

MAX-PLANCK-INSTITUT FÜR PLASMAPHYSIK
GARCHING BEI MÜNCHEN

DIFFERENZENSCHEMATA FÜR HYPERBOLISCHE SYSTEME
(DIFFERENCE SCHEMES FOR HYPERBOLIC SYSTEMS)

W. Höhn

IPP 6/109

July 1972

Die nachstehende Arbeit wurde im Rahmen des Vertrages zwischen dem Max-Planck-Institut für Plasmaphysik und der Europäischen Atomgemeinschaft über die Zusammenarbeit auf dem Gebiete der Plasmaphysik durchgeführt.

Abstract

In the case of non-strict hyperbolic differential systems the corresponding difference schemes usually involve amplification matrices increasing like polynomials. For this weak stability an equivalence theorem, similar to that of LAX, is established for initial value problems with time-dependent linear solution operators.

Difference schemes for hyperbolic systems using leap-frog differences - one of which is stable in the widest possible range - are studied and their stability or weak stability is investigated by direct evaluation of powers of certain normal forms of their amplification matrices.

Some numerical results demonstrate the practical usefulness of these schemes.

Inhaltsverzeichnis

Einführung	1
I. Ein Äquivalenzsatz	4
II. Ein 2-Schritt-Differenzenschema	17
A. Konsistenz	18
B. Stabilität im strikt hyperbolischen Falle	20
C. Stabilitätsbereich und Abhängigkeitsgebiet	27
III. Schwache Stabilität	33
IV. Ausweitung auf ähnliche Schemata	66
V. Numerische Resultate	69
Literaturverzeichnis	75

Einführung

Zur numerischen Behandlung von Anfangswertproblemen werden in der Praxis gern Differenzenschemata herangezogen. Diese Methode sichert meist einen verhältnismäßig geringen Aufwand, etwa durch ihre Schnelligkeit oder bequeme Handhabung beim Rechnen.

Für den Beweis der Konvergenz von Lösungen des Differenzenschemas zu denen des entsprechenden Anfangswertproblems hat sich im linearen Falle die Theorie von LAX und RICHTMYER durchgesetzt (s. z. B. [5], [8], [1]), deren Kernstück, der Äquivalenzsatz ([8] S45, 127; [1] S62), einen Zusammenhang zwischen Stabilität und Konvergenz herstellt. Im Konvergenzbegriff sind verallgemeinerte Lösungen des Anfangswertproblems mit einbezogen.

Um in der Praxis sinnvoll mit Differenzenschemata umgehen zu können, genügt aber nicht nur ihre Konvergenz allein, vielmehr benötigt man eine ausreichende Geschwindigkeit der Konvergenz, da merkliche Verkleinerung der Schrittweite h vor allem bei Problemen mit mehreren Raumdimensionen einen enormen Mehrbedarf an Speicherplatz bedeutet, oder aber das Verfahren durch Rundungsfehler ohne Erfolg bleibt.

Die Geschwindigkeit der Konvergenz läßt sich den praktischen Verhältnissen angemessen am einfachsten durch den Exponenten τ im Ausdruck $K(u) \cdot h^\tau$ angeben, den man als Schranke für den Fehler beweisen muß. Die Konstante $K(u) \in \mathbb{R}^+$ ist abhängig von der zu approximierenden exakten Lösung. Damit ist die unangenehmste Tatsache für die Behandlung verallgemeinerter Lösungen bereits angesprochen, denn diese lassen sich im allgemeinen nur durch Grenzwerte von Folgenfolgen $(u_\ell)_{\ell \in \mathbb{N}}$ gewinnen. Daß für jedes Folgeelement der Fehler

durch $K(u_\epsilon)h^2$ abgeschätzt werden kann, ist für den Grenzwert nutzlos da die Folge $(K(u_\epsilon))_{\epsilon \in \mathcal{N}}$ möglicherweise stark divergiert; man kann lediglich mit den bekannten Methoden überhaupt Konvergenz beweisen.

Unter diesem Gesichtspunkt scheint es gerechtfertigt, derartige Geschwindigkeitsuntersuchungen auf "klassische" Lösungen zu beschränken, solche Lösungen, die einer Konsistenzbedingung genügen. Gegenüber der LAX-RICHTMYER Theorie hat man dann die Begriffe Konvergenz und Konsistenz nicht mehr für eine Gesamtheit verallgemeinerter Lösungen zu betrachten sondern für eine einzige, feste Lösung, den Begriff Stabilität dagegen sehr wohl. Mit den so veränderten Definitionen wird im Abschnitt I ein Äquivalenzsatz auf analoge Weise bewiesen. Als Verallgemeinerungen berücksichtigt er zeitabhängige Probleme, schwache Stabilität und Geschwindigkeit der Konvergenz.

Im II. Kapitel werden diese Ergebnisse für ein von 2. Ordnung konsistentes Differenzenschema für zunächst strikt hyperbolische Systeme 1. Ordnung mit konstanten Koeffizienten verwendet, in welchen alle Ableitungen durch "Laubfrosch"-Differenzen ersetzt sind. Der Stabilitätsbereich des Schemas wird berechnet, und es wird gezeigt, daß er der größte überhaupt mögliche ist.

Schwach stabil wird dieses Schema für nicht mehr strikte sondern gewöhnlich hyperbolische Systeme. Dieser Fall wird im III. Abschnitt anhand einer für die Praxis brauchbaren Normalform für nicht mehr diagonalisierbare Amplifikationsmatrizen untersucht und die Stabilitätskonstante berechnet.

Im IV. Kapitel werden diese Ergebnisse auf andere 2-Schritt Verfahren verallgemeinert; als Beispiel wird der Stabilitätsbereich eines von 3. Ordnung konsistenten Schemas abgeschätzt.

Den Abschluß bilden numerische Resultate.

I. Ein Äquivalenzsatz

In einem BANACH-Raum $(B, \|\cdot\|)$, sei ein lineares Anfangswertproblem gegeben:

$$(1) \quad \frac{d}{dt} u(t) = A(t)u(t) \quad \text{für } t \in [0, T] \quad ; \quad u(0) = u_0$$

mit einer Schar linearer Operatoren $A(t): B' \rightarrow B$ wobei $B' \subset B, u(t) \in B'$ für $t \in [0, T]$, und zwar mit einem $T \in \mathbb{R}^+$.

Für alle u_0 aus einem $D, D \subset B$, existiere eine eindeutige Lösung $u(t)$ von (1). Für $t \in [0, T]$ sei $\{E(s)_t; s \in [t, T]\}$ eine Schar von Operatoren $E(s)_t: D_t \rightarrow B$ definiert auf:

$$D_t := \{u_t \in B; \text{es gibt ein } u_0 \in D \text{ mit } u(t) = u_t\}$$

durch:

$$E(s)_t u_t := u(s+t) \quad \text{für alle } s, t, s+t \in [0, T]$$

u_0 ist darin dem jeweiligen u_t zugehörig, $u(t)$ die Lösung von (1) zu u_0 .

Entsprechend zur Linearität von (1) ist $E(s)_t$ linear.

Falls $A(t)$ von der Zeit unabhängig ist, erfüllen die Lösungsoperatoren die Halbgruppeneigenschaft (S.[1], S.2,7), im zeitabhängigen Fall gilt das i.a. nicht mehr, sondern dagegen nur (auf D_t):

$$E(r)_s E(s)_t = E(r+s)_t \quad \text{für alle } r, s, t, r+s+t \in [0, T]$$

Definition:

Seien $k \in \mathbb{N}; B^k := B \times B \times \dots \times B = \bigotimes_{v=1}^k B$ und $\{C(h)_t; h \in [0, T]\}$ eine Schar linearer und stetiger Operatoren $C(h)_t: B^k \rightarrow B$ für $t \in [0, T]$.

$$(2) \quad U_h^{n+1} = C(h)_{n,h} [U_h^n, U_h^{n-1}, \dots, U_h^{n-k+1}]$$

mit vorgegebenen $U_h^m, U_h^{m-1}, U_h^{m-2}, \dots, U_h^{m-k+1} \in B, m \in \mathbb{N}_0, m \geq k-1$

ist ein lineares k -Schritt Differenzenverfahren (-Schema) mit h als Schrittweite und $n \in \mathbb{N}_0, n \geq m, n \cdot h \in [0, T]$ als Variable. $(\mathbb{N}_0 := \mathbb{N} \cup \{0\})$

Im folgenden wird für jedes $u_0 \in D = D_0$ die Lösung des Anfangswertproblems (1) mit $u(t) = E(t)_0 u_0$ bezeichnet, mit $U_h^{m,n}$ die des Differenzenschemas (2) zu den Anfangswerten $U_h^{m,m}, U_h^{m,m-1}, \dots, U_h^{m,m-k+1} \in B$.

Der Parameter m ermöglicht den Beginn der Approximation mit dem Schema nicht nur zur Zeit $t = 0$ sondern für beliebige $t \in [0, T]$.

Das Schema (2) läßt sich übersichtlicher als Operation in B^k

schreiben, dessen Elemente eine " \sim " bekommen.

$$\tilde{u} := (u_k, u_{k-1}, \dots, u_1), u_v \in B \text{ für } v=1, 2, \dots, k$$

Die u_v sind die Komponenten von \tilde{u} .

Für $t \in [0, T]$ seien $\tilde{C}_1(h)_t : B^k \rightarrow B^k$ definiert durch:

$$\tilde{C}_1(h)_t \tilde{u} := (C(h)_t [u_k, u_{k-1}, \dots, u_1], u_k, u_{k-1}, \dots, u_2)$$

für alle $\tilde{u} \in B^k$, solange $h, t \in [0, T]$,

und es seien $\tilde{U}_h^{m,m} := (U_h^{m,m}, U_h^{m,m-1}, \dots, U_h^{m,m-k+1}) \in B^k$.

Dann hat man (2) in der äquivalenten Form:

$$(3) \quad \tilde{U}_h^{m,n+1} = \tilde{C}_1(h)_{n \cdot h} \tilde{U}_h^{m,n} \text{ mit vorgegebenen } \tilde{U}_h^{m,m} \in B^k \text{ und } m \geq k-1, n \geq m, n \cdot h \leq T$$

Für Vergleiche wird $E(t)_t$ auch umgeschrieben:

Für $t \in [0, T]$ sei $D_t^k := \bigotimes_{v=1}^k D_t$ und für jedes $u_0 \in D$ wird vereinbart:

$$\tilde{u}(t) := (u(t), u(t), \dots, u(t)) \in B^k; \tilde{u}_0 := (u_0, u_0, \dots, u_0) \in B^k$$

Falls dann $h, t + (k-1) \cdot h \in [0, T]$ so sei $\tilde{J}_t^h : D_t^k \rightarrow B^k$

definiert durch :

$$\tilde{J}_t^h \tilde{u}(t) := (E((k-1)h)_t u(t), E((k-2)h)_t u(t), \dots, u(t))$$

für alle $u_0 \in D$

Falls weiter $s, s+t, s+t-(k-1)h \in [0, T]$ so wird $\tilde{E}(s)_t : \tilde{D}_{t,h} \rightarrow B^k$

definiert durch:

$$\tilde{E}_h(s)_t \tilde{u} := (E(s)_t u_k, E(s)_{t-h}, \dots, E(s)_{t-(k-1)h} u_1)$$

$$\text{für alle } \tilde{u} \in \tilde{D}_{t,h} := \left\{ \tilde{u} \in B^k : u_v \in D_{t-(k-v)h} \text{ für } v = 1, 2, \dots, k \right\}$$

Zur Abkürzung wird noch $\tilde{u}_{t,h}(s) := \tilde{E}_h(s-(k-1)h)_{t+(k-1)h} \tilde{J}_t^h \tilde{u}(t)$

für $u_0 \in D$ gesetzt, falls auch $s \geq (k-1)h$. Oder anders ge-

schrieben und für alle $t, s \in [0, T]$ mit $t+s, t+s-(k-1)h \in [0, T]$ gültig:

$$\tilde{u}_{t,h}(s) = (u(t+s), u(t+s-h), \dots, u(t+s-(k-1)h)) \text{ für } u_0 \in D.$$

$\tilde{C}(h)_t, \tilde{J}_t^h$ und $\tilde{E}_h(s)_t$ sind offenbar wieder linear.

Sei $\|\cdot\|_{\mathbb{R}^k}$ irgendeine Norm des \mathbb{R}^k . Dann ist die Abbildung

$\|\cdot\| : B^k \rightarrow \mathbb{R} \cup \{0\}$ definiert durch

$$\|\tilde{u}\| := \left\| (\|u_k\|, \|u_{k-1}\|, \dots, \|u_1\|) \right\|_{\mathbb{R}^k} \text{ für jedes } \tilde{u} \in B^k,$$

$\tilde{u} := (u_k, u_{k-1}, \dots, u_1)$, eine Norm in B^k . Da der \mathbb{R}^k endliche

Dimension besitzt, sind alle derartigen Normen äquivalent, insbesondere

also äquivalent zu der durch $\|\tilde{u}\|_1 := \sum_{v=1}^k \|u_v\|$ definierten Norm

in B^k . Sei $\|\cdot\|$ im folgenden eine solche Norm, die sich von der Norm

in B schreibtechnisch also nur durch das Argument unterscheidet, dann ist $(B^k, \|\cdot\|)$ ein BANACH-Raum.

Nach den Schreibkonventionen kommen jetzt die die Approximation betreffenden Begriffe.

Def.: (Konsistenz)

Das Schema (3) (bzw. (2)) heißt für ein $u_0 \in D$ Konsistent von der Ordnung $p, p \in \mathbb{R}^+$, wenn für ein $H_1 \in (0, T]$ und jedes $t \in [0, T]$ gilt:

$$(4) \quad \left\| \left[\tilde{C}(h)_t - \tilde{E}_h(h)_t \right] \tilde{E}_h(t - (k-1)h) \int_0^h \tilde{u}_0 \right\| \leq M(u_0) \cdot h^{p+1}$$

für alle $h \in (0, H_1)$; $t+h, t-(k-1)h \in [0, T]$

mit einer (von u_0 abhängigen) Konstanten $M(u_0) \in \mathbb{R}^+$.

Die Relation (4) läßt sich kürzer schreiben:

$$(4)^1 \quad \left\| \left[\tilde{C}(h)_t - \tilde{E}_h(h)_t \right] \tilde{u}_{0,h}(t) \right\| \leq M(u_0) \cdot h^{p+1}$$

Def.: (schwache Stabilität)

Das Schema (3) heißt schwach stabil vom Grade $\beta, \beta \in \mathbb{R}^+ \cup \{0\}$, wenn es ein $N \in \mathbb{R}^+$ und ein $H_2 \in (0, T]$ so gibt, daß für alle $n \in \mathbb{N}_0, (n \geq k-1), h \in (0, H_2)$ mit $n \cdot h \in [0, T]$ gilt:

$$(5) \quad \left\| \prod_{\nu=1}^n \tilde{C}(h)_{(n-\nu)h} \right\| \leq N \cdot h^{-\beta} \quad \text{für alle } \mu \in \mathbb{N}, k-1 \leq \mu \leq n$$

Im Falle $\beta = 0$ nennt man es auch stabil.

Anmerkung:

Im zeitunabhängigen Fall hat man statt (5) einfach

$$\left\| \tilde{C}(h)^n \right\| \leq N \cdot h^{-\beta}, \quad \text{der Index ist weggelassen,}$$

also eine von einem Parameter, nämlich μ , weniger abhängende Relation.

Mittels (5) wird verlangt, daß sich alle bei Anwendung des Differenzschemas zu einer Zeit $v \cdot h$: $0 \leq v \cdot h \leq n \cdot h$ entstandenen Fehler bis zur Zeit $n \cdot h$ nicht zu stark, d.h. nur polynomartig, nicht exponentiell, vergrößern.

Def.:

Seien $t \in [0, T]$, $H \in (0, T]$, F_{st}^H ist eine Menge von Paaren

dreier Folgen $(h_j)_{j \in \mathbb{N}}$, $(m_j)_{j \in \mathbb{N}}$ und $(n_j)_{j \in \mathbb{N}}$:

$$F_{st}^H := \left\{ [(h_j), (m_j), (n_j)] ; h_j \in (0, H), m_j, n_j \in \mathbb{N}_0, k-1 \leq m_j \leq n_j, \right. \\ \left. h_j(m_{j+1}) \in [0, T] \text{ für alle } j \in \mathbb{N}, \text{ und } \begin{cases} \lim_{j \rightarrow \infty} h_j = 0, \lim_{j \rightarrow \infty} m_j \cdot h_j = s, \\ \lim_{j \rightarrow \infty} n_j \cdot h_j = t \end{cases} \right\}$$

Def.: (Konvergenz)

Die Schar von Lösungen $\{U_h^{m,n}; k-1 \leq m \leq n, h, n \cdot h \in [0, T]\}$ des

Differenzschemas (2) heißt von der Ordnung τ konvergent (gegen die

Lösung $U(t)$ von (1)), wenn für jedes Paar (s, t) , $s, t \in [0, T]$, $s \leq t$, und

alle dazugehörigen $[(h_j), (m_j), (n_j)] \in F_{st}^T$ ein $U_0 \in \mathbb{R}^+$ und ein

$j_0(U_0, s, t) \in \mathbb{N}$ existieren, so daß

$$(6) \quad \left\| \left(\prod_{v=1}^{m_j - m_j} \tilde{C}(h_j)_{(m_j - v)h_j} \right) \tilde{U}_{h_j}^{m_j, m_j} - \tilde{E}_{h_j, ((m_j - k + 1)h_j)} \int_0^{h_j} \tilde{u}_0 \right\| \leq K(U_0) \cdot h_j^\tau$$

für alle $j \geq j_0(U_0, s, t)$

Falls der laufende Index des Produkts einmal kleiner als der oben

stehende Grenzindex sein sollte, so wird für das Produkt der Ein-

heitsoperator auf B^k vereinbart.

Links in (6) steht $\tilde{U}_{h_j}^{m_j, m_j}$, rechts $\tilde{U}_{m_j, h_j, h_j}((m_j - m_j)h_j) = \tilde{U}_{0, h_j}^{(m_j, h_j)}$

übersichtlicher geschrieben hat man daher statt (6):

$$(6)^1 \quad \left\| \tilde{U}_{h_j}^{m_j, m_j} - \tilde{U}_{m_j, h_j}((m_j - m_j)h_j) \right\| \leq K(U_0) \cdot h_j^\tau$$

In Übereinstimmung zur LAX-RICHTMYER-Formulierung ist hier die schwache Stabilität als eine Eigenschaft des Lösungsoperators des Differenzenschemas definiert, während dagegen Konsistenz und Konvergenz, abgesehen von der Ordnung, hier schwächere Eigenschaften sind, da sie sich nur auf ein $u_0 \in D$ beziehen.

Damit Lösungen von (2) gegen die von (1) konvergieren, wird man selbstverständlich voraussetzen, daß die Anfangswerte von (2) die exakte Lösung approximieren, und daß die Geschwindigkeit der Konvergenz am Anfang nicht kleiner ist als die gewünschte.

Def:

Die Schar $\{U_h^{mn}; k-1 \leq m \leq n, h, m \cdot h \in [0, T]\}$ von Lösungen von (2) approximiert am Anfang die Lösung $u(t)$ von (1) von der Ordnung μ , $\mu \in \mathbb{R}^+$, wenn für ein $H_3 \in (0, T]$ gilt:

$$\|U_h^{m-m-v} - u((m-v)h)\| \leq a'(u_0) \cdot h^{\mu}$$

für $v=0, 1, \dots, k-1$, alle $m \in \mathbb{N}_0, m \geq k-1, m \cdot h \in [0, T]$,
und alle $h \in (0, H_3)$ mit einem $a'(u_0) \in \mathbb{R}^+$

oder äquivalent dazu:

$$(7) \quad \|\tilde{U}_h^{mm} - \tilde{u}_{mh, h}(0)\| \leq a(u_0) \cdot h^{\mu} \quad \text{für alle } m \in \mathbb{N}_0 \text{ mit } m \geq k-1,$$

$$m \cdot h \in [0, T] \text{ und alle } h \in (0, H_3) \quad \text{mit einem } a(u_0) \in \mathbb{R}^+$$

Damit läßt sich für Lösungen von (1), die der Konsistenzbedingung (4) genügen, analog zur LAX'schen Vorgehensweise ein Äquivalenzsatz beweisen:

Satz 1: (Äquivalenzsatz)

Seien $\rho, \tau \in \mathbb{R}^+, \rho \geq \tau$

Für ein $u_0 \in D$ sei das Schema (2) konsistent von der

Ordnung ρ . Dann gilt:

Schwache Stabilität vom Grade $\rho - \tau$



Jede Schar von Lösungen von (2) $\{U_h^{mn}; k-1 \leq m \leq n; h, n \cdot h \in [0, T]\}$,

die am Anfang die Lösung $u(t)$ von (1) von der Ordnung ρ approximiert, konvergiert von der Ordnung τ gegen $u(t)$.

Bemerkung:

Die im Satz garantierte Konvergenz ist eine stärkere Eigenschaft als die, die man in der Praxis normalerweise benötigt. Man braucht da nämlich Konvergenz nur für Scharen mit $m = k-1 = \text{const}$, jedenfalls im Rahmen einer Theorie die Rundungsfehler ausschließt. Im zeitunabhängigen Falle verschwinden diese Unterschiede; ein ähnliches Verhalten wie bei der schwachen Stabilität.

Beweis des Satzes:

a) Das Verfahren sei konvergent.

Wegen der Linearität sowohl des Anfangswertproblems als auch des

Differenzenschemas kann man die Konvergenz von Lösungsscharen

$\{U_h^{mn}; k-1 \leq m \leq n; h, n \cdot h \in [0, T]\}$ gegen $u(t)$ auf den linearen Teilraum D^0 von B , $D^0 := \{c \cdot u_0, c \in \mathbb{R}\}$ erweitern. Benötigt wird das in

diesem Beweisteil nur für das Nullelement $0 \in B$, zu dem als Anfangswert die Konstante 0, $t \in [0, T]$, als Lösung von (1) gehört.

Seien $\tilde{0} := (0, 0, 0, \dots, 0) \in B^k$ und $\{O_h^{mn}; k-1 \leq m \leq n; h, n \cdot h \in [0, T]\}$ eine Schar von Lösungen von (2) die am Anfang 0 von

der Ordnung ρ approximiert. Sei \tilde{P}_h^{mn} die Lösung von (3) mit dem Anfangswert

$\tilde{P}_h^{mn} = \int_{(m-k+1)h}^h u_0$. \tilde{P}_h^{mn} stimmt dann am Anfang genau mit $u(t)$ überein, also sicher von der Ordnung ρ . $\tilde{P}_h^{mn} + O_h^{mn}$ approximiert

am Anfang $u(t)$ von der Ordnung ρ . Die Elemente ohne " \sim " sind darin

wieder die zugehörigen Lösungen von (2). Nach Voraussetzung folgt Konvergenz für beide Scharen von der Ordnung τ , d.h.:

für jedes Paar (s,t) , $j, t \in [0, T]$, $s \leq t$, und alle dazugehörigen $[(h_j), (m_j), (n_j)] \in F_{st}^T$ existieren ein $K(u_0) \in \mathbb{R}^+$ und ein $j_0(u_0, s, t) \in \mathbb{N}$, so daß

$$\begin{aligned} & \left\| \tilde{P}_{h_j}^{m_j, n_j} - \tilde{u}_{0, h_j}(m_j, h_j) \right\| \leq K(u_0) \cdot h_j^\tau \quad \text{und} \\ & \left\| \tilde{P}_{h_j}^{m_j, n_j} + \tilde{O}_{h_j}^{m_j, n_j} - \tilde{u}_{0, h_j}(m_j, h_j) \right\| \leq K(u_0) \cdot h_j^\tau \\ \Rightarrow & \left\| \tilde{O}_{h_j}^{m_j, n_j} - \tilde{0} \right\| = \left\| \tilde{O}_{h_j}^{m_j, n_j} \right\| \leq 2K(u_0) \cdot h_j^\tau \\ (8) \quad & \Rightarrow \quad \text{für } j \geq j_0(u_0, s, t) \end{aligned}$$

$\{O^{m, n}; k-1 \leq m \leq n; h, n \cdot h \in [0, T]\}$ konvergiert gegen 0 von der Ordnung τ .

Sei $\tilde{V} \in \mathbb{B}^k$ beliebig gewählt. Angenommen die Menge

$$\left\{ h^{\rho-\tau} \cdot \left\| \left(\prod_{v=1}^{m-n} \tilde{C}(h)_{(m-n)h} \right) \tilde{V} \right\|; h \in (0, H_2), n, m \in \mathbb{N}_0, k-1 \leq m \leq n, n \cdot h \in [0, T] \right\}$$

ist für kein $H_2 \in (0, T]$ beschränkt. Dann gibt es Folgen

$$(h_j), (m_j), (n_j), j \in \mathbb{N} \text{ mit } h_j \in (0, T], m_j, n_j \in \mathbb{N}_0, k-1 \leq m_j \leq n_j, n_j \cdot h_j \in [0, T]$$

für $j \in \mathbb{N}$; mit $\lim_{j \rightarrow \infty} h_j = 0$

$$\text{und } h_j^{\rho-\tau} \cdot \left\| \left(\prod_{v=1}^{n_j-m_j} \tilde{C}(h_j)_{(n_j-m_j)h_j} \right) \tilde{V} \right\| \xrightarrow{j \rightarrow \infty} \infty$$

und darin gibt es Teilfolgen $(h_{j_\ell}), (m_{j_\ell}), (n_{j_\ell}), j_\ell \in \mathbb{N}$

$j_{\ell+1} > j_\ell$ für $\ell \in \mathbb{N}$ (ist Index von j) mit

$$\lim_{\ell \rightarrow \infty} m_{j_\ell} h_{j_\ell} = s, \lim_{\ell \rightarrow \infty} n_{j_\ell} h_{j_\ell} = t \quad \text{für ein Paar } (s, t) \text{ mit } s, t \in [0, T], s \leq t$$

und

$$W_{j_\ell} := \left(\prod_{v=1}^{n_{j_\ell}-m_{j_\ell}} \tilde{C}(h_{j_\ell})_{(n_{j_\ell}-v)h_{j_\ell}} \right) \tilde{V} \neq \tilde{0} \quad \text{für alle } \ell \in \mathbb{N}$$

Ohne Einschränkung der Allgemeinheit erfüllen bereits die Folgen

$(h_j), (m_j), (n_j)$ diese Eigenschaften.

Sei $\tilde{V}_j := \frac{h_j^{\rho} \tilde{V}}{(h_j^{\rho-\tau} \|\tilde{W}_j\|)^{\frac{1}{2}}}$ für $j \in \mathcal{N}$.

Da der Nenner nach unten beschränkt ist, approximiert \tilde{V}_j das Element \tilde{O} mit h_j^{ρ} (d.h. die Relation (7) ist für die Folgen $(h_j), (m_j), (n_j)$ erfüllt). Wegen (8) gilt für alle $j \in \mathcal{N}, j \geq j_0$

mit einem $j_0 \in \mathcal{N}$:

$$\begin{aligned} 2K(u_0) \cdot h_j^{\tau} &\geq \left\| \left(\prod_{v=1}^{m_j-m_j} \tilde{C}(h_j)_{(m_j-v)h_j} \right) \tilde{V}_j \right\| = \\ &= \frac{h_j^{\rho} \left\| \left(\prod_{v=1}^{m_j-m_j} \tilde{C}(h_j)_{(m_j-v)h_j} \right) \tilde{V} \right\|}{\left\{ h_j^{\rho-\tau} \left\| \left(\prod_{v=1}^{m_j-m_j} \tilde{C}(h_j)_{(m_j-v)h_j} \right) \tilde{W}_j \right\| \right\}^{\frac{1}{2}}} = \\ &= h_j^{\rho-\frac{\rho-\tau}{2}} \|\tilde{W}_j\|^{\frac{1}{2}} = h_j^{\tau} \cdot (h_j^{\rho-\tau} \|\tilde{W}_j\|)^{\frac{1}{2}} =: h_j^{\tau} \cdot B_j \end{aligned}$$

$\Rightarrow B_j \leq 2K(u_0)$ für genügend großes j , mit einer Folge $(B_j)_{j \in \mathcal{N}}$, die $\lim_{j \rightarrow \infty} B_j = \infty$ erfüllt und das ist ein Widerspruch.

Die Annahme war also falsch; es gibt daher ein $H_2 \in (0, T]$, so daß

$$h^{\rho-\tau} \left\| \left(\prod_{v=1}^{n-m} \tilde{C}(h)_{(n-v)h} \right) \tilde{V} \right\| \leq K_n(\tilde{V})$$

für $h \in (0, H_2)$; $n, m \in \mathcal{N}_0, k-1 \leq m \leq n$, $n, h \in [0, T]$ mit einem $K_n(\tilde{V}) \in \mathbb{R}^+$.

Nach dem Prinzip von der gleichmäßigen Beschränktheit, angewandt auf die Menge von Operatoren in B^k

$$M := \left\{ h^{\rho-\tau} \prod_{v=1}^{n-m} \tilde{C}(h)_{(n-v)h} \mid h \in (0, H_2), n, m \in \mathcal{N}_0, k-1 \leq m \leq n, n, h \in [0, T] \right\}$$

gibt sein $N \in \mathbb{R}^+$, so daß

$$h^{\rho-\tau} \left\| \prod_{v=1}^{n-m} \tilde{C}(h)_{(n-v)h} \right\| \leq N \text{ für alle } h^{\rho} \prod_{v=1}^{n-m} \tilde{C}(h)_{(n-v)h} \in M,$$

d.h. schwache Stabilität vom Grade $\rho-\tau$.

Anmerkungen:

Die Konsistenz ging nicht ein in diesem Beweisteil.

Hier ist die Erfahrung ausgenutzt, daß im Beweise für den Äquivalenzsatz von LAX ([8], S.127) im entsprechenden Teil Konvergenz nur in einer Umgebung von $\tilde{O} \in \beta^k$ benötigt wird.

Hätte man genaueres Übereinstimmen von $\{U_h^{m,n}; k \rightarrow \infty, m \leq n, h_1, n, h \in [0, T]\}$ am Anfang mit $u(t)$ gefordert, etwa von der Ordnung $\rho + \varepsilon, \varepsilon > 0$, so könnte man auf diese Weise nicht mehr die Notwendigkeit der schwachen Stabilität vom Grade $\rho - \tau$ für Konvergenz der Ordnung τ beweisen.

b) Das Verfahren sei schwach stabil vom Grade $\rho - \tau$.

Seien $s, t \in [0, T], s \leq t$ und $[(h_s), (m_s), (n_s)] \in \bar{F}_{s,t}^T$ ein Folgentripel. u_0 erfülle die Konsistenzbedingung (4).

$\{U_h^{m,n}; k \rightarrow \infty, m \leq n; h_1, n, h \in [0, T]\}$ sei eine Schar von Lösungen von (2), die am Anfang $u(t)$ von der Ordnung ρ approximiert. Dann gilt für jedes $l \in \mathbb{N}_0$ mit $l < n_j - m_j$

$$\begin{aligned} \tilde{U}_{m_j, h_j, h_j}((l+1) \cdot h_j) &= \tilde{E}_{h_j}(h_j)_{(m_j+l)h_j} \tilde{U}_{m_j, h_j, h_j}(l h_j) \\ \tilde{U}_{h_j}^{m_j, m_j+l+1} &= \tilde{C}_1(h_j)_{(m_j+l)h_j} U_{h_j}^{m_j, m_j+l} \end{aligned}$$

Daraus erhält man für den Fehler

$$\tilde{Z}_{h_j}^{m_j, m_j+l} := \tilde{U}_{h_j}^{m_j, m_j+l} - \tilde{U}_{m_j, h_j, h_j}(l h_j)$$

definiert für $l \in \mathbb{N}_0, l < n_j - m_j$

die rekursive Relation:

$$\begin{aligned} \tilde{Z}_{h_j}^{m_j, m_j+l+n} &= \tilde{U}_{h_j}^{m_j, m_j+l+n} - \tilde{U}_{m_j, h_j, h_j}((l+1)h_j) = \\ &= [\tilde{C}(h_j)_{(m_j+l)h_j} - \tilde{E}_{h_j}(h_j)_{(m_j+l)h_j}] \tilde{U}_{m_j, h_j, h_j}(l h_j) + \\ &\quad + \tilde{C}_1(h_j)_{(m_j+l)h_j} [\tilde{U}_{h_j}^{m_j, m_j+l} - \tilde{U}_{m_j, h_j, h_j}(l h_j)] \\ &= [\tilde{C}(h_j)_{(m_j+l)h_j} - \tilde{E}_{h_j}(h_j)_{(m_j+l)h_j}] \tilde{U}_{m_j, h_j, h_j}(l h_j) + \tilde{C}_1(h_j)_{(m_j+l)h_j} \tilde{Z}_{h_j}^{m_j, m_j+l} \end{aligned}$$

Die eindeutige Lösung dieser (Operatoren-) Differenzgleichung

zum Anfangswert $\tilde{z}_{h_j}^{m_j, m_j}$ ist gegeben durch:

$$(9) \quad \begin{aligned} & \tilde{z}_{h_j}^{m_j, m_j+l} = \\ & = \sum_{\mu=0}^{l-1} \left(\prod_{v=\mu}^{l-1} \tilde{C}(h_j)_{(m_j+l-v)h_j} \right) \left[\tilde{C}(h_j)_{(m_j+l-1-\mu)h_j} - \tilde{E}_{h_j}(h_j)_{(m_j+l-1-\mu)h_j} \right] \cdot \\ & \quad + \prod_{v=1}^l \tilde{C}(h_j)_{(m_j+l-v)h_j} \tilde{z}_{h_j}^{m_j, m_j} \quad \cdot \tilde{U}_{0, h_j}((m_j+l-1-\mu)h_j) \\ & \quad \text{für } l \leq m_j - m_j \end{aligned}$$

Beweis durch vollständige Induktion:

Für $l=0$ ist die Aussage trivial.

Der Induktionsschritt läßt sich direkt verifizieren:

$$\begin{aligned} & \left[\tilde{C}(h_j)_{(m_j+l)h_j} - \tilde{E}_{h_j}(h_j)_{(m_j+l)h_j} \right] \tilde{U}_{m_j, h_j, h_j}(lh_j) + \tilde{C}(h_j)_{(m_j+l)h_j} \tilde{z}_{h_j}^{m_j, m_j+l} = \\ & = [\dots] \tilde{U}_{m_j, h_j, h_j}(lh_j) + \sum_{\mu=0}^{l-1} \left(\tilde{C}(h_j)_{(m_j+l)h_j} \prod_{v=\mu}^{l-1} \tilde{C}(h_j)_{(m_j+l-v)h_j} \right) \cdot \\ & \quad \cdot \left[\tilde{C}(h_j)_{(m_j+l-1-\mu)h_j} - \tilde{E}_{h_j}(h_j)_{(m_j+l-1-\mu)h_j} \right] \cdot \tilde{U}_{0, h_j}((m_j+l-1-\mu)h_j) \\ & \quad + \tilde{C}(h_j)_{(m_j+l)h_j} \prod_{v=1}^l \tilde{C}(h_j)_{(m_j+l-v)h_j} \tilde{z}_{h_j}^{m_j, m_j} \quad \Rightarrow \\ \text{Es ist:} \end{aligned}$$

$$\tilde{C}(h_j)_{(m_j+l)h_j} \prod_{v=1}^l \tilde{C}(h_j)_{(m_j+l-v)h_j} = \prod_{v=0}^l \tilde{C}(h_j)_{(m_j+l-v)h_j}$$

Summiert man 1 bis l statt von 0 bis $l-1$ so steht da:

$$\begin{aligned} & \Rightarrow \left[\tilde{C}(h_j)_{(m_j+l)h_j} - \tilde{E}_{h_j}(h_j)_{(m_j+l)h_j} \right] \tilde{U}_{m_j, h_j, h_j}(lh_j) + \\ & + \sum_{\mu=1}^l \left(\prod_{v=\mu}^l \tilde{C}(h_j)_{(m_j+l+1-v)h_j} \right) \left[\tilde{C}(h_j)_{(m_j+l-\mu)h_j} - \tilde{E}_{h_j}(h_j)_{(m_j+l-\mu)h_j} \right] \cdot \\ & \quad \cdot \tilde{U}_{0, h_j}((m_j+l-\mu)h_j) \\ & + \prod_{v=1}^{l+1} \tilde{C}(h_j)_{(m_j+l+1-v)h_j} \tilde{z}_{h_j}^{m_j, m_j} \quad \Rightarrow \end{aligned}$$

also nach Indizesverschiebung von μ (bzw. ν im hinteren Produkt) und damit zwangsläufig für ν (in den Produkten der Summanden).

$$\Rightarrow \sum_{h_j}^{m_j, m_j+l+1} ; \text{ damit ist (9) bewiesen.}$$

Anmerkung:

Diese Herleitung begründet das Auftreten der LAX'schen Summe im Äquivalenzbeweis ([8], S.46, (3.10)), der Fehler der Anfangswerte wird in ihr berücksichtigt und die LAX'sche Methode ist durch sie auf zeitabhängige Probleme ausdehnbar.

Fortsetzung des Beweisteils b):

Aus (9) folgt für $l \in \mathbb{N}_0, l \leq m_j - m_j$:

$$\begin{aligned} \left\| \sum_{h_j}^{m_j, m_j+l} \right\| &\leq \sum_{\mu=0}^{l-1} \left\| \prod_{\nu=1}^{\mu} \tilde{C}(h_j)_{(m_j+l-\nu)h_j} \right\| \cdot M(u_0) \cdot h_j^{\rho+1} \\ &+ \left\| \prod_{\nu=1}^l \tilde{C}(h_j)_{(m_j+l-\nu)h_j} \right\| \cdot \left\| \sum_{h_j}^{m_j, m_j} \right\| \leq \end{aligned}$$

wegen der Beschränktheit der Operatoren und der Konsistenzbedingung

(4)¹, und zwar für $h_j < H_1$.

$$\leq M(u_0) \cdot N \cdot h_j^{\rho+1} \cdot h_j^{\tau-\rho} \sum_{\mu=0}^{l-1} 1 + a(u_0) \cdot N \cdot h_j^{\tau-\rho} \cdot h_j^{\rho} \Rightarrow$$

wegen schwacher Stabilität (5), und der Approximation am Anfang

(s. (7)), die Voraussetzung des Satzes und $\sum_{h_j}^{m_j, m_j} = \tilde{U}_{h_j}^{m_j, m_j} - \tilde{U}_{m_j, h_j, h_j}^{(0)}$.

$$= h_j^{\tau} \cdot N \cdot [M(u_0) \cdot h_j \cdot l + a(u_0)]$$

$$\leq h_j^{\tau} \cdot N \cdot [M(u_0) \cdot T + a(u_0)]$$

für $h_j < H_1$

und das ist sogar mehr als in der Konvergenzbedingung (6) gefordert

wird. Q.e.d.

Anmerkungen:

Daß der Fehler (s. letzte Beweiszeile) bei festgehaltener Schrittweite linear in T ansteigt, paßt zu der praktischen Erfahrung wachsender Fehler, wenn man ein Anfangswertproblem über längere Zeiten näherungsweise berechnet.

Weniger genau stimmende Anfangswerte, etwa von der Ordnung $\rho - \varepsilon$, $0 < \varepsilon < \tau$, hätten langsamere Konvergenz, und zwar von der Ordnung $\tau - \varepsilon$ bedingt.

In der Praxis ist ohnehin meist die Ordnung $\rho + 1$ erforderlich.

Da dieser Beweis für eine bestimmte Lösung von (1) geführt wurde,

werden von (1) nicht globale Eigenschaften der Operatorenschar

$\{E(t)_0; t \in [0, T]\}$ wie etwa dichter Definitionsbereich in B

oder gleichmäßige Beschränktheit gefordert. Dagegen gehen die

"Glattheitseigenschaften" der gerade gesuchten Lösung des Anfangs-

wertproblems im vollen Umfang ein in der Konsistenzbedingung und

in den Koeffizienten vor h_j^{τ} .

Nur Konvergenz verallgemeinerter Lösungen ohne Angabe der Ge-

schwindigkeit hätte man im zeitabhängigen Fall analog zum LAX'schen

Beweise liefern können.

II. Untersuchung des 2-Schritt Schemas

$$(10) \quad U^{m+1} - U^{m-1} = \sum_{v=1}^N \frac{A_v}{R_v} (T_v - T_v^{-1}) U^m$$

Das Schema soll zur näherungsweise Berechnung der Lösungen des Anfangswertproblems

$$(11) \quad \frac{\partial u(t, x)}{\partial t} = \sum_{v=1}^N A_v \frac{\partial u(t, x)}{\partial x_v} \quad \text{für } t \in [0, T] ; \quad x \in \mathbb{R}^N ;$$

$$u(0, \cdot) = u_0(\cdot)$$

für reellwertige Funktionen $u: \mathbb{R} \times \mathbb{R}^N \rightarrow \mathbb{R}^m$ dienen. Die Matrizen $A_v, v \in \mathcal{N}$, sind konstant, N ist die Raumdimension. Als Norm wählt man, um später mit FOURIER-Transformationen arbeiten zu können, die L_2 -Norm. Sei

$$I := \bigotimes_{v=1}^N [a_v, b_v] \quad ; \quad a_\ell, b_\ell \in \mathbb{R}, a_\ell < b_\ell \quad \text{für } \ell = 1, 2, \dots, N$$

und im folgenden $B = L_2(I)$ d.h. $v(\cdot)$ bzgl. I periodisch, und

$$\|v(\cdot)\| := \left(\int_I (v(x), v(x)) dx \right)^{\frac{1}{2}} \quad \text{für alle } v(\cdot) \in B$$

Mit (\cdot, \cdot) ist das Skalarprodukt zweier Vektoren $p, q \in \mathbb{R}^\ell, \ell \in \mathcal{N}$,

gemeint: $(p, q) := \sum_{i=1}^\ell p_i q_i$; p_i, q_i sind die Komponenten von p bzw. q .

Zu einem $u_0(\cdot) \in B$ existiere eine eindeutige Lösung $u(t) \in B$ für $t \in [0, T]$ von (11).

Die T_v in (10) sind Translationsoperatoren um den Betrag $R_v \cdot h, R_v \in \mathbb{R}^+$,

für $v = 1, 2, \dots, N$. Seien $y, y' \in \mathbb{R}^{N+1}$, $y := (x_0, x_1, \dots, x_N)$,

$y'_v := (x_0, x_1, \dots, x_{v-1}, x_v + R_v h, x_{v+1}, \dots, x_N)$ wird mit x_0 identifiziert,

$R_0 := 1$. T_v ist in y für alle in y und y_v definierten Funktionen $v(\cdot)$

erklärt durch:

$$T_v v(y) = v(y'_v) \quad , \quad \text{und zwar für } v = 0, 1, 2, \dots, N$$

A. Konsistenz.

Von $u_0(\cdot)$ wird gefordert, daß einige partielle Ableitungen 3. Ordnung von $u(t, \cdot)$ existieren, beschränkt sind und zu B gehören, und zwar: für $i = 0, 1, 2, \dots, N$ wird verlangt

$$\frac{\partial^3 u(t, \cdot)}{\partial x_i^3} \in B \quad \text{für } t \in [0, T]$$

und für $(t, x) \in [0, T] \times I_0$, I_0 eine offene Menge: ($I_0 \supset I$)

$$\left| \frac{\partial^3 u(t, x)}{\partial x_i^3} \right| \leq C_i \quad \text{mit gewissen Konstanten } C_i \in \mathbb{R}^+$$

Aus dem TAYLOR'schen Satz folgt dann für $(t, x) \in [h, T-h] \times I$

und für $v = 0, 1, 2, \dots, N$:

$$\begin{aligned} \left| (T_v - T_v^{-1}) u(t, x) - 2R_v h \frac{\partial u(t, x)}{\partial x_v} \right| &\leq \frac{R_v^3 h^3}{3!} \left(\left| \frac{\partial^3 u(z^1)}{\partial x_v^3} \right| + \left| \frac{\partial^3 u(z^2)}{\partial x_v^3} \right| \right) \\ &\leq \frac{R_v^3}{3!} h^3 P_1 \quad \text{mit einem } P_1 \in \mathbb{R}^+ \end{aligned}$$

wobei $z^1 := (x_0, x_1, \dots, x_{v-1}, x_v', x_{v+1}, \dots, x_N)$ $x_v' \in [x_v, x_v + R_v h]$

und $z^2 := (x_0, x_1, \dots, x_{v-1}, x_v'', x_{v+1}, \dots, x_N)$ $x_v'' \in [x_v - R_v h, x_v]$

\Rightarrow
(12)

$$\left\| (T_v - T_v^{-1}) u(t, \cdot) - 2R_v h \frac{\partial u(t, \cdot)}{\partial x_v} \right\| \leq h^3 \frac{R_v^3 P_1}{6} \leq h^3 P_2$$

für $t \in [h, T-h]$ und $v=0,1,2,\dots,N$ mit einem $P_2 \in \mathbb{R}^+$

Schreibt man (10) als System in $\begin{pmatrix} U^{n+1} \\ U^n \end{pmatrix}$ d.h.

$$(13) \quad \begin{pmatrix} U^{n+1} \\ U^n \end{pmatrix} = \begin{pmatrix} \sum_{v=1}^N \frac{A_v}{R_v} (T_v - T_v^{-1}) & I \\ I & 0 \end{pmatrix} \begin{pmatrix} U^n \\ U^{n-1} \end{pmatrix}$$

(genauer sind es ortsabhängige $U^m(\cdot) \in B$, I ist die Einheitsmatrix) und bezeichnet wieder die zugehörigen in Abschnitt I erklärten Lösungsoperatoren mit $\tilde{C}(h)$ bzw. $\tilde{E}(h)$ (der zeitabhängige Index ist weggelassen), dann ist für alle h, t mit $h, t, t+2h \in [0, T]$:

$$\begin{aligned} & \left\| [\tilde{C}(h) - \tilde{E}(h)] \begin{pmatrix} u(t+h, \cdot) \\ u(t, \cdot) \end{pmatrix} \right\|_1 = \\ & = \left\| \sum_v \frac{A_v}{R_v} (T_v - T_v^{-1}) u(t+h, \cdot) + u(t, \cdot) - u(t+2h, \cdot) \right\| + \underbrace{\|u(t+h, \cdot) - u(t+h, \cdot)\|}_{=0} \\ & = \left\| \sum_v \frac{A_v}{R_v} (T_v - T_v^{-1}) u(t+h, \cdot) - 2h \sum_v A_v \frac{\partial}{\partial x_v} u(t+h, \cdot) - \right. \\ & \qquad \qquad \qquad \left. u(t+2h, \cdot) - u(t, \cdot) + 2h \frac{\partial u(t+h, \cdot)}{\partial t} \right\| \Rightarrow \end{aligned}$$

wegen der Differentialgleichung (11). Sei $\|A\|$ die Spektralnorm einer quadratischen Matrix A . Sie ist verträglich mit der EUKLID'ischen Vektornorm.

$$\begin{aligned} & \rightarrow \leq \sum_{v=1}^N \frac{\|A_v\|}{R_v} \cdot \left\| (T_v - T_v^{-1}) u(t+h, \cdot) - 2h R_v \frac{\partial u(t+h, \cdot)}{\partial x_v} \right\| \\ & \qquad \qquad \qquad + \left\| (T_0 - T_0^{-1}) u(t+h, \cdot) - 2h \frac{\partial u(t+h, \cdot)}{\partial t} \right\| \\ & \leq h^3 \cdot P_2 \cdot \left(\sum_v \frac{\|A_v\|}{R_v} + 1 \right) \quad \text{für } t, h \in [0, T] \text{ mit } t+2h \in [0, T] \end{aligned}$$

wegen (12). Das ist genau die zu erfüllende Bedingung

(4)' ($k=2$). \Rightarrow Konsistenz von der Ordnung 2 für u_0 .

B. Stabilität im strikt hyperbolischen Falle

Mit Hilfe von (10) läßt sich U^n in den Punkten eines Gitters

$$(14) \quad G := \{x \in \mathbb{R}^N; x = h \cdot (j_1 R_1, \dots, j_N R_N); j_1, \dots, j_N \in \mathbb{Z}\}$$

berechnen, wenn man nur U^0 und U^1 auf diesem Gitter kennt.

Stabilität kann man mit Hilfe einer FOURIER-Transformation ermitteln (vgl. [1] S. 79).

Die Funktionen $e^{i(k,x)}$, $k \in \tilde{G}$, bilden für

$$\tilde{G} := \{k \in \mathbb{R}^N; k = 2\pi \left(\frac{j_1}{a_1 - b_1}, \dots, \frac{j_N}{a_N - b_N} \right); j_1, \dots, j_N \in \mathbb{Z}\}$$

ein vollständiges Orthogonalsystem in $L_2(I)$. Nach dem Satz von

RIESZ-FISCHER ist $U^n \in B$ daher eindeutig darstellbar als

$$(15) \quad U^n(x) = \frac{1}{\sqrt{|I|}} \sum_{k \in \tilde{G}} V^n(k) e^{i(k,x)}$$

Für die FOURIER-Koeffizienten gilt umgekehrt:

$$V^n(k) = \frac{1}{\sqrt{|I|}} \int_I U^n(x) e^{-i(k,x)} dx$$

Aus der Vollständigkeit des Orthonormalsystems folgt die PARSEVAL'sche Gleichung:

$$(16) \quad \|U^n\|^2 = \int_I |U^n(x)|^2 dx = \sum_{k \in \tilde{G}} |V^n(k)|^2 =: \|V^n\|_{\tilde{G}}^2$$

Setzt man (15) in (10) ein, so liefert ein Koeffizientenvergleich vor $e^{i(k,x)}$

$$(17) \quad \begin{aligned} V^{n+1}(k) - V^n(k) &= \sum_{\nu=n}^N \frac{A_\nu}{R_\nu} (e^{i k_\nu R_\nu h} - e^{-i k_\nu R_\nu h}) V^n(k) \\ &= 2i V^n(k) \sum \frac{A_\nu}{R_\nu} \sin k_\nu R_\nu h =: 2i V^n(k) A(k) \end{aligned}$$

Durch die FOURIER-Transformation hat man ein Differenzenschema in $V^n(k)$ erhalten, in das nur noch Differenzen in eine Richtung eingehen, das aber abhängig von einem Parameter ist.

Die naheliegendste Vereinfachung von (17) ist die Transformation von $A(k)$ auf eine geeignete Normalform.

Sei (11) strikt hyperbolisch, dann ist $A(k)$ diagonalisierbar; sei etwa

$$T^{-1}(k) A(k) T(k) = \Lambda(k) \quad \text{für } k \in \mathbb{R}^n$$

mit nicht singulären Matrizen $T(k)$ und mit Diagonalmatrizen $\Lambda(k)$ in denen die Eigenwerte von $A(k)$ stehen.

Sei $w^n(k) := T^{-1}(k) V^n(k)$

Um die Existenz von $\|w^0\|_{\tilde{G}}, \|w^1\|_{\tilde{G}}$ (20) zu wissen, muß man $\|T^{-1}(k)\| \leq C_{-1} = \text{const.}$ fordern, desgleichen braucht man für die andere Richtung $\|T(k)\| \leq C_1 = \text{const.}$, $C_1, C_{-1} \in \mathbb{R}^+$ d.h., um aus der Beschränktheit von $\|w^n\|_{\tilde{G}}$ die von $\|V^n\|_{\tilde{G}}$ zu folgern.

Diese Voraussetzung ist bedingt durch die Parameterabhängigkeit von (17). Im Falle HERMITE'scher Matrizen ($\Rightarrow A(k)$ HERMITESch) gibt es sogar unitäre Matrizen $T(k)$, die $A(k)$ auf Diagonalform transformieren; die Voraussetzung ist dann also erfüllt.

$w^n(k)$ erfüllt das Differenzenschema (ohne Parameter k):

$$w^{n+1} - w^{n-1} = 2i \Lambda w^n \quad ; \quad \Lambda(k) = \begin{pmatrix} \lambda_1(k) & & 0 \\ & \lambda_2(k) & \\ 0 & & \ddots \\ & & & \lambda_m(k) \end{pmatrix}$$

welches in die Komponenten w_l von w zerfällt:

(18) $w_l^{n+1} - w_l^{n-1} = 2i \lambda_l w_l^n, \quad l=1, \dots, m$

Diese Differenzengleichung läßt sich bei Vorgabe von w_l^0, w_l^1 exakt lösen. Die charakteristische Gleichung von (18) ist (der Index l wird unterdrückt):

$$\lambda^2 - 2i\lambda - 1 = 0$$

$$\Rightarrow \lambda_{\pm} = i\lambda \pm \sqrt{1-\lambda^2} =: \lambda_{\pm}$$

Für $\lambda_+ = \lambda_- =: \lambda^*$ d.h. $\lambda^* = \pm i$ hat man die Eigenlösungen $(\pm i)^n$ und $n(\pm i)^n$, $|\omega^n|$ bleibt daher in diesem Falle im allgemeinen nicht beschränkt.

Für $\lambda_+ \neq \lambda_-$ sind λ_+^n und λ_-^n die beiden Eigenlösungen von (18). Für vorgegebene Anfangswerte ω^0, ω^1 lautet daher die Lösung von (18):

$$\omega^n = a\lambda_+^n + b\lambda_-^n$$

wobei

$$a = \frac{\omega^1 - \omega^0\lambda_-}{\lambda_+ - \lambda_-}; \quad b = -\frac{\omega^1 - \omega^0\lambda_+}{\lambda_+ - \lambda_-}; \quad \lambda_+ - \lambda_- = 2\sqrt{1-\lambda^2} \neq 0$$

Offensichtlich ist für die Beschränktheit von $|\omega^n|$ im allgemeinen diejenige von λ_+^n und λ_-^n notwendig. Da λ und ω^n noch von einem Parameter k abhängen, ist zu erwarten, daß außerdem Nenner von a und b : $\lambda_+ - \lambda_- = 2\sqrt{1-\lambda^2}$ betragsmäßig nach unten durch eine von k unabhängige konstante $\varepsilon' > 0$ abzuschätzen ist: $|\lambda_+ - \lambda_-| \geq 2\varepsilon' > 0$

Zunächst hat man

$$|\lambda_+^n|, |\lambda_-^n| \leq \text{const. zu erfüllen.}$$

$$\Leftrightarrow |\lambda_+|, |\lambda_-| \leq 1$$

Wegen $\lambda_+ \cdot \lambda_- = -1$ ist das gleichbedeutend mit

$$|\lambda_+| = |\lambda_-| = 1$$

Es gibt daher einen Zweig des natürlichen Logarithmus, so daß $\log(-i\lambda_+)$ und $\log(-i\lambda_-)$ definiert sind.

$$-i \log(-i\lambda_{\pm}) = -i \log(\lambda \pm \sqrt{\lambda^2 - 1}) = \arccos \lambda$$

darin ist der durch die linke Seite definierte Zweig des \arccos zu nehmen.

$$-i\lambda_{\pm} = e^{i \arccos \lambda}$$

Es soll sein $|\lambda_+| = |\lambda_-| = 1$

$$\Leftrightarrow \arccos \lambda = b$$

mit $b \in \mathbb{R}$, b beliebig.

$$\Leftrightarrow \lambda = \cos b \Rightarrow \lambda \in \mathbb{R}$$

(19) $\Leftrightarrow -1 \leq \lambda \leq 1$

D.h. $A(\kappa)$ darf für alle Werte von κ nur reelle Eigenwerte haben.
Das ist genau die Hyperbolicität von (11).

Um in κ gleichmäßige Beschränktheit der Lösungen von (18) zu haben wird außer (19) vorausgesetzt

$$|\lambda_+ - \lambda_-| = 2\sqrt{1-\lambda^2} \geq 2\varepsilon' > 0, \quad 0 < \varepsilon' < 1$$

$$\Leftrightarrow |\lambda| \leq \sqrt{1-\varepsilon'^2} =: 1-\varepsilon \Rightarrow 0 < \varepsilon < 1$$

(20) d.h. $-1+\varepsilon \leq \lambda \leq 1-\varepsilon$

(20) umfaßt offenbar (19).

Faßt man Beweisteile und Voraussetzungen zusammen, so erhält man die "Stabilität" von U^n :

U^0, U^1 seien vorgegeben, wegen (16) sind damit $\|V^0\|_{G_1}^2, \|V^1\|_{G_1}^2$ sowie $\|w^0\|_{G_1}^2, \|w^1\|_{G_1}^2$ bestimmt, wegen

$$|w^x(\kappa)|^2 = |T^{-x}(\kappa)V^x(\kappa)|^2 \leq \|T^{-x}(\kappa)\|^2 |V^x(\kappa)|^2 \leq C_{-1}^2 \|V^x(\kappa)\|^2$$

für $x=0,1$

Es ist $w^n(\kappa) = \oplus(\kappa) [(\Lambda_+^n(\kappa) - \Lambda_-^n(\kappa)) \omega^1(\kappa) + (\Lambda_+^{n-1}(\kappa) - \Lambda_-^{n-1}(\kappa)) \omega^0(\kappa)]$

mit $m \times m$ Diagonalmatrizen $\oplus(\kappa), \Lambda_+(\kappa), \Lambda_-(\kappa)$ definiert

durch:

$$\oplus(\kappa) := \begin{pmatrix} [\lambda_{n+}(\kappa) - \lambda_{n-}(\kappa)]^{-1} & & \\ & \ddots & \\ & & [\lambda_{2+}(\kappa) - \lambda_{2-}(\kappa)]^{-1} \\ & & & \ddots \\ & & & & [\lambda_{m+}(\kappa) - \lambda_{m-}(\kappa)]^{-1} \end{pmatrix}; \quad \Lambda_{\pm}(\kappa) := \begin{pmatrix} \lambda_{n\pm}(\kappa) & & \\ & \ddots & \\ & & \lambda_{e\pm}(\kappa) \\ & & & \ddots \\ & & & & \lambda_{m\pm}(\kappa) \end{pmatrix}$$

$$\Rightarrow |\omega^n(k)| \leq \|\oplus(k)\| \cdot [(\|\mathcal{L}_+^n(k)\| + \|\mathcal{L}_-^n(k)\|) \cdot |\omega^1(k)| + (\|\mathcal{L}_+^{n-1}(k)\| + \|\mathcal{L}_-^{n-1}(k)\|) |\omega^0(k)|] \leq$$

hier ist wieder die Spektralnorm gemeint.

$$\leq \frac{1}{\min_{\ell=1, \dots, m} |\lambda_{\ell+} - \lambda_{\ell-}|} \left[(\max_{\ell=1, \dots, m} |\lambda_{\ell+}|^n + \max_{\ell=1, \dots, m} |\lambda_{\ell-}|^n) |\omega^1(k)| + (\max_{\ell=1, \dots, m} |\lambda_{\ell+}|^{n-1} + \max_{\ell=1, \dots, m} |\lambda_{\ell-}|^{n-1}) |\omega^0(k)| \right] =$$

$$= \frac{2}{\min_{\ell=1, \dots, m} |\lambda_{\ell+} - \lambda_{\ell-}|} [|\omega^1(k)| + |\omega^0(k)|] \leq$$

$$\leq \frac{1}{\varepsilon^1} [|\omega^1(k)| + |\omega^0(k)|] =$$

$$= \frac{1}{\sqrt{1 - (\varepsilon^1)^2}} [|\omega^1(k)| + |\omega^0(k)|]$$

$$\Rightarrow |\omega^n(k)|^2 \leq \frac{2}{\varepsilon^{1,2}} [|\omega^1(k)|^2 + |\omega^0(k)|^2]$$

$$\Rightarrow \|\omega^n\|_{\mathcal{G}}^2 \leq \frac{2}{\varepsilon^{1,2}} [\|\omega^1\|_{\mathcal{G}}^2 + \|\omega^0\|_{\mathcal{G}}^2]$$

$$\Rightarrow \|U^n\|^2 = \|V^n\|_{\mathcal{G}}^2 \leq \frac{2C^2 C_{-1}^2}{\varepsilon^{1,2}} [\|V^1\|_{\mathcal{G}}^2 + \|V^0\|_{\mathcal{G}}^2] = \frac{2C^2 C_{-1}^2}{\varepsilon^{1,2}} [\|U^1\|^2 + \|U^0\|^2]$$

d.h. Stabilität

Für alle Matrizen $D \in \mathbb{R}^m \times \mathbb{R}^m$ sei $|D| := \max\{\lambda_{\max}, \lambda_{\min}\}$
 Eigenwert von D }

Dann ist bewiesen der

Satz 2:

Das Differenzenschema

$$(10) \quad U^{n+1} - U^{n-1} = \sum_{v=n}^N \frac{A_v}{R_v} (T_v - T_v^{-1}) U^n$$

ist stabil, wenn die folgenden drei Aussagen gelten:

a) Das Differentialgleichungssystem

$$(11) \quad u_t(t, x) = \sum_{v=n}^N A_v u_{x_v}(t, x) \quad \text{für } 0 \leq t \leq T \text{ und } x \in I$$

ist strikt hyperbolisch; d.h.:

Es gibt ein $\mu > 0$, so daß für alle $\xi \in W_\mu(0) := \{ \xi \in \mathbb{R}^N; |\xi_v| \leq \mu, v=1, \dots, N \}$ die Matrix $B(\xi) := \sum_{v=1}^N \frac{A_v}{R_v} \xi_v$ mit einer Matrix $S(\xi)$ diagonalisierbar ist, d.h. $S^{-1}(\xi)$ existiert und $S^{-1}(\xi) B(\xi) S(\xi)$ ist eine Diagonalmatrix.

b) Es gibt Konstanten $C, C_1 \in \mathbb{R}^+$, so daß

$$\|S^{-1}(\xi)\| \leq C_1; \|S(\xi)\| \leq C \text{ für alle } \xi \in W_\mu(0)$$

c) (eigentliche Stabilitätsbedingung)

Es gibt ein $\eta > 0$, so daß

$$(21) \quad |B(\xi)| \leq 1 - \eta \quad \text{für alle } \xi \in W_\mu(0)$$

Bemerkung 1:

Daß hier statt $W_\mu(0) = \{ \xi \in \mathbb{R}^N; \xi_v = \sin(k_v R_v t) \text{ für } v=1, 2, \dots, N; k \in \mathbb{R}^N \}$ der Würfel $W_\mu(0)$ gesetzt wurde, ist erlaubt wegen $W_\mu(0) = \{ \xi = \frac{z}{\mu}; z \in W_\mu(0) \}$, und weil jede $m \times m$ Matrix $a \cdot M$ für $a \in \mathbb{R}$ diagonalisierbar ist durch eine Matrix S , wenn das für M gilt.

Bemerkung 2:

Falls (11) nicht hyperbolisch oder $|A(k)| \geq 1 + \varepsilon$ für ein $\varepsilon \in \mathbb{R}^+$ und für alle $k \in \bar{N} \subset \mathbb{R}^N$, wobei \bar{N} nicht eine Nullmenge ist, so wächst $\|V^n\|_{\mathcal{G}}$ und damit $\|U^n\|_F$ exponentiell mit $(1 + \varepsilon)^n$, d.h. es gibt ein $b \in \mathbb{R}^+$, so daß

$$\|U^n\| \geq b (1 + \varepsilon)^n$$

Bemerkung 3:

Im Falle HERMITEScher Systeme, d.h. für A_v HERMITESCH ($v=1, 2, \dots, N$) hat man nur die Bedingung c) zu erfüllen, denn $\sum_{v=1}^N \frac{A_v}{R_v} \xi_v$ ist dann für

alle $f \in \mathbb{R}^N$ HERMITESch.

Man kann noch die Aussage des Satzes abschwächen durch:

$$\begin{aligned} f \in W_1(0) : |B(f)| &= \left| \sum_{v=1}^N \frac{A_v}{R_v} f_v \right| = \left| \max_{\substack{x \in \mathbb{R}^N \\ |x|=1}} \sum_{v=1}^N \frac{(A_v x)_v}{R_v} f_v \right| \leq \\ &\leq \sum_{v=1}^N \frac{|f_v|}{R_v} \left| \max_{\substack{x \in \mathbb{R}^N \\ |x|=1}} (A_v x)_v \right| = \sum_{v=1}^N \frac{|f_v|}{R_v} |A_v| \leq \sum_{v=1}^N \frac{|A_v|}{R_v} \end{aligned}$$

wegen der Extremaleigenschaft HERMITEScher Formen.

\Rightarrow Stabilität herrscht für

$$\sum_{v=1}^N \frac{|A_v|}{R_v} \leq 1 - \varepsilon$$

Bemerkung 4:

Unter den Voraussetzungen des Satzes existiert insbesondere zu einem wie unter II.A. vorausgesetzten u_0 eine eindeutige Lösung von (11).

Jede am Anfang von der Ordnung 2 approximierende Schar von Lösungen von (10) konvergiert daher von derselben Ordnung gegen $u(t)$. Die Wahl von $T \in \mathbb{R}^+$ ist hier willkürlich.

C. Stabilitätsbereich und Abhängigkeitsgebiet

Es läßt sich zeigen, daß das Schema (10) in der Geometrie seines Gitters das Schema mit dem größten überhaupt möglichen Stabilitätsbereich ist. Diese Eigenschaft ist von praktischer Bedeutung, da erfahrungsgemäß für die größten Zeitschrittweiten, für die das Schema noch stabil ist, der Fehler am kleinsten wird.

In ihrer klassischen Arbeit [2] haben COURANT, FRIEDRICHS und LEWY Instabilität unter dem Gesichtspunkt von Abhängigkeitsfragen behandelt. Die Lösung $u(t, x)$ des Anfangswertproblems (11) ist an der Stelle (t, x) ; $t \in \mathbb{R}_+^+$ eine eindeutige Funktion von Punkten $(0, x^*)$ aus kompakten Mengen in der Hyperebene $t = 0$. Der Durchschnitt aller dieser Mengen ist das Abhängigkeitsgebiet $D(t, x)$ des Punktes (t, x) , s. LAX [4]. Für einen Punkt $(t_n, x_n) \in [0, T] \times G =: G_n^T$, $t_n = \tau_n h$ gibt es eine endliche Menge $\hat{D}_n(t_n, x_n)$ in der Anfangsebene $t = 0$, von der der Wert der Lösung von (10) in (t_n, x_n) (abgesehen von Punkten in der zweiten Schicht) allein abhängt. Die konvexe Hülle $D_n(t_n, x_n)$ von $\hat{D}_n(t_n, x_n)$ bezeichnet man als das Abhängigkeitsgebiet des Differenzenschemas im Gitter für den Punkt (t_n, x_n) , kurz "Gitterabhängigkeitsgebiet". Wenn die Zeitschrittweite so groß im Verhältnis zu den Raumschrittweiten gewählt wird, daß $F_n(t_n, x_n) := D(t_n, x_n) \setminus D_n(t_n, x_n) \neq \emptyset$, so kann man den Anfangswert $u_0(\cdot)$ in $F_n(t_n, x_n)$ nachträglich so verändern, daß die zugehörige exakte Lösung $u(t, x)$ in (t_n, x_n) einen anderen Wert bekommt, während sich am Wert der Lösung des Differenzenschemas u_j^m an dieser Stelle nichts ändert ($j = (j_1, \dots, j_n)$).

Bei Verfeinerung des Gitters kann daher U_j^n nicht gegen $u(t, x)$ konvergieren, da es unendlich viele exakte Lösungen von (11) gibt, die mit U_j^0, U_j^1 für alle möglichen j beliebig genau übereinstimmen.

Die Forderung $D(t_n, x_n) \subset D_n(t_n, x_n)$ für alle $(t_n, x_n) \in G_n^T$ bedeutet daher eine unmittelbar einleitende, "natürliche" Bedingung für Stabilität eines Schemas. Was über sie hinausgeht wird man als ein Maß für Ungereimtheiten betrachten dürfen, die durch die Konstruktion des Differenzschemas, nicht durch das zu approximierende Anfangswertproblem, herbeigeführt sind.

Da das System (11) linear und homogen mit konstanten Koeffizienten ist, genügt es das Abhängigkeitsgebiet D des Punktes $(1, 0) \in \mathbb{R}^{N+1}$ zu betrachten. Dazu werden solche G_n^T gewählt, die diesen Punkt enthalten; das zugehörige Gitterabhängigkeitsgebiet sei D_n . Wegen der Konvexität von D_n ist

$$(22) \quad D \subset D_n \iff H \subset D_n$$

worin H die konvexe Hülle von D ist.

D_n läßt sich leicht bestimmen: zur Berechnung von U_j^{n+1} werden im Schema (10) die Werte $T_v U_j^n, T_v^{-1} U_j^n, U_j^{n-1}$ für $v=1, 2, \dots, N$ benötigt. Stellt man sich G_n^T noch um eine Schicht $\{-h\} \times G$ erweitert vor und läßt das Schema (10) in ihr beginnen, so ist D_n der N -dimensionale Rhombus mit den Eckpunkten

$$\begin{aligned} & (R_n, 0, \dots, 0), (-R_n, 0, \dots, 0) \\ & (0, R_n, 0, \dots, 0), (0, -R_n, 0, \dots, 0) \\ & \vdots \\ & (0, 0, \dots, 0, R_n), (0, 0, \dots, 0, -R_n) \end{aligned}$$

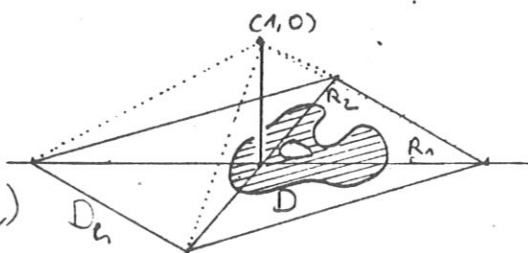


Fig. 1

d.h. die konvexe Hülle der aus diesen Punkten bestehenden Menge (s. Figur 1).

H ist von LAX für das Problem (11) in [4] durch die Trägerfunktion "f" bestimmt worden, die für eine beliebige, beschränkte Menge $M \subset \mathbb{R}^N$ so definiert ist:

$$f_M(k) := \sup_{x \in M} k \cdot x \quad \text{für alle } k \in \mathbb{R}^N$$

Mit $k \cdot x$ ist das Skalarprodukt $k \cdot x = \sum_{i=1}^N k_i \cdot x_i$ gemeint.

Da M beschränkt, existiert $f_M(k)$ immer.

f_M ist durch die konvexe Hülle von M eindeutig bestimmt und umge-

kehrt. Für zwei beschränkte konvexe Mengen M, M' gilt genau dann

$$M' \subset M$$

wenn $f_{M'}(k) \leq f_M(k)$ für alle $k \in \mathbb{R}^N$

Die Notwendigkeit dieser Bedingung ist anhand der Definition von f_M unmittelbar einzusehen. Daß sie auch hinreichend ist, folgt genauso wie die eindeutige Bestimmtheit aus der Theorie der konvexen Mengen (s. LAX [4]).

Da D_k und H konvex sind, lautet die COURANT-FRIEDRICHS-LEWY-(CFL) Bedingung (22):

$$(23) \quad f_H(k) \leq f_{D_k}(k) \quad \text{für alle } k \in \mathbb{R}^N$$

Wegen der Homogenität von f_M für nicht-negative Parameter:

$$f_M(a \cdot k) = a \cdot f_M(k) \quad \text{für } a \in \mathbb{R}^+ \cup \{0\}, k \in \mathbb{R}^N$$

ist (23) äquivalent zu

$$(231) \quad f_H(k) \leq f_{D_k}(k) \quad \text{für alle } k \in P$$

mit einer Menge P , die die Bedingung

$$\mathbb{R}^N = \{k' \in \mathbb{R}^N; \text{es gibt } k \in P, b \in \mathbb{R}^+ \cup \{0\} \text{ mit } k' = b \cdot k\} \text{ erfüllt.}$$

Eine solche Menge P , die von jeder Halbgeraden mit 0 als Anfangspunkt geschnitten wird, ist zum Beispiel der Rand ∂Q_R des N -dimensionalen Quaders Q_R :

$$Q_R := \left\{ k \in \mathbb{R}^N; |k_v| \leq \frac{1}{R_v} \text{ für } v=1,2,\dots,N \right\}; \quad R := (R_1, R_2, \dots, R_N) \\ (R_v > 0 \text{ für } v=1,\dots,N)$$

Darauf läßt sich f_{D_Q} bestimmen:

Behauptung:

$$f_{D_Q}(k) = 1 \quad \text{für alle } k \in \partial Q_R$$

Beweis:

Wegen der Symmetrie des Rhombus D_Q und des Quaders Q_R gibt es zu jedem $x \in D_Q$ ein $x^+ \in D_Q$ mit $x_v^+ = |x_v|$, zu jedem $k \in \partial Q_R$ ein $k^+ \in \partial Q_R$ mit $k_v^+ = |k_v|$ und für $y \in D_Q, k \in \partial Q_R$ ein $\tilde{y} \in D_Q$ mit $|\tilde{y}_v| = |y_v|$ und $\tilde{y}_v \cdot k_v \geq 0$ für $v=1,2,\dots,N$.

$$\Rightarrow y \cdot k \leq \tilde{y} \cdot k = y^+ \cdot k^+$$

$$\Rightarrow f_{D_Q}(k) = f_{D_Q^+}(k^+) \quad \text{für } k \in \partial Q_R$$

$$\text{mit } D_Q^+ := \{x \in D_Q, x_v \geq 0 \text{ für } v=1,2,\dots,N\}$$

oder anders geschrieben:

$$f_{D_Q}(k) = f_{D_Q^+}(k^+) \quad \text{für } k^+ \in \partial Q_R^+$$

mit analogen ∂Q_R^+ :

$$\partial Q_R^+ := \{k \in \partial Q_R, k_v \geq 0 \text{ für } v=1,2,\dots,N\}$$

D_Q^+ ist das von der Hyperebene $\sum_{v=1}^N \frac{x_v}{R_v} = 1$ und den Koordinatenebenen

$x_x = 0, x=1,2,\dots,N$ erzeugte Simplex:

$$D_Q^+ = \left\{ x \in \mathbb{R}^N, \sum_{v=1}^N \frac{x_v}{R_v} \leq 1, x_x \geq 0 \text{ für } x=1,2,\dots,N \right\}$$

Sei $k \in \partial Q_R^+$ fest. Dann gibt es ein $l \in \mathcal{N}$, $l \leq \mathcal{N}$ mit

$$k_l = \frac{1}{R_l}.$$

Für $x \in D_k^+$ ist:

$$\begin{aligned} x_l &\leq R_l \cdot \left(1 - \sum_{\substack{v=1 \\ v \neq l}}^{\mathcal{N}} \frac{x_v}{R_v}\right) \\ \Rightarrow k \cdot x &= \sum_{v=1}^{\mathcal{N}} k_v x_v \leq \sum_{\substack{v=1 \\ v \neq l}}^{\mathcal{N}} k_v x_v + \frac{1}{R_l} \cdot R_l \left(1 - \sum_{\substack{v=1 \\ v \neq l}}^{\mathcal{N}} \frac{x_v}{R_v}\right) \\ &= \sum_{\substack{v=1 \\ v \neq l}}^{\mathcal{N}} \left(k_v - \frac{1}{R_v}\right) x_v + 1 \\ &\leq 1 \quad \text{da } k_v \leq \frac{1}{R_v}, x_v \geq 0 \text{ für } v=1, 2, \dots, \mathcal{N} \end{aligned}$$

Für $x_{R_l} := (0, \dots, 0, R_l, 0, \dots, 0) \in D_k^+$, d.h. die l -te Komponente

ist größer als Null, ist aber $k \cdot x_{R_l} = 1$

\Rightarrow Für alle $k \in \partial Q_R^+$ existiert sogar $\max_{x \in D_k^+} k \cdot x$ und ist gleich 1.

Es ist also

$$f_{D_k^+}(k^+) = f_{D_k^+}(k) = 1 \text{ für alle } k \in \partial Q_R \text{ q.e.d.}$$

LAX hat nun in [4] die Trägerfunktion f_H berechnet zu:

$$f_H(k) = \lambda_{\max} \left(\sum_{v=1}^{\mathcal{N}} k_v \cdot A_v \right)$$

worin $\lambda_{\max}(M)$ den größten Eigenwert der $m \times m$ Matrix M bezeichnet.

Wegen (23)' ist daher die CFL-Bedingung äquivalent zu:

$$\lambda_{\max} \left(\sum k_v \cdot A_v \right) \leq 1 \quad \text{für alle } k \in \partial Q_R$$

$$\Leftrightarrow \lambda_{\max} \left(\sum k_v \cdot \frac{A_v}{R_v} \right) \leq 1 \quad \text{für alle } k \in \partial W_1 (= \partial W_1(0))$$

mit einem Würfel $W_1 := \{k \in \mathbb{R}^{\mathcal{N}}; |k_v| \leq 1 \text{ für } v=1, 2, \dots, \mathcal{N}\}$

$$\Leftrightarrow \lambda_{\max} \left(\sum k_v \cdot \frac{A_v}{R_v} \right) \leq 1 \quad \text{für alle } k \in W_1 \rightarrow$$

Denn für $k' \in W_n$ gibt es ein $b \in [0, 1]$ und
ein $k \in \partial W_n$ mit $k' = b \cdot k$.

Wegen $\lambda_{\max}(-B) = -\lambda_{\min}(B)$ und

$k \in W_n \Leftrightarrow -k \in W_n$ folgt:

$$\Rightarrow \lambda_{\min} \left(\sum_{v=1}^N k_v \frac{A_v}{R_v} \right) \geq -1$$

für $k \in W_n$

Man hat damit bewiesen:

Satz 3: Das Schema (10) erfüllt genau dann die CFL-Bedingung für
(11), wenn

$$\left| \sum_{v=1}^N k_v \frac{A_v}{R_v} \right| \leq 1 \quad \text{für alle } k \in W_n$$

Dies ist bis auf ein beliebig klein wählbares $\zeta > 0$ genau (21), die
Stabilitätsbedingung c) des Satzes 1. Aus den Ergebnissen des fdgen-
den Abschnitts wird man ersehen, daß das ζ die Folge schwacher d.h.
polynomartiger Instabilitäten, im Gegensatz zu den exponentiellen, ist.

III. Schwache Stabilität

In diesem Abschnitt wird der Fall der nicht nur strikten Hyperbolizität von (11) untersucht. (10) ist dann nicht mehr wie in II. stabil sondern i.a. schwach stabil.

Man wählt wieder die Form (13) des Schemas (10):

$$(13) \quad \tilde{U}^{n+1} = \begin{pmatrix} \sum_{v=1}^N \frac{A_v}{R_v} (T_v - T_v^{-1}) & I \\ I & 0 \end{pmatrix} \tilde{U}^n \quad \text{mit} \quad \tilde{U}^n := \begin{pmatrix} U^{n+1} \\ U^n \end{pmatrix} \\ =: \tilde{C}(h) \tilde{U}^n$$

$\tilde{U}^n = \tilde{U}^n(\cdot)$ ist darin eine Funktion abhängig von $x \in \mathbb{R}^N$.

Sei jetzt einfach $h \in \mathbb{R}^+$ ohne die Beschränkung $h \in [0, T]$.

Mit der Argumentation in (15)-(17) erhält man für die Transformierte:

$$\tilde{V}^{n+1} = \underbrace{\begin{pmatrix} 2i \sum_{v=1}^N \frac{A_v}{R_v} \sin k_v R_v h & I \\ I & 0 \end{pmatrix}}_{=: \tilde{A}_h(k)} \tilde{V}^n(k) \quad \text{für } k \in \mathbb{R}^N \\ \text{mit } \tilde{V}^n(k) = \begin{pmatrix} V^{n+1}(k) \\ V^n(k) \end{pmatrix}$$

Um schwache Stabilität zu beweisen hat man daher die Potenzen der $2m \times 2m$

Matrix $\tilde{A}_h(k)$ zu betrachten. Dazu bietet sich die Wahl einer

Normalform an. Hier wird vorläufig die von JORDAN genommen, und zwar

seien:

$$A_h(k) := \sum_{v=1}^N \frac{A_v}{R_v} \sin k_v R_v h,$$

$$J_h(k) = T_h^{-1}(k) A_h(k) T_h(k) \quad \text{mit regulären Matrizen } T_h(k) \\ \text{für } k \in \mathbb{R}^N$$

$J_h(k)$ ist eine $m \times m$ Diagonalübermatrix:

$$J_h(k) := \text{Diag}(J_{h,1}(k), J_{h,2}(k), \dots, J_{h,r}(k))$$

mit $m_v \times m_v$ Matrizen $J_{h,v}(k)$ für $v = 1, 2, \dots, r$,
und $\sum_{v=1}^r m_v = m$.

Die $J_{h,v}(k)$ sind wieder Diagonalübermatrizen.

$$J_{h,v}(k) := \text{Diag}(J_{h,v,1}(k), J_{h,v,2}(k), \dots, J_{h,v,s}(k))$$

mit $m_{\nu\mu} \times m_{\nu\mu}$ Matrizen $J_{h,v,\mu}(k)$ für $\mu = 1, 2, \dots, s$,
und $\sum_{\mu=1}^s m_{\nu\mu} = m_v$. $J_{h,v,\mu}(k)$ hat schließlich die Gestalt:

$$J_{h,v,\mu}(k) = \lambda_{h,v}(k) I_{m_{\nu\mu}} + N_{m_{\nu\mu}}$$

Die $\lambda_{h,1}(k), \dots, \lambda_{h,r}(k)$ sind die r verschiedenen Eigenwerte
von $A_h(k)$; $I_{m_{\nu\mu}}$ ist die $m_{\nu\mu} \times m_{\nu\mu}$ Einheits-
matrix und $N_{m_{\nu\mu}}$ mit derselben Zeilenlänge quadratisch und
aus lauter Nullen außer $m_{\nu\mu} - 1$ Einsen in der ersten oberen
Nebendiagonale bestehend:

$$I_{m_{\nu\mu}} := \underbrace{\begin{pmatrix} 1 & & & 0 \\ & \ddots & & \\ & & \ddots & \\ 0 & & & 1 \end{pmatrix}}_{m_{\nu\mu}}; \quad N_{m_{\nu\mu}} := \underbrace{\begin{pmatrix} 0 & 1 & & 0 \\ & 0 & \ddots & \\ & & \ddots & 1 \\ 0 & & & 0 \end{pmatrix}}_{m_{\nu\mu}}$$

$N_{m_{\nu\mu}}$ ist also nilpotent.

$$\text{Dann gilt: } \begin{pmatrix} T_h^{-1}(k) & 0 \\ 0 & T_h^{-1}(k) \end{pmatrix} \tilde{A}_h(k) \begin{pmatrix} T_h(k) & 0 \\ 0 & T_h(k) \end{pmatrix} =$$

$$= \begin{pmatrix} T_h^{-1}(k) & 0 \\ 0 & T_h^{-1}(k) \end{pmatrix} \begin{pmatrix} 2i A_h(k) & I \\ I & 0 \end{pmatrix} \begin{pmatrix} T_h(k) & 0 \\ 0 & T_h(k) \end{pmatrix} =$$

$$= \begin{pmatrix} 2i J_h(k) & I \\ I & 0 \end{pmatrix} =: \tilde{J}_h(k)$$

Das Verhalten von $\tilde{A}_i(\lambda)^m$ ist daher bis auf Transformationsmatrizen durch das von $\tilde{J}_h(\lambda)^m$ gegeben. Die Multiplikation der

Untermatrizen gebildet aus

$\tilde{J}_e(\lambda)$ durch Wegnahme

der Zeilen von der 1.

bis $\sum_{\alpha=1}^{\nu} m_{\alpha} + \sum_{\alpha=1}^{k-1} m_{\nu+\alpha} t_{\alpha}$,

der $\sum_{\alpha=1}^{\nu} m_{\alpha} + \sum_{\alpha=1}^{k-1} m_{\nu+\alpha+1} t_{\alpha}$

bis zur

$$\sum_{\alpha=1}^{\nu} m_{\alpha} + \sum_{\alpha=1}^{k-1} m_{\nu+\alpha} + m t_{\alpha}$$

und von der

$$\sum_{\alpha=1}^{\nu} m_{\alpha} + \sum_{\alpha=1}^{k-1} m_{\nu+\alpha+1} + m t_{\alpha}$$

bis zur $2m$ ten Zeile

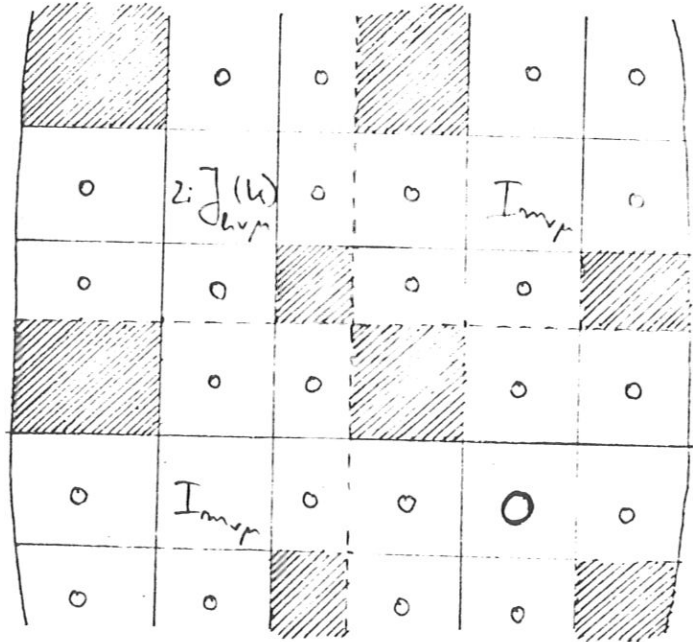


Fig. 2

und derselben Spalten, d.h. die Multiplikation der Matrizen

$$\begin{pmatrix} 2i J_{h, \nu, \nu}(\lambda) & I_{m_{\nu, \nu}} \\ I_{m_{\nu, \nu}} & 0 \end{pmatrix}$$

ist voneinander unabhängig, wegen des Schachbrettaufbaus von $\tilde{J}_h(\lambda)$ (s. Fig. 2). Man kann sich daher auf die Potenzen dieser Matrizen beschränken.

Hilfssatz:

Das charakteristische und das Minimalpolynom der $2p \times 2p$ Matrix

($p \in \mathbb{N}$)

$$\tilde{J} := \begin{pmatrix} 2i J_{\lambda} & I_p \\ I_p & 0 \end{pmatrix} \text{ mit } J_{\lambda} := \lambda I_p + N_p = \begin{pmatrix} \lambda & & & 0 \\ & \lambda & & 0 \\ & & \ddots & \\ 0 & & & \lambda \end{pmatrix} \text{ für } \lambda \in \mathbb{C}$$

sind gegeben durch $(z^2 - 2i\lambda z - 1)^p$.

Korollar:

Für $\lambda \neq \pm 1$ ist die JORDAN'sche Normalform \tilde{J}_0 von \tilde{J} gegeben

durch $\begin{pmatrix} J_{\lambda_+} & 0 \\ 0 & J_{\lambda_-} \end{pmatrix}$ mit $\lambda_{\pm} := i\lambda \pm \sqrt{1-\lambda^2}$
(bis auf Vertauschung von $J_{\lambda_+}, J_{\lambda_-}$),

und für $\lambda = \pm 1$ durch $\tilde{J}_0 = \pm iI_{2p} + N_{2p}$

Die Lösungen der charakteristischen Gleichung $(z^2 - 2i\lambda z - 1)^p = 0$ sind nämlich λ_+ und λ_- , und zwar jeweils von der Vielfachheit p.

$$\lambda_+ \lambda_- = -1, \quad \lambda_+ + \lambda_- = 2i\lambda$$

Da das Minimalpolynom mit den charakteristischen Polynom zusammenfällt ist die Länge der JORDAN-Matrizen dann p bzw. 2p. Für den Beweis des Hilfssatzes ist nützlich das

Lemma:

$$D_{2p} := \text{Det} \begin{pmatrix} aI_p + bN_p & I_p \\ I_p & cI_p \end{pmatrix} = \text{für } a, b, c \in \mathbb{C}, p \in \mathbb{N}$$

Beweis des Lemmas durch vollständige Induktion über p:

Für $p = 1$ ist $D_2 = \text{Det} \begin{pmatrix} a & 1 \\ 1 & c \end{pmatrix} = (ac-1)^1$

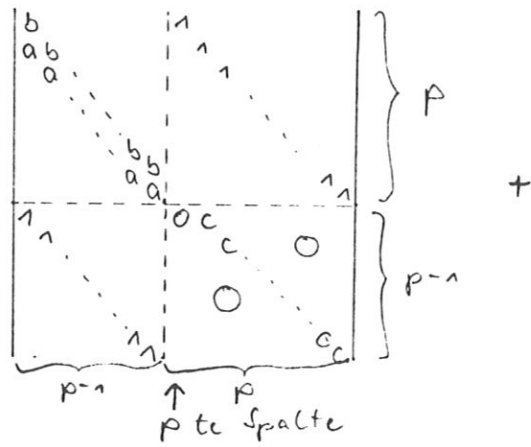
Die Behauptung sei für $p-1$ bewiesen:

Sei $|D| := \text{Det}(D)$ für jede quadratische Matrix D.

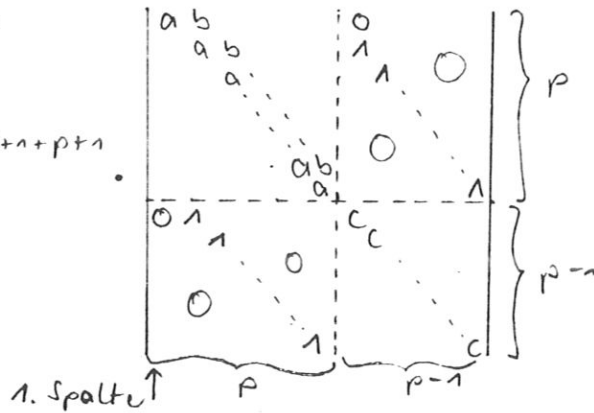
$$D_{2p} = \begin{vmatrix} ab & & & & 1 & & & & 0 \\ ab & & & & 1 & & & & 0 \\ ab & & & & & & & & \\ 0 & & & & & & & & \\ 0 & & b & & & & & & \\ & & ab & & & & & & \\ & & ab & & & & & & \\ 1 & & & & & & & & \\ 1 & & & & & & & & \\ 0 & & & & c & & & & \\ & & & & c & & & & \\ & & & & c & & & & \\ 0 & & & & & & & & \\ & & & & & & & & \\ & & & & 1 & & & & \\ & & & & & & & & \\ & & & & & & & & c \\ & & & & & & & & c \end{vmatrix} \begin{matrix} \leftarrow \\ \\ \\ \\ \\ \end{matrix}$$

Entwicklung nach der $p+1$ ten Zeile

$$D_{2p} = (-1)^{p+1+p+1} \cdot$$

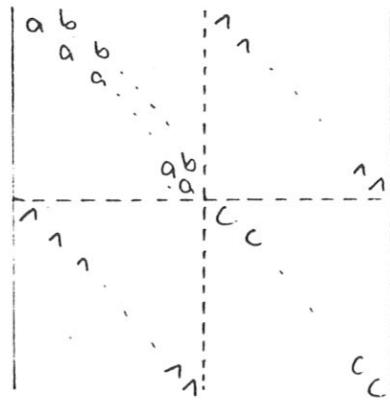


$$+ C \cdot (-1)^{p+1+p+1} \cdot$$



Die entstandenen Determinanten lassen sich nach der p-ten bzw. 1. Spalte entwickeln, wobei dieselbe Restdeterminante übrig bleibt.

$$D_{2p} = [(-1)^p \cdot (-1)^{p+1} + C \cdot a] \cdot$$



$$= (ac - 1) \cdot D_{2p-2}$$

$$= (ac - 1)^p$$

nach Induktionsvoraussetzung.

g.e.d. (Lemma)

Anmerkung: b geht nicht ein in die Determinante.

Beweis des Hilfssatzes:

Setzt man $a = 2i\lambda - b$, $b = 2i$, $c = -b$, so hat man die Determinante von

$$\left(\begin{array}{c|c} 2iJ_\lambda - bI_p & I_p \\ \hline I_p & -bI_p \end{array} \right) \text{ zu } (b^2 - 2i\lambda b - 1)^p \text{ bestimmt.}$$

Damit ist die Gestalt des charakteristischen Polynoms geklärt.

Das Minimalpolynom muß im Falle $\lambda_+ \neq \lambda_-$, d.h. $\lambda \neq \pm 1$, den Faktor $b^2 - 2i\lambda b - 1$ enthalten, da das charakteristische Polynom dann diese beiden verschiedenen Nullstellen besitzt.

Es ist:

$$\tilde{J}^2 = \left(\begin{array}{c|c} (2iJ_\lambda)^2 + I_p & 2iJ_\lambda \\ \hline 2iJ_\lambda & I_p \end{array} \right)$$

und weiter

$$\tilde{J}^2 - 2i\lambda\tilde{J} - I_p = \left(\begin{array}{c|c} (2iJ_\lambda)^2 - (2i)^2\lambda J_\lambda & 2iJ_\lambda - 2i\lambda I_p \\ \hline 2iJ_\lambda - 2i\lambda I_p & 0 \end{array} \right)$$

Für J_λ hat man: $J_\lambda^2 = \lambda^2 I_p + 2\lambda N_p + N_p^2$

$$\Rightarrow J_\lambda^2 - 2\lambda J_\lambda = \lambda^2 I_p + 2\lambda N_p + N_p^2 - \lambda^2 I_p - 2\lambda N_p = N_p^2 = N_p(\lambda + N_p) = N_p J_\lambda$$

$$\Rightarrow \tilde{J}^2 - 2i\lambda\tilde{J} - I_p = \begin{pmatrix} -4N_p J_\lambda & 2iN_p \\ 2iN_p & 0 \end{pmatrix} = 2i \begin{pmatrix} N_p & 0 \\ 0 & N_p \end{pmatrix} \begin{pmatrix} 2iJ_\lambda & I_p \\ I_p & 0 \end{pmatrix} = 2i \begin{pmatrix} N_p & 0 \\ 0 & N_p \end{pmatrix} \tilde{J}$$

Die beiden letzten Matrizen darf man vertauschen, denn N_p kommutiert mit $J_\lambda = \lambda I_p + N_p$. \tilde{J} ist eine nicht-singuläre Matrix, denn sie besitzt nur die Eigenwerte λ_+, λ_- mit $\lambda_+ \lambda_- = -1 \Rightarrow$

$$\lambda_+, \lambda_- \neq 0$$

N_p dagegen ist nilpotent, und zwar ist $\alpha = p$ in $N_p^\alpha = 0$

der kleinste Exponent, der diese Relation erfüllt. Daher verschwindet

$$(\tilde{J}^2 - 2i\lambda \tilde{J} - I_{2p})^\alpha$$

zuerst für $\alpha = p$, d.h.:

$$(\lambda^2 - 2i\lambda - 1)^p$$

ist auch das Minimalpolynom für $\lambda \neq \pm 1$

Im Falle $\lambda = +1$ oder -1 hat das charakteristische Polynom

die Gestalt:

$$(\lambda \mp i)^{2p}$$

Man weiß daher nur, daß das Minimalpolynom den linearen Faktor $\lambda \mp i$ enthält und irgendeine Potenz mit einem Exponenten $\leq 2p$ davon ist.

Der Fall gerader Exponenten läßt sich wie oben behandeln:

$$(\tilde{J} \mp i I_{2p})^{2\alpha} = (\tilde{J}^2 \mp 2i\tilde{J} - I_{2p})^\alpha = (2i)^\alpha \begin{pmatrix} N_p & 0 \\ 0 & N_p \end{pmatrix}^\alpha \tilde{J}^\alpha$$

und das verschwindet zuerst für $\alpha = p$. Falls

$$(\tilde{J} \mp i I_{2p})^\tau = 0$$

für ein ungerades $\tau \in \mathbb{N}$, $\tau < 2p$,

so müßte insbesondere $(\tilde{J} \mp i I_{2p})^{\tau p - 1}$ gleich der Nullmatrix sein; für

jene gilt:

$$\begin{aligned} (\tilde{J} \mp i I_{2p})^{2p-1} &= (\tilde{J} \mp i I_{2p})^{2(p-1)} (\tilde{J} \mp i I_{2p}) = \\ &= (2i)^{p-1} \tilde{J}^{p-1} \begin{pmatrix} N_p^{p-1} & 0 \\ 0 & N_p^{p-1} \end{pmatrix} \begin{pmatrix} 2i\tilde{J} \mp i I_p & I_p \\ I_p & \mp i I_p \end{pmatrix} = \end{aligned}$$

$$= (2i)^{p-1} \tilde{J}^{p-1} \begin{pmatrix} 0 & 1 & | & 0 \\ 0 & & | & 0 \end{pmatrix} \begin{pmatrix} \pm i I_p + 2i N_p & I_p \\ I_p & \mp i I_p \end{pmatrix} =$$

\uparrow \uparrow
 p te $2p$ te Spalte

$$= (2i)^{p-1} \tilde{J}^{p-1} \begin{pmatrix} 0 & \pm i & | & 0 & 1 \\ 0 & 1 & | & 0 & \mp i \end{pmatrix} \begin{matrix} \leftarrow 1. \text{ Zeile} \\ \leftarrow (p+1) \text{te} \end{matrix}$$

Der 2ρ -Vektor $\begin{matrix} \lambda \rightarrow \\ \rho+1 \rightarrow \end{matrix} \begin{pmatrix} \pm i \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \pm i \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ \pm i \\ 0 \\ \vdots \\ 0 \end{pmatrix}$ ist ein Eigenvektor von \tilde{J} zum Eigen-

wert $\pm i$:
$$\left(\begin{array}{c|c|c} \pm 2i I_\rho + 2i N_\rho & I_\rho & \\ \hline I_\rho & 0 & \end{array} \right) \begin{pmatrix} \pm i \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} -1 \\ 0 \\ \vdots \\ 0 \\ \pm i \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \pm i \begin{pmatrix} \pm i \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\Rightarrow (\tilde{J} \mp i I_{2\rho})^{2\rho-1} = (2i)^{\rho-1} (\pm i)^{\rho-1} \begin{pmatrix} 0 & \pm i & & 0 \\ & 0 & & 0 \\ & & & 0 \\ & & & 0 \\ 0 & & & 1 \\ & & & 0 \\ & & & 0 \\ & & & 0 \end{pmatrix} \neq 0$$

\Rightarrow

$(\tilde{J} \mp i I_{2\rho})^{\tau}$ verschwindet also nicht für ungerade $\tau < 2\rho$.

Damit ist der Hilfssatz bewiesen.

Die Potenzen von \tilde{J} sind bis auf eine nicht-singuläre Transformationsmatrix \tilde{T} gegeben durch \tilde{J}_0^n ; aus

$$\tilde{J} = \tilde{T}^{-1} \tilde{J}_0 \tilde{T} \quad \text{folgt} \quad \tilde{J}^n = \tilde{T}^{-1} \tilde{J}_0^n \tilde{T} \quad \text{für } n \in \mathbb{N}$$

Polynomartiges Anwachsen der Komponenten der Matrix $\tilde{A}_\mu(\kappa)$

läßt sich daher für feste Werte der Parameter recht leicht kontrollieren.

Zunächst hat man von den Eigenwerten λ von $A_\mu(\kappa)$ vorauszusetzen

$$|\lambda| \leq 1, \quad \lambda \in \mathbb{R},$$

da sonst (s.II.B.) für die Eigenwerte λ_+, λ_- von $\tilde{A}_\mu(\kappa)$ gilt:

$$|\lambda_+|, |\lambda_-| > 1, \quad ,$$

so daß zum Beispiel die Elemente in der Hauptdiagonale von \tilde{J}_0^n , die gleich λ_+^n bzw. λ_-^n sind, exponentiell divergieren. (Hier geht auch die Hyperbolizität wieder ein: $\lambda \in \mathbb{R}$).

Für JORDAN Matrizen $\mu I_p + N_p, \mu \in \mathbb{C}$, gilt nämlich:

$$(24) \quad (\mu I_p + N_p)^n = \sum_{v=0}^n \binom{n}{v} \mu^{n-v} N_p^v = \sum_{v=0}^{\min(n, p-1)} \binom{n}{v} \mu^{n-v} N_p^v$$

für $\mu \in \mathbb{C}, p \in \mathbb{N}$

da N_p nilpotent ist mit $N_p^p = 0$ ($N_p^0 := I_p$).

Dann hat man $A_L(k)$ nach Eigenwerten $+1$ oder -1 zu untersuchen;

sei q deren größter Elementarteilerexponent im Minimalpolynom.

Aus dem Korollar des Hilfssatzes weiß man, daß die JORDAN'sche Normal-

form der zugehörigen Matrix \tilde{J} die Matrix $\pm I_{2q} + N_{2q}$ ist, wegen (24)

enthält \tilde{J}_0^n daher in der rechten oberen Ecke ein Element $(\pm i)^{n-2q+1} \binom{n}{2q-1}$,
für $n \geq 2q-1$

Das ist ein Polynom in n vom Grade $2q-1$, also dasjenige mit dem höchsten Grad.

Es bleiben noch die Eigenwerte λ von $A_L(k)$ mit $-1 < \lambda < 1$.

Sei p ihr größter Elementarteilerexponent, das zugehörige \tilde{J} hat

jetzt die JORDAN'sche Normalform

$$\begin{pmatrix} J_{\lambda_+} & 0 \\ 0 & J_{\lambda_-} \end{pmatrix} \quad \text{mit} \quad J_{\lambda_{\pm}} = \lambda_{\pm} I_p + N_p$$

$J_{\lambda_{\pm}}^n$ hat in der rechten oberen Ecke ein Element $(\lambda_{\pm})^{n-p+1} \cdot \binom{n}{p-1}$

für $n \geq p-1$, das ist das Polynom in n mit dem höchsten Grad unter allen Elementen von \tilde{J}^n . Da $|\lambda_{\pm}| = 1$, ändert die $n-p+1$ te Potenz von λ_{\pm} daran nichts, d.h. λ_{\pm}^n konvergiert nicht gegen 0 und divergiert nicht gegen ∞ , sondern bleibt betragsmäßig nach oben und nach unten beschränkt.

Da man durch die Wahl größerer Gitterverhältnisse R_{ν} die Eigenwerte von $A_L(k)$ verkleinern kann, sind die ersten beiden der drei Fälle leicht auszuschließen, etwa durch Auferlegung einer Bedingung (21). Der dritte Fall ist durch das Problem (11) unweigerlich bedingt. Es wird daher i.a.

Elemente in $\tilde{A}_L(\kappa)$ geben, die Polynome vom Grade $p-1$ in n sind, mit dem obigen ρ . Von diesen Elementen hat man schwache Stabilität desselben Grades $p-1$ zu erwarten, wegen $L^{1-p} = \frac{n^{p-1}}{t^{p-1}}$ für alle in (5) zugelassenen $h, n, t=nh$.

Diese Argumentation vermag zwar eine plausible Erklärung für das zu erwartende Verhalten der Potenzen der Amplifikationsmatrix $\tilde{A}_L(\kappa)$ geben, mathematisch ist sie schlicht unbefriedigend, da die Parameter unberücksichtigt geblieben sind. Die Schwäche dieser Vorgehensweise liegt in der Unkenntnis der Transformationsmatrizen

$$\tilde{T}_L(\kappa) := \begin{pmatrix} T_L(\kappa) & 0 \\ 0 & T_L(\kappa) \end{pmatrix} .$$

Die in solchen Situationen häufig gestellte Forderung der gleichmäßigen Beschränktheit dieser Matrizen in einer geeigneten Norm, und zwar gleichmäßig beschränkt in K , ist eine reichlich realitätsfremde Auflage, wenn sich der maximale Elementarteilerexponent für einen Eigenwert auch wirklich ändert. Die Eigenwerte von $\tilde{A}_L(\kappa)$ selbst kann man noch als stetig von K abhängig auffassen, da die Koeffizienten des sie zu Nullstellen habenden charakteristischen Polynoms stetige Funktionen in K sind. An einer solchen Sprungstelle κ^* der Elementarteilerexponenten wird man aber wohl kaum eine Schar von derartigen Matrizen $T_L(\kappa)$

$$\left. \begin{array}{l} \{ T_L(\kappa) \cdot T_L^{-1}(\kappa) \text{ existiert, } \gamma_L(\kappa) = T_L^{-1}(\kappa) A_L(\kappa) T_L(\kappa) \text{ für} \\ \text{alle } \kappa \text{ aus einer Umgebung von } \kappa^* \} \end{array} \right\} \text{ finden,}$$

die zusammen mit der analogen Schar für $T_L^{-1}(\kappa)$ dort stetig oder auch nur gleichmäßig beschränkt ist in der Norm, wenn so etwas nicht überhaupt grundsätzlich unmöglich ist.

Wählt man dagegen eine gegenüber der von JORDAN etwas veränderte Normalform, nämlich eine solche, in der statt der $J_{\lambda, \nu, \mu}(\kappa)$ Matrizen der Form

$$L_{\lambda, \nu, \mu}(\kappa) := \lambda_{\lambda, \nu}(\kappa) I_{m_{\nu}} + c_{\lambda, \nu, \mu}(\kappa) N_{m_{\nu}}$$

so lassen sich damit in der Praxis brauchbare Kriterien für die schwache Stabilität von $\tilde{A}_e(\kappa)$ entwickeln. Ein Beispiel wird noch später gegeben, welches auch der Leistungsfähigkeit des Schemas (10) angemessen ist, das ja in der Regel von 2. Ordnung konsistent ist, also nur schwache Stabilität für Grade < 2 kompensieren kann.

Zunächst ist das auch tatsächlich eine Normalform, denn

$$L_{\lambda} := \lambda I_p + c N_p \quad \text{mit} \quad \lambda, c \in \mathbb{C}, \quad p \in \mathbb{N} \quad \text{hat das}$$

$c \neq 0$

charakteristische und das Minimalpolynom $(\lambda - \lambda)^p$, und ist daher für $c \neq 0$ der JORDAN Matrix $\lambda I_p + N_p$ ähnlich. Denn für das charakteristische Polynom sieht man das sofort, und der Ausdruck

$$(\lambda I_p + c N_p - \lambda I_p)^{\alpha} = c^{\alpha} N_p^{\alpha} \text{ verschwindet zuerst für } \alpha = p \text{ falls } c \neq 0$$

Der nächste Schritt besteht in der direkten Ausrechnung der Potenzen von

$$\tilde{L} := \begin{pmatrix} \lambda & L_{\lambda} & I_p \\ I_p & \lambda & 0 \end{pmatrix}$$

Dazu ist einiger rechnerischer Aufwand erforderlich.

Lemma: Sei $p \in \mathbb{N}$.

Die Potenzen der $2p \times 2p$ Matrix $\tilde{K} := \begin{pmatrix} X_p & I_p \\ I_p & 0 \end{pmatrix}$ mit einer $p \times p$ Matrix $X_p \in \mathbb{C}^p \times \mathbb{C}^p$ sind gegeben durch:

$$(25) \quad \tilde{\mathcal{K}}^n = \begin{pmatrix} \mathcal{K}_n & \mathcal{K}_{n-1} \\ \mathcal{J}_{n-1} & \mathcal{J}_{n-2} \end{pmatrix} \quad \text{für } n \in \mathbb{N}$$

Darin sind die \mathcal{K}_v $p \times p$ Matrizen, und zwar:

$$\mathcal{K}_\ell := \begin{cases} \sum_{\substack{v=0 \\ 2v \leq \ell}}^{\lfloor \ell/2 \rfloor} \binom{\ell-v}{v} X_p^{\ell-2v} & \text{für } \ell \in \mathbb{N} \cup \{0\} \\ 0 & \text{für } \ell = -1 \end{cases}$$

Die Summation ist dabei über alle ganzen Zahlen von 0 bis zu der zu führen, die gerade noch kleiner oder gleich $\frac{\ell}{2}$ ist.

Beweis durch vollständige Induktion über n :

Die Behauptung (25) ist offenbar richtig für $n=1$ ($X_p^0 := I_p$).

Der Index p wird jetzt weggelassen. Für $n=2$ hat man:

$$\mathcal{K}_2 = I + X^2, \quad \mathcal{K}_1 = X, \quad \mathcal{K}_0 = I \quad \text{und es ist auch}$$

$$\tilde{\mathcal{K}}^2 = \begin{pmatrix} X^2 + I & X \\ X & I \end{pmatrix}$$

Sei (25) für $v \in \mathbb{N}$, $v \leq n$ bewiesen. Dann ist:

$$\tilde{\mathcal{K}}^{n+1} = \tilde{\mathcal{K}} \cdot \tilde{\mathcal{K}}^n = \begin{pmatrix} X & I \\ I & 0 \end{pmatrix} \cdot \begin{pmatrix} \mathcal{K}_n & \mathcal{K}_{n-1} \\ \mathcal{K}_{n-1} & \mathcal{K}_{n-2} \end{pmatrix} \quad \text{und das soll sein}$$

$$\text{gleich} \quad \stackrel{!}{=} \begin{pmatrix} \mathcal{K}_{n+1} & \mathcal{K}_n \\ \mathcal{K}_n & \mathcal{K}_{n-1} \end{pmatrix}$$

In der 2. Zeile stehen rechts und links dieselben Matrizen. Man hat nur noch zu beweisen:

$$\mathcal{K}_{n+1} = X \mathcal{K}_n + \mathcal{K}_{n-1} \quad \text{und noch einmal}$$

$$\mathcal{K}_n = X \mathcal{K}_{n-1} + \mathcal{K}_{n-2}$$

$$X \mathcal{K}_n + \mathcal{K}_{n-1} = \sum_{\substack{v=0 \\ 2v \leq n}} \binom{n-v}{v} X^{n-2v+1} + \sum_{\substack{v=0 \\ 2v \leq n-1}} \binom{n-1-v}{v} X^{n-2v} \cong$$

Indexverschiebung in der 2. Summe:

$$\cong \sum_{\substack{v=0 \\ 2v \leq n}} \binom{n-v}{v} X^{n-2v+1} + \sum_{\substack{v=1 \\ 2v \leq n+1}} \binom{n-v}{v-1} X^{n+1-2v} =$$

$$= X^{n+1} + \sum_{\substack{v=1 \\ 2v \leq n}} \binom{n+1-v}{v} X^{n+1-2v} + \sum_{n < 2v \leq n+1} \binom{n-v}{v-1} X^{n+1-2v} =$$

$$\text{wegen } \binom{n-v}{v} + \binom{n-v}{v-1} = \binom{n+1-v}{v}$$

$$= \sum_{\substack{v=0 \\ 2v \leq n+1}} \binom{n+1-v}{v} X^{n+1-2v} = \mathcal{K}_{n+1}$$

denn

$$\binom{n - \frac{n+1}{2}}{\frac{n+1}{2} - 1} = \binom{\frac{n}{2} - \frac{1}{2}}{\frac{1}{2} - \frac{1}{2}} = 1 = \binom{n+1 - \frac{n+1}{2}}{\frac{n+1}{2}} \text{ falls } n \text{ ungerade}$$

g. e. d. (Lemma)

Das Lemma gestattet einen Einblick in die Struktur von $\tilde{\mathcal{L}}^m$:

$$\tilde{\mathcal{L}}^m = \begin{pmatrix} \mathcal{K}_n & \mathcal{K}_{n-1} \\ \mathcal{K}_{n-1} & \mathcal{K}_{n-2} \end{pmatrix} \quad \text{für } m \in \mathbb{N}$$


mit
$$\mathcal{K}_\ell := \sum_{\substack{v=0 \\ 2v \leq \ell}} \binom{\ell-v}{v} (z_i)^{\ell-2v} L_\lambda^{\ell-2v} \quad \text{für } \ell \in \mathcal{N}_0$$

und $\mathcal{K}_\ell = 0$ für $\ell = -1$

Sei $\ell \in \mathcal{N}_0$. Es ist $L_\lambda = 2I_p + cN_p$.

$$\begin{aligned} \Rightarrow \mathcal{K}_\ell &= \sum_{\substack{v=0 \\ 2v \leq \ell}} \binom{\ell-v}{v} (z_i)^{\ell-2v} (\lambda I_p + cN_p)^{\ell-2v} = \\ &= \sum_{\substack{v=0 \\ 2v \leq \ell}} \binom{\ell-v}{v} (z_i)^{\ell-2v} \sum_{\alpha=0}^{\min(\ell-2v, p-1)} \binom{\ell-2v-\alpha}{\alpha} \lambda^{\ell-2v-\alpha} c^\alpha N^\alpha = \\ &= \sum_{\alpha=0}^{p-1} c^\alpha N^\alpha \sum_{\substack{v=0 \\ 2v \leq \ell-\alpha}} \binom{\ell-v}{v} \binom{\ell-2v-\alpha}{\alpha} (z_i)^{\ell-2v} \lambda^{\ell-2v-\alpha} \end{aligned}$$

Für die 2. Summe wird darin 0 vereinbart, falls $\ell - \alpha < 0$

\mathcal{K}_ℓ ist also eine Dreiecksmatrix: . In jeder oberen Nebendiagonale sind alle Elemente gleich. Sei $k_{\ell\alpha}$ ein Element in der α ten ($\alpha = 0, 1, 2, \dots, p-1$; für $\alpha = 0$ hat man die Hauptdiagonale), oberen Nebendiagonalen, dann ist:

$$k_{\ell\alpha} = (z_i c)^\alpha \sum_{\substack{v=0 \\ 2v \leq \ell-\alpha}} \binom{\ell-v}{v} \binom{\ell-2v-\alpha}{\alpha} (z_i \lambda)^{\ell-2v-\alpha}$$

An dieser Darstellung kann man allerdings das Verhalten der Elemente von \tilde{L}^n noch nicht ablesen; es bedarf noch einer Umschreibung, die die Eigenwerte λ_+, λ_- einführt.

Transformiert man \tilde{L} auf seine JORDAN'sche Normalform und berechnet mit dem Ergebnis die Potenzen von \tilde{L} , so gewinnt man für $p=2, 3$ überschaubare Ausdrücke für $k_{\ell\alpha}$. Von $p=4$ ab wird das Verfahren recht umfangreich; insbesondere ist für allgemeine $p \in \mathcal{N}$ zwar noch die Matrix \tilde{T} zu berechnen, die \tilde{L} auf JORDAN'sche

Normalform transformiert, die Inversion von \tilde{T} aber ist recht kompliziert. Die für $p=2, 3$ erhaltenen Formeln für $K_{e,x}$ lassen glücklicherweise eine rein formal begründete Extension für größere p zu. Dies zur Motivation des folgenden

Hilfssatzes:

(26) Für $n, x, n-x \in \mathbb{N} \setminus \{0\}$ und a, b aus einem kommutativen Ring mit Einselement gilt:

$$\begin{aligned}
 \text{(i)} \quad & (a-b)^{2x+1} \sum_{\substack{v=0 \\ 2v \leq n-x}}^{\binom{n-v}{v} \binom{n-2v}{x}} (a+b)^{n-x-2v} (ab)^v = \\
 \text{(ii)} \quad & = (a-b)^{2x+1} \sum_{\mu=0}^{n-x} \binom{n-\mu}{x} \binom{x+\mu}{x} a^{n-x-\mu} b^{\mu} = \\
 \text{(iii)} \quad & = \sum_{\tau=0}^x \binom{n+x+1}{\tau} \binom{n-\tau}{x-\tau} [a^{n+1+x-2\tau} - b^{n+1+x-2\tau}] (ab)^{\tau}
 \end{aligned}$$

Die Summen in (i) und (ii) sind ebenfalls gleich.

Beweis durch Doppelinduktion über n und x .

Wenn man sich anhand der Matrizenmultiplikation klargemacht hat, wie der Induktionsschritt auszusehen hat, besteht der Beweis nur noch aus Rechnerei.

Für den Induktionsanfang hat man zunächst den Fall $x=0$

zu beweisen:

$$\begin{aligned}
 (a-b) \sum_{\substack{v=0 \\ 2v \leq n}}^{\binom{n-v}{v}} (a+b)^{n-2v} (ab)^v &= (a-b) \sum_{\mu=0}^n a^{n-\mu} b^{\mu} = \\
 &= a^{n+1} - b^{n+1}
 \end{aligned}$$

Die 2. Gleichheit ist die bekannte Darstellung der geometrischen Reihe, die der beiden Summen (bezeichnet mit M_n) folgt wie im Lemma durch Induktion über n :

$$\begin{aligned} \text{für } n=0 \quad \text{ist} \quad & : 1 = 1 \\ n=1 \quad & : a+b = a+b \end{aligned}$$

Induktionsschritt: $M_{n+1} = (a+b)M_n - abM_{n-1}$

(i) Beh.: M_n erfüllt den Induktionsschritt für $n \in \mathbb{N}$

Beweis:

$$\begin{aligned} & (a+b)M_n - abM_{n-1} = \\ &= (a+b) \sum_{\substack{v=0 \\ 2v \leq n}} \binom{n-v}{v} (a+b)^{n-2v} (-ab)^v - ab \sum_{\substack{v=0 \\ 2v \leq n-1}} \binom{n-1-v}{v} (a+b)^{n-1-2v} (-ab)^v = \end{aligned}$$

Indexverschiebung wie oben

$$\begin{aligned} &= (a+b)^{n+1} + \sum_{\substack{v=1 \\ 2v \leq n}} \left[\binom{n-v}{v} + \binom{n-v}{v-1} \right] (a+b)^{n+1-2v} (-ab)^v \\ & \quad + \sum_{n \leq v < n+1} \binom{n-v}{v-n} (a+b)^{n+1-2v} (-ab)^v = \\ &= M_{n+1} \end{aligned}$$

(ii) Beh.: M_n erfüllt den Induktionsschritt für $n \in \mathbb{N}$:

Beweis:

$$\begin{aligned} & (a+b)M_n - abM_{n-1} = \\ &= (a+b) \sum_{\mu=0}^n a^{n-\mu} b^\mu - ab \sum_{\mu=0}^{n-1} a^{n-1-\mu} b^\mu = \\ &= a^{n+1} + b^{n+1} + 2 \sum_{\mu=0}^{n-1} ab \cdot a^{n-1-\mu} b^\mu - ab \sum_{\mu=0}^{n-1} a^{n-1-\mu} b^\mu = \\ &= \sum_{\mu=0}^{n+1} a^{n+1-\mu} b^\mu = M_{n+1} \end{aligned}$$

q. e. d. (Induktionsanfang 1. Teil)

Man braucht weiter für den Induktionsanfang (26) für die Parameter $\alpha = n$ und $\alpha = n-1$. Der Fall $\alpha = n$ (ist leicht einzusehen, in (i) und (ii) ergeben die Summen beide den Wert 1 und in (iii) steht die Binomialentwicklung von $(a-b)^{2n+1}$.

Sei $\alpha = n-1$, $n \in \mathbb{N}$. In (i), (ii) stehen die Ausdrücke

$$(i) \quad (a-b)^{2n-1} \cdot n \cdot (a+b) =$$

$$(ii) \quad (a-b)^{2n-1} \cdot [na + nb]$$

die beiden Summen sind also wieder gleich.

Die Übereinstimmung mit (iii) sieht man so ein:

$$\begin{aligned} n(a-b)^{2n-1}(a+b) &= n \sum_{v=0}^{2n-1} \binom{2n-1}{v} a^{2n-1-v} (-b)^v (a+b) = \\ &= n \sum_{v=0}^{2n-1} \binom{2n-1}{v} a^{2n-v} (-b)^v - n \sum_{v=0}^{2n-1} \binom{2n-1}{v} a^{2n-1-v} (-b)^{v+1} = \\ &= n(a^{2n} + b^{2n}) + n \sum_{v=n}^{2n-1} \left[\binom{2n-1}{v} - \binom{2n-1}{v-1} \right] a^{2n-v} (-b)^v = \end{aligned}$$

nach Indexverschiebung in der 2. Summe.

$$\begin{aligned} &= \sum_{v=0}^{2n} \binom{2n}{v} (n-v) a^{2n-v} (-b)^v, \text{ denn } n \cdot [\dots] = \binom{2n}{v} (n-v) \\ &= \sum_{\tau=0}^{n-1} \binom{2n}{\tau} (n-\tau) [a^{2n-2\tau} - b^{2n-2\tau}] (-ab)^\tau \end{aligned}$$

Für $\alpha = 0, n-1, n$ ist der Hilfssatz also richtig, das ist der Induktionsanfang.

Induktionsschritt:

Der Ausdruck unter (i), (ii), (iii) im Hilfssatz wird jetzt jeweils

mit $M_{n,x}$ bezeichnet für alle $n, x, n-x \in \mathbb{N}_0$. Für

$n, x, n-x \in \mathbb{N}$ wird dann in allen drei Fällen der Induktionsschritt

$$(27) \quad M_{n+1,x} = (a+b)M_{n,x} - abM_{n-1,x} + (a-b)^2 M_{n,x-1} \text{ bewiesen,}$$

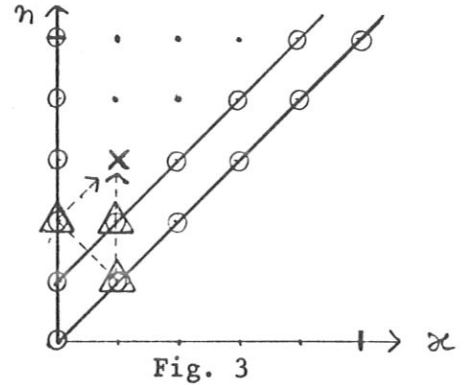
aus welchem dann die Gleichheit von

(i), (ii), (iii) in (26) folgt. In (i), (ii)

sind die Zusammenhänge einfacher, dort

bezeichne $N_{n,x}$ jeweils die Summe.

Sei also $n, x, n-x \in \mathbb{N}$.



(i) Es muß (27) gezeigt werden. Wegen

$$M_{n,x} = (a-b)^{2x+1} N_{n,x} \quad \text{folgt das aus}$$

$$(28) \quad N_{n+1,x} = (a+b)N_{n,x} - abN_{n-1,x} + N_{n,x-1}$$

Es ist

$$\begin{aligned} & (a+b)N_{n,x} - abN_{n-1,x} + N_{n,x-1} = \\ & = \sum_{\substack{v=0 \\ 2v \leq n-x}} \binom{n-v}{v} \binom{n-2v}{x} (a+b)^{n-x-2v+1} (-ab)^v + \\ & + \sum_{\substack{v=0 \\ 2v \leq n-1-x}} \binom{n-1-v}{v} \binom{n-1-2v}{x} (a+b)^{n-1-x-2v} (-ab)^{v+1} + \\ & + \sum_{\substack{v=0 \\ 2v \leq n-x+1}} \binom{n-v}{v} \binom{n-2v}{x-1} (a+b)^{n-x-2v+1} (-ab)^v \quad \Rightarrow \end{aligned}$$

Indexverschiebung in der 2. Summe:

$$\begin{aligned}
 &\Rightarrow \dots + \\
 &+ \sum_{\substack{v=1 \\ 2v \leq n+1-\alpha}} \binom{m-v}{v-1} \binom{m-2v+1}{\alpha} (a+b)^{m-\alpha-2v+1} (-ab)^v + \\
 &+ \dots = \\
 &= \left[\binom{n}{\alpha} + \binom{n}{\alpha-1} \right] (a+b)^{n-\alpha+1} + \sum_{n-\alpha < 2v \leq n-\alpha+1} \left[\binom{m-v}{v-1} \binom{m-2v+1}{\alpha} + \binom{m-v}{v} \binom{m-2v}{\alpha-1} \right] \cdot \\
 &\quad \cdot (a+b)^{m-\alpha-2v+1} (-ab)^v + \\
 &+ \sum_{\substack{v=1 \\ 2v \leq n-\alpha}} \left[\binom{m-v}{v} \binom{m-2v}{\alpha} + \binom{m-v}{v-1} \binom{m-2v+1}{\alpha} + \binom{m-v}{v} \binom{m-2v}{\alpha-1} \right] (a+b)^{m-\alpha-2v+1} (-ab)^v = \\
 &= \binom{m+1}{\alpha} (a+b)^{n-\alpha+1} + \sum_{n-\alpha < 2v \leq n-\alpha+1} \underbrace{\left[\binom{m-v}{v-1} \binom{m-n+\alpha-1+1}{\alpha} + \binom{m-v}{v} \binom{m-n+\alpha-1}{\alpha-1} \right]}_{= \binom{m+1-v}{v}} \cdot (a+b)^{m-\alpha-2v+1} (-ab)^v + \\
 &+ \sum_{\substack{v=1 \\ 2v \leq n-\alpha}} \binom{m+1-v}{v} \binom{m+1-2v}{\alpha} (a+b)^{m+1-\alpha-2v} (-ab)^v \Rightarrow
 \end{aligned}$$

Weil in der Randsumme allenfalls ein Term mit $v = \frac{n-\alpha+1}{2}$

beiträgt, und wegen:

$$\begin{aligned}
 &\binom{m-v}{v} \binom{m-2v}{\alpha} + \binom{m-v}{v-1} \binom{m-2v+1}{\alpha} + \binom{m-v}{v} \binom{m-2v}{\alpha-1} = \\
 &= \left[\binom{m-v}{v} + \binom{m-v}{v-1} \right] \binom{m-2v+1}{\alpha} = \binom{m+1-v}{v} \binom{m+1-2v}{\alpha}
 \end{aligned}$$

$$\begin{aligned}
 &\Rightarrow \sum_{\substack{v=0 \\ 2v \leq n+1-\alpha}} \binom{m+1-v}{v} \binom{m+1-2v}{\alpha} (a+b)^{m+1-\alpha-2v} (-ab)^v = N_{m+1, \alpha}
 \end{aligned}$$

Die Relation (27) ist also erfüllt.

(ii)

Wieder genügt es statt (27) die einfachere Relation (28) zu zeigen.

Es ist:

$$\begin{aligned}
 & (a+b) N_{n,x} - a b N_{n-1,x} + N_{n,x-1} = \\
 & = \sum_{\mu=0}^{n-x} \binom{n-\mu}{x} \binom{x+\mu}{x} a^{n-x-\mu} b^{\mu} (a+b) + \\
 & - \sum_{\mu=0}^{n-x-1} \binom{n-1-\mu}{x} \binom{x+\mu}{x} a^{n-1-\mu-x} b^{\mu+1} + \\
 & + \sum_{\mu=0}^{n+1-x} \binom{n-\mu}{x-1} \binom{x-1+\mu}{x-1} a^{n+1-x-\mu} b^{\mu} = \\
 & = \sum_{\mu=0}^{n-x} \binom{n-\mu}{x} \binom{x+\mu}{x} a^{n+1-x-\mu} b^{\mu} + \sum_{\mu=0}^{n-x} \binom{n-\mu}{x} \binom{x+\mu}{x} a^{n-x-\mu} b^{\mu+1} + \\
 & - \dots + \\
 & + \dots = \\
 & = \sum_{\mu=0}^{n-x} \left[\binom{n-\mu}{x} \binom{x+\mu}{x} + \binom{n-\mu}{x-1} \binom{x-1+\mu}{x-1} \right] a^{n+1-x-\mu} b^{\mu} + \binom{n}{x-1} b^{n+1-x} + \\
 & + \sum_{\mu=0}^{n-1-x} \underbrace{\left[\binom{n-\mu}{x} - \binom{n-1-\mu}{x} \right]}_{= \binom{n-1-\mu}{x-1}} \binom{x+\mu}{x} a^{n-x-\mu} b^{\mu+1} + \binom{n}{x} b^{n+1-x} = \\
 & = \dots + \binom{n+1}{x} b^{n+1-x} + \\
 & + \sum_{\mu=1}^{n-x} \binom{n-\mu}{x-1} \binom{x+\mu-1}{x} a^{n+1-x-\mu} b^{\mu} \cong
 \end{aligned}$$

nach Indexverschiebung

$$\begin{aligned} &\Rightarrow \left[\binom{m}{x} + \binom{m}{x-1} \right] a^{m+1-x} + \binom{m+1}{x} b^{m+1-x} + \\ &+ \sum_{\mu=1}^{m-x} \left[\binom{m-\mu}{x} \binom{x+\mu}{x} + \binom{m-\mu}{x-1} \binom{x-1+\mu}{x-1} + \binom{m-\mu}{x-1} \binom{x+\mu-1}{x} \right] a^{m+1-x-\mu} b^{\mu} = \\ &= \sum_{\mu=0}^{m+1-x} \binom{m+1-\mu}{x} \binom{x+\mu}{x} a^{m+1-x-\mu} b^{\mu} = M_{m+1, x} \end{aligned}$$

wegen

$$\begin{aligned} &\binom{m-\mu}{x} \binom{x+\mu}{x} + \binom{m-\mu}{x-1} \left[\binom{x-1+\mu}{x-1} + \binom{x-1+\mu}{x} \right] = \\ &= \left[\binom{m-\mu}{x} + \binom{m-\mu}{x-1} \right] \binom{x+\mu}{x} = \binom{m+1-\mu}{x} \binom{x+\mu}{x} \end{aligned}$$

Hier ist (27) also auch erfüllt.

(iii)

In diesem Falle zeigt man am besten (27) direkt:

$$\begin{aligned} &(a+b) M_{m, x} - ab M_{m-1, x} + (a-b)^2 M_{m, x-1} = \\ &= \sum_{\tau=0}^x \binom{m+x+1}{\tau} \binom{m-\tau}{x-\tau} [a^{m+1+x-2\tau} - b^{m+1+x-2\tau}] (-ab)^{\tau} (a+b) + \\ &+ \sum_{\tau=0}^x \binom{m+x}{\tau} \binom{m-1-\tau}{x-\tau} [a^{m+x-2\tau} - b^{m+x-2\tau}] (-ab)^{\tau+1} + \\ &+ \sum_{\tau=0}^{x-1} \binom{m+x}{\tau} \binom{m-\tau}{x-1-\tau} [a^{m+x-2\tau} - b^{m+x-2\tau}] (-ab)^{\tau} (a-b)^2 \Rightarrow \end{aligned}$$

Es ist

$$\begin{aligned}
 & [a^{n+1+\alpha-2\tau} - b^{n+1+\alpha-2\tau}] (-ab)^\tau (a+b) = \\
 & = (-1)^\tau [a^{n+2+\alpha-2\tau} b^\tau + a^{n+1+\alpha-\tau} b^{\tau+1} - b^{n+1+\alpha-\tau} a^{\tau+1} + \\
 & \quad - b^{n+2+\alpha-\tau} a^\tau] \\
 & [a^{n+\alpha-2\tau} - b^{n+\alpha-2\tau}] (-ab)^\tau \cdot (a^2 - 2ab + b^2) = \\
 & = (-1)^\tau [a^{n+2+\alpha-\tau} b^\tau - 2a^{n+1+\alpha-\tau} b^{\tau+1} + a^{n+\alpha-\tau} b^{\tau+2} + \\
 & \quad - b^{n+\alpha-\tau} a^{\tau+2} + 2b^{n+1+\alpha-\tau} a^{\tau+1} - b^{n+2+\alpha-\tau} a^\tau]
 \end{aligned}$$

Setzt man in die Summen ein und ordnet nach Potenzen von a und b,

so steht da:

$$\begin{aligned}
 \Rightarrow & \sum_{\tau=0}^{\alpha} \left\{ (-1)^\tau \binom{n+\alpha+1}{\tau} \binom{n-\tau}{\alpha-\tau} [a^{n+2+\alpha-\tau} b^\tau - b^{n+2+\alpha-\tau} a^\tau] + \right. \\
 & \left. + [(-1)^\tau \binom{n+\alpha+1}{\tau} \binom{n-\tau}{\alpha-\tau} + (-1)^{\tau+1} \binom{n+\alpha}{\tau} \binom{n-1-\tau}{\alpha-\tau}] [a^{n+1+\alpha-\tau} b^{\tau+1} - b^{n+1+\alpha-\tau} a^{\tau+1}] \right\} + \\
 & + \sum_{\tau=0}^{\alpha-1} (-1)^\tau \binom{n+\alpha}{\tau} \binom{n-\tau}{\alpha-1-\tau} \cdot [a^{n+2+\alpha-\tau} b^\tau - 2a^{n+1+\alpha-\tau} b^{\tau+1} + \\
 & \quad + a^{n+\alpha-\tau} b^{\tau+2} - b^{n+\alpha-\tau} a^{\tau+2} + 2b^{n+1+\alpha-\tau} a^{\tau+1} - b^{n+2+\alpha-\tau} a^\tau] \Rightarrow
 \end{aligned}$$

Jetzt werden für die einzelnen Summationsterme die Indizes verschoben, so daß a und b nur noch in Monomen der Form

$a^{n+2+\alpha-\tau} b^\tau$ und $b^{n+2+\alpha-\tau} a^\tau$ auftreten.

$$\begin{aligned}
 \Rightarrow & \sum_{\tau=0}^{\alpha} (-1)^\tau \binom{n+\alpha+1}{\tau} \binom{n-\tau}{\alpha-\tau} [a^{n+2+\alpha-\tau} b^\tau - b^{n+2+\alpha-\tau} a^\tau] + \\
 & + \sum_{\tau=1}^{\alpha+1} [(-1)^{\tau-1} \binom{n+\alpha+1}{\tau-1} \binom{n+1-\tau}{\alpha+1-\tau} + (-1)^\tau \binom{n+\alpha}{\tau-1} \binom{n-\tau}{\alpha+1-\tau}] \cdot [a^{n+2+\alpha-\tau} b^\tau - b^{n+2+\alpha-\tau} a^\tau] + \\
 & + \sum_{\tau=0}^{\alpha-1} (-1)^\tau \binom{n+\alpha}{\tau} \binom{n-\tau}{\alpha-1-\tau} [a^{n+2+\alpha-\tau} b^\tau - b^{n+2+\alpha-\tau} a^\tau] + \\
 & + \sum_{\tau=1}^{\alpha} (-1)^\tau \binom{n+\alpha}{\tau-1} \binom{n+1-\tau}{\alpha-\tau} \cdot 2 [a^{n+2+\alpha-\tau} b^\tau - b^{n+2+\alpha-\tau} a^\tau] + \\
 & + \sum_{\tau=2}^{\alpha+1} (-1)^\tau \binom{n+\alpha}{\tau-2} \binom{n+2-\tau}{\alpha+1-\tau} [a^{n+2+\alpha-\tau} b^\tau - b^{n+2+\alpha-\tau} a^\tau] \Rightarrow
 \end{aligned}$$

$$\begin{aligned}
 &= \left[\binom{n}{x} + \binom{n}{x-1} \right] \cdot [a^{n+2+x} - b^{n+2+x}] + \\
 &+ \left[-(n+x+1) \binom{n-1}{x-1} + \binom{n+1}{x} - \binom{n-1}{x} - (n+x) \binom{n-1}{x-2} + 2 \binom{n}{x-1} \right] \cdot [a^{n+1+x} b - b^{n+1+x} a] + \\
 &+ (-1)^x \left[\binom{n+x+1}{x} - (n+1+x) \binom{n+x+1}{x-1} + (n-x) \binom{n+x}{x-1} + 2 \binom{n+x}{x-1} + \right. \\
 &\quad \left. + (n+2-x) \binom{n+x}{x-2} \right] \cdot [a^{n+2} b^x - b^{n+2} a^x] + \\
 &+ (-1)^x \underbrace{\left[\binom{n+x+1}{x} - \binom{n+x}{x} - \binom{n+x}{x-1} \right]}_{=0} \cdot [a^{n+1} b^{x+1} - b^{n+1} a^{x+1}] + \\
 &+ \sum_{\tau=2}^{x-1} (-1)^\tau \left[\binom{n+x+1}{\tau} \binom{n-\tau}{x-\tau} - \binom{n+x+1}{\tau-1} \binom{n+1-\tau}{x+1-\tau} + \binom{n+x}{\tau-1} \binom{n-\tau}{x+1-\tau} + \right. \\
 &\quad \left. + \binom{n+x}{\tau} \binom{n-\tau}{x-1-\tau} + 2 \binom{n+x}{\tau-1} \binom{n+1-\tau}{x-\tau} + \binom{n+x}{\tau-2} \binom{n+2-\tau}{x+1-\tau} \right] \cdot [a^{n+2+x-\tau} b^\tau - b^{n+2+x-\tau} a^\tau]
 \end{aligned}$$

Darin wird vereinbart: $\binom{l}{x-2} = 0$ für $l \in \mathbb{N}_0$, falls $x=1$.

Die Randterme lassen sich ohne Umwege in die gewünschte Form bringen:

$$\binom{n}{x} + \binom{n}{x-1} = \binom{n+1}{x}$$

$$\begin{aligned}
 &-(n+x+1) \binom{n-1}{x-1} + \binom{n+1}{x} - (n+x) \binom{n-1}{x-2} - 2 \binom{n}{x-1} = \\
 &= -(n+x) \left[\binom{n-1}{x-1} + \binom{n-1}{x-2} \right] + \binom{n+1}{x} - \binom{n-1}{x-1} - \binom{n-1}{x} - 2 \binom{n}{x-1} = \\
 &= -(n+x) \binom{n}{x-1} + \binom{n+1}{x} \underbrace{- \binom{n}{x-1}}_{-\binom{n+1}{x}} = -(n+1+x) \binom{n}{x-1}
 \end{aligned}$$

$$\begin{aligned}
 &\binom{n+x+1}{x} - (n+1-x) \binom{n+x+1}{x-1} + (n-x) \binom{n+x}{x-1} + 2 \binom{n+x}{x-1} + (n+2-x) \binom{n+x}{x-2} = \\
 &= \binom{n+x+1}{x} - \binom{n+x+1}{x-1} + 2 \binom{n+x}{x-1} + 2 \binom{n+x}{x-2} + (n-x) \underbrace{\left[-\binom{n+x+1}{x-1} + \binom{n+x}{x-1} + \binom{n+x}{x-2} \right]}_{=0} = \\
 &= \binom{n+x+1}{x} - \binom{n+x+1}{x-1} + 2 \binom{n+x+1}{x-1} = \\
 &= \binom{n+2+x}{x}
 \end{aligned}$$

Die in der Summe auftretenden Binominalkoeffizienten zerlegt man dagegen am besten:

$$\begin{aligned}
 & \binom{n+\alpha+1}{\tau} \binom{n+\alpha+1}{\alpha-\tau} - \binom{n+\alpha+1}{\tau-1} \binom{n+\alpha+1}{\alpha-\tau} + \binom{n+\alpha}{\tau-1} \binom{n-\tau}{\alpha+1-\tau} + \\
 & + \binom{n+\alpha}{\tau} \binom{n-\tau}{\alpha-1-\tau} + 2 \binom{n+\alpha}{\tau-1} \binom{n+\alpha}{\alpha-\tau} + \binom{n+\alpha}{\tau-2} \binom{n+\alpha}{\alpha+1-\tau} = \\
 & = \left[\binom{n+\alpha}{\tau} + \binom{n+\alpha}{\tau-1} \right] \left[\binom{n-\tau}{\alpha-\tau} - \left[\binom{n+\alpha}{\tau-1} + \binom{n+\alpha}{\tau-2} \right] \cdot \left[\binom{n-\tau}{\alpha+1-\tau} + \binom{n-\tau}{\alpha-\tau} \right] \right] + \\
 & + \binom{n+\alpha}{\tau-1} \binom{n-\tau}{\alpha+1-\tau} + \binom{n+\alpha}{\tau} \binom{n-\tau}{\alpha-1-\tau} + 2 \binom{n+\alpha}{\tau-1} \binom{n+\alpha}{\alpha-\tau} + \\
 & + \binom{n+\alpha}{\tau-2} \left[\binom{n-\tau}{\alpha-1-\tau} + 2 \binom{n-\tau}{\alpha-\tau} + \binom{n-\tau}{\alpha+1-\tau} \right] =
 \end{aligned}$$

Die "Hälfte" vom ersten Term, der ganze zweite Term, der dritte und die "Hälfte" vom sechsten Term heben sich heraus. Es bleibt übrig:

$$\begin{aligned}
 & = \binom{n+\alpha}{\tau} \left[\binom{n-\tau}{\alpha-\tau} + \binom{n-\tau}{\alpha-1-\tau} \right] + 2 \binom{n+\alpha}{\tau-1} \binom{n+\alpha}{\alpha-\tau} + \\
 & + \binom{n+\alpha}{\tau-2} \left[\binom{n-\tau}{\alpha-1-\tau} + \binom{n-\tau}{\alpha-\tau} \right] = \\
 & = \left[\binom{n+\alpha}{\tau} + 2 \binom{n+\alpha}{\tau-1} + \binom{n+\alpha}{\tau-2} \right] \cdot \binom{n+\alpha}{\alpha-\tau} = \\
 & = \binom{n+\alpha}{\tau} \binom{n+\alpha}{\alpha-\tau}
 \end{aligned}$$

Es folgt:

$$\begin{aligned}
 & \stackrel{\geq}{=} \sum_{\tau=0}^{\alpha} (-1)^{\tau} \binom{n+\alpha}{\tau} \binom{n+\alpha}{\alpha-\tau} \left[a^{n+\alpha-\tau} b^{\tau} - b^{n+\alpha-\tau} a^{\tau} \right] = \\
 & = M_{n+1} x
 \end{aligned}$$

Damit ist der Hilfssatz bewiesen.

Mit seiner Hilfe läßt sich die Potenzbeschränktheit von

$$\tilde{L} = \begin{pmatrix} 2i\lambda I_{\rho} + 2iN_{\rho} & I_{\rho} \\ \vdots & \vdots \\ I_{\rho} & 0 \end{pmatrix} \quad \text{vollständig überblicken.}$$

Mit $a = \lambda_+$, $b = \lambda_- \Rightarrow a + b = 2i\lambda$, $-ab = 1$ kann man direkt in ihn eingehen und erhält für $k_{n\alpha}$ die Darstellung

$$(29) \quad (i) \quad k_{n\alpha} = (2ic)^x \sum_{\substack{v=0 \\ 2v \leq n-\alpha}}^1 \binom{n-v}{v} \binom{n-2v}{\alpha} (2i\lambda)^{n-\alpha-2v} =$$

$$(ii) \quad = (2ic)^x \sum_{\mu=0}^{n-\alpha} \binom{n-\mu}{\alpha} \binom{\alpha+\mu}{\alpha} \lambda_+^{n-\alpha-\mu} \lambda_-^{\mu} =$$

$$(iii) \quad = \frac{(2ic)^x}{(\lambda_+ - \lambda_-)^{2\alpha+1}} \sum_{\tau=0}^{\alpha} \binom{n+\alpha+1}{\tau} \binom{n-\tau}{\alpha-\tau} \left[\lambda_+^{n+\alpha-\tau} - \lambda_-^{n+\alpha-\tau} \right]$$

wovon die letzte Beziehung (iii) nur für $\lambda_+ \neq \lambda_-$ gilt.

Aus ihr sieht man im Falle $|\lambda| > 1$, d.h. $|\lambda_+|$ oder $|\lambda_-| > 1$,

wegen $|\lambda_+| \cdot |\lambda_-| = 1$ sofort das exponentielle Anwachsen aller

nicht verschwindenden Elemente von \sum^n bedingt durch den führenden

Term $n^x \lambda_+^{n+\alpha+1}$ bzw. $n^x \lambda_-^{n+\alpha+1}$, der für

genügend großes n alle anderen dem Betrage nach majorisiert.

Falls $\lambda_- = \lambda_+$ d.h. $\lambda = +1$ oder $\lambda = -1$, erhält

man aus der Darstellung (ii) für $k_{n\alpha}$:

$$k_{n\alpha} = (\pm 1)^{n-\alpha} (2c)^x i^n \sum_{\mu=0}^{n-\alpha} \binom{n-\mu}{\alpha} \binom{\alpha+\mu}{\alpha}$$

$$= (\pm 1)^{n-\alpha} (2c)^x i^n \binom{n+\alpha+1}{2\alpha+1}$$

Die darin verwandte elementare Reaktion zwischen Binominalkoeffizienten

beweist man am einfachsten durch Koeffizientenvergleich vor dem Monom

$a^x b^x$ in der Summe $\sum_{\mu=0}^{n-x} (a+1)^{n-\mu} (b+1)^{\mu+x}$ und in einem anderen Ausdruck

für sie, den man über Umformungen mittels der Darstellung von geometrischen Reihen erhält. Man hat dann noch das allgemeine Additionstheorem für Binominalkoeffizienten anzuwenden.

In der allgemeineren Form:

$$\binom{n+p}{m} = \sum_{k=0}^m \binom{n-k}{n-m} \binom{p+k-1}{k}$$

findet man die obige Relation

z.B. in [9] auf S. 8 (3b).

Für genügend große n wachsen daher die Elemente $k_{n\alpha}$ dem Betrage nach wie $n^{2\alpha+1}$, abgesehen vom konstanten Faktor $2^\alpha |c|^\alpha$. Das ist das Ergebnis, welches aus den Überlegungen anhand der JORDAN'schen Normalform zu erwarten war.

Es bleibt der hauptsächlich interessierende Fall $-1 < \alpha < 1$ zu untersuchen, in welchem $\alpha_+ \neq \alpha_-$, $|\alpha_+| = |\alpha_-| = 1$. Man schätzt

$$|k_{n\alpha}| \text{ mittels (29) (iii) ab, indem } |\alpha_+^k - \alpha_-^k| \leq 2 \text{ gesetzt wird.}$$

Die noch verbleibende, unhandliche Summe gestattet auch eine bequeme Umformung.

Lemma: Für $n, \alpha, n-\alpha \in \mathbb{N}_0$ gilt:

$$(30) \quad \sum_{\tau=0}^{\alpha} \binom{n+\alpha+1}{\tau} \binom{n-\tau}{\alpha-\tau} = \frac{2^{\alpha} \alpha!}{\alpha!} \prod_{\tau=0}^{\alpha} (n+\alpha-2\tau)$$

Für das Produkt wird der Wert 1 vereinbart, falls $\alpha = 0$, dann stimmt (30) offenbar.

Beweis:

Setzt man in (26) (i) und (iii) für a und b die Zahlen $a = 1$, $b = -1$ ein, so steht da:

$$2^{2\alpha+1} \sum_{\substack{v=0 \\ 2v \leq n-\alpha}}^{\alpha} \binom{n-v}{v} \binom{n-2v}{\alpha} \cdot 0^{n-\alpha-2v} = \sum_{\tau=0}^{\alpha} \binom{n+\alpha+1}{\tau} \binom{n-\tau}{\alpha-\tau} [1 - (-1)^{n+\alpha-2\tau}]$$

Darin verschwinden beide Seiten, falls $n-x$ oder dazu äquivalent $n+x$ ungerade ist. Für $n-x$ gerade, und das ist dasselbe wie $n+x$ gerade, verschwindet in der linken Summe nur der Term mit $v = \frac{n-x}{2}$ nicht, während rechts $1 - (-1)^{n+1+x-2\tau} = 2$ ist, d.h.:

$$2^{2x+1} \binom{n - \frac{n-x}{2}}{\frac{n-x}{2}} \binom{n-n+x}{x} = 2 \sum_{\tau=0}^x \binom{n+x+1}{\tau} \binom{n-\tau}{x-\tau} =$$

$$= 2^{2x+1} \binom{\frac{n+x}{2}}{\frac{n-x}{2}} = \frac{2^{2x+1}}{x!} \prod_{v=0}^{x-1} \left(\frac{n+x}{2} - v\right) = \frac{2^{x+1}}{x!} \prod_{v=0}^{x-1} (n+x-2v)$$

Hält man x fest und betrachtet die rechte und die linke Seite von (30) als Polynome in n , so stimmen diese Polynome im Grade, bei beiden x , und für unendlich viele Werte von n überein, und zwar für $n \in \mathbb{N}_0$ mit $n+x$ gerade. Ihre Differenz ist daher ein Polynom höchstens vom Grade x , das an unendlich vielen Stellen verschwindet, also das Nullpolynom. Damit ist (30) auch für $n+x$ ungerade bewiesen.

Für k_{nx} hat man daher im Falle $-1 < \lambda < 1$:

$$(31) \quad |k_{nx}| \leq \frac{|2c| x \cdot 2^{x+1} \prod_{v=0}^{x-1} (n+x-2v)}{|\lambda_+ - \lambda_-|^{2x+1} \cdot x!} \leq \frac{2 \cdot |4c|^x}{|\lambda_+ - \lambda_-|^{2x+1} \cdot x!} \cdot (n+1)^x$$

wegen $(n+1+a)(n+1-a) = (n+1)^2 - a^2 \leq (n+1)^2$ für $a \in \mathbb{R}$. Der Grad des Polynoms, welches die Elemente von k_n dem Betrage nach abschätzt, ist hier um $x+1$ kleiner als für $\lambda_+ = \lambda_-$.

Um aus (31) schwache Stabilität vom Grade x abzuleiten, hat man zunächst eine Norm $\|\cdot\|$ für die Matrizen \tilde{L}^n zu wählen; hier wird wieder die sich anbietende Spektralnorm genommen.

Für jede $l \times l$ Matrix $Z, l \in \mathbb{N}$, also $Z = (z_{\mu\nu})_{\mu, \nu=1, \dots, l}$ mit $z_{\mu\nu} \in \mathbb{C}$ für alle $\mu, \nu = 1, \dots, l$ sei Z_+ definiert durch:

$$Z_+ := (|z_{\mu\nu}|)_{\mu, \nu=1, \dots, l}$$

und wie üblich die komplex konjugierte Matrix \bar{Z} sowie die Transponierte Z^T durch:

$$\bar{Z} := (\bar{z}_{\mu\nu})_{\mu, \nu=1, \dots, l} \quad ; \quad Z^T := (z_{\nu\mu})_{\mu, \nu=1, \dots, l}$$

Dann ist (s.z.B. [10] S.175) allgemein:

$$|Z| \leq |Z_+| \quad , \text{ und weiter:}$$

$$|Z \bar{Z}^T| \leq |(Z \cdot \bar{Z}^T)_+| \leq |Z_+ (\bar{Z}^T)_+| = |Z_+ (Z_+)^T|$$

d.h.:

$$\|Z\| \leq \|Z_+\|$$

Mit derselben Argumentation weiß man von zwei $l \times l$ Matrizen X, Y mit $X = X_+, Y - X = (Y - X)_+$ die Eigenschaft:

$$\|X\| \leq \|Y\|$$

Damit läßt sich die Spektralnorm von

$$L^n = \begin{pmatrix} \mathcal{K}_n & \mathcal{K}_{n-1} \\ \mathcal{K}_{n-1} & \mathcal{K}_{n-2} \end{pmatrix} \quad \text{für } n \geq 2 \text{ abschätzen.}$$

Zunächst wird eine zusätzliche Bedingung an die Eigenwerte λ_+, λ_- gestellt, die später wie in $\mathbb{I}B$ gleichmäßige Beschränktheit der

Norm in K garantiert, und zwar

$$(20) \quad -1 + \varepsilon \leq \lambda \leq 1 - \varepsilon \quad \text{für ein } \varepsilon \in \mathbb{R}^+, \varepsilon < 1$$

In (31) kann die rechte Seite für größeres λ noch kleiner werden, falls η klein genug. Damit die unschöne Korrektur "für genügend großes η " nicht mitgeschleppt werden muß, wird ein $d \in \mathbb{R}^+ \cup \{0\}$ eingeführt, und zwar sei $\varepsilon' \in \mathbb{R}^+, 0 < \varepsilon' < 1$ und

$$\Rightarrow \quad \sqrt{1 - \varepsilon'^2} := 1 - \varepsilon$$

$$|\lambda_+ - \lambda_-| = 2 \sqrt{1 - \lambda^2} \geq 2 \sqrt{1 - (1 - \varepsilon)^2} = 2 \sqrt{\varepsilon'^2} = 2 \varepsilon'$$

sowie

$$d := \max [|\lambda|, \varepsilon'^2(p-1)] \Rightarrow \frac{d}{\varepsilon'^2(p-1)}^{(n+1)} \geq 0 \quad \text{für } n \geq 0, p \geq 2$$

$$(32) \quad \Rightarrow \quad |k_{n\alpha}| \leq \frac{d^{p-1}}{(p-1)! \varepsilon'^2 p^{-1}} \cdot (n+1)^{p-1} =: m_{n+} \quad \text{für } \alpha = 0, 1, \dots, p-1$$

und das gilt auch für $p=1$ wegen (29)ⁱⁱⁱ, d.h. für $p \in \mathbb{N}, n \in \mathbb{N}_0$.
(Es wird $0^0 = 1$ vereinbart).

Seien $M_n, M_{n+} \in \mathbb{R}^{2p} \times \mathbb{R}^{2p}$ definiert durch

$$M_1 := \left(\underbrace{\begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix}}_{2p} \right) \quad M_{n+} := m_{n+} \cdot M_1$$

Wegen (32) folgt zunächst für $n \geq 2$

$$M_{n+} - (\tilde{L}^n)_+ = (M_{n+} - (\tilde{L}^n)_+)_+$$

Das gilt auch für $n=1$ und $n=0$, denn

$$m_{n+1}, m_{n+1} \geq 1 \quad (\text{für alle } p \in \mathcal{N})$$

$$\Rightarrow \|\tilde{L}^n\| \leq \|M_{n+1}\| = m_{n+1} \|M_1\| = 2^p m_{n+1}$$

M_1 hat offenbar 2^p und sonst die Null 2^{p-1} mal zum Eigenwert; aus der besonderen Gestalt dieser Matrix folgt dann der Wert 2^p für ihre Spektralnorm.

$$(33) \quad \Rightarrow \|\tilde{L}^n\| \leq 2^p \frac{d^{p-1}}{(p-1)! \varepsilon^{1-2p-1}} \quad \text{für } \begin{array}{l} p \in \mathcal{N}, n \in \mathcal{N}_0 \\ \lambda \in \mathbb{R}, |\lambda| \leq 1 - \varepsilon \\ \varepsilon' = \sqrt{1 - (1 - \varepsilon)^2} \end{array}$$

Diese Abschätzung wird zu der der Potenzen von $\tilde{A}_\varepsilon(k)$ dienen.

Der Parameterbereich der Matrizen

$A_\varepsilon(k) = \sum_{v=1}^N \frac{A_v}{R_v} \sin \varepsilon R_v k_v$ ist zwar der \mathbb{R}^N für k , \mathbb{R}^+ für ε , aber die beschränkte Funktion \sin reduziert ihn mit der einfachen Transformation $\xi \in \mathbb{R}^N$, $\xi_v := \sin \varepsilon R_v k_v$ für $k \in \mathbb{R}^N$ auf W_1 , d.h. die Matrizen

$$B(\xi) := A_\varepsilon(k) \quad \text{besitzen kompakten}$$

Parameterbereich.

Läßt man in der JORDAN'schen Normalform $J_\varepsilon(k)$ von $A_\varepsilon(k)$ statt der $J_{\alpha, \nu, \mu}(k)$ die Matrizen $L_{\alpha, \nu, \mu}(k) =: L_{\nu, \mu}(\xi)$ zu, mit $C_{\alpha, \nu, \mu}(k) \in \mathbb{C}$ in der ersten oberen Nebendiagonale, so hat man eine andere Normalform, sie heiße $L(\xi) =: L_\varepsilon(k)$, mit den Eigenwerten $\lambda(\xi)$ in der Hauptdiagonalen.

Satz 4.

a) Das Differentialgleichungssystem

$$(11) \quad u_t(t, x) = \sum_{v=1}^N A_v u_{x_v}(t, x) \quad \text{für } 0 \leq t \leq T$$

sei hyperbolisch, d.h. die Eigenwerte $\lambda(\xi)$ der Matrix

$$B(\xi) := \sum_{v=1}^N \frac{A_v}{R_v} \xi_v \quad \text{für } \xi \in W_1$$

sind reell, diese erfüllen weiter für ein $\varepsilon \in \mathbb{R}^+, 0 < \varepsilon < 1$ die Bedingung

$$(20) \quad -1 + \varepsilon \leq \lambda(\xi) \leq 1 - \varepsilon \quad \text{für } \xi \in W_1.$$

$B(\xi)$ besitzt für jedes ξ eine Normalform $L(\xi)$, $\xi \in W_1$.

Sei $\varepsilon' := \sqrt{1 - (1 - \varepsilon)^2}$.

b) Unter den für $\xi \in W_1$ auftretenden Elementarteilerexponenten von $B(\xi)$ sei $p, p \in \mathbb{N}$, der größte.

c) Die Elemente $c(\xi)$ in der ersten oberen Nebendiagonalen besitzen eine Schranke $d \in \mathbb{R}^+ \cup \{0\}$

$$|c(\xi)| \leq d \quad \text{für alle } \xi \in W_1 \\ \text{mit } d \geq \varepsilon'^2(p-1)$$

d) Es gebe eine Schar von regulären Transformationsmatrizen $T(\xi)$ mit $L(\xi) = T^{-1}(\xi) B(\xi) T(\xi)$ und konstanten $C, C_{-1} \in \mathbb{R}^+$ so daß

$$\|T(\xi)\| \leq C, \|T^{-1}(\xi)\| \leq C_{-1} \quad \text{für alle } \xi \in W_1$$

Dann ist das Differenzenschema

$$(10) \quad U^{n+1} - U^{n-1} = \sum_{v=1}^N \frac{A_v}{R_v} (T_v - T_v^{-1}) U^n$$

schwach stabil vom Grade $p - 1$. Für seine Amplifikationsmatrix gilt:

Beweis:

Seien

$$\tilde{B}(\xi) := \begin{pmatrix} 2; B(\xi) & I_m \\ I_m & 0 \end{pmatrix}, \quad \tilde{T}(\xi) := \begin{pmatrix} T(\xi) & 0 \\ 0 & T(\xi) \end{pmatrix},$$

$$\tilde{L}(\xi) := \begin{pmatrix} 2; L(\xi) & I_m \\ I_m & 0 \end{pmatrix}, \quad \tilde{L}_{\nu, \mu}(\xi) := \begin{pmatrix} 2; L_{\nu, \mu}(\xi) & I_{m_{\nu, \mu}} \\ I_{m_{\nu, \mu}} & 0 \end{pmatrix}$$

Dann ist

$$\begin{aligned} \|\tilde{B}(\xi)^{-1}\| &= \|\tilde{T}(\xi)^{-1} \tilde{L}(\xi)^{-1} T(\xi)\| \leq \|\tilde{T}(\xi)^{-1}\| \cdot \|\tilde{L}(\xi)^{-1}\| \cdot \|T(\xi)\| \leq \\ &\leq C \cdot C^{-1} \|\tilde{L}(\xi)^{-1}\| \end{aligned}$$

Es ist weiter:

$$\|\tilde{L}(\xi)^{-1}\| = |\overline{\tilde{L}(\xi)^{-1}} \cdot \tilde{L}(\xi)^{-1}|$$

$\tilde{L}(\xi)^{-1}, \overline{\tilde{L}(\xi)^{-1}}$ haben beide denselben Aufbau wie $\tilde{Y}_h(\xi)$ in

Fig. 2, die Multiplikation der durch Streichen von Zeilen und Spalten

entstehenden Untermatrizen $\tilde{L}_{\nu, \mu}(\xi)$ bzw. $\overline{\tilde{L}_{\nu, \mu}(\xi)}$

ist daher jeweils für ein Indexpaar ν, μ von den anderen unabhängig.

Daher hat auch die Matrix $\tilde{L}(\xi)^{-1} \cdot \overline{\tilde{L}(\xi)^{-1}} - 2I_{2m}$

dieselbe Struktur. Das charakteristische Polynom von $\tilde{L}(\xi)^{-1} \cdot \overline{\tilde{L}(\xi)^{-1}}$

zerfällt deshalb in die charakteristischen Polynome der

$$L_{\nu, \mu}(\xi)^{-1} \cdot \overline{L_{\nu, \mu}(\xi)^{-1}}$$

Mit anderen Worten: die Spektralnorm von $\tilde{L}(\xi)^{-1}$ ist kleiner oder

gleich der größten Spektralnorm der $\tilde{L}_{\nu, \mu}(\xi)^{-1}$. Wegen (33) und

aus der Definition von d in b) sowie p in c) folgt

$$\|\tilde{L}(\xi)^{-1}\| \leq 2p \frac{d^{p-1}}{(p-1)! \varepsilon^{1/2 p-1}} \cdot (n+1)^{p-1}$$

für $\xi \in W_n, n \in \mathbb{N}_0, -1+\varepsilon \leq \lambda \leq 1-\varepsilon$

Damit ist (34) bewiesen.

Bem. 1

Falls die Amplifikationsmatrix $\tilde{A}_L(u)$ ein Minimalpolynom mit höchstens quadratischen Faktoren besitzt, und sich für sie die Bedingungen c) und d) erfüllen lassen, so folgt Konvergenz von der Ordnung 1 aus dem Äquivalenzsatz für solche $u_0(\cdot)$, wie sie in II A vorausgesetzt sind.

Bem. 2

Mittels apriori-Abschätzungen der Lösungen von (11) lassen sich oft die dritten Ableitungen von $u(\cdot, \cdot)$ die (11) auch lösen, eingrenzen. In solchen Fällen kann man die Konstante P_2 in (12) und damit $M(u_0)$ in (4) berechnen. Verwendet man im Beweis des Äquivalenzsatzes noch (34), so erhält man direkte Fehlerschranken für die Näherungslösung.

Beispiel: $u_t = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} u_x + \begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix} u_y$

Es ist: $B(\mathcal{F}) = \begin{pmatrix} s_1 + s_2 & s_1 - s_2 \\ 0 & s_1 + s_2 \end{pmatrix}$. Sei $s_v = \frac{f_v}{R_v}, v=1, 2$.
Auf JORDAN'sche Normalform transformiert z.B. die Diagonalmatrix $\begin{pmatrix} s_1 - s_2 & 0 \\ 0 & 1 \end{pmatrix}$:

$$\frac{1}{s_1 - s_2} \begin{pmatrix} 1 & 0 \\ 0 & s_1 - s_2 \end{pmatrix} \begin{pmatrix} s_1 + s_2 & s_1 - s_2 \\ 0 & s_1 + s_2 \end{pmatrix} \begin{pmatrix} s_1 - s_2 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} s_1 + s_2 & 1 \\ 0 & s_1 + s_2 \end{pmatrix},$$

die Transformation ist so also nur für $s_1 \neq s_2$ möglich. Wählt man dagegen $\frac{1}{R_1} = \frac{1}{R_2} = \frac{1}{2} - \frac{\epsilon}{2}, \alpha = 1, p = 2$, so sind die Voraussetzungen des Satzes erfüllt.

IV Erweiterung auf ähnliche Schemata

Der Satz 4 zusammen mit der Kenntnis von Instabilitäten für den Fall, in welchem seine Voraussetzung a) nicht gegeben ist, lassen sich ohne Schwierigkeiten auf alle Schemata anwenden, die den Aufbau

$$(35) \quad U^{n+1} - U^{n-1} = \sum_{\nu=1}^N \sum_{\alpha=-L}^L A_{\nu\alpha}(h) T_{\nu}^{\alpha} U^n \quad \text{besitzen.}$$

$h \in \mathcal{N}_0$

Statt $A_{\nu\alpha}(h)$ hat man dann die Matrizen

$$\sum_{\nu=1}^N \sum_{\alpha=-L}^L A_{\nu\alpha}(h) e^{i\alpha h R_{\nu}} K_{\nu} \quad \text{zu untersuchen.}$$

Genausogut aber auch auf Schemata der Form:

$$(36) \quad U^{n+1} - e^{i\pi\varphi} U^{n-1} = \sum_{\nu=1}^N \sum_{\alpha=-L}^L A_{\nu\alpha}(h) T_{\nu}^{\alpha} U^n; \quad \varphi \in \mathbb{R}$$

mit anderen Matrizen $A_{\nu\alpha}(h)$, denn (36) geht durch die Transformation

$$U^n =: e^{i\pi\frac{\varphi}{2}} V^n \quad \text{in (35) über.}$$

Anders ausgedrückt: die beiden $2m \times 2m$ Matrizen

$$\begin{pmatrix} A & I_m \\ e^{i\pi\varphi} I_m & 0 \end{pmatrix} \quad \text{und} \quad \begin{pmatrix} e^{-i\pi\frac{\varphi}{2}} A & I_m \\ I_m & 0 \end{pmatrix} \quad \text{sind bis auf}$$

einen komplexen Faktor vom Betrage 1 unitär kongruent wegen

$$\begin{pmatrix} e^{i\pi\frac{\varphi}{2}} I_m & 0 \\ 0 & I_m \end{pmatrix} \begin{pmatrix} A & I_m \\ e^{i\pi\varphi} I_m & 0 \end{pmatrix} \begin{pmatrix} e^{-i\pi\frac{\varphi}{2}} I_m & 0 \\ 0 & I_m \end{pmatrix} = e^{i\pi\frac{\varphi}{2}} \begin{pmatrix} e^{-i\pi\frac{\varphi}{2}} A & I_m \\ I_m & 0 \end{pmatrix}$$

Damit hat man eine große Auswahl zur Konstruktion stabiler, hochkonsistenter 2-Schritt Schemata.

Das Beispiel für hyperbolische Systeme mit konstanten Koeffizienten

$$(37) \quad U^{n+1} + U^{n-1} = \left\{ 2 + \sum_{v=1}^N \frac{A_v^2}{R_v} (T_v - 2 + T_v^{-1}) + \sum_{\substack{v < l \\ v, l = 1}}^N (A_v A_l + A_l A_v) \frac{T_v - T_v^{-1}}{2R_v} \cdot \frac{T_l - T_l^{-1}}{2R_l} \right\} U^n$$

gehört zur Klasse (36). Für ausreichend Δt differenzierbare Lösungen des Systems (11) gilt:

$$U_{tt} = \sum_{v,l=1}^N A_v A_l U_{x_v x_l}$$

und für sie ist in (37) der Abschneidefehler ein Term $O(\Delta t^4)$.

Die Methode, bei einem zeitabhängigen System 1. Ordnung zum iterierten überzugehen, garantiert das Auftreten der zweiten Ableitungen in der Zeit, die sich besonders gut durch zentrierte Differenzen annähern lassen. Ersetzt man die übrigen Ableitungen durch ebenso genaue Differenzen so erhält man Konsistenz 3. Ordnung für genügend glatte Lösungen.

Die Stabilitätsuntersuchung wird für (37) nur im HERMITESchen Falle durchgeführt. Sei $\bar{A}_v^T = A_v$ für $v = 1, \dots, N$. Dann ist

$$2A_n(k) := 2I_m - 4 \sum_{v=1}^N \frac{A_v^2}{R_v} \sin^2(k R_v \frac{\Delta t}{2}) + \sum_{v < l}^N \frac{A_v A_l + A_l A_v}{R_v R_l} \sin(k R_v k_v) \cdot \sin(k R_l k_l) \Rightarrow$$

HERMITESch und für ihre Eigenwerte $2\lambda_n(k)$ hat man die Relation (20) zu etablieren.

$$\begin{aligned} &\Rightarrow 2I_m - 4 \sum_{v=1}^N \frac{A_v^2}{R_v} \sin^2(k R_v \frac{\Delta t}{2}) - \left(\sum_{v=1}^N \frac{A_v}{R_v} \sin k R_v k_v \right)^2 \\ &\quad + \sum_{v=1}^N \frac{A_v^2}{R_v} \sin^2 k R_v k_v = \\ &= 2I_m - 4 \sum_{v=1}^N \frac{A_v^2}{R_v} \sin^2(k R_v k_v) - \left(\sum_{v=1}^N \frac{A_v}{R_v} \sin k R_v k_v \right)^2 \Rightarrow \end{aligned}$$

wegen: $4 \min^2 d - \min^2 2d = 4 \min^2 d$ für $d \in \mathbb{R}$

$$\stackrel{\geq}{=} 2I_m - 2F_n(k)$$

Da $F_n(k)$ HERMITESch ist, sind seine Eigenwerte in der Menge

$$L := \left\{ (y, F_n(k)y), y \in \mathbb{R}^m, |y| = 1 \right\} \text{ zu finden.}$$

Nun sind sowohl $\frac{A_v^2}{R_v^2} \sin^2 k R_v \frac{k_v}{2}$ für $v=1, \dots, N$
als auch $\left(\sum_v \frac{A_v}{R_v} \sin k R_v k_v \right)^2$ positiv semidefinit.

\Rightarrow Für $r \in L$ gilt: $r \geq 0$

$$\Rightarrow (|A_n(k)| \leq 1 - \varepsilon \iff |F_n(k)| \leq 2 - \varepsilon) \text{ für } 0 < \varepsilon < 1$$

Es gilt:

$$|F_n(k)| \leq 4 \sum \frac{|A_v|^2}{R_v^2} + \left(\sum \frac{|A_v|}{R_v} \right)^2$$

wählt man $R_v \geq |A_v| \sqrt{\frac{4N + N^2}{2 - \varepsilon}}$ so folgt:

$$|F_n(k)| \leq 4N \frac{2 - \varepsilon}{4N + N^2} + N^2 \frac{2 - \varepsilon}{4N + N^2} = 2 - \varepsilon$$

Man hat also das Ergebnis:

Für HERMITESche Systeme ist (37) stabil, falls

$$R_v \geq |A_v| \cdot \sqrt{\frac{4N + N^2}{2 - \varepsilon}} \text{ gewählt wird.}$$

V. Numerische Resultate

In diesem Abschnitt werden die beiden Schemata (10) und (37), im folgenden ZENT und DZENT genannt, mit dem von 2. Ordnung konsistenten 1-Schritt Schema (3.3) aus [6], im folgenden LAXW, in ihrer Genauigkeit für drei hyperbolische Systeme verglichen. Das LAX-WENDROFF Schema LAXW erhält man aus

$$u(t+h) = u(t) + h u_t(t) + \frac{h^2}{2} u_{tt}(t) + O(h^3)$$

worin alle auftretenden Zeitableitungen mittels (11) durch Raumableitungen ersetzt werden und dann $\frac{\partial}{\partial x_i}$ gegen $\frac{T_i - T_{i-1}}{2R_i h} =: \Delta_i$,

$$\frac{\partial^2}{\partial x_i \partial x_\ell} \quad \text{gegen } \Delta_i \Delta_\ell \quad (\text{für } i \neq \ell) \quad \text{und} \quad \frac{\partial^2}{\partial x_i^2} \quad \text{gegen} \quad \frac{T_i - 2 + T_{i-1}}{R_i^2 h^2}$$

ausgewechselt werden. LAXW benötigt als 1-Schritt Verfahren etwa halb so viel Speicherplatz wie ZENT und DZENT, dagegen etwas mehr Rechenaufwand als DZENT und erheblich mehr als ZENT.

Bei den folgenden Gleichungen, die ersten zwei inhomogene, symmetrische Systeme mit konstanten bzw. mit orts- und zeitabhängigen Koeffizienten, die dritte quasilinear, wurden für jeden Zeitschritt die maximal auftretenden Fehler der Näherungslösungen zur jeweilig gewählten exakten Lösung niedergeschrieben.

Die Anfangswerte für LAXW in der ersten Zeitschicht ($t=0$) und für ZENT und DZENT in den ersten beiden Schichten wurden exakt vorgegeben. Um die Wirkung von Anfangsfehlern zu testen, wurden für ZENT und DZENT in einem weiteren Experiment in der ersten Schicht die exakten und in der zweiten die von LAXW berechneten Werte genommen. Die zugehörigen Fehler sind unter ZENTL bzw. DZENTL aufgeführt. Auf den Rändern des Raumgitters sind in allen Fällen die Werte der exakten Lösung genommen.

Die Raumgitter wurden in den drei Beispielen festgehalten und die Zeitschrittweite h bis an die Stabilitätsgrenze variiert.

1. Beispiel

$$u_t(t, x, y) = \begin{pmatrix} 2 & 1 \\ 1 & -1 \end{pmatrix} u_x(t, x, y) + \begin{pmatrix} 3 & 1 \\ 1 & -2 \end{pmatrix} u_y(t, x, y) + \begin{pmatrix} (2 \cos x + 3 \cos y) \cdot \cos t \\ (\sin x + 2 \sin y) \cdot \sin t \end{pmatrix}$$

Exakte Lösung

$$u(t, x, y) = \begin{pmatrix} (\sin x + \sin y) \cdot \cos t \\ (\cos x + \cos y) \cdot \sin t \end{pmatrix}$$

In diesem wie auch im folgenden Beispiel wurde als Raumgitter die diskrete Menge

$$\left\{ (x, y) \in \mathbb{R}^2; x = m \cdot \frac{2\pi}{100}, y = n \cdot \frac{2\pi}{100}, m, n \in \mathbb{N}_0, m, n \leq 10 \right\}$$

gewählt.

Die Indizes in den Tabellen geben die Komponenten der Vektoren an.

In den Tabellen auf Seite ⁷³ zeigt sich fast kein Unterschied bei den Fehlern von ZENT und ZENTL. Die Genauigkeit des LAX-WENDROFF-Schemas reicht also aus, um die Anfangswerte zu berechnen. Zwischen DZENT und DZENTL gibt es schon größere Differenzen, in der ersten Komponente störend größere. Dieselben Unterschiede traten bei allen Rechnungen auf, ZENTL und DZENTL sind daher in den folgenden Tabellen nicht aufgeführt.

ZENT ist deutlich schlechter als LAXW und LAXW erheblich schlechter als DZENT in diesem Beispiel. Bemerkenswert sind die kleinen Fehler von DZENT für größtmögliche Schrittweiten h . Während für $h = 0.015$ LAXW und ZENT schon instabil sind, liefert DZENT besonders gute Ergebnisse. Für $h = 0.02$ ist auch DZENT instabil. Nahe der Stabili-

tätsgrenze ist das Verfahren also sehr genau. Daß ZENT noch vor DZENT instabil wird, mag an der Inhomogenität des Problems oder am "Festhalteeffekt" der Randwerte liegen.

2. Beispiel

$$u_t = \begin{pmatrix} t & y \\ y & t \end{pmatrix} u_x + \begin{pmatrix} t & x \\ x & t \end{pmatrix} u_y + (y+x) \cdot \begin{pmatrix} \sin(x+y+t^2) \\ -\cos(x+y+t^2) \end{pmatrix}$$

Exakte Lösung

$$u(t, x, y) = \begin{pmatrix} \sin(x+y+t^2) \\ \cos(x+y+t^2) \end{pmatrix}$$

Nach der Tabelle auf Seite 74 fällt DZENT in diesem Beispiel wegen der Größe seiner Fehler aus dem Rahmen. ZENT ist klar besser als LAXW.

Nach [3] und [7] (s.auch [8] S.119) läßt sich die Stabilität von LAXW unter anderem für dieses Beispiel beweisen. Das dabei verwandte KREISSsche Matrix Theorem läßt sich sicher nicht auf ZENT und DZENT anwenden, da die Amplifikationsmatrizen dieser beiden Schemata im Stabilitätsbereich nur Eigenwerte des Absolutbetrages 1. haben, sie beide also niemals dissipativ sind (s.[8] s. 109). Die für Stabilität nicht notwendige Dissipativitätsbedingung ist also auch kein Garant für die Qualität eines Verfahrens.

Für $h_1 = 0.02$ wurde ZENT, für $h_1 = 0.025$ wurden auch LAXW und DZENT instabil.

3. Beispiel

$$u_t = u u_x$$

Exakte Lösung: $u(t, x) = \frac{x}{1-t}$; $t \neq 1$

Ein Vergleich der beiden Tabellen liefert für um einen Faktor 10

veränderte Schrittweite stark veränderte Fehler. DZENT ist für

$h=0.005$ das schlechteste, für $h=0.05$ das beste

Verfahren. Umgekehrt verhält es sich mit ZENT, das im Falle

$h=0.005$ die insgesamt kleinsten Fehler aufweist. Das Ergebnis legt dar, wie sehr im quasilinearen Falle die Qualität von Differenzenverfahren nicht nur vom gerade zu behandelnden System sondern auch von der Wahl der Schrittweite abhängt.

Maximale Fehler in 10^{-3} Einheit

Tabellen zum Beispiel 1

$h = 0,01$	LAXW1	ZENT1	ZENTL1	DZENT1	DZENTL1	LAXW2	ZENT2	ZENTL2	DZENT2	DZENTL2
$t = 0,1$	0,58	0,99	1,06	0,02	0,29	0,15	0,26	0,27	0,00	0,11
0,2	1,07	2,08	2,11	0,03	0,13	0,26	0,52	0,52	0,01	0,04
0,3	1,05	2,19	2,15	0,01	0,36	0,30	0,57	0,58	0,02	0,12
0,4	1,01	2,06	2,02	0,00	0,26	0,30	0,86	0,87	0,05	0,06
0,5	0,96	1,92	1,95	0,01	0,18	0,32	0,95	0,94	0,13	0,17
0,6	0,90	2,19	2,20	0,02	0,29	0,33	1,03	1,02	0,13	0,14
0,9	0,70	3,05	2,99	0,04	0,11	0,46	1,55	1,54	0,05	0,12

$h = 0,015$	LAXW	ZENT	DZENT	LAXW	ZENT	DZENT
$t = 0,15$	0,70	1,56	0,003	0,17	0,40	0,004
0,3	0,78	18,88	0,003	0,23	4,61	0,011
0,45	langsam	45763,05	0,009	langsam	10803,35	0,031
0,6	instabil	schnell	0,017	instabil	schnell	0,043
0,9		instabil	0,006		instabil	0,031
1,2			0,008			0,016
1,5	430,44		0,004	103,73		0,012

Maximale Fehler in 10^{-3} Einheit

Tabelle zum Beispiel 2

h = 0.015	LAXW	ZENT	DZENT	LAX	ZENT	DZENT
t = 0.3	0.87	0.24	10.79	2.21	0.19	38.03
0.6	3.60	0.36	32.60	2.85	0.31	81.21
0.9	12.88	0.34	57.64	4.26	1.31	87.72
1.2	13.70	1.31	39.52	7.60	1.69	23.06
1.5	4.06	2.19	19.98	12.97	1.64	55.90

Tabellen zum Beispiel 3

h = 0.005	LAXW	ZENT	DZENT
t = 0.1	0.02	0.01	0.15
0.2	0.06	0.02	0.65
0.3	0.10	0.03	1.47
0.4	0.16	0.05	2.59
0.5	0.24	0.09	3.96
0.6	0.36	0.30	158.83

h = 0.05	LAXW	ZENT	DZENT
t = 0.1	0.15	0.28	0.02
0.2	0.53	0.64	0.10
0.3	1.14	1.35	0.29
0.4	2.18	2.54	0.68
0.5	4.23	4.87	1.45
0.6	8.37	9.79	3.08
0.7	18.95	22.90	7.59
0.8	56.86	70.38	25.33
0.9	353.62	489.92	229.50

Literaturverzeichnis

- [1] ANSORGE, R. und HASS, R.: Konvergenz von Differenzenverfahren für lineare und nichtlineare Anfangswertaufgaben. Lecture notes in mathematics. Springer 1970, Heidelberg.
- [2] COURANT, R., FRIEDRICHS, K.O., und LEWY, H.: Über die partiellen Differenzgleichungen der mathematischen Physik. Math. Ann., vol.100 (1928), p.32
- [3] KREISS, H.O.: On difference approximations of the dissipative type for hyperbolic differential equations. Comm.Pure Appl.Math., vol.17(1964), p.335
- [4] LAX, P.D.: Differential equations, difference equations and matrix theory. Comm.Pure Appl.Math., vol.11(1958), p.175
- [5] LAX, P.D. and RICHTMYER, R.D.: Survey of the stability of linear finite difference equations. Comm.Pure Appl. Math., vol. 9 (1956), p.267
- [6] LAX, P.D. and WENDROFF, B.: Difference schemes for hyperbolic equations with high order of accuracy. Comm.Pure Appl. Math., vol.17(1964), p.381
- [7] PARLETT, B.N.: Accuracy and dissipation of difference schemes. Comm. Pure Appl. Math., vol. 19 (1966), p.111
- [8] RICHTMYER, R.D. and MORTON, K.W.: Difference methods for initial value problems. Interscience Publishers 1967 2nd ed., New York
- [9] RIORDAN, J.: Combinatorial Identities. WILEY 1968, New York
- [10] GRÖBNER, W.: Matrizenrechnung. B.I.1966, Mannheim

This IPP report is intended for internal use.

IPP reports express the views of the authors at the time of writing and do not necessarily reflect the opinions of the Max-Planck-Institut für Plasmaphysik or the final opinion of the authors on the subject.

Neither the Max-Planck-Institut für Plasmaphysik, nor the Euratom Commission, nor any person acting on behalf of either of these:

1. Gives any guarantee as to the accuracy and completeness of the information contained in this report, or that the use of any information, apparatus, method or process disclosed therein may not constitute an infringement of privately owned rights; or
2. Assumes any liability for damage resulting from the use of any information, apparatus, method or process disclosed in this report.