin{pmatrix} -s \\ j \end{pmatrix} ds$
j	0	1	2	3	4	5	6	7	8
κ_j	2	0	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{29}{90}$	$\frac{14}{45}$	$\frac{1139}{3780}$		

■ Beispiele:

$$k = 1 : u_{n+1} = u_{n-1} + 2hf_n \ (\widehat{=} k = 2) \Rightarrow \text{Explizite Mittelpunktsregel},$$

 $k = 3 : u_{n+1} = u_{n-1} + h\left[\frac{7}{3}f_n - \frac{2}{3}f_{n-1} + \frac{1}{3}f_{n-2}\right].$

- 1. Nyström-Formeln sind explizite MSV!
- 2. Konsistenzordnung = k (siehe Pkt. 3.4.2).
- 3. Stabilität: 0-stabil (siehe Pkt. 3.4.3).

3.4.1.4 Die Milne-Simpson-Formeln

- Idee: Interpolationsquadratur mit
 - Stützstellen t_i • Stützwerten $f_i = f(t_i, u_i)$ $i = n-k+1, \dots, n-1, n, n+1$

aus Newton-Darstellung des Interpolationspolynoms mit Hilfe dividierter Differenzen ⇒

Hilfe dividierter Differenzen
$$\Rightarrow$$

$$u(t_{n+1}) = u(t_{n-1}) + \int_{t_{n-1}}^{t_{n+1}} \underbrace{f(t, u(t))}_{t_{n-1}} dt$$

$$f_{n-k+1} \qquad f_{n-1} \qquad f_{n} \qquad f_{n+1} \qquad p^*(t_n + sh) = \sum_{j=0}^k (-1)^j \binom{-s+1}{j} \nabla^j f_{n+1}$$

$$mit \nabla f_{n+1} = f_{n+1} - f_n, \nabla^{j+1} = \nabla^j \nabla, \nabla^0 = I,$$

$$\binom{x}{j} = \frac{x(x-1) \cdot \ldots \cdot (x-j+1)}{1 \cdot 2 \cdot \ldots \cdot j} = \frac{x!}{(x-j)!j!}$$

$$u_{n+1} = u_{n-1} + h \int_{-1}^{1} p^* (t_n + sh) ds$$

■ Resultat: Milne-Simpson-Formeln (k = 0, 1, 2, ...)

u_n .	₊₁ =	= <i>u</i> _n _	1 +	$h\sum_{j=1}^{k}$	$\sum_{j=0}^{k} \kappa_{j}^{*} \nabla^{j}$	f_{n+1} r	$\text{mit } \kappa_j^* =$	$(-1)^j \int_{-1}^+$	$\binom{1}{s}\binom{-s+1}{j}ds$
j	0	1	2	3	4	5	6	7	8
κ_j^*	2	-2	$\frac{1}{3}$	0	$-\frac{1}{90}$	$-\frac{1}{90}$	$-\frac{37}{3780}$		

■ Beispiele:

 $k=0:u_{n+1}=u_{n-1}+2hf_{n+1} \Rightarrow \text{Impliziter Euler mit doppeltem Schritt},$

 $k = 1 : u_{n+1} = u_{n-1} + 2hf_n \Rightarrow \text{Explizite Mittelpunktsregel } (2h),$

$$k = 2: u_{n+1} = u_{n-1} + h\left[\frac{1}{3}f_{n+1} + \frac{4}{3}f_n + \frac{1}{3}f_{n-1}\right] \Rightarrow \text{Simpson-Regel.}$$

- 1. Milne-Simpson-Formeln sind <u>implizite</u> MSV! \longrightarrow Für hinr. kleines h ist Auflösbarkeit nach u_{n+1} durch Fixpunktiteration garantiert: Startnäherung $u_{n+1}^{(0)} := u_n$ bzw. $u_{n+1}^{(0)} := \text{Nyström-N\"{a}herung}.$
- 2. Konsistenzordnung = k + 1 (siehe Pkt. 3.4.2).
- 3. Stabilität: 0-stabil (siehe Pkt. 3.4.3).

3.4.1.5Die BDF-Verfahren (Backward Differencing Formula)

- Idee: Numerische Differentiation des Interpolationspolynoms mit

 - Stützstellen t_i Stützwerten u_i $i = n k + 1, \dots, n 1, n, n + 1$

aus Newton-Darstellung des Interpolationspolynoms mit Hilfe dividierter Differenzen ⇒

$$q \in P_k$$

$$u_{n-k+1}$$

$$u_{n-1}$$

$$u_n$$

$$t$$

$$t_{n-k+1}$$

$$t_{n-k+1}$$

 $q'(t_{n+1}) = f_{n+1} := f(t_{n+1}, u_{n+1})$

$$\underbrace{u'(t_{n+1})}_{u_{n-k+1}} = f(t_{n+1}, u(t_{n+1}))$$
Hilfe dividierter Differenzen \Rightarrow

$$q(t_n + sh) = \sum_{j=0}^k (-1)^j {\binom{-s+1}{j}} \nabla^j u_{n+1}$$

$$\min \nabla u_{n+1} = u_{n+1} - u_n, \nabla^{j+1} = \nabla^j \nabla, \nabla^0 = I,$$

$$\binom{x}{j} = \frac{x(x-1) \cdot \dots \cdot (x-j+1)}{1 \cdot 2 \cdot \dots \cdot j} = \frac{x!}{(x-j)!j!}$$

Resultat: BDF-Verfahren $(k = \emptyset, 1, 2, ...)$

$$\frac{\sum_{j=0}^{k} \delta_{j}^{*} \nabla_{j} u_{n+1} = h f_{n+1} \text{ mit } \delta_{j}^{*} = (-1)^{j} \frac{d}{ds} \binom{-s+1}{j} \Big|_{s=1}}{\delta_{0}^{*} = 0, \quad \delta_{j}^{*} = \frac{1}{j} \text{ für } j \ge 1}$$

■ Beispiele:

 $k = 1 : u_{n+1} - u_n = h f_{n+1} \implies \text{Impliziter Euler},$

$$k = 2 : \frac{3}{2}u_{n+1} - 2u_n + \frac{1}{2}u_{n-1} = hf_{n+1},$$

$$k = 3 : \frac{11}{6}u_{n+1} - 3u_n + \frac{3}{2}u_{n-1} - \frac{1}{3}u_{n-2} = hf_{n+1}.$$

- 1. BDF-Methoden sind implizite MSV!
- 2. Konsistenzordnung = k (siehe Pkt. 3.4.2).
- 3. Stabilität: 0-stabil für $k \le 6$, nicht 0-stabil für $k \ge 7$ (siehe Pkt. 3.4.3).

3.4.2 Konsistenz linearer MSV

■ Lineare MSV:

Alle bisher betrachteten MSV sind von der Form

(53)
$$\begin{cases} n+1 & n-k+1 \\ \uparrow & \uparrow & \uparrow \\ \alpha_k \mathbf{u_{n+k}} + \alpha_{k-1} u_{n+k-1} + \ldots + \alpha_0 u_n = h \sum_{j=0}^k \beta_j f_{n+j}, \ n = 0, 1, \ldots, m-k, \\ \text{mit idealisierter Startphase (o. B. d. Allg.):} & \downarrow \\ u_j = u(t_j), \ j = 0, 1, \ldots, k-1. & f(t_{n+j}, u_{n+j}) \end{cases}$$

MSV der Form (53) heißen <u>lineare MSV</u> (k-Schrittverfahren).

Durch eine Normierungsbedingung

(54)
$$\alpha_k = 1$$
 oder $\sum_{j=0}^k \beta_j = 1$

wird das MSV eindeutig festgelegt. Betrachten im folgenden nur lineare MSV!

■ Lokaler (Diskretisierungs-)Fehler, Konsistenz, Konsistenzordnung:

(55)
$$u(s) : u'(s) = f(s, u(s)), s \ge t$$

 $u(t) \text{ geg.}$

(56)
$$u_h(t+kh): \sum_{j=0}^k \alpha_j u_h(t+jh) = h \sum_{j=0}^k \beta_j f(t+jh, u_h(t+jh))$$

 $u_h(t+jh) = u(t+jh), \quad j = 0, 1, \dots, k-1$

• **Definition 3.27:** (Konsistenz von MSV)

Ein MSV heißt konsistent, falls
$$u(t+kh)-u_h(t+kh)\longrightarrow 0 \text{ für } h\to 0 \quad \forall t.$$
 Falls
$$u(t+kh)-u_h(t+kh)=O(h^{p+1}) \text{ für } h\to 0 \quad \forall t,$$
 dann ist die Konsistenzordnung $p.$

Dabei wird $u(\cdot)$ immer als hinreichend glatt vorausgesetzt.

• Lemma 3.28: (Darstellung des lokalen Fehlers)

 $\underline{\text{Vor.:}}\ f$ sei stetig differenzierbar.

Bh.: Dann $\exists \delta \in (0,1)$:

(57)
$$u(t+kh) - u_h(t+kh) = [\alpha_k I - h\beta_k f_u(t+kh,\nu)]^{-1} L(u;t,h),$$

$$\text{mit } \nu = \nu(\delta) = u_h(t+kh) + \delta[u(t+kh) - u_h(t+kh)],$$

(58)
$$L(u;t,h) := \sum_{j=0}^{k} [\alpha_j u(t+jh) - h\beta_j u'(t+jh)].$$

Im Fall von Systemen ist (57) komponentenweise zu verstehen mit $\delta_i \in (0,1)$ und $\nu_i = \nu(\delta_i), i = \overline{1,N}$.

Beweis:

Aus (56) folgt sofort die Beziehung

(59)
$$\sum_{j=0}^{k-1} [\alpha_{j} u(t+jh) - h\beta_{j} \underbrace{f(t+jh, u(t+jh))}_{u'(t+jh)}] + \alpha_{k} u_{h}(t+kh) - h\beta_{k} f(t+kh, u_{h}(t+kh)) = 0.$$

Aus (58) und (59) ergibt sich dann

$$L(u, t, h) = \alpha_k u(t + kh) - h\beta_k f(t + kh, u(t + kh)) + \sum_{j=0}^{k-1} [\alpha_j u(t + jh) - h\beta_j u'(t + jh)] = \underbrace{\sum_{j=0}^{(59)} [\alpha_k u_h(t + kh) - h\beta_k f(t + kh, u_h(t + kh))]}_{(59)}$$

$$\stackrel{(59)}{=} -[\alpha_k u_h (t+kh) - h\beta_k f(t+kh, u_h (t+kh))]$$

$$= \alpha_k [u(t+kh) - u_h(t+kh)] - h\beta_k [f(t+kh, u(t+kh)) - f(t+kh, u_h(t+kh))]$$

 $\stackrel{\text{[16] (Satz I.5.5)}}{=} \left[\alpha_k I - h\beta_k f_u(t+kh,\nu)\right] \left(u(t+kh) - u_h(t+kh)\right).$

• Bemerkung: Für explizite MSV (
$$\beta_k = 0$$
) und Normierungsbed. ($\alpha_k = 1$): $\Rightarrow u(t+kh) - u_h(t+kh) = L(u,t,h)$, sonst $= O(L(u;t,h))$.

■ Näherungsgleichung und lokaler Abschneidefehler:

• Schreiben MSV (53) als Operatorgleichung in der Form

(60) Ges.
$$u_h \in X_h : F_h(u_h) = 0 \text{ in } Y_h$$

mit
$$X_h = C^1(I_h)$$
, $Y_h = C(I_h)$ und

(61)
$$\begin{cases} F_{h}(v_{h})(t_{i}+kh) := \frac{1}{h} \sum_{j=0}^{k} \alpha_{j} v_{h}(t_{i}+jh) - \sum_{j=0}^{k} \beta_{j} f(t_{i}+jh, v_{h}(t_{i}+jh)), \\ i = \overline{0, m-k} \\ \text{und idealisierter Startphase} \\ F_{h}(v_{h})(t_{0}+jh) := v_{h}(t_{0}+jh) - u(t_{0}+jh), \quad j = 0, 1, \dots, k-1. \end{cases}$$

Dabei ist es jetzt zweckmäßig, (o. B. d. Allg.) von der Normierungsbed.

(62)
$$\sum_{j=0}^{k} b_j = 1 \quad \text{auszugehen.}$$

- Für den <u>lokalen Abschneidefehler</u>
 - $\tau_h := F_h(u)$, wobei u Lsg. von (1) ist, (63)

gilt offenbar die Beziehung

Der damit verbundene Begriff der Konsistenzordnung p

(65)
$$\|\tau_h\|_{Y_h=C(I_h)} = O(h^p)$$

entspricht wegen Lemma 3.28 genau dem in Def. 3.27 eingeführten Begriff der Konsistenzordnung über die Kleinheit des lokalen Fehlers.

■ Zur Bestimmung der Konsistenzordnung linearer MSV:

• Ordnen den Koeffizienten eines linearen MSV Polynome zu:

(66)
$$\rho(z) = \alpha_k z^k + \alpha_{k-1} z^{k-1} + \ldots + \alpha_0; \quad \text{(charak. Polynom)}$$
$$\sigma(z) = \beta_k z^k + \beta_{k-1} z^{k-1} + \ldots + \beta_0.$$

• Satz 3.29:

Vor.: f sei hinreichend oft stetig differenzierbar.

Bh.: Ein MSV der Form (53) besitzt genau dann die

Konsistenzordnung
$$p$$
, falls
$$\begin{cases} \sum_{j=0}^{k} \alpha_j = 0, \\ \sum_{j=0}^{k} \alpha_j j^q = q \sum_{j=0}^{k} \beta_j j^{q-1}, \quad q = 1, 2, \dots, p. \end{cases}$$

Beweis:

$$L(u,t,h) = \sum_{j=0}^{k} \left[\alpha_{j} \underbrace{u(t+jh)}_{-} - h \beta_{j} \underbrace{u'(t+jh)}_{-} \right]^{\text{Taylor}} \stackrel{\text{Taylor}}{=}$$

$$\sum_{q=0}^{p} \frac{j^{q}}{q!} u^{(q)}(t) h^{q} + O(h^{p+1}) \qquad \sum_{r=0}^{p-1} \frac{j^{r}}{r!} u^{(r+1)}(t) h^{r} + O(h^{p})$$

$$= u(t) \sum_{j=0}^{k} \alpha_{j} + \sum_{q=1}^{p} \frac{u^{(q)}(t)}{q!} h^{q} \sum_{j=0}^{k} \alpha_{j} j^{q} - \sum_{q\equiv r+1=1}^{p} \frac{u^{(q)}(t)q}{(q-1)!q} h^{q} \sum_{j=0}^{k} \beta_{j} j^{q-1} + O(h^{p+1}) =$$

$$= u(t) \sum_{j=0}^{k} \alpha_{j} + \sum_{q=1}^{p} \frac{h^{q}}{q!} u^{(q)}(t) \underbrace{\left[\sum_{j=0}^{k} \alpha_{j} j^{q} - q \sum_{j=0}^{k} \beta_{j} j^{q-1}\right]}_{=0 \text{ für } q=1,2,...,p} + O(h^{p+1}) =$$

$$= O(h^{p+1})$$

$$\mathbf{q.e.d.}$$

• Bemerkung:

Für den Fall p = 1 (Konsistenzordnung 1) lautet (67):

(68)
$$\begin{cases} \rho(1) = 0 & \left[\sum_{j=0}^{k} \alpha_j = 0\right], \\ \rho'(1) = \sigma(1) & \left[\sum_{j=0}^{k} j\alpha_j = \sum_{j=0}^{k} \beta_j\right]. \end{cases}$$

• Folgerung 3.30:

Die höchste erreichbare Konsistenzordnung eines k-Schrittverfahrens = 2k.

Beweis: mms.

• Folgerung 3.31: (Konsistenzordnung der MSV aus Pkt. 3.4.1)

MSV (k-Schritt)	Konsistenzordnung	MSV ist exakt für $u'(t) = qt^{q-1}$
Adams-Bashforth Adams-Moulton	$ \begin{vmatrix} k \\ k+1 \end{vmatrix}$	$q = 0, 1, \dots, k$ $q = 0, 1, \dots, k + 1$
Nyström	k	$q = 0, 1, \dots, k$
Milne–Simpson BDF	$\begin{vmatrix} k+1\\k \end{vmatrix}$	$q = 0, 1, \dots, k + 1$ $q = 0, 1, \dots, k$

<u>Beweis:</u> Unter Benutzung von (57) aus Lemma 3.28 und der Exaktheit des MSV für $u'(t) = qt^{q-1}$, d.h. für $u(t) = t^q + c$ erhalten wir sofort die Identität

temma 3.28, (57) $0 \stackrel{\downarrow}{=} [\alpha_{k}I - 0](u(t + kh) - u_{h}(t + kh)) = L(t^{q} + c, t, h) =$ $= \sum_{j=0}^{k} [\alpha_{j}[(t + jh)^{q} + c] - h\beta_{j}q(t + jh)^{q-1}] =$ $= c \sum_{j=0}^{k} \alpha_{j} + \sum_{j=0}^{k} [\alpha_{j}(t + jh)^{q} - h\beta_{j}q(t + jh)^{q-1}]$ $= h^{q} \sum_{j=0}^{k} [\alpha_{j}j^{q} - q\beta_{j}j^{q-1}]$ $\uparrow \qquad \qquad \downarrow = 0$ q.e.d.

3.4.3 Stabilität linearer MSV

■ Im Gegensatz zu den RK-Formeln (Satz 3.25: f – Lipschitz-stetig \Rightarrow Stabilität) ist die Untersuchung der Stabilität für MSV komplizierter!

$$\ddot{\mathbf{U}}$$
 3.10 Man zeige zunächst, daß das explizite Zweischrittverfahren $(k=2)$

$$(*) \qquad u_{n+2} + 4u_{n+1} - 5u_n = h(4f_{n+1} + 2f_n)$$

die höchstmögliche Konsistenzordnung p=3 hat!

Wendet man dieses Verfahren auf das AWP

$$u'(t) = f(t, u(t)) := u(t), \quad t \ge 0; \text{ Lsg. sei } \underline{u(t) = e^t}$$

mit der exakten Startphase $\underline{u(0)} = 1$ und $\underline{u(h)} = e^h$ an, so erhält man völlig <u>unbrauchbare</u> Resultate, auch und insbesondere für $h \to 0$, obwohl f(t,u) := u Lipschitz-stetig mit der Lipschitz-Konstanten L = 1 ist!

Bestätigen Sie diese Behauptung durch numerische Experimente, und begründen Sie die Behauptung theoretisch!

<u>Hinweis:</u> • <u>theoretisch:</u> (*) mit Ansatz: $u_j = \zeta^j$ $\Rightarrow \text{Charakteristisches Polynom: } \zeta^2 + 4(1-h)\zeta - (5+2h) = 0$ $\Rightarrow \text{Wurzeln: } \zeta_1(h) = 1 + h + O(h^2), \quad \zeta_2(h) = -5 + O(h)$ $u_n = a_1 \zeta_1^n(h) + a_2 \zeta_2^n(h)$

■ 0-Stabilität von MSV:

- Untersuchen Stabilität von linearen MSV zunächst für die triviale RS $f \equiv 0$ (" $h \rightarrow 0$ "):
 - (69) $\alpha_k u_{n+k} + \alpha_{k-1} u_{n+k-1} + \ldots + \alpha_0 u_n = 0$

als Näherungsgleichung für die triviale Dgl.

- (70) u'(t) = 0.
- Der nächste Satz liefert die allgem. Lsg. der Differenzengleichung (69):

Satz 3.32:

<u>Vor.:</u> Seien $\zeta_1, \zeta_2, \ldots, \zeta_l$ die Wurzeln von $\rho(z) = \sum\limits_{j=0}^k \alpha_j z^j$ mit den Vielfachheiten m_1, m_2, \ldots, m_l $(\sum\limits_{j=0}^l m_j = k, \ l \leq k).$

<u>Bh.:</u> Dann ist die allgemeine Lösung von (69) durch

(71)
$$u_n = p_1(n)\zeta_1^n + p_2(n)\zeta_2^n + \ldots + p_l(n)\zeta_l^n$$
gegeben, wobei $p_j(\cdot)$ bel. Polynome vom Grad $m_j - 1$ sind.

Beweis: siehe auch [5] S. 328:

- o l=k, d.h. $\rho(z)$ hat nur einfache Nullstellen $(m_j=1\ \forall j)$ $\Rightarrow u_n=a_1\zeta_1^n+\ldots+a_k\zeta_k^n$ ist offenbar allgem. Lsg. von (69)!
- o l < k, d.h. $\rho(z)$ hat mehrfache Nullstellen $(\exists m_j > 1)$:
 - 1) Sei ζ eine Nullstelle der Vielfachheit m. Dann sind $u_n = \binom{n}{i} \zeta^n$ für jedes $i = 0, 1, \ldots, m-1$ spezielle Lsg. von (71).

Tatsächlich, wegen der Identität $\binom{n+j}{i} = \sum_{l=0}^{i} \binom{n}{i-l} \binom{j}{l}$ (mms) gilt:

$$\sum_{j=0}^{k} \alpha_{j} u_{n+j} = \sum_{j=0}^{k} \alpha_{j} \binom{n+j}{i} \zeta^{n+j} = \zeta^{n} \sum_{l=0}^{i} \binom{n}{i-l} \sum_{j=0}^{k} \alpha_{j} \binom{j}{l} \zeta^{j} =$$

$$= \zeta^{n} \sum_{l=0}^{i} \binom{n}{i-l} l! \zeta^{l} p^{(l)}(\zeta) = 0$$

2) Superposition:
$$u_n = \sum_{j=1}^l \sum_{i=0}^{m_j-1} a_{ij} \binom{n}{i} \zeta_j^n = \sum_{j=1}^l p_j(n) \zeta_j^n$$
 mit frei wählbaren Koeffizienten $\{a_{ij}\}_{j=\overline{1,l};\ i=\overline{0,m_j-1}}$.

• Durch die Vorgabe der k Startwerte

$$u_0, u_1, \ldots, u_{k-1}$$

wird die Lsg. (71) von (69) eindeutig definiert.

• Damit Störungen δ_i in der Startphase, d.h. $u_i + \delta_i$, $i = \overline{0, k-1}$, nicht zu unbeschränkten Lösungen führen, muß für jede Nullstelle ζ von $\rho(\cdot)$ offenbar gelten:

(72)
$$\begin{cases} \bullet & |\zeta| \leq 1 \text{, falls } \zeta - \text{ einfache Nullstelle,} \\ \bullet & |\zeta| < 1 \text{, falls } \zeta - \text{ mehrfache Nullstelle.} \end{cases}$$

Das führt auf die folgende Definition:

Definition 3.33: (0-Stabilität)

Das MSV (53) heißt <u>0-stabil</u>, falls jede Nullstelle ζ des charakteristischen Polynoms $\rho(\cdot)$ die Bedingung (72) erfüllt.

- Beispiele:
 - 1. <u>Adams–Bashforth und Adams–Moulton:</u> <u>0</u>—stabil, da $\rho(z) = z^k z^{k-1} \Rightarrow \zeta_1 = 0 \ (k-1)$ —fach; $\zeta_2 = 1$ (einfach).
 - 2. Nyström und Milne–Simpson–Formeln: 0–stabil, da $\rho(z)=z^k-z^{k-2}\Rightarrow \zeta_1=0\;(k-2)\text{-fach}\,;\,\zeta_{2,3}=\pm 1\;(\text{einfach}).$
 - 3. <u>BDF-Formeln:</u> $k \leq 6 \Rightarrow 0$ -stabil; $k \geq 7 \Rightarrow \underline{\text{nicht}} 0$ -stabil! Analyse komplizierter: $\sum_{j=1}^{k} \frac{1}{j} \nabla_j u_{n+1} \mapsto \rho(z) = \dots$? siehe [5] S. 329 ff.
- Die erste Dahlquist-Barriere (1956):

Für die höchste, erreichbare Ordnung p eines 0-stabilen k-Schrittverfahrens gilt:

$$(73) p \leq \begin{cases} k+2, & \text{falls } k-\text{ungerade}, \\ k+1, & \text{falls } k-\text{gerade}, \\ k, & \text{falls } \beta_k/\alpha_k \leq 0 \text{ (also insbes. für } \underline{\text{explizite MSV}}). \end{cases}$$

Beweis: siehe [5] S. 332 ff.

3.4.4 Konvergenz linearer MSV

■ Idee: [Butcher, 1966]

■ Schreiben MSV in $I\!\!R^N$ als ESV in $(I\!\!R^N)^k$:

 \bullet Sei $f(\cdot,\cdot)$ stetig und im 2. Arg. Lipschitz–stetig, d.h.

$$(74) |f(t,v) - f(t,w)| \le L|v - w|, \quad \forall t \in I, \quad \forall v, w \in \mathbb{R}^N.$$

Dann gibt es für bel. (fix. !) Werte

$$t_n, u_n, u_{n+1}, \ldots, u_{n+k-1}$$

und <u>hinreichend kleinen h</u> eine eindeutig bestimmte Lösung

$$u_{n+k} = g(t_n; u_n, u_{n+1}, \dots, u_{n+k-1}; h)$$

des MSV(53)

$$\sum_{j=0}^{k-1} \alpha'_{j} u_{n+j} + \mathbf{u}_{n+k} = \underbrace{h \sum_{j=0}^{k-1} \beta'_{j} f_{n+j} + h \beta'_{k} f(t_{n} + kh, \mathbf{u}_{n+k})}_{h \psi};$$

$$\alpha'_j = \frac{\alpha_j}{\alpha_k}, \quad \beta'_j = \frac{\beta_j}{\beta_k}.$$

<u>Beweis:</u> Banachscher Fixpunktsatz angewandt auf

$$\mathbf{u}_{n+k} = h\beta_k' f(t_n + kh, \mathbf{u}_{n+k}) + \text{Rest !} \quad (\text{mms})$$

Damit läßt sich die Funktion ψ definieren:

$$\psi(t_n; u_n, u_{n+1}, \dots, u_{n+k-1}; h) := \sum_{j=0}^{k-1} \beta'_j f_{n+j} + \beta'_k f(t_n + kh, \underbrace{g(t_n; u_n, u_{n+1}, \dots, u_{n+k-1}; h)}_{=u_{n+k}})$$

Damit erhält das MSV (53) die Form:

(75)
$$u_{n+k} = -\sum_{j=0}^{k-1} \alpha'_j u_{n+j} + h \psi(t_n; \underbrace{u_n, u_{n+1}, \dots, u_{n+k-1}}_{=:U_n}; h).$$

Definieren

$$U_n = \begin{bmatrix} u_{n+k-1} \\ u_{n+k-2} \\ \vdots \\ u_n \end{bmatrix}, A = \begin{bmatrix} -\alpha'_{k-1}I & -\alpha'_{k-2}I & \dots & -\alpha'_0I \\ I & \mathbf{O} & \dots & \mathbf{O} \\ & \ddots & \ddots & \\ & & I & \mathbf{O} \end{bmatrix}, e_1 = \begin{bmatrix} I \\ \mathbf{O} \\ \vdots \\ \mathbf{O} \end{bmatrix}$$

und können damit das MSV (75) \equiv (53) als $\underline{\text{ESV}}$

(76)
$$U_{n+1} = A U_n + h \Phi(t_n, U_n, h), \quad n = \overline{0, m - k}$$
AB: U_0 geg. (ESV aus Pkt. 3.3: $A = I$!)

im Produktraum $(I\!\!R^N)^k = I\!\!R^N \times \ldots \times I\!\!R^N$ schreiben, wobei

$$\Phi(t, U, h) = \psi(t, U, h) e_1,$$

$$\mathbf{X}_h = C^1(\tilde{I}_h, (\mathbb{R}^N)^k), \quad \mathbf{Y}_h = C(\tilde{I}_h, (\mathbb{R}^N)^k), \quad \tilde{I}_h = \{t_0, \dots, t_{m-k+1}\}.$$

■ Konsistenz: MSV(53) = (75) $\hat{=}$ ESV(76)

Satz 3.34:

Aus der Konsistenz(-ordnung) des MSV (53) = (75) gemäß Def. 3.27 folgt die Konsistenz(-ordnung) des zum MSV äquivalenten ESV (76) und umgekehrt.

Beweis:

- Sei u(s) die exakte Lösung des AWP (1) und (77) $U(s) = [u(s+(k-1)h), u(s+(k-2)h), \dots, u(s)]^T$
- Die Näherungslösung an der Stelle t+h, die man aus dem ESV (76) mit dem Startwert U(t) gewinnt, sei $U_h(t+h)$!

• Dann gilt offenbar:

$$U(t+h) - U_h(t+h) = \begin{bmatrix} u(t+kh) \\ u(t+(k-1)h) \\ \vdots \\ u(t+h) \end{bmatrix} - \begin{bmatrix} u_h(t+kh) \\ u(t+(k-1)h) \\ \vdots \\ u(t+h) \end{bmatrix} = \begin{bmatrix} u(t+kh) - u_h(t+kh) \\ \vdots \\ u(t+h) \end{bmatrix}$$

Daraus folgt:

(78)
$$||U(t+h) - U_h(t+h)||_{\mathbf{X}_h} = ||u(t+kh) - u_h(t+kh)||_{X_h} = O(h^{p+1}).$$

■ Stabilität: 0-Stabilität des MSV + L-Stetigkeit von $f \Rightarrow$ Stabilität v. (76)

• <u>Lemma 3.35:</u>

Vor.: Das MSV sei 0-stabil (Def. 3.33).

Bh.: Dann gibt es eine Vektornorm in $(\mathbb{R}^N)^k$, sodaß für die zugeordnete Matrixnorm gilt:

$$||A|| \le 1.$$

Beweis: siehe [5] S. 343 f.

- Bemerkung: Lineares MSV ist 0-stabil, genau dann, wenn $\exists c = \text{const.} \geq 0$: $||A^n|| \leq c \quad \forall n \in \mathbb{N}$.
- Satz 3.36: (Stabilität des ESV (76) i. S. der Def. 3.21)

Vor.: 1) MSV
$$(53) = (75)$$
 sei 0-stabil.

2)
$$f(\cdot,\cdot)$$
 erfülle die Lipschitz-Bedingung
$$|f(t,v)-f(t,w)| \leq L|v-w| \quad \forall v,w \in I\!\!R^N \quad \forall t \in I.$$

Bh.: Dann ist das zum MSV zugeordnete ESV (76) stabil i. S. der Def. 3.21.

Beweisskizze:

o Aus der Lipschitz-Stetigkeit von $f(\cdot, \cdot)$ folgt (mms):

$$\|\Phi(t,V,h)-\Phi(t,W,h)\|\leq L^*\|V-W\|\quad\forall V,W\in(I\!\!R^N)^k\ \forall t,\ \forall h\leq H.$$

o Unter Benutzung von <u>Lemma 3.35</u> ($\Rightarrow ||A|| \le 1$) zeigen wir zunächst Stabilität bzgl. der AB (Startphase) analog zu <u>Lemma 3.18</u>:

$$Z_{j} = U_{j}^{(77)} - V_{j}^{(76)}$$

$$Z_{n+1} = AZ_{n} + h[\Phi(t_{n}, U_{n}, h) - \Phi(t_{n}, V_{n}, h)]$$

$$\|Z_{n+1}\| \leq \|A\| \|Z_{n}\| + L^{*}h \|Z_{n}\| \leq (1 + L^{*}h) \|Z_{n}\|$$

$$\leq e^{Lt_{n+1}}$$
usw.

o Daraus folgt analog zu <u>Lemma 3.19</u> und <u>Satz 3.20</u> die Stabilität:

$$||V_h - W_h||_{\mathbf{X}_h} \le C ||F_h(V_h) - F_h(W_h)||_{\mathbf{Y}_h}.$$

q.e.d.

■ Konvergenz: Stabilität + Konsistenz = Konvergenz

• Satz 3.37: (Konvergenz des ESV (76))

Vor.: 1) MSV habe Konsistenzordnung p.

- 2) f sei entsprechend oft differenzierbar.
- 3) f sei Lipschitz-stetig im 2. Arg.
- 4) MSV sei 0-stabil.

Bh.: Dann ist das ESV (76) konvergent von der Ordnung p.

<u>Beweis:</u> analog zum Beweis von <u>Satz 3.23</u> unter Benutzung der Sätze 3.34 und 3.36.

q.e.d.

• Folgerung 3.38: (Konvergenz des MSV)

Unter den Voraussetzungen von Satz 3.37 ist das MSV (53) konvergent von der Ordnung p.

Beweis: MSV (53) \equiv (75) = ESV (76). q.e.d.

3.5 Steife Differentialgleichung

■ Btr. Beispiel E 6 aus P IX :

(80)
$$\begin{cases} u'(t) = -50(u(t) - \cos(t)), & t \in [0, T] \\ \underline{AB}: & u(0) = 0 \end{cases}$$

Das <u>explizite</u> und das <u>implizite</u> Euler-Verfahren verhalten sich völlig unterschiedlich (siehe P IX):

- expliziter Euler: sinnvolle Resultate erst für $h = \frac{1}{40}, \frac{1}{30} \le \frac{2}{50}$!
- ullet impliziter Euler: keine Probleme auch für größere h!,

obwohl beide Verfahren

- Konsistenzordnung 1 haben und
- stabil im Sinne der Def. 3.21 sind.

Offenbar reicht der bisherige Stabilitätsbegriff nicht aus, um das Verhalten der Verfahren für größere ("endliche") Schrittweiten h zur Lösung steifer AWP (1), d.h. große Lipschitz-Konstante L bzw. große Kondition der Jacobi-Matrix $J=f_u(t,u(t))$, zu erklären:

- Lemma 3.18: $(1 + \stackrel{?}{L}h_j) \le e^{Lh_j} \Rightarrow e^{Lh_j} \cdot \ldots \cdot e^{Lh_0} = e^{Lt_{j+1}} \le e^{LT}$,
- Beweis Satz 3.25: h hinreichend (?) klein.
- Zur Untersuchung der Stabilitätsverhalten eines Verfahrens mit "endlicher" Schrittweite wird das Verhalten des Verfahrens für das $\underline{\text{Test-}}$ problem (81) mit $\lambda \in \mathcal{C}$ betrachtet:

$$\lambda = \lambda_1 + i\lambda_2, \, \lambda_1 = \operatorname{Re} \lambda, \, \lambda_2 = \operatorname{Im} \lambda$$

$$\downarrow \qquad \qquad \downarrow$$

$$z(t) = \tilde{u}(t) - u(t) = \delta e^{\lambda t} = \delta e^{\operatorname{Re} \lambda \cdot t} e^{i \cdot \operatorname{Im} \lambda \cdot t}$$

$$\Rightarrow \qquad \qquad |z(t)| = \delta e^{\operatorname{Re} \lambda \cdot t} \leq \delta, \quad \text{falls } \operatorname{Re} \lambda \leq 0$$

<u>o.K.</u>, aber bekommt man so Aussagen über den allgemeinen Fall u'(t) = f(t, u(t))?

■ Motivation für das Testproblem (81):

(1)
$$\bullet$$
 $\begin{cases} u'(t) = f(t, u(t)) \\ u(0) = u_0 \end{cases} \Rightarrow u(t)$ exakte Lösung

(83)
$$\bullet \quad \tilde{u}'(t) = f(t, \tilde{u}(t)) \\ \tilde{u}(0) = u_0 + \delta$$
 $\Rightarrow \tilde{u}(t)$ gestörte Lösung

• Linearisieren gestörtes Problem (83) in der Näherung der exakten Lösung u(t):

$$\Rightarrow \tilde{u}'(t) = \underbrace{f(t, u(t))}_{= u'(t)} + f_u(t, u(t)) \left(\tilde{u}(t) - u(t) \right) + o(\tilde{u}(t) - u(t))$$

 \Rightarrow Erhalten also für den Fehler $z(t) = \tilde{u}(t) - u(t)$:

$$z'(t) = \underbrace{f_u(t, u(t))}_{\text{Einfrieren der Jacobi-Matrix}} z(t) + o\left((z(t))\right), \quad z(0) = \delta$$

$$J = f_u(t_*, u(t_*)) = [\dots]_{N \times N}, \text{ z.B. für } t_* = 0$$

Dann ist die Dgl.

(84)
$$z'(t) = J z(t)$$
 (vgl. Kap. 2: $J = M_h^{-0.5} K_h M_h^{-0.5}$)

ein Modell für die Fehlerausbreitung.

• Sei J diagonalisierbar, d.h.

○ ∃ vollst. System v. EV :
$$v_1, ..., v_N$$

○ EW : $\lambda_1, ..., \lambda_N$.

Ansatz:

$$z(t) = \sum_{i=1}^{N} \eta_i(t) v_i$$

Eingesetzt in (84) ergibt das:

$$z'(t) \equiv \sum_{i=1}^{N} \eta'_i(t) v_i = Jz(t) \equiv \sum_{i=1}^{N} \eta_i(t) \lambda_i v_i,$$

d.h.
$$\eta'_i(t) = \lambda_i \eta_i(t), \quad \forall \ i = \overline{1, N}$$
 $\hat{=}$ (81)

• Schlußfolgerung:

Verhalten des Integrationsverfahrens (ESV) für (81) mit $\lambda \in \sigma(f_u(t, u(t))) = \{ EW \text{ der Jacobi-Matrix} \}$ gibt Aufschluß über Stabilität des Verfahrens!

- Wenden nun <u>expliziten</u> und <u>impliziten Euler</u> auf unser Testproblem (81) an:
 - expliziter Euler:

$$u_{n+1} = u_n + h\lambda u_n \equiv R(h\lambda)u_n$$
, mit
 $R(z) = 1 + z$

• impliziter Euler:

$$u_{n+1} = u_n + h\lambda u_{n+1}$$
, also
$$u_{n+1} = \frac{1}{1 - h\lambda} u_n \equiv R(h\lambda) u_n$$
, mit
$$R(z) = \frac{1}{1 - z}$$

■ Wendet man nun ein ESV auf ein "stabiles" AWP (81) an,

 \longrightarrow d.h. mit $Re\;\lambda \le 0$ (Störungen in den AW vergrößern sich nicht !), dann sollte das ESV

$$u_{n+1} = R(h\lambda)u_n$$

die gleiche Eigenschaft haben, d.h.

$$|R(h\lambda)| \le 1$$
,

falls $\lambda \in \mathcal{C}$: $Re \lambda \leq 0$ und $\lambda \in \sigma (f_u(t, u(t)))$.

Diese Stabilitätsinterpretation führt zum Begriff der <u>A-Stabilität</u> (Pkt. 3.5.1)!

3.5.1 A-Stabilität (\longrightarrow ESV)

- Stabilitätsbereich:
 - Ein ESV zur Lösung des Testproblems

$$(81) u'(t) = \lambda u(t)$$

habe die Form

$$(85) u_{n+1} = R(h\lambda)u_n$$

mit der Stabilitätsfunktion R(z).

• Definition 3.40: (Stabilitätsbereich)

Die Menge

$$S = \{ z \in \mathcal{C} : |R(z)| \le 1 \}$$

heißt Stabilitätsbereich des ESV.

• Bemerkung:

- 1. Für $z = h\lambda \in S$ produziert ESV (81) beschränkte (d.h. stabile) Lösungen !
- 2. Fordern daher, Schrittweite h so zu wählen, daß

$$h \lambda \in S$$

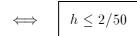
für alle "relevanten" (†) $\lambda: Re \ \lambda \leq 0$.

• Beispiel: Expliziter Euler: R(z) = 1 + z

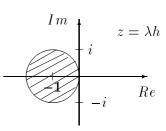
$$\Rightarrow S = \{ z \in \mathcal{C} : |z - (-1)| \le 1 \}$$

Für das Bsp. (80) aus $\boxed{\text{P IX}} \Rightarrow \lambda = -50 = Re \lambda$

 \Rightarrow Stabilitätsforderung: $-50h \in S,$ d.h. $-1 \leq -50h + 1 \leq 1$



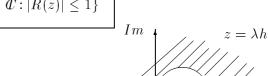
 $h \le 2/50$ (vgl. mit num. Resultaten!)



■ A-Stabilität für ESV

- Ideal: $C^- := \{z : Re \ z \le 0\} \subset S$
 - $\Rightarrow \forall \lambda \in \mathcal{C}: Re \lambda \leq 0$ hätte man Stabilität im obigen Sinne unabhängig von der Schrittweite h!
- **Definition 3.41:** (A–Stabilität für ESV)

Ein ESV heißt <u>A-stabil</u>, falls $\mathcal{C}^- \subset S := \{z \in \mathcal{C} : |R(z)| \le 1\}$



Re

• Beispiel: Impliziter Euler: R(z) = 1/(1-z)

$$\Rightarrow S = \{z \in \mathcal{C} : |1 - z|^{-1} \le 1\} =$$

$$= \{z \in \mathcal{C} : 1 \le |z - 1|\}$$

- $\Rightarrow~\mathcal{C}^-\subset S,$ d.h. impliziter Euler ist A—stabil !
 - o d.h. impliziter Euler produziert beschränkte Näherungen, falls Testproblem (81) beschränkte Lösung hat!
 - o d.h. Fehler im AW werden nicht vergrößert, falls $Re \lambda \leq 0$!

■ A-Stabilität von Runge-Kutta-Verfahren:

• Wenden allgemeines RK-Verfahren (o. B. d. Allg.: $N=1,\ h_n=h)$

$$\begin{array}{c|c} c & A \\ \hline & b \\ \hline \\ g_i & = u_n + h \sum_{j=1}^l b_j f(t_n + c_j h, g_j) \\ g_i & = u_n + h \sum_{j=1}^l a_{ij} f(t_n + c_j h, g_j), \quad i = \overline{1, l} \end{array}$$

auf das Testproblem (81) $u'(t) = \underline{\lambda} \, u(t) \equiv \underline{f}$ an:

(86)
$$\begin{cases} g = u_n e + h \lambda A g, & e = (1, 1, ..., 1)^T, \text{ i. allgem. } (1 \to I), \\ u_{n+1} = u_n + h \lambda b^T g, \end{cases}$$

$$\Rightarrow (I - h\lambda A)g = u_n e, \text{ d.h. } g = (I - h\lambda A)^{-1}(u_n e).$$

Also gilt die Darstellung

(87)
$$u_{n+1} = u_n + h\lambda b^T (I - h\lambda A)^{-1} e u_n = R(h\lambda) u_n$$

mit der Stabilitätsfunktion

(88)
$$R(z) = 1 + zb^{T}(I - zA)^{-1}e.$$

- Anwendung der Cramerschen Regel auf das zu (86) äquivalente System ($z=\lambda h$)

$$\begin{bmatrix} (I - zA) & 0 \\ -zb^T & 1 \end{bmatrix} \begin{bmatrix} g \\ \mathbf{u}_{n+1} \end{bmatrix} = u_n \begin{bmatrix} e \\ 1 \end{bmatrix}, \quad z = \lambda h,$$

liefert eine weitere Darstellung der Stabilitätsfunktion

(88)'
$$R(z) = \frac{\det \begin{pmatrix} I - zA & e \\ -zb^T & 1 \end{pmatrix}}{\det (I - zA)} = \frac{\det (I - zA + zeb^T)}{\det (I - zA)}$$
$$\det \begin{pmatrix} I - zA & e \\ -zb^T & 1 \end{pmatrix} = \det \begin{pmatrix} I - zA + zeb^T & 0 \\ -zb^T & 1 \end{pmatrix} = \det (I - zA + zeb^T)$$

• Betr. explizite RK-Verfahren, d.h.

$$A = \begin{bmatrix} 0 \\ a_{21} & 0 \\ a_{31} & a_{32} & 0 \\ \vdots & \ddots & \ddots \\ a_{l1} & \dots & a_{l-1} & 0 \end{bmatrix} = \text{echte untere Dreiecksmatrix.}$$

$$\implies \det (I - zA) = \det \begin{bmatrix} 1 & \mathbf{O} \\ & \ddots \\ x & 1 \end{bmatrix} = 1.$$

$$\Longrightarrow |R(z)| = |\det{(I - zA + zeb^T)}| \to \infty \text{ für } |z| \to \infty \text{ !}$$

$$\Longrightarrow$$
 Explizite RK-Formeln können nicht A -stabil sein !

⇒ Damit gewinnen implizite RK-Formeln für steife Dgl. an Bedeutung !

- Beispiele A-stabiler impliziter RK-Formeln:
 - 1. Impliziter Euler (↑)

2. Implizite Mittelpunktsregel:
$$R(z) = 1 + \frac{z}{1 - \frac{z}{2}} = \frac{2 + z}{2 - z}$$

$$f = \lambda u : u_{n+1} = u_n + h\lambda g \text{ mit } g = u_n + \frac{h}{2}\lambda g$$

$$u_{n+1} = u_n + h\lambda \left(1 - \frac{h\lambda}{2}\right)^{-1} u_n = \left(1 + h\lambda \left(1 - \frac{h\lambda}{2}\right)^{-1}\right) u_n$$

$$\implies S := \{z \in \mathcal{C} : |R(z)| \le 1\} = \mathcal{C}^- \text{ (mms)}$$

3. Implizite Trapezregel (CN-Schema):

Implizite Trapezregel (CN-Schema):
$$R(z) = \frac{2+z}{2-z}$$

$$\uparrow$$

$$f = \lambda u : u_{n+1} = u_n + \frac{h\lambda}{2}(u_n + u_{n+1}) \Rightarrow u_{n+1} = \frac{(1+\frac{z}{2})}{(1-\frac{z}{2})}u_n, \quad z = \lambda h$$

$$\Longrightarrow S = \mathbb{C}^- \quad (\text{mms})$$

4. Implizite *l*-stufige RK-Formeln vom Gauß-Typ (mms).

3.5.2 L-Stabilität

■ Motivation

• Für die implizite Mittelpunktsregel und für die implizite Trapezregel (CN-Schema) gilt z.B.

$$S := \{ z \in \mathcal{C} : |R(z)| \le 1 \} = \mathcal{C}^{-1}$$

$$\implies |R(iy)| = 1 \quad \forall \ y \in \mathbb{R}^1$$

$$\Longrightarrow \lim_{z\to -\infty} R(z) = \lim_{z\to +\infty} R(z) = \lim_{y\to \infty} R(iy), \, \mathrm{da} \ R(\cdot,\cdot) - \mathrm{rationale} \ \mathrm{Funktion}$$

$$\Longrightarrow z \in {\rm I\!\!\!C} \colon Re \ z << 0 \Rightarrow |R(z)| \stackrel{\approx}{\scriptscriptstyle <} 1,$$

während exakte Lsg. e^z $(z = \lambda t)$ sehr klein ist.

⇒ Steife Komponenten werden langsam gedämpft (↑)

■ **Definition 3.42:** (*L*-Stabilität)

Ein ESV mit der Stabilitätsfunktion R(z) heißt <u>L-stabil</u>, falls es A-stabil ist und

$$\lim_{z \to \infty} R(z) = 0.$$

 \blacksquare Beispiel: Implizites Euler-Verfahren ist L-stabil, da

$$\lim_{z \to \infty} R(z) = \lim_{z \to \infty} \frac{1}{1 - z} = 0.$$

Impl. MP-Regel und impl. TR sind zwar A-stabil, aber nicht L-stabil (mms)!

■ <u>Satz 3.43:</u>

<u>Vor.:</u> Für eine implizite, l-stufige RK-Formel $\frac{c \mid A}{b}$ mit regulärer Matrix A sei eine der folgenden beiden Bedingungen erfüllt:

(89)
$$a_{lj} = b_j, \quad j = 1, 2, ..., l \quad \text{bzw.}$$

(90)
$$a_{i1} = b_1, \quad i = 1, 2, \dots, l.$$

Bh.: Dann gilt:

$$\lim_{z \to \infty} R(z) = 0.$$

Beweis:

• Aus Darstellung (88) folgt:

$$R(z) = 1 + zb^{T}(I - zA)^{-1}e = 1 + zb^{T}((-zA)(-\frac{1}{z}A^{-1} + I))^{-1}e$$

$$= 1 - zb^{T}\left(I - \frac{1}{z}A^{-1}\right)^{-1}\frac{1}{z}A^{-1}e$$

$$\Rightarrow \lim_{z \to \infty} R(z) = 1 - b^{T}A^{-1}e.$$

- <u>Bed. (89)</u>: $A^T e_l = b$ $\Rightarrow 1 - b^T A^{-1} e = 1 - e_l^T A A^{-1} e = 1 - 1 = 0.$
- <u>Bed. (90):</u> $A e_1 = b_1 e$ $\frac{1}{b_1} e_1 = A^{-1} e$ $(b_1 \neq 0, \text{ da } A \text{ reg. !})$ $\Rightarrow 1 - b^T A^{-1} e = 1 - b^T \frac{1}{b_1} e_1 = 1 - \frac{b_1}{b_1} = 1 - 1 = 0.$

q.e.d.

3.5.3 B-Stabilität

- Zur Behandlung <u>stark nichtlinearer steifer Probleme</u> benötigt man Verfahren mit noch stärkeren Stabilitätseigenschaften, wie etwa die <u>B-Stabilität</u>:
 - Dazu betr. wir eine Klasse nichtlinearer Testprobleme

(91)
$$\begin{cases} u'(t) = f(t, u(t)) \text{ mit der Eigenschaft} \\ (f(t, v) - f(t, w), v - w) \leq 0 \quad \forall \ v, w \in \mathbb{R}^N \quad \forall \ t \in I. \end{cases}$$

Solche Dgl.-Systeme nennt man dissipativ.

• Seien v(t) und w(t) zwei Lösungen des dissipativen Systems (91). Dann gilt:

(93)
$$||v(t) - w(t)|| < ||v(t_0) - w(t_0)|| \quad \forall t > t_0,$$

d.h. insbesondere, daß Fehler im Startpunkt $t_0=0$ nicht zunehmen.

Beweis: (mms)

Btr.
$$m(t) := ||v(t) - w(t)||^2 \Rightarrow \text{monoton fallend, da}$$

 $m'(t) = 2(v'(t) - w'(t), v(t) - w(t)) =$
 $= 2(f(t, v(t)) - f(t, w(t)), v(t) - w(t)) < 0 \#$

- Für ein B-stabiles Verfahren (ESV) fordert man nun das gleiche qualitative Verhalten!
- **Definition 3.44:** (*B*–Stabilität)

Ein ESV $u_{n+1} = u_n + h_n \varphi(t_n, u_n, t_n)$ heißt <u>B-stabil</u>, falls aus der Dissipativität

$$(94) (f(t,v) - f(t,w), v - w) < 0 \forall v, w \in \mathbb{R}^N \forall t \in I$$

folgt, daß

$$(95) ||v_1 - w_1|| \le ||v_0 - w_0||$$

für alle Schrittweiten $h_{(=)}>0$ gilt, mit

$$v_1 = v_0 + h\varphi(t_0, v_0, h)$$
 und $w_1 = w_0 + h\varphi(t_0, w_0, h)$.

■ <u>Satz 3.45:</u> Aus der *B*-Stabilität folgt die *A*-Stabilität.

Beweis: (mms).

■ Die folgenden 2 Bedingungen (algebraische Stabilitätsbed.) garantieren die B-Stabilität von (impliziten) RK-Formeln:

1.
$$b_i \geq 0 \quad \forall \ i = \overline{1, l}$$

2. Matrix $M = [m_{ij} := b_i a_{ij} + b_j a_{ji} - b_i b_j]$ ist positiv semi-definit.

Beweis: siehe [5] S. 193 f. #

- Beispiele: B-stabiler RK-Formeln (mms)
 - 1. Impliziter Euler: $\begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array}$
 - 2. Implizite Mittelpunktsregel: $\begin{array}{c|c} 1/2 & 1/2 \\ \hline & 1 \end{array}$
 - 3. Implizite RK-Formeln vom Gauß-Typ (siehe [5] S. 194 ff auch weiter B-stabile RK-Formeln).

3.5.4 A-Stabilität von MSV

- Wendet man <u>lineares MSV</u>
 - (53) $\alpha_k u_{n+k} + \alpha_{k-1} u_{n+k-1} + \ldots + \alpha_0 u_n = h(\beta_k f_{n+k} + \ldots + \beta_0 f_n)$ auf das Testproblem
 - (81) $u'(t) = \overline{\lambda u(t)} =: f(t, u(t)) \equiv f(u(t))$

an, so erhält man die Differenzengleichung

(96) $(\alpha_k - \mu \beta_k) u_{n+k} + (\alpha_{k-1} - \mu \beta_{k-1}) u_{n+k-1} + \ldots + (\alpha_0 - \mu \beta_0) u_n = 0$ mit $\mu = h\lambda$.

Die charakteristische Gleichung der Differenzengleichung (96) lautet:

(97) $\rho(z) - \mu \sigma(z) = 0.$

Im <u>Stabilitätsbereich</u> S werden nun diejenigen Werte $\mu \in \mathcal{C}$ zusammengefaßt, die stabile (d.h. beschr.) Lösungen der Differenzengleichung (96) zur Folge haben:

(98) $S:=\{\mu\in \mathcal{C}: \text{ Alle Nullstellen } \zeta \text{ von } \rho(z)-\mu\,\sigma(z) \text{ erfüllen die Bed. (72)}, \\ \text{d.h. } |\zeta|\leq 1, \text{ falls einfach und } |\zeta|<1, \text{ falls mehrfach } \}$

Um das Testproblem stabil zu lösen, muß gelten:

- (99) $h \lambda \in S$.
- Bemerkung: MSV (53) ist 0-stabil, genau dann, wenn $0 \in S$.
- <u>Definition 3.46:</u> (A-Stabilität von MSV)

Ein lineares MSV heißt <u>A-stabil</u>, falls ${\mathfrak C}^-\subset S.$

■ Satz 3.47: (2. Dahlquist – Barriere)

Ein A-stabiles lineares MSV muß die Ordnung

$$p \leq 2$$

haben.

Beweis: siehe Literatur [1], Satz 7.36, S. 318-321. #

■ Beispiel: Impl. TR = CN = Adams-Moulton (k = 1)Die implizite Trapezregel (CRANK-NICOLSON)

$$\begin{split} u_{n+1} - u_n &= h\left(\frac{1}{2}f_{n+1} + \frac{1}{2}f_n\right) \equiv \mu\left(\frac{1}{2}u_{n+1} + \frac{1}{2}u_n\right) \\ \left(1 - \mu\frac{1}{2}\right)z + \left(-1 - \mu\frac{1}{2}\right) &= 0 \qquad \text{(charakteristische Gl.)} \\ \zeta &= \frac{1 + \frac{\mu}{2}}{1 - \frac{\mu}{2}} = \frac{2 + \mu}{2 - \mu} \qquad - \text{einfache Nullstelle} \\ |\zeta| &= \left|\frac{2 + \mu}{2 - \mu}\right| \leq 1 \text{ gdw. } \mu \in \mathcal{C}^- = S, \end{split}$$

ist A-stabil (\uparrow o.k.) und hat unter allen A-stabilen MSV mit der Konsistenzordnung $p \leq 2$ den <u>kleinsten führenden Fehlerterm</u>: $\Rightarrow \underline{\text{lineare MSV sind für steife Dgl.}}$ uninteressant!

Kapitel 4

Numerische Behandlung hyperbolischer ARWA

- 4.1 Dreischichtige Differenzenschemata für hyperbolische ARWA
- 4.1.1 Allgemeine und kanonische Form dreischichtiger Schemata
 - Allgemeine Form:

(1)
$$C_1 v^{j+1} + C_0 v^j + C_{-1} v^{j-1} = \tau \varphi^j, \quad j = 1, 2, \dots, m-1,$$

$$AB: v^0, v^1 \text{ geg.},$$

wobei
$$v=v^j:\omega_h\longrightarrow I\!\!R^1$$
 – Gitterfkt. auf j –ter Zeitschicht,
$$\varphi=\varphi^j:\omega_h\longrightarrow I\!\!R^1$$
 – rechte Seite auf j –ter Zeitschicht,
$$C_{-1},C_0,C_1$$
 – lineare Operatoren (RB–eingearbeitet),
$$\tau$$
 – Zeitschritt ($m\tau=T$), h – Ortsdikretisierungsparameter,
$$\omega\equiv\omega_h$$
 – Gitter für Ω (vgl. Pkt. 2.2.1).

■ Entsteht z.B. bei der Diskretisierung:

- hyperbolischer ARWA:
$$\frac{\partial^2 u}{\partial t^2} \longmapsto u_{\bar{t}t} \quad (x: \text{FDM, FEM})$$

$$\begin{vmatrix} \frac{\partial^2 u}{\partial t^2} + L \, u(x,t) = f(x,t), & (x,t) \in Q_T = \Omega \times \mathbf{T}, & \mathbf{T} = (0,T); \\ + & \text{RB:} & 1. - 4. \text{ Art;} \\ + & \text{AB:} & u(x,0) = u_0(x), \frac{\partial u}{\partial t}(x,0) = u_1(x), & x \in \bar{\Omega}; \\ L-\text{elliptischer Differentialausdruck (siehe [17] Nu II).} \end{vmatrix}$$

- parabolische ARWA:
$$\frac{\partial u}{\partial t} \longrightarrow u_{\stackrel{\circ}{t}}$$
 (x: FDM, FEM) (siehe Kap. 2, insbesondere Pkt. 2.2).

■ Kanonische Form:

(2)
$$Bv_{0}^{j} + \tau^{2}Rv_{\bar{t}t}^{j} + Av^{j} = \varphi^{j}, \quad j = 1, 2, \dots, m-1,$$

$$AB: v^{0}, v^{1} \text{ geg}.$$

■ Beziehung zwischen kanonischer und allgemeiner Form:

$$(2) \Leftrightarrow B \frac{v^{j+1} - v^{j-1}}{2\tau} + \tau^2 R \frac{v^{j+1} - 2v^j - v^{j-1}}{\tau^2} + Av^j = \varphi^j \mid \cdot \tau$$

$$\Leftrightarrow \underbrace{\left(\frac{1}{2}B + \tau R\right)}_{=C_1} v^{j+1} + \underbrace{\left(\tau A - 2\tau R\right)}_{=C_0} v^j + \underbrace{\left(-\frac{1}{2}B + \tau R\right)}_{=C_{-1}} v^{j-1} = \tau \varphi^j$$

(3)
$$C_{1} = \frac{1}{2}B + \tau R , C_{0} = \tau (A - 2R), C_{-1} = -\frac{1}{2}B + \tau R$$
$$B = C_{1} - C_{-1}, R = \frac{1}{2\tau}(C_{1} + C_{-1}), A = \frac{1}{\tau}(C_{1} + C_{0} + C_{-1})$$

4.1.2 Ein allgemeines Stabilitätsresultat für dreischichtige Schemata

■ Standardvoraussetzungen an (2):

$$H$$
 – (diskreter) reeller Hilbert-Raum (\rightarrow Grundraum)
i. S. Raum von Gitterfkt.: $\|\cdot\|$, (\cdot, \cdot) .

(4)
$$\begin{cases} 1) & B, R, A : H \mapsto H - \text{lineare Operatoren;} \\ 2) & A = A^*, R = R^*; \\ 3) & A > 0, B + 2\tau R > 0 \quad (\Rightarrow 2C_1 > 0 \Rightarrow \exists C_1^{-1}!); \\ 4) & B, R, A - \text{zeitunabhängig.} \end{cases}$$

■ <u>Idee:</u> Rückführung des 3-schichtigen DS in H auf ein 2-schichtiges DS in $H \times H$ und Anwendung der Resultate aus Pkt. 2.2!

Sei m – ungerade, $m \geq 3$:

■ Bezeichnungen:

• Vorwärtige Differenzen:

$$V_t \equiv V_t^j := \frac{1}{2\tau} (\hat{V} - V) \equiv \frac{1}{\bar{\tau}} (\hat{V} - V) \text{ mit } \bar{\tau} = 2\tau,$$

• Skalarprodukt und Norm im Faktorraum $\mathcal{H} = H \times H$:

$$(Y,Z) = \left(\begin{pmatrix} y^1 \\ y^2 \end{pmatrix}, \begin{pmatrix} z^1 \\ z^2 \end{pmatrix} \right) := (y^1, z^1) + (y^2, z^2) \quad \forall \, Y, Z \in \mathcal{H} ,$$
$$\|Y\|^2 = (Y,Y) = \|y^1\|^2 + \|y^2\|^2 \quad \forall \, Y = \begin{pmatrix} y^1 \\ y^2 \end{pmatrix} \in \mathcal{H} .$$

• Energetisches Skalarprodukt:

$$\|Y\|_{\boldsymbol{\mathcal{A}}}^2 \; := (\boldsymbol{\mathcal{A}} \; Y, Y) \; \text{für} \; \boldsymbol{\mathcal{A}} \; : \boldsymbol{\mathcal{H}} \; \mapsto \boldsymbol{\mathcal{H}} \; - \text{linear, s.a., p.d.}$$

■ Resultat:

■ <u>Satz 4.1:</u>

<u>Vor:</u> 1. Standardvoraussetzungen (4): 1) - 4).

- 2. $R > \frac{1}{4}A$.
- 3. B > 0.

 $\begin{array}{ll} \underline{\mathbf{Bh.:}} & \mathrm{Dann} \ \mathrm{ist} \ \mathrm{das} \ 3\text{--schichtige Schema} \ (1) \equiv (2) \\ & (\Leftrightarrow 2\text{--schichtiges Schema} \ (5)) \ \underline{\mathrm{stabil}} \ (\mathrm{bzgl.} \ \underline{\mathbf{AB}} \ \mathrm{u.} \ \mathrm{RS}) \\ & \mathrm{und} \ \forall \ j: (2j+1)\tau \leq T \ \mathrm{gilt} \ \mathrm{die} \ \mathrm{A-priori-Absch\"{a}tzung} : \end{array}$

(6)
$$\|V^{j+1}\|_{\mathcal{A}} \le \|V^j\|_{\mathcal{A}} + T \max_{1 \le k \le j} \|\mathcal{L}^{-1}\Phi^k\|_{\mathcal{A}}$$

Beweis: folgt aus Satz 2.16 angewandt auf das 2-schichtige DS (5):

• Standardvor. 1) - 4) an 2-schichtige DS aus Kap. 2, Pkt. 2.2.2:

- 1) \mathcal{A} , \mathcal{L} : $\mathcal{H} \longrightarrow \mathcal{H}$ lineare Operatoren: o.k.
- 2) $\mathcal{A} = \mathcal{A}^* > 0$ folgt sofort aus 2) + 3) von (4) und $R > \frac{1}{4}A$: $\mathcal{A} = \mathcal{A}^*$ (mms) unter Benutzung von 2): o.k.

 \mathcal{A} p.d.:

$$(\mathcal{A} \ V, V) = \left(\begin{bmatrix} 2R & A - 2R \\ A - 2R & 2R \end{bmatrix} \begin{bmatrix} v^1 \\ v^2 \end{bmatrix}, \begin{bmatrix} v^1 \\ v^2 \end{bmatrix} \right) =$$

$$= (2Rv^1, v^1) + (2Rv^2, v^2) + \underline{((A - 2R)v^2, v^1) + ((A - 2R)v^1, v^2)}$$

$$= (2Rv^1, v^1) + (2Rv^2, v^2) + ((A - 2R)(v^2 - v^1), v^1 - v^2) +$$

$$+ ((A - 2R)v^1, v^1) + ((A - 2R)v^2, v^2)$$

$$= \underline{(Av^1, v^1) + (Av^2, v^2) + ((2R - A)(v^2 - v^1), v^2 - v^1)} =$$

$$= \underline{\frac{1}{2}(A(v^1 + v^2), v^1 + v^2) + \underline{\frac{1}{2}(A(v^2 - v^1), v^2 - v^1)}$$

$$= \underline{\frac{1}{2}(A(v^1 + v^2), v^1 + v^2) + ((2R - \frac{1}{2}A)(v^2 - v^1), v^2 - v^1) > 0$$

$$= \underline{\frac{1}{2}(A(v^1 + v^2), v^1 + v^2) + ((2R - \frac{1}{2}A)(v^2 - v^1), v^2 - v^1) > 0}$$

$$= \underline{\frac{1}{2}(A(v^1 + v^2), v^1 + v^2) + ((2R - \frac{1}{2}A)(v^2 - v^1), v^2 - v^1) > 0}$$

- 3) $\mathcal{L} > 0$: Folgt aus der noch zu beweisenden Stabilitätsbedingung $\mathcal{L} \geq \frac{\bar{\tau}}{2} \mathcal{A} = \tau \mathcal{A} > 0.$ 4) \mathcal{A} , \mathcal{L} – seien zeitunabhängig folgt aus (4) 4).
- Stabilitätsbed. $\mathcal{L} \geq \frac{\bar{\tau}}{2} \mathcal{A} = \tau \mathcal{A}$ 2-schichtiges DS:

$$\begin{split} &((\mathcal{L} - \tau \mathcal{A})V, V) = \\ &= \left(\left[\begin{pmatrix} B + 2\tau R & 2\tau (A - 2R) \\ \mathbf{O} & B + 2\tau R \end{pmatrix} - \tau \begin{pmatrix} 2R & A - 2R \\ A - 2R & 2R \end{pmatrix} \right] \begin{bmatrix} v^1 \\ v^2 \end{bmatrix}, \begin{bmatrix} v^1 \\ v^2 \end{bmatrix} \right) \\ &= \left(\begin{bmatrix} B & \tau (A - 2R) \\ -\tau (A - 2R) & B \end{bmatrix} \begin{bmatrix} v^1 \\ v^2 \end{bmatrix}, \begin{bmatrix} v^1 \\ v^2 \end{bmatrix} \right) = \\ &= (Bv^1, v^1) + (Bv^2, v^2) \ge 0, \text{ da } B \ge 0. \end{split}$$

q.e.d.

- 4.1.3 <u>Beispiel:</u> Gewichtete dreischichtige Differenzenschemata für die Saitenschwingungsgleichung
 - Betr. ARWA für Saitenschwingungsgleichung:

(7)
$$\frac{\partial^{2} u}{\partial t^{2}} - a^{2} \frac{\partial^{2} u}{\partial x^{2}} = f(x, t), \ x \in \Omega = (0, 1), \ t \in \mathbf{T} = (0, T),$$

$$u(x, 0) = u_{0}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x, 0) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1],$$

$$\frac{\partial u}{\partial t}(x) = u_{1}(x), \ x \in \bar{\Omega} = [0, 1$$

- AB: $1. \ v_i^0 = u_0(x_1), \ i = \overline{0,n}$ $2. \ u(x_i,\tau) = u(x_i,0) + \int_0^\tau \frac{\partial u(x_i,t)}{\partial t} \, dt$ $\approx u_0(x_i) + \tau u_1(x_i)$ $t_j = j\tau$ Also: $v_i^0(x) = u_1(x), x \in \overline{\omega}_h \text{ bzw. unter}$ $Benutzung \frac{\partial^2 u(x,0)}{\partial t^2} = \dots O(\tau^3)$ Bem.: $\Omega = (0,l)$ T = (0,T) $T' = T \cdot a/l$ $T = \frac{T}{m}$ $T = \frac{T}{m}$
- Gewichtetes, 3-schichtiges DS für (7): $(0 \le \sigma \le 1)$

$$v_{\bar{t}t}^{j} - \sigma a^{2} v_{\bar{x}x}^{j+1} - (1 - 2\sigma) a^{2} v_{\bar{x}x}^{j} - \sigma a^{2} v_{\bar{x}x}^{j-1} = \varphi^{j}(x), \ x \in \omega_{h}, \ j = \overline{1, m - 1}$$

$$\underline{AB:} \quad v_{i}^{0} = u_{0}(x_{i}), \quad v_{i}^{1} = u_{0}(x_{i}) + \tau u_{1}(x_{i}), i = \overline{0, n} \text{ bzw. } i = \overline{1, n - 1}$$

$$\underline{RB:} \quad v_{0}^{j} = v_{n}^{j} = 0, \quad j = \overline{0, m}$$

■ Formulierung von (8) als Operatorgleichung:

Def. diskreten Hilbert-Raum

$$H = L_2(\omega_h) = L_2^0(\bar{\omega}_h) \ni v : \bar{\omega}_h \to \mathbb{R}^1 : v|_{\gamma_h = \bar{\omega}_h \setminus \omega_h} = 0$$

und den Differenzenoperator

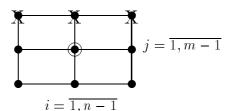
$$\bar{A}v := -a^2 v_{\bar{x}x}, \quad v \in H.$$

Dann läßt sich (8) als dreischichtiges Schema in allgemeiner und kanonischer Form schreiben:

$$\begin{array}{c|c} (9) & & \\ & \widehat{=} \\ & \text{allgem.} \\ \text{Form (1)} \end{array} \\ \hline + \text{AB: } v^0, v^1 \text{ geg.} \\ \hline \end{array} \\ \hline \begin{array}{c} \left(\frac{1}{\tau}I + \tau\sigma\bar{A}\right) v^{j+1} + \underbrace{\left(-\frac{2}{\tau} + (1-2\sigma)\tau\bar{A}\right)}_{=:C_0} v^j + \underbrace{\left(\frac{1}{\tau}I + \tau\sigma\bar{A}\right)}_{=:C_{-1}} v^{j-1} = \tau\varphi \\ \hline \end{array}$$

(10)

$$\stackrel{\triangle}{=} (I + \tau^2 \sigma \bar{A}) v_{\bar{t}t} + \bar{A}v = \varphi$$
+ AB: v^0, v^1 geg.



■ <u>Resultat:</u> Approximation + Stabilität ⇒ diskr. Konvergenz

1. Approximation (Konsistenz): (mms)

$$\psi = \psi(x,t) = u_{\bar{t}t}^j - \sigma a^2 u_{\bar{x}x}^{j+1} - (1-2\sigma)a^2 u_{\bar{x}x}^j - \sigma a^2 u_{\bar{x}x}^{j-1} - \varphi^j(x,t) = .?.$$

- Methode: Taylor-Entwicklung im Pkt. $(x, t) = (x_i, t_j)$ (\rightarrow siehe Pkt. 2.1.2 für 2-schichtige DS)

$$\begin{array}{ccc} - & \underline{\text{Resultate:}} & \psi = O(\tau^2 + h^2) & \forall \sigma \in [0,1] \\ \\ \psi = O(\tau^2 + h^4) & \text{falls} & \sigma = \sigma_* := \bar{\sigma} - \frac{h^2}{12\tau^2} \in [0,1], \\ \\ \varphi = f - \sigma_* \tau^2 \frac{\partial^2 f}{\partial x^2}, \end{array}$$

$$\bar{\sigma} \neq \bar{\sigma}(h, \tau)$$
 bel. Konst. (Stabil.) $(a^2 = 1)$

2. Stabilität:

- Methode: a) Fourier-Analyse nach Eigenfkt. von \bar{A} (\rightarrow siehe Pkt. 2.1.3 für 2-schichtige DS !)
 - b) Allgemeine Stabilitätstheorie aus Pkt. 4.1.2: <u>Satz 4.1:</u> \Rightarrow Voraussetzungen überprüfen! 1. Standardvor. (4): 1) - 4) o.k. 2. $R \equiv \frac{1}{\tau^2}I + \sigma \bar{A} > \frac{1}{4}A \equiv \frac{1}{4}\bar{A}$!

$$3. B \equiv \mathbf{O} \geq 0 \text{ o.k.}$$

- Resultate: 2.?

$$\sigma \ge \frac{1}{4} \Longrightarrow \underline{\text{unbedingt stabil}} : \Rightarrow \text{A-priori-Abschätzung (6)}.$$

 $\sigma = 0 \Longrightarrow \underline{\text{bedingt stabil}} \colon \Rightarrow \text{A-priori-Abschätzung (6)}.$

$$(11) \quad \frac{\text{CFL-Bed.}}{h} \quad \frac{\frac{\tau}{h}}{\frac{1}{a}} \leq \frac{1}{a} \quad \text{Firedrichs und Levi}$$

Tatsächlich:

$$\begin{array}{l} \frac{1}{4}\,A = \frac{1}{4}\,a^2\,\Lambda \ \leq \ \frac{1}{4}\,a^2\,\lambda_{\max}\left(\Lambda\right)I \ < \ \frac{1}{4}\,a^2\,\frac{4}{h^2}\,I \ \leq \ \frac{1}{\tau^2}\,I \ ! \ (\sigma = 0) \\ \uparrow \qquad \qquad \qquad \uparrow \qquad \qquad \downarrow \uparrow \\ \Lambda v := -v_{\bar{x}x} \qquad \text{Pkt. 2.1.3: } \lambda_{\max}(\Lambda) < \frac{4}{h^2} \qquad \frac{a^2}{h^2} \leq \frac{1}{\tau^2} \end{array}$$

- $\ddot{\mathbf{U}}$ **4.1** Untersuchen Sie die Stabilität des $O(\tau^2+h^4)$ –Schemas!

4.2 Zurückführung auf AWA mittels Semidiskretisierung

■ Btr. zunächst <u>hyperbolische ARWA</u> in <u>klassischer Formulierung</u> (vgl. [17] Nu II, Pkt. 1.2):

Ges.
$$u(x,t) \in X = C^{2,2}(Q_T) \cap C^{0,1}(\bar{\Omega} \times [0,T)) \cap C^{1,0}(\Omega \cup \Gamma_2 \cup \Gamma_3 \times (0,T)):$$

$$\frac{\partial^2 u}{\partial t^2} - \sum_{i,j=1}^d \frac{\partial}{\partial x_i} \left(a_{ij}(x,t) \frac{\partial u}{\partial x_j} \right) + a(x,t) u(x,t) = f(x,t),$$
ellipt. Anteil
$$\forall (x,t) \in Q_T = \Omega \times (0,T)$$

$$+ \underline{RB}: \qquad \text{o. B. d. Allg.}$$

$$u(x,t) = g_1(x,t) \stackrel{\downarrow}{=} 0, \quad x \in \Gamma_1$$

$$\frac{\partial u}{\partial N} := \sum_{i,j=1}^d a_{ij}(x,t) \frac{\partial u}{\partial x_j} n_i = g_2(x,t), \quad x \in \Gamma_2$$

$$\frac{\partial u}{\partial N} + \kappa(x,t) u(x,t) = g_3(x,t), \quad x \in \Gamma_3$$

$$+ \underline{AB}: u(x,0) = u_0(x), \quad \frac{\partial u}{\partial t}(x,0) = u_1(x), \quad x \in \bar{\Omega}$$

■ Linienvariationsformulierung (vgl. Pkt. 2.3.1):

(13) Ges.
$$u \in C(\bar{T}, V_0)$$
 mit $\dot{u} \in L_{\infty} (T, V_0) \cap C(\bar{T}, L_2(\Omega))$ und $\ddot{u} \in L_{\infty} (T, L_2(\Omega))$:
$$(\ddot{u}(t), v)_0 + a(t; u(t), v) = \langle F(t), v \rangle \quad \forall v \in V_0, \quad \forall \text{ f.\"{u}. } t \in T$$

$$+ \text{AB:} \quad u(0) = u_0 \text{ in } V_0,$$

$$\dot{u}(0) = u_1 \text{ in } L_2(\Omega)$$

$$\text{mit } V_0 := \{ v \in V = H^1(\Omega) : v = 0 \text{ auf } \Gamma_1 \}, \quad F \in C(T, V_0^*).$$

Bemerkung: Existenz-, Eindeutigkeits- und Regularitätsaussagen siehe [10] Method of Rothe in Evolution Equations.

Teubner-Verlag, Leipzig 1985.

■ Galerkin-FEM-Semidiskretisierung von (13):

Ansatz:
$$u_h(x,t) = \sum_{i \in \omega_h} u^{(i)}(t) p^{(i)}(x) \in V_{0h} = \text{span } \{p^{(i)} : i \in \omega_h\} \subset V_0$$
 (14)

 $(13)_h$

$$\begin{aligned} \text{Ges. } u_h(x,t) \colon & \frac{d^2}{dt^2}(u_h,v_h)_0 + a(t;u_h,v_h) = < F, v_h > & \forall \ v_h \in V_{0h} \ \ \forall \ \text{f.\"{u}.} \ t \in \textbf{\textit{T}} \\ & + \text{AB: } \ \ \text{z.B. } L_2\text{-Projektion:} \\ & (u_h(\cdot,0),v_h)_0 = (u_0(\cdot),v_h)_0 \ \ \forall \ v_h \in V_{0h}, \\ & (\dot{u}_h(\cdot,0),v_h)_0 = (u_1(\cdot),v_h)_0 \ \ \forall \ v_h \in V_{0h}. \end{aligned}$$



 $(\underline{13})_h$

Ges.
$$\underline{u}_h(t) := [u^{(i)}(t)]_{i \in \omega_h}$$
: $M_h \underline{\ddot{u}}_h(t) + K_h(t)\underline{u}_h(t) = \underline{f}_h(t), \quad t \in \mathbf{T}$
 $+ \text{AB}$: $\underline{u}_h(0) = M_h^{-1}\underline{u}_{0h},$
 $\dot{u}_h(0) = M_h^{-1}\underline{u}_{1h}.$

System gewöhnlicher Dgl. 2. Ordnung

Ges.
$$y(t) := M_h^{0.5} \underline{u}_h(t) : I \equiv \overline{T} \longmapsto I\!\!R^{N=N_h} :$$

 $y''(t) = f(t, y(t)) := -M_h^{-0.5} K_h M_h^{-0.5} y(t) + M_h^{-0.5} \underline{f}_h(t), \quad t \in I$
 $+ \text{AB: } y(0) = y_0 := M_h^{-0.5} \underline{u}_{0h}, \quad y'(0) = v_0 := M_h^{-0.5} \underline{u}_{1h}.$

$$\bigoplus_{13)'_h} U := \begin{bmatrix} y \\ v \end{bmatrix} \text{ mit } v(t) = y'(t)$$

■ Volldiskretisierung mittels gew. 3-schichtiger DS:

Ges.
$$\underline{v}^{j} = \underline{v}_{h}^{j} := [v^{(i)}(t_{j})]_{i \in \omega_{h}} \in \mathbb{R}^{N_{h}}:$$

$$M\underline{v}_{\overline{t}t}^{j} + \sigma K(t_{j+1})\underline{v}^{j+1} + (1 - 2\sigma)K(t_{j})\underline{v}^{j} + \sigma K(t_{j-1})\underline{v}^{j-1} = \underline{\varphi}^{j}$$

$$+ AB: \underline{v}^{0}, \ \underline{v}^{1} \text{ geg.}$$

$$\text{mit } M = M_{h}, \ K = K_{h}(t), \ \underline{\varphi}^{j} \stackrel{\text{z.B.}}{=} \underline{f}_{h}(t_{j}).$$

$$\frac{\text{Allgemeine Form:}}{K \neq K(t)} \frac{C_1 \underline{v}^{j+1} + C_0 \underline{v}^j + C_{-1} \underline{v}^{j-1} = \tau \underline{\varphi}^j}{\text{mit } C_1 = \frac{1}{\tau} M + \tau \sigma K, \ C_0 = -\frac{2}{\tau} M + (1 - 2\sigma) K, \ C_{-1} = \frac{1}{\tau} M + \tau \sigma K}$$

Kanonische Form:
$$B=C_1-C_{-1}=\mathbf{0}$$
, $R=\frac{1}{2\tau}(C_1+C_{-1})=\frac{1}{\tau^2}M+\sigma K$, $A=\frac{1}{\tau}(C_1+C_0+C_{-1})=K$

$$A = \frac{1}{\tau}(C_1 + C_0 + C_{-1}) = K$$

$$(M + \tau^2 \sigma K) \underline{v}_{\bar{t}t} + K\underline{v} = \varphi$$

$$+ \underline{AB}: \underline{v}^0, \underline{v}^1 \text{ geg.}$$

Stabilität: $R \equiv \frac{1}{\tau^2}M + \sigma K > \frac{1}{4}A \equiv \frac{1}{4}K$

- $\sigma \ge \frac{1}{4} \Rightarrow \underline{\text{unbedingt stabil !}}$ $\sigma = 0 \Rightarrow \underline{\text{bedingt stabil: CFL-Bedingung !}}$

$$\frac{4}{\tau^{2}} > \max_{\underline{v}} \frac{(K\underline{v},\underline{v})}{(M\underline{v},\underline{v})} = \lambda_{\max}(M^{-1}K) = \lambda_{\max}(M^{-0.5}KM^{-0.5}) = \frac{c^{2}}{h^{2}}$$

$$\uparrow \qquad \qquad \uparrow \qquad \qquad \uparrow$$

$$\Rightarrow \frac{\tau^{2} < \frac{4}{\lambda_{\max}(M_{h}^{-1}K_{h})}}{\frac{\tau}{h} < \frac{2}{c}} \qquad \qquad = \text{CFL-Bed. !}$$

$$\Rightarrow \tau < 2/\sqrt{\lambda_{\max}(M_{h}^{-1}K_{h})}.$$

- Volldiskretisierung mittels A-stabiler Einschrittformeln angewandt auf (13)_h: U'(t) = F(t, U(t)), U(o) geg.:
 - z.B. <u>implizite Trapez–Regel</u> (A–stabil \Rightarrow unbedingt stabil):

$$U_{j+1} = U_j + \frac{\tau}{2} (F_j + F_{j+1}), \quad U_0 = \begin{bmatrix} y_0 \\ v_0 \end{bmatrix} \text{ geg.}$$

$$\begin{cases} y_{j+1} = y_j + \frac{\tau}{2} (v_j + v_{j+1}) \\ v_{j+1} = v_j + \frac{\tau}{2} (f(t_j, y_j) + f(t_{j+1}, y_{j+1})) \end{cases}$$

$$\underline{u}_{h}^{j+1} = \underline{u}_{h}^{j} + \frac{\tau}{2} (\underline{u}_{h}^{j} + \underline{u}_{h}^{j+1})$$

$$M_{h} \underline{u}_{h}^{j+1} = M_{h} \underline{u}_{h}^{j} + \frac{\tau}{2} (\underline{f}_{h}^{j} + \underline{f}_{h}^{j+1} - K_{h} \underline{u}_{h}^{j} - K_{h} \underline{u}_{h}^{j+1})$$

$$\underline{AB:} \quad \underline{u}_{h}^{0} := M_{h}^{-1} \underline{u}_{0h}, \quad \underline{u}_{h}^{0} := M_{h}^{-1} \underline{u}_{1h}$$

$$M_h \, \underline{\ddot{u}}_h^{j+1} + \underbrace{C_h}\underline{\dot{u}}_h^{j+1} + K_h \, \underline{u}_h^{j+1} = \underline{f}_h^{j+1}$$

$$\underline{u}_h^{j+1} = \underline{u}_h^j + \tau \, \underline{\dot{u}}_h^j + \frac{\tau^2}{2} \left\{ (1 - 2\beta) \, \underline{\ddot{u}}_h^j + 2\beta \, \underline{\ddot{u}}_h^{j+1} \right\}$$

$$\underline{\dot{u}}_h^{j+1} = \underline{\dot{u}}_h^j + \tau \, \left\{ (1 - \gamma) \, \underline{\ddot{u}}_h^j + \gamma \, \underline{\ddot{u}}_h^{j+1} \right\}$$

Literaturverzeichnis

- [1] P. Deuflhard and F. Bornemann. Numerische Mathematik II: Integration gewöhnlicher Differentialgleichungen. de Gruyter Lehrbuch, Berlin, New York, 1994.
- [2] H. Gajewski, K. Gröger, and K. Zacherias. *Nichtlineare Operatorgleichungen und Operatordifferentialgleichungen*. Akademie-Verlag, Berlin, 1974.
- [3] R. Grigorieff. Numerik gewöhnlicher Differentialgleichungen I. B.G. Teubner, Stuttgart, 1972.
- [4] C. Großmann and H.-G. Roos. Numerik partieller Differentialgleichungen. B.G. Teubner, Stuttgart, 1992.
- [5] E. Hairer, S. Nørsett, and G. Wanner. Solving Ordinary Differential Equations I: Nonstiff Problems. Springer Verlag, Berlin, Heidelberg, 1987.
- [6] E. Hairer and G. Wanner. Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems. Springer Verlag, Berlin, Heidelberg, 1991.
- [7] P. J. P. J. R. Dormand. A Family of embedded Runge-Kutta formulae. J. Comp. Appl. Math. 6, 19 – 26, 1980.
- [8] P. J. P. J. R. Dormand. Higher order embedded Runge-Kutta formulae. J. Comp. Appl. Math. 7, 67 75, 1981.
- [9] M. Jung and U. Langer. Skriptum zur Vorlesung FEM (Eine Einführung für Ingenieurstudenten). TU Chemnitz (Fakultät für Mathematik) und Johannes Kepler Universität (Institut für Mathematik), Chemnitz und Linz, 1993.
- [10] J. Kacur. Method of Rothe in Evolution Equation. B.G. Teubner, Leipzig, 1985.
- [11] N. Kikuchi. Finite Element Method in Mechanics. Cambridge University Press, Cambridge, 1986.
- [12] V. Korneev and U. Langer. Approximate Solution of Plastic Flow Theory Problems. B.G. Teubner, Leipzig, 1984.
- [13] D. Kröner. Numerical schemes for conservation laws. B.G. Teubner, Stuttgart, 1996.
- [14] O. Ladyschenskaja. Aufgaben der Mathematischen Physik. Nauka, Moskau, 1973. (in Russisch).

- [15] U. Langer. Skriptum zur Vorlesung MULTIGRID METHODEN. Johannes Kepler Universität, Institut für Mathematik, Linz, 1996.
- [16] U. Langer. Skriptum zur Vorlesung NUMERIK I (Operatorgleichungen). Johannes Kepler Universität, Institut für Mathematik, Linz, 1996.
- [17] U. Langer. Skriptum zur Vorlesung NUMERIK II (RWA). Johannes Kepler Universität, Institut für Mathematik, Linz, 1996.
- [18] K. Rektorys. The Method of Discretization in Time and PDEs. Dordrecht, Bosten, 1982.
- [19] H. Stetter. Analysis of Discretization Methods for Ordinary Differential Equations. Springer Verlag, Berlin, Heidelberg, New York, 1973.
- [20] J. Thomas. Numerical Partial Differential Equations: Finite Difference Methods. Springer, New York, 1995.
- [21] V. Thomée. Galerkin Finite Element Methods for Parabolic Problems. Springer Verlag, Berlin, Heidelberg, New York, Tokyo, 1984.
- [22] E. Zeidler. Vorlesungen über nichtlineare Funktionalanalysis I: Fixpunktsätze. Teubner-Texte zur Mathematik. B.G. Teubner, Leipzig, 1976.
- [23] E. Zeidler. Vorlesungen über nichtlineare Funktionalanalysis II: Monotone Operatoren, volume 9 of Teubner-Texte zur Mathematik. B.G. Teubner, Leipzig, 1977.

1

Kapitel 5

Praktikum

${\bf ``Zeit integrations methoden''}$

(Gundolf Haase)

zur Vorlesung

"Numerik III (Anfangs- und Anfangsrandwertaufgaben)"

PRAKTIKUM

 ${\bf ``Zeit integrations methoden''}$

(Gundolf Haase)

zur Vorlesung

"Numerik III (Anfangs- und Anfangsrandwertaufgaben)"

P 0 Einführungspraktikum¹

5.1 Einführung

- Im Praktikum "Zeitintegrationsverfahren" werden
 - Übungsaufgaben $\begin{bmatrix} \ddot{\mathbf{U}} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \ddot{\mathbf{U}} & \mathbf{22} \end{bmatrix}$ gelöst,
 - kleinere Programmieraufgaben mit numerischen Experimenten $oxed{E}$ $oxed{1}$ $oxed{-}$ $oxed{E}$ $oxed{9}$ durchgeführt und
 - von einem Team (in der Regel 2 Studenten) ähnlich zum Praktikum in Numerik II
 eine größere Praktikumsaufgabe P x gelöst.
- Für die größere Praktikumsaufgabe stehen die parabolischen bzw. hyperbolischen Analoga zu den im Praktikum zur Vorlesung Numerik II vorgestellten Aufgaben zur Wahl:

¹Für jedes Praktikum Px stehen 45 Minuten zur Verfügung.

Px	RWA	Numerik	II	Numerik	III	parabol.	hyperbol.	Bem.
		Team		Team		ARWA	ARWA	
1	Kolben					X		
3	Schellbach						X	nicht-
								linear
4	Courant						X	
7	Kanal 1					X		
9	Kammer						X	
11	Wärmeaus-					X		
	tauscher							
12	TWD					X		
13	E-Magnet						X	
14	Chip					X		
15	Draht					X		Konvek-
								tion
16	Tetraeder					X		3D-FEM

Abrechnung : 1. Vortrag (15 - 20 min)

2. Belegarbeit

• Bewertung: $\ddot{\mathrm{U}}:40~\%$

 $\mathrm{E}:20~\%$

P:40 %

PI

Praktikum I

5.2Numerische Behandlung parabolischer ARWA

5.2.1Differenzenverfahren

5.2.1.1Ein spezielles diskretes Eigenwertproblem

-	-						
stetig	diskret						
Ges. $u \in C^2(0,\ell) \cap C[0,\ell] : u \not\equiv 0 \text{ und } \lambda \in \mathbf{R}$	Ges. $v : \overline{\omega}_h \mapsto \mathbf{R}^1 : v \not\equiv 0 \text{ und } \lambda \in \mathbf{R},$						
:	$x \in \omega_h = \{x_i = ih : i = \overline{1, n-1} :$						
$-u''(x) = \lambda u(x), x \in \Omega = (0,6.2.1)$	$-v_{\overline{x}x}(x) = \lambda v(x) \qquad (1_h)$ $v_0 = v_n = 0 \qquad h = \ell/n$						
$u(0) = u(\ell) = 0$	$v_0 = v_n = 0 \qquad h = \ell/n$						
Homog. lin. Dgl. mit konst. Koeff.!	Matrixeigenwertproblem						
$-u'' - \lambda u = 0$	$\begin{bmatrix} 2 & -1 & \end{bmatrix} \begin{bmatrix} v_1 \end{bmatrix} \begin{bmatrix} v_1 \end{bmatrix}$						
	$\begin{bmatrix} \frac{1}{h^2} & 2 & -1 & & & \\ -1 & 2 & -1 & & 0 \\ & \ddots & \ddots & \ddots & \\ 0 & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_{n-2} \\ v_{n-1} \end{bmatrix} = \lambda \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_{n-2} \\ v_{n-1} \end{bmatrix}$						
	$\left \frac{1}{12} \right \cdots \left \frac{1}{12} \right = \lambda \vdots \left \frac{1}{12} \right $						
	$\begin{bmatrix} h^2 & 0 & -1 & 2 & -1 & v_{n-2} \end{bmatrix}$						
	-1 2 $\begin{vmatrix} v_{n-1} \\ v_{n-1} \end{vmatrix}$						
1. Eigenfunktionen							

$$u_k(x) = \sqrt{\frac{2}{\ell}} \sin \frac{k\pi x}{\ell}, \quad x \in [0, \ell]$$

$$\mu_k(x) = \sqrt{\frac{2}{\ell}} \sin \frac{k\pi x}{\ell}, \quad x \in \overline{\omega}_h$$

$$k = 1, 2, \dots, n - 1$$

$$2. \text{ Eigenwerte}$$

$$\lambda_k = \left(\frac{k\pi}{\ell}\right)^2 \qquad k = 1, 2, \dots \qquad \qquad \lambda_k = \frac{4}{h^2} \sin^2 \frac{k\pi h}{2\ell} \quad k = 1, 2, \dots, n-1$$

$$0 < \left(\frac{\pi}{\ell}\right)^2 = \lambda_1 < \lambda_2 < \lambda_3 < \longrightarrow \infty \qquad \qquad \frac{8}{\ell^2} \leq \lambda_1 < \lambda_2 < \dots < \lambda_k < \dots < \lambda_{n-1} < \frac{4}{h^2} \longrightarrow \infty$$

$$3. \text{ Orthonormalität der Eigenfunktionen}$$

$$(u_k, u_m)_{L_2(\Omega)} := \int_0^\ell u_k(x) u_m(x) dx = \delta_{k,m} \qquad (\mu_k, \mu_m)_h := \sum_{x \in \omega_h} h \mu_k(x) \mu_m(x) = \delta_{k,m}$$

$$k, m = 1, 2, \dots$$

$$k, m = 1, 2, \dots, n-1$$

4. Ableitungen (bzw. Differenzen) der Eigenfunktionen

$$u_k'(x) = \sqrt{\lambda_k} \sqrt{\frac{2}{\ell}} \cos \frac{k\pi x}{\ell}, \quad x \in [0, \ell] \qquad \mu_{k, \overline{x}}(x) = \sqrt{\lambda_k} \sqrt{\frac{2}{\ell}} \cos \frac{k\pi (x - h/2)}{\ell}, \quad x \in \overline{\omega}_h$$

$$k = 1, 2, \dots, n - 1$$
5. Orthogonalität der Ableitungen (bzw. Differenzen) der Eigenfunktionen

$$\int_{0}^{\ell} u_{k}'(x)u_{m}'(x)dx = \lambda_{k}\delta_{k,m} \qquad (\mu_{k,\overline{x}}, \mu_{m,\overline{x}}]_{h} = \sum_{x \in \vec{\omega}_{h}} h\mu_{k,\overline{x}}(x)\mu_{m,\overline{x}}(x) = \lambda_{k}\delta_{k,m}$$

$$k = 1, 2, \dots, n - 1$$

$$\vec{\omega}_{h} = \{x_{i} = ih : i = \overline{1,n}\}$$
6. Über die Vollständigkeit des Systems der Eigenfunktionen

Sei $f \in L_2(\Omega)$, Dann gilt für die Fourierreihe:

$$f(x) \quad \text{``=``} \quad \sum_{k=1}^{\infty} f_k \, u_k(x)$$

$$\downarrow \qquad \qquad \downarrow \qquad \qquad \downarrow \qquad \qquad \downarrow \qquad$$

- d.h. $\| f \sum_{k=1}^{m} f_k u_k(x) \|_{L_2} \xrightarrow{m \to \infty} 0$
- i.S. $\overset{\circ}{W}_{2}^{1}(\Omega)$, falls $f \in \overset{\circ}{W}_{2}^{1}(\Omega)$.
- i.S. $C(\overline{\Omega})$, falls $f \in \mathring{W}_2^1(\Omega) \cap C^2(\overline{\Omega})$.

 $\Omega = (0,\ell)$. $\forall f(\cdot) : \omega_h \longmapsto \mathbf{R}^1$ -Gitterfunktionen gilt :

7. Die Parseval'sche Gleichung

$$\| f \|_{L_{2}(0,\ell)}^{2} := \int_{0}^{\ell} f(x)^{2} dx = \sum_{k=1}^{\infty} f_{k}^{2} \qquad \| f \|_{L_{2}(\omega_{h})}^{2} = \sum_{k=1}^{n-1} f_{k}^{2}$$

- Man löse das diskrete EWP (1_h) bzw. man beweise, daß die angegebenen Eigenfunktionen (1.) und Eigenwerte (2.) richtig sind!
 - O Hinweis: Ansatz $v(x) = c \cdot \sin(\alpha x)$ oder exp-Ansatz.
- U2 | Man beweise die Ungleichungen : a) $\left(\frac{2k}{\ell}\right)^2 < \lambda_k < \left(\frac{\pi k}{\ell}\right)^2$ $\forall k = \overline{1, n-1}$

- b) $\frac{8}{\ell^2} < \lambda_1 < \lambda_2 < \ldots < \lambda_k < \ldots < \lambda_{n-1} < \frac{4}{h^2}$!
- Ü3 Man zeige, daß die Eigenfunktionen $\overline{\mu}_k = c \cdot \sin \frac{k\pi}{\ell} x$, $x \in \omega_h$, $k = \overline{1, n-1}$ für $c = \sqrt{2/\ell}$ bzgl. des Skalarproduktes $(\cdot, \cdot)_h$ orthonormal sind!
- Ü4* Man zeige die Beziehungen 4. und 5. für die rückwärtigen Differenzen der diskreten Eigenfunktionen!

5.2.1.2 Konsistenz- und Stabilitätsuntersuchungen

Betrachten das instationäre 1D-Wärmeleitproblem

Ges.
$$u(x,t)$$
: $\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = f(x,t), \quad x \in (0,1), t \in \mathbf{T} = (0,T)$
+ RB: $u(0,t) = g_0(t), \quad t \in \mathbf{T}$
 $u(1,t) = g_1(t), \quad t \in \mathbf{T}$
+ AB: $u(x,0) = u_0(x), \quad x \in [0,1]$ (5.2.2)

und approximieren es durch das Leapfrog-Schema (RICHARDSON-Schema):

Ges.
$$v: \overline{\omega}_{h\tau} \longmapsto \mathbf{R}^1:$$

$$v_{\stackrel{\circ}{v}} - v_{\overline{x}x} = \varphi(x,t) , \qquad (x,t) \in \overline{\omega}_{h\tau} \quad \text{bzw.}$$

$$\frac{v_i^{j+1} - v_i^{j-1}}{2\tau} - \frac{v_{i-1}^j - 2v_i^j + v_{i+1}^j}{h^2} = \varphi_i^j, \qquad i=\overline{1,n-1}, j=\overline{1,m-1}$$

$$+ \text{RB}: v_0^j = g_0(t_j) , \qquad v_n^j = g_1(t_j) \qquad j=\overline{1,m}$$

$$+ \text{AB}: v_i^0 = u_0(x_i) \qquad i=\overline{0,n}$$

$$+ \text{geeignetes (?) Einschrittverfahren zur Bestimmung von } v_i^1, i=\overline{1,n} \quad ??$$

$$\text{mit } \varphi_i^j = f(x_i, t_j) . \qquad (5.2.3)$$

- Man untersuche die lokale Approximationsordnung (= Konsistenzordnung) des <u>Leapfrog-Schemas</u> (5.2.3) für (5.2.2) und gebe die führenden Terme des lokalen Approximationsfehlers (= lokaler Abschneidefehler = local truncation error) $\psi = .$?. $h^p + .$?. $\tau^q + 0(h^p + \tau^q) = O(h^p + \tau^q)$ an !

 Welches Einschrittverfahren schlagen Sie zur Bestimmung von v_i^1 vor ?
- Ü6 Untersuchen Sie die Stabilität (im v. Neumannschen Sinne) des Leapfrog-Schemas mit homogener rechter Seite!
 - Hinweis:
 - <u>Ansatz</u> : $v_s^{\jmath}=(e^{\alpha\tau})^{\jmath}\,e^{i\lambda sh}$, $(i^2=-1)$ (Fortpflanzung der harmonischen Anfangsstörung)
 - <u>Stabilitätskriterium</u>: $|e^{\alpha \tau}| \leq 1$
 - Allgemeines Stabilitätskriterium : $|e^{\alpha\tau}| \leq 1 + c\tau$ mit $c = \text{const.} \neq c(\tau, h)$ (exponentielles Wachstum ist zugelassen !)

Untersuchen Sie das <u>DU-FORT-FRANKEL-Schema</u> (= Modifikation des Leapfrog-Schemas (5.2.3) mit $v_i^j \longrightarrow \frac{1}{2}(v_i^{j+1} + v_i^{j-1})$):

Ges.
$$v: \overline{\omega}_{h\tau} \longmapsto \mathbf{R}^1:$$

$$\frac{v_i^{\jmath+1} - v_i^{\jmath-1}}{2\tau} - \frac{v_{i-1}^{\jmath} - (v_i^{\jmath+1} + v_i^{\jmath-1}) + v_{i+1}^{\jmath}}{h^2} = \varphi_i^{\jmath}, \quad i = \overline{1, n-1}, j = \overline{1, m-1}$$

$$+ \operatorname{RB}: v_0^{\jmath} = g_0(t_j), \quad v_n^{\jmath} = g_1(t_j) \quad j = \overline{1, m}$$

$$+ \operatorname{AB}: v_i^0 = u_0(x_i) \quad i = \overline{0, n}$$

$$+ \operatorname{geeignetes}(?) \operatorname{Einschrittverfahren zur Bestimmung von } v_i^1, \quad i = \overline{1, n} \quad !?$$

$$\operatorname{mit} \varphi_i^{\jmath} = f(x_i, t_j). \quad (5.2.4)$$

zur Diskretisierung von (5.2.2) auf Konsistenz (lokaler Approximationsfehler) und Stabilität (im v. Neumannschen Sinne).

PIII

Praktikum III

Numerische Experimente 5.2.1.3

E1 | Die parabolische ARWA

Ges.:
$$u(x,t)$$
 mit $x \in [a,b]$, $t \in [0,T]$:
$$\frac{\partial u(x,t)}{\partial t} - \frac{\partial^2 u(x,t)}{\partial x^2} = f(x,t) := t^{s-1} (s(x-a)(x-b) - 2t) \quad (5.2.5)$$

$$+ AB: \quad u(x,0) = x \quad \forall x \in [a,b]$$

$$+ RB: \quad u(a,t) = a \quad \forall t \in (0,T]$$

$$u(b,t) = b \quad \forall t \in (0,T]$$

besitzt für $s \ge 1$ die exakte Lösung

$$u(x,t) = x + (x-a)(x-b)t^{s} . (5.2.6)$$

Man diskretisiere (5.2.5) mit dem σ -gewichteten Differenzenschema auf einem gleichmäßigen Gitter der Form $\overline{\omega}_{h\tau} = \overline{\omega}_h \times \overline{\omega}_{\tau}$, $\overline{\omega}_h = \{x_i = x_0 + \imath h : \imath = \overline{0,n}\}$, $\overline{\omega}_{\tau} = \{t_j = \jmath \tau : \jmath = \overline{0,m}\}$ mit $x_0 = a, h = (b - a)/n, \tau = T/m$.

Implementieren Sie die Familie von Differenzenschemata und führen Sie Rechnungen für verschiedene

- Ausgangsdaten: a = 0 , b = 1 , T = 1, s = 2
- $\sigma = 1 , \sigma = \frac{1}{2} , \sigma = \sigma_* = \frac{1}{2} \frac{h^2}{12\tau} , \sigma = 0$ 1) $h = \frac{1}{10} , \tau = \frac{1}{10}$ 2) $h = \frac{1}{10} , \tau = \frac{1}{200}$ 3) $h = \frac{1}{4} , \tau = \frac{1}{16}$ - Gewichte:
- Diskretisierungen :

durch. Beachten Sie, daß die rechte Seite entsprechend der Parameterwahl diskretisiert werden sollte (Satz 2.10). Vergleichen Sie die diskrete Lösung mit der exakten Lösung (5.2.6)!

E2

Man diskretisiere die parabolische ARWA (vgl. E1

$$\begin{array}{rclcrcl} & \text{Ges.:} & u(x,t) & \text{mit} & x \in [0,1] \;, \; t \in [0,1] \;: \\ & \frac{\partial u(x,t)}{\partial t} - \frac{\partial^2 u(x,t)}{\partial x^2} & = & t \left(2x(x-1) - 2t \right) & x \in (0,1) \; t \in (0,1] \\ & + \text{AB:} & u(x,0) & = & x & \forall x \in [0,1] \\ & + \text{RB:} & u(0,t) & = & 0 & \forall t \in (0,T] \\ & & u(1,t) & = & 1 & \forall t \in (0,T] \\ & \text{L\"osung:} & u(x,t) & = & x + x(x-1) \, t^2 \end{array}$$

mit dem <u>DU-FORT-FRANKEL-Schema</u> auf einem gleichmäßigen Gitter der Form $\overline{\omega}_{h\tau}=\overline{\omega}_h\times\overline{\omega}_\tau$, $\overline{\omega}_h=\{x_i=x_0+\imath h\,:\,\imath=\overline{0,n}\}\;,\;\;\overline{\omega}_\tau=\{t_\jmath=\jmath\tau\,:\,\jmath=\overline{0,m}\}\;\;$ mit $x_0=a,\;h=(b-a)/n,$ $\tau = T/m \text{ (vgl. } \ddot{\text{U}}7 \text{)}.$

Wählen Sie ein geeignetes Einschrittverfahren im Startschritt. Implementieren Sie das Differenzenschema und führen Sie Testrechnungen für folgende Fälle durch:

1)
$$h = \frac{1}{10}$$
, $\tau = \frac{1}{10}$

- 2) Wählen Sie h und τ so, daß der Fehler in der Größenordnung 10^{-2} liegt bei möglichst geringen Aufwand !
- E3 Die parabolische ARWA

Ges.:
$$u(x,t)$$
 mit $x \in [0,1]$, $t \in [0,1]$:
 $5 \frac{\partial u(x,t)}{\partial t} - \frac{\partial^2 u(x,t)}{\partial x^2} = f(x,t) = 10t - x$ (5.2.7)
 $+ AB$: $u(x,0) = x^3 \quad \forall x \in [0,1]$
 $+ RB$: $u(0,t) = t^2 \quad \forall t \in (0,1]$
 $u(1,t) = t^2 + t + 1 \quad \forall t \in (0,1]$

besitzt die exakte Lösung

$$u(x,t) = x^3 + xt + t^2 . (5.2.8)$$

Man diskretisiere (5.2.7) mit dem σ -gewichteten Differenzenschema auf einem gleich mäßigen Gitter der Form $\overline{\omega}_{h\tau}=\overline{\omega}_h\times\overline{\omega}_{\tau}$, $\overline{\omega}_h=\{x_i=x_0+\imath h:\imath=\overline{0,n}\}$, $\overline{\omega}_{\tau}=\{t_j=\jmath \tau:\jmath=\overline{0,m}\}$ mit $x_0=a,\,h=(b-a)/n,\,\tau=T/m$.

Wie sieht die L_2 -Stabilitätsbedingung aus ?

Implementieren Sie die Familie von Differenzenschemata und führen Sie Rechnungen für verschiedene

– Gewichte:
$$\sigma=1$$
 , $\sigma=\frac{1}{2}$, $\sigma=\sigma_*=\frac{1}{2}$ -? , $\sigma=0$

- Diskretisierungen : 1)
$$h = \frac{1}{10}$$
, $\tau = \frac{1}{10}$

1)
$$h = \frac{1}{10}, \tau = \frac{1}{10}$$

2) Fehler 10^{-2} für $\sigma = 0, h = \dots, \tau = \dots$

durch. Beachten Sie, daß die rechte Seite entsprechend der Parameterwahl diskretisiert werden sollte (Satz 2.10). Vergleichen Sie die diskrete Lösung mit der exakten Lösung (5.2.8)!

E4 Randbedingungen der Form $\frac{\partial u}{\partial x}(0,t) = f(t)$

können bei der numerischen Behandlung der homogenen Wärmeleitgleichung

$$\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = 0$$

durch die Approximation

$$u(0,t+\tau) \approx u(0,t) + \frac{2\tau}{h^2} [u(h,t) - u(0,t) - h \cdot f(t)] , \text{d.h.}$$

$$v_0^{j+1} = v_0^j + \frac{2\tau}{h^2} [v_1^j - v_0^j - h \cdot f(t_j)]$$

$$j=0,1,\dots,m-1$$

berücksichtigt werden (Warum?).

Man diskutiere auf dieser Basis die ARWA

Ges.
$$u(x,t)$$
: $\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = 0$, in $Q_T = \underbrace{(0,1)}_{\Omega} \times \underbrace{(0,1)}_{\overline{\mathbf{T}}}$

$$+ RB : \frac{\partial u(0,t)}{\partial x} = 10^6 t, \quad t \in \overline{\mathbf{T}}$$

$$u(1,t) = 0, \quad t \in \overline{\mathbf{T}}$$

$$+ AB : \quad u(x,0) = 0, \quad x \in \overline{\Omega}$$

mit dem expliziten Euler-Verfahren ($\sigma=0$) unter Verwendung der Schrittweiten

- a) h = 0.2 , $\tau = 0.01$
- b) h = 0.1 , $\tau = 0.01$

bis $T=0.08 \ (m=8)$. Vergleichen Sie die numerischen Ergebnisse!

Stabilität des σ -gewichteten Differenzenschemas in der Energienorm

Man zeige, daß für beliebige Gitterfunktionen

$$u, v : \overline{\omega}_h = \{x_i = x_0 + ih : i = \overline{0, n}, h = \frac{b-a}{n}\} \mapsto \mathbf{R}^1$$

mit $x_0 = a, x_n = b, b > a$, die folgenden Formeln gelten :

a) Die diskrete Produktregel:

$$(u \cdot v)_{x,i} = u_{x,i}v_{i+1} + u_i v_{x,i}$$
 (5.2.9)

b) Die Formel der partiellen Summation (= diskrete Formel der part. Integration):

$$(v_x, u) = -(v, u_{\overline{x}}] + u_n v_n - u_0 v_1$$
 (5.2.10)

mit
$$(u,v) := \sum_{i=1}^{n-1} h v_i u_i$$
 , $(u,v] := \sum_{i=1}^n h v_i u_i$

- Man zeige, daß für beliebige Gitterfunktionen $u, v: \overline{\omega}_h \mapsto \mathbf{R}^1$ mit $u_0 = u_n = 0$, $v_0 = v_n = 0$ die folgenden Formeln gelten :
 - a) Die diskrete Greensche Formel:

$$-(u_{\overline{x}x}, v) = (u_{\overline{x}}, v_{\overline{x}}] \qquad \forall u, v \in \overset{\circ}{W} \,_{2}^{1}(\omega_{h})$$
 (5.2.11)

d.h. für u = v gilt also insbesondere

$$-(u_{\overline{x}x}, u) = ||u_{\overline{x}}||^2 := \sum_{i=1}^n h u_{\overline{x},i}^2$$

b) Die diskrete Friedrichs-Ungleichung:

$$\parallel u \parallel \leq c_F \parallel u_{\overline{x}} \mid] \tag{5.2.12}$$

mit
$$||u||^2 = \sum_{i=1}^{n-1} h u_i^2 = (u, u)$$
 und $c_F = \frac{b-a}{\sqrt{2}}$.
Hinweis: $u_i = \sum_{j=1}^{n} u_{\overline{x},j} h = \sum_{j=1}^{n} (u_j - u_{j-1}) = u_i - u_0 = u_i$

c) Die diskrete C-Einbettung (1D-Fall!!):

$$\| u \|_{C(\omega_h)} = \max_{i=\overline{1,n-1}} |u_i| \le \sqrt{b-a} \| u_{\overline{x}} |$$
 (5.2.13)

- $\ddot{\text{U}}10^*$ Man gebe die optimale (kleinste !) Konstante c_F in der diskreten Friedrichsungleichung (5.2.12) an !
- $\ddot{\mathrm{U}}11$ Man zeige, daß das σ -gewichtete Differenzenschema (Fehlerschema)

$$z_{t,i}^{j} - \sigma z_{\overline{x}x,i}^{j+1} - (1 - \sigma) z_{\overline{x}x,i}^{j} = \psi_{i}^{j} , \quad i = \overline{1, n-1}, j = \overline{0, m-1}$$
(5.2.14)
$$RB : \qquad z_{0}^{j} = z_{n}^{j} = 0 \qquad \forall j = \overline{1, m}$$

$$AB : \qquad z_{i}^{0} = w_{i}^{0} \qquad \forall i = \overline{0, n} \quad w^{0} \text{ (ist geg. Störung der AB)}$$

im energetischen Raum H_A ($H=L_2^0(\overline{\omega}_h)=L_2(\omega_h)$, $\overline{\omega}_h=\{x_i=\imath h:\imath=\overline{1,n}\}$ $A=\overline{A}$, $\overline{A}v:=-v_{\overline{x}x}$) stabil bzgl. AB und RS ist und daß die a-priori-Abschätzungen

a)
$$\|z_{\overline{x}}^{j+1}\| \le \|z_{\overline{x}}^{0}\| + T \max_{k=0,j} \|\psi_{\overline{x}}^{k}\|$$

$$\mathrm{b}) \quad \parallel z^{\jmath+1} \parallel_{C(\omega_h)} \ := \ \max_{\imath = \overline{1, n-1}} |z_{\imath}^{\jmath+1}| \quad \leq \quad \parallel z_{\overline{x}}^{0}| \Big] \ + \ T \ \max_{k = \overline{0, \jmath}} \ \parallel \psi_{\overline{x}}^{k}| \Big]$$

gelten, falls $1 \geq \sigma \geq \sigma_0 := \frac{1}{2} - \frac{h^2}{4\tau}$ gilt.

<u>Hinweis</u>: 1) Satz 2.16 aus der Vorlesung anwenden!

- 2) Beziehungen (5.2.11) und (5.2.13) benutzen!
- 3) Zeigen, daß

$$\|(I+\sigma\tau A)^{-1}\| \stackrel{!}{=} \text{Spektral norm } ([I+\sigma\tau A]^{-1}) \le 1$$

5.2.2 FEM-Galerkin-Verfahren für parabolische ARWA

Ü12 | Man beweise den <u>Satz 2.8 von Picard und Lindelöf</u>:

Vor.: Betrachten AWA für folgendes System gewöhnlicher Differentialgleichungen

$$\frac{dy(t)}{dt} + A(t) \underline{y}(t) = \underline{f}(t) \qquad \text{f.\"{u.}} \ t \in \mathbf{T} = (0, T) ,$$

$$\underline{y}(0) = \underline{y}_0$$
 (5.2.15)

mit gesuchter Vektorfunktion $\underline{y} = \underline{y}(t) = (y_1(t), \dots, y_N(t))^T$, gegebenen Anfangswerten $\underline{y}_0 = (y_{01}, \dots, y_{0N})^T \in \mathbf{R}^N$, und gegebenen Daten $A = A(t) = [a_{ki}(t)]_{k,i=\overline{1,N}} : a_{ki}(\cdot) \in L_{\infty}(0,T)$ sowie $\underline{f}(t) = (f_1(t), \dots, f_N(t))^T \in [L_2(0,T)]^N$.

 $\underline{\mathbf{Bh.:}} \ \exists ! \, \underline{y}(t) \in \left[W_2^1(0,T) \right]^N \ : (5.2.15)$

Hinweise:

1) Übergang zur äquivalenten Igl.: $\int\limits_0^t \left(\ \right) d\xi$, d.h.

$$\underline{y}(t) = -\int_{0}^{t} A(\xi) \underline{y}(\xi) d\xi + \int_{0}^{t} \underline{f}(\xi) d\xi + \underline{y}_{0}$$

$$y = B y \quad (= \text{Fixpunktgleichung})$$

$$\begin{array}{ll} \text{mit} & B \,:\, \mathcal{X} \,=\, \left[C(\overline{T})\right]^N \,\mapsto\, \left[W_2^1(0,T)\right]^N \subset \mathcal{X} \\ (\mathcal{X} \,\text{-}\, \text{Banachraum} \stackrel{?}{\Longrightarrow} \|\cdot\|_{\mathcal{X}}). \end{array}$$

2) Wenden Sie auf die Fixpunktgleichung

Ges.:
$$\underline{y} \in \mathcal{X}$$
 : $\underline{y} = B \underline{y}$ in \mathcal{X}

den verallgemeinerten Banachschen Fixpunktsatz I.2.5 an (vgl. auch Bsp. I.2.4)!

Ü13 Man zeige am Beispiel linearer Dreieckselemente, daß die Massenmatrix M_h gut konditioniert ist, falls die Triangulation im Sinne der Definition II.4.3 (vgl. Ü.II.4.4) regulär ist, d.h. $\kappa(M_h) = \lambda_{max}(M_h)/\lambda_{min}(M_h) = O(1)$!

Hinweise:

1) Definition der Massenmatrix:

$$(M_h \underline{u}_h, \underline{v}_h) = \int_{\Omega} u_h(x) v_h(x) dx \qquad \forall \underline{u}_h, \underline{v}_h \in \mathbf{R}^{N_h} ,$$

$$\mathbf{R}^{N_h} \ni \underline{u}_h, \underline{v}_h \longleftrightarrow u_h, v_h \in \mathbf{V}_{0h}$$
 (FE-Raum),

wobei (\cdot,\cdot) das Euklidische Skalarprodukt im R^{N_h} ist.

2) Wdh. Def. II.4.3 (Reguläre Triangulation):

Wir nennen die Familie $\{\tau_h\}_{h\in\Theta}$ der Triangulation $\tau_h=\{\delta_r:r\in\mathbf{R}_h\}$ regulär, falls positive, h-unabhängige Konstanten $\underline{c}_1,\overline{c}_1,c_2,c_3=\mathrm{const.}>0$ existieren :

$$\frac{c_1 h^d}{\|J_{\delta_r}\|} \leq \overline{c}_1 h^d \qquad \forall \xi \in \overline{\triangle} \ \forall r \in \mathbf{R}_h \ \forall h \in \Theta$$

$$\|J_{\delta_r}\| := \sqrt{\lambda_{max} (J_{\delta_r}^T J_{\delta_r})} \leq c_2 h \qquad \forall \xi \in \overline{\triangle} \ \forall r \in \mathbf{R}_h \ \forall h \in \Theta$$

$$\|J_{\delta_r}^{-T}\| := \sqrt{\lambda_{max} (J_{\delta_r}^{-1} J_{\delta_r}^{-T})} \leq c_3 h^{-1} \qquad \forall \xi \in \overline{\delta_r} \ \forall r \in \mathbf{R}_h \ \forall h \in \Theta ,$$

$$\text{wobei } J_{\delta_r} = \frac{dx_{\delta_r}(\xi)}{d\xi} \ , \ J_{\delta_r}^{-1} = \frac{d\xi_{\delta_r}(x)}{dx} \ , \ \delta_r \ \stackrel{\xi = \xi_{\delta_r}(x)}{\longrightarrow} \ \triangle \ \text{und} \ \delta_r \ \stackrel{x = x_{\delta_r}(\xi)}{\longleftarrow} \ \triangle \ .$$

- 3) Zu zeigen ist also : $\exists \gamma_1, \gamma_2 = \text{const.} > 0$: $\gamma_i \neq \gamma_i(h)$: $\gamma_1 h^d \leq \lambda_{min}(M_h)$ und $\lambda_{max}(M_h) \leq \gamma_2 h^d$. Gehen Sie analog zum Beweis von Satz II.4.4 vor, in dem die Eigenwerte der Steifigkeitsmatrix K_h abgeschätzt werden.
- 4) Geben sie für das Beispiel der gleichmäßigen Triangularisierung des Einheitsquadrates γ_1 und γ_2 an.



Ü14 Man zeige, daß die Stabilitätsbedingung

$$B = M_h + \tau \sigma K_h \geq \frac{\tau}{2} A = \frac{\tau}{2} K_h$$

i.S. $(B\underline{u},\underline{u}) \geq \frac{\tau}{2}(A\underline{u},\underline{u}) \quad \forall \underline{u} \in \mathbf{R}^N$ für das Schema $(\underline{27})_{h\tau}$ aus der Vorlesung NuIII, Pkt. 2.3.2, äquivalent ist zur Bedingung

$$\sigma \geq \frac{1}{2} - \frac{1}{\tau \lambda_{max}} ,$$

wobei $\lambda_{max} = \text{Maximaler Eigenwert des verallgemeinerten Eigenwertproblems}$ $K_h \underline{u}_h = \lambda M_h \underline{u}_h,$

 $M_h = \text{Massenmatrix},$

 K_h = Steifigkeitsmatrix ist.

Die Steifigkeitsmatrix wurde hier als t-unabhängig vorausgesetzt!

 $\mathbf{P} \ \mathbf{VI}$

Praktikum VI

5.2.3 Konsultation zur Praktikumsaufgabe Px

- 5.2.3.1 Konsultation zu den parabolischen ARWA
 - Theorie : siehe Vorlesung

5.2.3.2 Hyperbolische ARWA

Ü15 Geben Sie analog zu den parabolischen ARWA die Linienformulierung der hyperbolischen ARWA der nichtgedämpften Schwingungen einer fest eingespannten Membran

Ges.
$$u(x,t)$$
:
$$\frac{\partial^2 u}{\partial t^2} - \triangle u(x,t) = f(x,t), \quad x \in \Omega \subset \mathbf{R}^2, t \in \mathbf{T} = (0,T)$$

$$+ \text{RB} : u(0,t) = 0, \quad \forall x \in \partial\Omega ; \forall t \in \mathbf{T}$$

$$+ \text{AB} : u(x,0) = u_0(x), \quad x \in \overline{\Omega} \quad \text{(Kompatibilität !)}$$

$$\frac{\partial u}{\partial t}(x,0) = u_1(x), \quad x \in \overline{\Omega} \quad .$$

an und leiten Sie mit der vertikalen Linienmethode (FE-Ortsdiskretisierung) die semidiskrete Ersatzaufgabe her!

 $\ddot{\mathrm{U}}16^*$ Betrachten Sie das System gewöhnlicher Differentialgleichungen 2. Ordnung

Ges.
$$u(t) = [u_1(t), ..., u_N(t)]^T \in [C^2(I)]^N$$
 :
$$M \ddot{u}(t) + C \dot{u}(t) + K u(t) = f(t) , t \in I = \overline{\mathbf{T}} = [0, T]$$
 mit geg. AB :
$$u(0) = u_0 , \dot{u}(0) = u_1 .$$

Hierbei sind $f \in [C(I)]^N$ die gegebene rechte Seite und

$$\begin{array}{ll} M,C,K & -(N\times N) \text{ - Matrizen :} \\ M=M^T \text{ p.d.} & -\text{Massenmatrix,} \\ C & -\text{D\"{a}mpfungsmatrix,} \\ K=K^T \text{ p.d.} & -\text{Steifigkeitsmatrix.} \end{array}$$

Schreiben Sie die folgenden Zeitintegrationsschemata auf und motivieren Sie diese:

- a) σ -gewichtetes dreischichtiges Differenzenschema (siehe Vorlesungsmanuskript)
- b) NEWMARK- β -Methode (Kikuchi[1], S. 165-169),
- c) Wilson-Θ-Methode (Kikuchi[1], S. 165-169),
- d) Rückführung auf ein System 1. Ordnung und Anwendung von Integrationsverfahren für Systeme 1. Ordnung soweit bekannt.

5.3 Numerische Behandlung von AWA für gewöhnliche Dgl.

5.3.1 Einschrittverfahren

Ü17 Der Oregonator wird durch das Reaktionsschema

$$\begin{array}{cccc} \operatorname{BrO}_3^- + \operatorname{Br}^- & \xrightarrow{K_1} & \operatorname{HBrO}_2 \\ \operatorname{HBrO}_2 + \operatorname{Br}^- & \xrightarrow{K_2} & \operatorname{P} \\ \operatorname{BrO}_3^- + \operatorname{HBrO}_2 & \xrightarrow{K_3} & \operatorname{2HBrO}_2 + \operatorname{Ce}(\operatorname{IV}) \\ & & & & \\ \operatorname{2HBrO}_2 & \xrightarrow{K_4} & \operatorname{P} \\ & & & & \\ \operatorname{Ce}(\operatorname{IV}) & \xrightarrow{K_5} & \operatorname{Br}^- \end{array}$$

mit gegebenen Reaktionsgeschwindigkeitskoeffizienten K_1, K_2, K_3, K_4, K_5 beschrieben. Stellen Sie das Differentialgleichungssystem für die gesuchten Konzentrationen

$$\begin{array}{rcl} c_{1}(t) & = & c_{\mathrm{BrO}_{3}^{-}}(t) \\ c_{2}(t) & = & c_{\mathrm{Br}^{-}}(t) \\ c_{3}(t) & = & c_{\mathrm{HBrO}_{2}}(t) \\ c_{4}(t) & = & c_{\mathrm{P}}(t) \\ c_{5}(t) & = & c_{\mathrm{Ce(IV)}}(t) \end{array}$$

auf, und stellen Sie sinnvolle Anfangsbedingungen zum Zeitpunkt $t=t_0=0$!

O Bemerkung:

Im Praktikum P IX (18.12.1996) werden wir das abgeleitete Dgl.-System numerisch integrieren und damit den Oregonator numerisch simulieren.

Ü18 Die Robertson-Reaktion wird durch das Reaktionsschema

$$\begin{array}{ccc} A & \xrightarrow{0.04} & B & \text{(langsame Reaktion)} \\ B+B & \xrightarrow{3.10^7} & C+B & \text{(sehr schnell)} \\ B+C & \xrightarrow{10^4} & A+C & \text{(schnell)} \end{array}$$

beschrieben. Stellen Sie das Dgl.-System zur Bestimmung der Konzentrationen $c_A(t)$, $c_B(t)$, $c_C(t)$ auf, und schreiben Sie geeignete Anfangsbedingungen zum Zeitpunkt $t = t_0 = 0$ vor!

O Bemerkung:

Im Praktikum P IX (18.12.1996) werden wir das abgeleitete Dgl.-System numerisch integrieren und damit die Robertson-Reaktion numerisch simulieren.

Ü 19 Man beweise die folgende Konvergenzaussage zum Eulerschen Polygonzugverfahren (EPZV) auf äquidistantem Gitter (Satz 3.4 aus der Vorlesung):

Vor.:

1.)
$$u(t): u'(t) = f(t, u(t)), t \in I = [0, T]$$
 $u(0) = g_0$

3.)
$$f: D \to \mathbb{R}^N$$
, wobei $D = I \times \mathbb{R}^N = \{(t, v) : t \in I = [0, T], | v - u_0 | < \infty = b\}$
4.) $| f(t, v) - f(t, \overline{v}) | \le L | v - \overline{v} | \forall (t, v), (t, \overline{v}) \in D.$

Dann gilt die Fehlerabschätzung

$$|u(t_j) - u_h(t_j)| \le e^{Lt_j} \left[|\delta| + \frac{\tau}{L} \right],$$

mit $\tau = \max_{j=1,\dots,m} |\tau_j|$ und den lokalen Abschneidefehlern (siehe Vorlesung)

$$\tau_j \equiv \tau_h(t_j) = \frac{1}{h} \left[\frac{u(t_j) - u(t_{j-1})}{h} - f((t_{j-1}, u(t_{j-1}))) \right]$$

O Hinweis:

a.) Schreiben Sie unter Benutzung der Darstellung (*)

$$u(t_{j+1}) = u(t_j) + hf(t_j, u(t_j)) + h\underbrace{\left[\frac{u(t_{j+1}) - u(t_j)}{h} - f((t_j, u(t_j))\right]}_{=:\tau_h(t_{j+1}) = \tau_{j+1}}$$

und des EPZV (**)

$$u_{j+1} = u_j + h f(t_j, u_j)$$

eine Rekursionsbeziehung ((*) - (**)!) für den Fehler

$$e_{i+1} = u(t_{i+1}) - u_h(t_{i+1}) = u(t_{i+1}) - u_{i+1}$$

auf, und schätzen Sie diese ab.

b.) Verwenden Sie dabei die elementare Beziehung

$$(1+hL)^{j+1} \le e^{(j+1)hL} = e^{Lt_{j+1}} \le e^{LT}$$

c.) Interpretieren Sie die zu beweisende Behauptung!

P VIII

Praktikum VIII

Ü20 | V

Verwendet man anstelle der Mittelpunktsregel die <u>Trapezregel</u> (TR) zur Berechnung

des Integrals $\int\limits_t^{t+h} f(s,u(s))\,ds$ so erhält man das Verfahren von Heun :

$$\int_{t}^{t+h} f(s, u(s)) ds \overset{\text{TR}}{\approx} \frac{h}{2} \left[f(t, u(t)) + f(t+h, u(t+h)) \right]$$

Damit ergibt sich für das Verfahren von HEUN das folgende Tableau:

$$\begin{array}{c|cccc}
0 & & & \\
1 & 1 & & \\
\hline
& 1/2 & 1/2 & \end{array} ,$$

d.h. das Verfahren von Heun ist eine 2-stufige Runge-Kutta-Formel und es gilt:

$$u(t+h) = u(t) + \frac{h}{2} [f(t,u) + f(t+h, u+hf(t,u))].$$

Man zeige durch Taylor-Entwicklung des lokalen Fehlers $u(t+h)-u_h(t+h)$, daß das Verfahren von Heun die Konsistenzordnung 2 besitzt !

 $\ddot{\mathrm{U}}21$

Für die autonome Differentialgleichung

$$\begin{cases} v'(t) &= g(v(t)) \\ v(0) &= 0 \end{cases}$$

sei eine ℓ -stufige Runge-Kutta-Formel durch das Tableau (Koeffizienten $\{c_i\}_{i=\overline{2,\ell}}$ kommen im Falle autonomer Dgl. nicht vor !)

gegeben und habe die Ordnung p.

Für eine allgemeine Differentialgleichung der Form

$$\begin{cases} u'(t) &= f(t, u(t)) \\ u(0) &= u_0 \end{cases}$$

ist dann eine ℓ -stufige Runge-Kutta-Formel der Ordnung p durch das Tableau

gegeben, falls

$$c_i = \sum_{j=1}^{i-1} a_{ij} , i = \overline{2,\ell} .$$
 (5.3.16)

gilt. Zu zeigen ist also, daß die Runge-Kutta-Formel (*) für die autonome Dgl. dieselben Näherungen liefert wie (**) für eine allg. Dgl., falls (5.3.16) gilt.

Ü22 Man zeige, daß das klassische Runge-Kutta-Verfahren

die Konsistenzordnung 4 besitzt!

P IX

Praktikum IX

5.3.2 Praktische Durchführung von Einschrittverfahren

 \times Man studiere den gleichnamigen Punkt 3.3.5 im Vorlesungsskriptum Numerik III.

5.3.3 Einfache numerische Experimente mit Einschrittverfahren

E5 Wir betrachten das Anfangswertproblem

 $\text{mit } B = 3 \text{ und } A = 1 \cdot (e^B - 1) .$

a) Man bestimme die exakte Lösung von (5.3.17) analytisch!

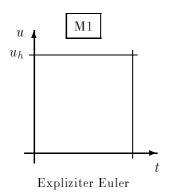
b) Man löse (5.3.17) mit den folgenden expliziten Runge-Kutta-Verfahren Mx unter Verwendung der Schrittweiten h=1, $\frac{1}{2}$, $\frac{1}{4}$, $\frac{1}{8}$ und $\frac{1}{16}$;

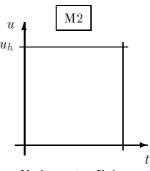
M1 Explizites Euler-Verfahren,

M2 Verbessertes Euler-Verfahren,

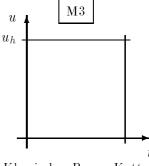
M3 "klassisches" Runge-Kutta-Verfahren der Ordnung 4 (siehe Ü22).

Man veranschauliche die Ergebnisse graphisch durch Vergleich der numerischen Lösungen für h=1, $\frac{1}{2}$, $\frac{1}{4}$, $\frac{1}{8}$, $\frac{1}{16}$ mit der analytischen Lösung:





Verbesserter Euler



Klassisches Runge-Kutta

c) Man trage in die folgende Tabelle $u_h(1)$ als Approximation von $u(1)=\dots$ für die betrachteten Verfahren Mx ein :

	h	\longrightarrow					
Mx \		1	1/2	1/4	1/8	1/16	u(1)
M1	(*)						
	(E)					_	
M2	(*)						
	(E)					_	
М3	(*)						
	(E)					_	
MxE							

wobei (*) ursprüngliches Verfahren

(E) Verbesserung der Werte durch globale Extrapolation (siehe S. 115/116 im Skriptum) nach der Formel

$$\hat{u}_h(t) = u_{h/2}(t) + \frac{u_{h/2}(t) - u_h(t)}{2^p - 1}$$

mit p = Ordnung des Verfahrens.

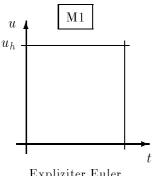
(MxE)Verfahren Mx für ein $x \in \{1, 2, 3\}$ mit lokaler Extrapolation (siehe S. 115/116 im Skriptum).

E6 Wir betrachten das Anfangswertproblem

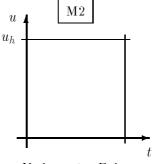
$$u'(t) = -50 (u(t) - \cos(t)) , t \in I = [0, 1.5] u(0) = 0 .$$
 (5.3.18)

- a) Man bestimme die exakte Lösung von (5.3.18) analytisch!
- b) Man löse (5.3.18) mit den folgenden Runge-Kutta-Verfahren | Mx | unter Verwendung der Schrittweiten $h = \frac{1}{20}$, $\frac{3}{80}$, $\frac{1}{30}$, $\frac{1}{40}$;
 - M1Explizites Euler-Verfahren,
 - M2Verbessertes Euler-Verfahren,
 - M3Implizites Euler-Verfahren.

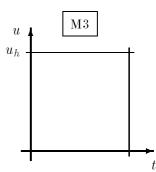
Man veranschauliche die Ergebnisse graphisch durch Vergleich der numerischen Lösungen für $h=\frac{1}{20}$, $\frac{3}{80}$, $\frac{1}{30}$, $\frac{1}{40}$ mit der analytischen Lösung :



Expliziter Euler



Verbesserter Euler



Impliziter Euler

- c) Man starte die Explizite Eulermethode bei $t_0=1/2$ mit der exakten Lösung $u(t_0)\equiv u(1/2)$ unter Verwendung der Schrittweite h=1/20 und stelle das Resultat wieder im Vergleich mit der analytischen Lösung graphisch dar. Ab welcher Schrittweite muß das Verfahren stabil werden?
- d) Lösen Sie (5.3.18) nochmals mit dem impliziten Euler-Verfahren und der Schrittweite h=1/2. Stellen Sie das Resultat wieder im Vergleich mit der analytischen Lösung graphisch dar.
- e) Man trage in die folgende Tabelle $u_h(1.5)$ als Approximation von $u(1.5) = \dots$ für die betrachteten Verfahren Mx ein :

h					
Mx \	1/20	3/80	1/30	1/40	u(1.5)
M1					
M2					
M3					

5.3.4 Numerische Lösung komplizierter Anfangswertprobleme

E7 Die Bahn eines Satelliten in der Ebene des Erde-Mond-Systems läßt sich durch das folgende System von Dgl. 2. Ordnung beschreiben (vgl. Bsp. 3.5 der Vorlesung Numerik III und [2, 3]):

$$y_1'' = y_1 + 2y_2' - (1 - \mu) \frac{y_1 + \mu}{D_1} - \mu \frac{y_1 - (1 - \mu)}{D_2} ,$$

$$y_2'' = y_2 - 2y_1' - (1 - \mu) \frac{y_2}{D_1} - \mu \frac{y_2}{D_2} , t \in [0, T] ,$$

$$+ AB : y_1(0) = 0.994$$

$$y_1'(0) = 0$$

$$y_2(0) = 0$$

$$y_2(0) = -2.00158510637908252240537862224 ,$$

$$mit D_1 = \left[(y_1 + \mu)^2 + y_2^2 \right]^{3/2} ,$$

$$D_2 = \left[(y_1 - (1 - \mu))^2 + y_2^2 \right]^{3/2} ,$$

$$\mu = 0.012277471$$

Für diese Daten ergibt sich eine periodische Lösung mit der Periode

$$t_{per} = t = \underline{17.06521656015796}25588917206249$$

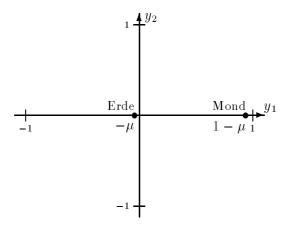
- 1) Man überführe (5.3.19) in eine äquivalente AWA für ein System gewöhnlicher Dgl. 1. Ordnung! Ist dieses System autonom?
- 2) Man löse dieses System numerisch mit den folgenden Integrationsverfahren:

M1 Klassisches Runge-Kutta-Verfahren mit dem Tableau

mit h = T/m, m = 6000 und m = ? (eigene Wahl).

M2 Verfahren eigener Wahl!

3) Man stelle die Lösungstrajektoren $(y_{1h}(t), y_{2h}(t))$, $t \in I_h = \{t_0, t_1, \ldots, t_m\}$, linear interpoliert, grafisch dar!



E8 Die Robertson-Reaktion [4] wird durch das Dgl.-System

$$c'_{A}(t) = -0.04 c_{A}(t) + 10^{4} c_{B}(t) c_{C}(t) ,$$

$$c'_{B}(t) = 0.04 c_{A}(t) - 3 \cdot 10^{7} c_{B}^{2}(t) - 10^{4} c_{B}(t) c_{C}(t) ,$$

$$c'_{C}(t) = 3 \cdot 10^{7} c_{B}^{2}(t) , t \in [0, T], T = 0.3$$

$$+ AB : c_{A}(0) = 1 c_{B}(0) = 0 c_{C}(0) = 0$$

$$(5.3.20)$$

beschrieben (vgl. \square). Aufgrund der Größenordnungsunterschiede in den Reaktionsgeschwindigkeitskoeffizienten ist zu erwarten, daß das Dgl.-System (5.3.20) steif ist. Lösen Sie die AWA (5.3.20) mit einem geeigneten Integrationsverfahren und stellen Sie den zeitlichen Verlauf der Konzentrationen $c_A(t)$, $c_B(t)$ und $c_C(t)$ (getrennt) grafisch dar!

E9 | Der <u>Oregonator</u> [2] wird durch das Dgl.-System

$$c'_{1}(t) = -k_{1}c_{1}c_{2} - k_{3}c_{1}c_{3}$$

$$c'_{2}(t) = -k_{1}c_{1}c_{2} - k_{2}c_{2}c_{3} + k_{5}c_{5}$$

$$c'_{3}(t) = k_{1}c_{1}c_{2} - k_{2}c_{2}c_{3} + k_{3}c_{1}c_{3} - 2k_{4}c_{3}^{2}$$

$$c'_{4}(t) = k_{2}c_{2}c_{3} + 2k_{4}c_{3}^{2}$$

$$c'_{5}(t) = k_{2}c_{2}c_{3} - k_{5}c_{5} \quad , t \in [0, T],$$

$$+ AB : c_{1}(0) = 1/2 \quad c_{2}(0) = 1/2 \quad c_{3}(0) = 0$$

$$c_{4}(0) = 0 \quad c_{5}(0) = 0,$$
wobei
$$k_{1} = 1.34 \quad k_{2} = 1.6 \cdot 10^{9} \quad k_{3} = 8.0 \cdot 10^{3}$$

$$k_{4} = 4.0 \cdot 10^{7} \quad k_{5} = 1.0$$

$$(5.3.21)$$

beschrieben (vgl. Ü17). Aufgrund der Größenordnungsunterschiede in den Reaktionsgeschwindigkeitskoeffizienten ist zu erwarten, daß das Dgl.-System (5.3.21) steif ist.

Lösen Sie die AWA (5.3.21) mit einem geeigneten (?) Integrationsverfahren und stellen Sie den zeitlichen Verlauf der Konzentrationen $c_1(t)$, $c_2(t)$, $c_3(t)$, $c_4(t)$ und $c_5(t)$ (getrennt) grafisch dar !

Literaturverzeichnis

- [1] N. Kikuchi. Finite Element Method in Mechanics. Cambridge University Press, 1986.
- [2] P. Deuflhard. Numerische Mathematik II: Integration gewöhnlicher Differentialgleichungen. de Gruyter, Berlin, 1994. ISBN 3-11-013977-5.
- [3] E. Hairer, S.P. Nørsett, and G. Wanner. Solving Ordinary Differential Equations 1: Non-stiff Problems. Springer, 1993.
- [4] E. Hairer and G. Wanner. Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Equations, volume 14 of Springer Series in Computational Mathematics. Springer, 1996. ISBN 3-540-60452-9.

Kapitel 6

Ergebnisse der numerischen Experimente

Für die Überlassung der folgenden Ergebnisgraphiken danken wir :

Gerald Fehringer, Franz Gruber, Ferdinand Kickinger, Roswitha Kroiss, Joachim Schöberl