

Differenzenverfahren für Partielle Differentialgleichungen

Wolf Hofmann

25. August 2005

Inhaltsverzeichnis

I	Parabolische Differentialgleichungen	1
§ 1	Die Wärmeleitungsgleichung	1
§ 2	Diskretisierung (einfachster Fall)	4
§ 3	Hilfsmittel aus der linearen Algebra	10
	Eigenschaften von A_h^0	12
	Eigenwerte von A_h^0	14
	Vergleich der Eigenwerte von kontinuierlicher und diskreter Aufgabe . . .	15
	Eigenwertschranken	16
	Skalarprodukte, Normen und Abschätzungen	17
§ 4	Stabilität (und „bessere“ Verfahren)	20
	Mittelung der Werte auf alter und neuer Zeitschicht	22
§ 5	Approximations- und Verfahrensfehler	31
§ 6	Spezielle Lösungsverfahren für lineare Gleichungssysteme	38
§ 7	Die Gleichung $u_t = \frac{\partial}{\partial x} \left(k(x) \frac{\partial u}{\partial x} \right) + f$	46
§ 8	Die allg. 1-dimensionale Wärmeleitungsgleichung	51
	Erweiterung des Stabilitätssatzes (4.2)	55
	Die Wärmeleitung mit zeitabhängigem Diffusionskoeffizienten	59
	Die Stabilität bzgl. der rechten Seite	61
§ 9	Gleichmäßige Stabilität und Konvergenz	64
	Oszillationsfreiheit	66

§ 10 Die mehrdimensionale Wärmeleitungsgleichung	67
Verfahren der Alternierenden Richtungen	72
Konvergenz der mehrdimensionalen Verfahren	77
II Elliptische Gleichungen	80
§ 11 Die Poissongleichung - Einleitung	80
§ 12 Die erste RWA für die Poissongleichung im Rechteck	82
Das diskrete Maximumprinzip	88
§ 13 Die 3. RWA für die Poissongleichung	93
§ 14 Die 1. RWA der Poissongl. in allgemeineren Gebieten	100
§ 15 Jacobi und Gauss-Seidel	110
Das Gauß-Seidel-Verfahren	115
Das symmetrische Gauß-Seidel-Verfahren	117
III Das Mehrgitterverfahren (MGV)	119
§ 16 Motivation und Grobstruktur	119
Zwei-Gitter-Verfahren	127
Mehrgitter-Verfahren	129
Volles Mehrgitter-Verfahren	133
Geschichtlicher Überblick	134
§ 17 Glättung, Restriktion, Prolongation	135
Glättungsiterationen	136
Parallelisierung	138
Prolongation und Restriktion	141
§ 18 Konvergenz des ZGV	146
Glättungseigenschaften des gedämpften Jacobi-Verfahrens	149
Das symmetrische Gauß-Seidel-Verfahren als Glättungsiteration	153
§ 19 Konvergenz des Mehrgitterverfahrens	157
Konstruktion der Iterationsmatrix M_l des MGV auf Level l	158

	Überlegungen zur Iterationszahl des einfachen MGW	165
§ 20	MGW für nichtlineare Probleme	169
	MGW zur Lösung von Integralgleichungen	172
IV	Hyperbolische Differentialgleichungen	174
§ 21	Die Wellengleichung	174
§ 22	Die Neumann'sche Stabilitätsanalyse	180
	Beispiele zur Stabilitätsanalyse	184
	Diskretisierungsfehler für die 1D-Wellengleichung	188
§ 23	Literatur	190

Vorbemerkung: Dieses Skript führt in die Theorie der Differenzen-Verfahren für Partielle Differentialgleichungen ein. Behandelt werden sollen hierbei die Wärmeleitungsgleichung, die Poissongleichung und die Wellengleichung. Die ersten Überlegungen behandeln den räumlich eindimensionalen Fall, der schon (fast) alle auftretenden Schwierigkeiten enthält. Später folgt eine Ausdehnung auf den räumlich zweidimensionalen Fall. Die Übertragung auf höhere Dimensionen ist dann fast nur Technik.

Es wird sich zeigen, daß zur numerischen Lösung der auftretenden Probleme eine ganze Menge Numerische Lineare Algebra nötig ist. Sie wird, soweit nicht aus der Vorlesung Numerik I+II für Studienanfänger bekannt, ebenfalls behandelt. Dies betrifft u.a. insbesondere eine relativ ausführliche Darstellung des Mehrgitter-Verfahrens, welches *das* Verfahren ist, das der numerischen Lösung der auftretenden linearen Gleichungssysteme angemessen ist.

Das Skript hat seinen Ursprung in einer Vorlesung über Numerische Behandlung von Partiellen Differentialgleichungen, die im Rahmen einer Gastvorlesung von Prof. Dr. G. Stoyan von der Elte-Universität Budapest gehalten worden ist. Aus seiner professionellen Erfahrung stammen auch viele Hinweise auf die Anwendbarkeit oder Nichtanwendbarkeit der verschiedenen Verfahren oder ihrer Varianten. Solche Hinweise sind in Lehrbüchern leider nur sehr sehr selten zu finden.

Dieses Skript enthält a) eine notwendigerweise beschränkte und subjektive Stoffauswahl aus dem Gebiet der Partiellen Differentialgleichungen (PDG) und b) – mit einiger Wahrscheinlichkeit – auch eine Reihe von Fehlern. Aus beiden Gründen ist es ungeeignet, ein Lehrbuch zu ersetzen. Sein Zweck ist es, den Hörer vom Zwang des Mitschreibens zu befreien. Es entbindet ihn nicht von der Notwendigkeit, den Stoff in Lehrbüchern nachzulesen und zu vertiefen und sich mit der notwendigen Referenzliteratur vertraut zu machen, die es ihm gestattet, Stoffgebiete nachzulesen, die nicht in der Vorlesung behandelt wurden.

Kapitel I

Parabolische Differentialgleichungen

§ 1 Die Wärmeleitungsgleichung

Unser Ziel ist die numerische Behandlung der allgemeinen Wärmeleitungsgleichung mittels Differenzenverfahren, die (eindimensional) wie folgt aussieht:

$$c_p \rho \frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left(k \frac{\partial u}{\partial x} \right) - v \frac{\partial u}{\partial x} - qu + f$$

- $u \hat{=}$ Temperatur
- $c_p \hat{=}$ Wärmekapazität
- $\rho \hat{=}$ Dichte
- $k \hat{=}$ Wärmeleitfähigkeit
- $v \hat{=}$ Geschwindigkeit
- $q \hat{=}$ Abbaurate
- $f \hat{=}$ Quellterm

Üblicherweise sind:

- $c_p, \rho > 0$ häufig konstant, möglicherweise abhängig von u
- $v = v(x, t)$ gegeben
- $f = f(x, t)$ oft nichtlinear
- $k = k(x, t)$ oft stückweise konstant (z.B. beim Übergang von einem Medium in ein anderes)

Nutzanwendung: Meteorologie, Luftverschmutzung, Bodenverschmutzung (Grundwasser).

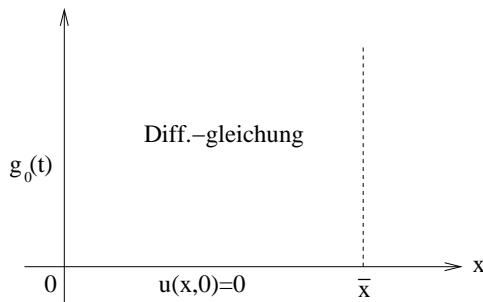
Ein typisches **Beispiel**: Aus einem Container (Tanker) strömt Gas aus.

$$\frac{\partial u}{\partial t} = D \frac{\partial^2 u}{\partial x^2} - v \frac{\partial u}{\partial x} - qu$$

Dabei sind

$u \hat{=}$ Gaskonzentration
 $D \hat{=}$ Diffusionskoeffizient
 $v \hat{=}$ Windgeschwindigkeit
 $q \hat{=}$ Abbauglied

Container im Nullpunkt: Randbedingungen in $x = 0$: $g_0(t) \hat{=}$ ausströmendes Gas
 Null-Anfangsbedingungen



Parameter in der Differentialgleichung:

v , meßbar

D , schwierig, Vorwissen nötig (was ist im Tank drin? Diffusion in der Luft?)

q , abhängig von den chemischen Eigenschaften des Gases

Randbedingungen in \bar{x} ? künstliche?

beliebt Neumann : $\frac{\partial u}{\partial x}(\bar{x}) = 0$
 (primitiv, klappt aber oft).

Problem z.B. Tschernobyl: q war nicht bekannt, man wußte nicht, was drinnen war.

Ein weiteres Anwendungsbeispiel: Das Börsenverhalten von Aktien (Modell von Black-Scholes 1973) zur Unterstützung von Kauf und Verkauf kann durch eine parabolische Differentialgleichung modelliert werden.

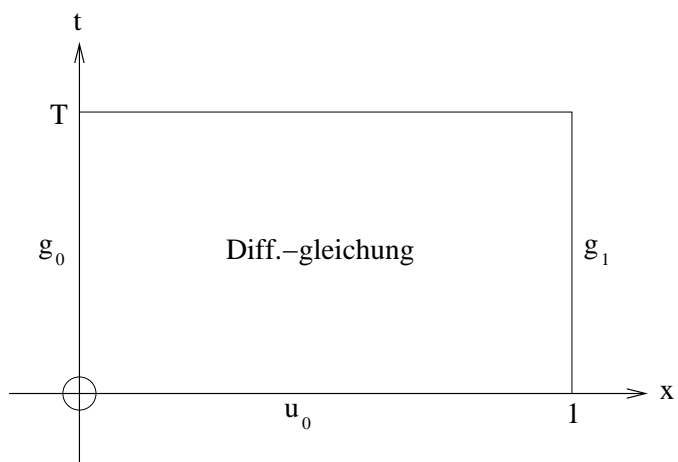
Wir beschränken uns in unseren Untersuchungen zunächst auf das einfachste Beispiel:

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + f, \quad 0 < x < 1, \quad 0 < t < T$$

$$t = 0 : u(x, 0) = u_0(x), \quad 0 \leq x \leq 1$$

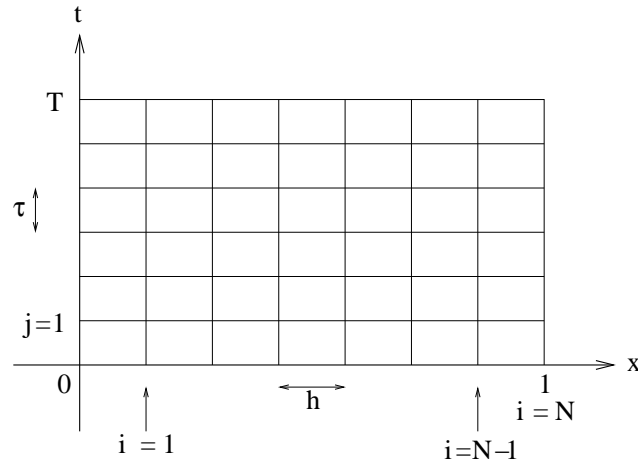
$$t > 0 : u(0, t) = g_0(t), \quad 0 \leq x \leq 1$$

$$u(1, t) = g_1(t).$$



§ 2 Diskretisierung (einfachster Fall)

Wir überziehen das Gebiet, in dem die Lösung berechnet werden soll, mit einem (nicht notwendigerweise quadratischen) Gitter. τ bzw. h sind Zeit- bzw. Ortsschrittweite.



$$\begin{aligned} \tau &= T/m, \quad m \geq 1 \\ h &= 1/N, \quad N \geq 2 \quad (\text{damit mindestens ein innerer Punkt existiert}) \end{aligned}$$

und definieren eine

Gitterfunktion $u_i^j = u(x_i, t_j)$, $x_i = i \cdot h$, $t_j = j \cdot \tau$; $i, j = 0, 1, 2, \dots$

und bezeichnen ihre Approximation mit

$$y_i^j \approx u(x_i, t_j), \quad x_i = i \cdot h, \quad t_j = j \cdot \tau; \quad i, j = 0, 1, 2, \dots$$

Eine *Zeitschicht* umfaßt alle Gitterpunkte für ein festes t (alle auf einer Linie).

Ersetzt man die Zeitableitung durch den

vorwärtsgenommenen Differenzenquotienten

$$(2.1) \quad \frac{\partial u(x_i, t_j)}{\partial t} \approx \frac{y_i^{j+1} - y_i^j}{\tau},$$

so kann man die Güte dieser Approximation abschätzen durch die Taylorentwicklung in Zeitrichtung (wir unterdrücken das x -Argument)

$$u(t_{j+1}) = u(t_j) + \tau \dot{u}(t_j) + \frac{\tau^2}{2} \ddot{u}(t_j) + 0(\tau^3) \quad (\text{Voraussetzung: } u \in C^2)$$

und erhält

$$(2.2) \quad \frac{u(t_{j+1}) - u(t_j)}{\tau} = \dot{u}(t_j) + \frac{\tau}{2} \ddot{u}(t_j) + 0(\tau^2),$$

eine (schlechte) Approximation der 1. Ordnung (lineare Konsistenzordnung).

Hier und im Folgenden bezeichnen wir mit \dot{u} Ableitungen nach der Zeit- und mit u' Ableitungen nach der Ortsvariablen (entsprechend für höhere Ableitungen mit mehr Punkten bzw. Strichen).

Eine bessere Approximationseigenschaft erhält man durch den

Zentrale Differenzenquotienten:

1. Ableitung:

Wir betrachten den Differenzenquotienten (2.1) als eine Approximation für die Ableitung auf der Zeitstufe $t_{j+1/2}$.

Wir betrachten die Taylorentwicklung in den Punkten $t_{j+\frac{1}{2}\pm\frac{1}{2}}$ an der Stelle $t_{j+\frac{1}{2}}$ (und unterdrücken die Ortsabhängigkeit in der Bezeichnung).

$$u(t_{j+\frac{1}{2}\pm\frac{1}{2}}) = u(t_{j+\frac{1}{2}}) \pm \frac{\tau}{2}\dot{u}(t_{j+\frac{1}{2}}) + \frac{1}{2!}\left(\frac{\tau}{2}\right)^2\ddot{u}(t_{j+\frac{1}{2}}) \pm \frac{\tau^3}{48}\ddot{\ddot{u}}_{\pm} \quad (\text{Vor: } u \in C^3).$$

Dabei werden mit $\ddot{\ddot{u}}_{\pm} = \ddot{\ddot{u}}(\xi_{\pm})$ Zwischenstellen bezeichnet.

Subtrahiert man die Darstellungen für t_{j+1} und t_j , so folgt

$$u(t_{j+1}) - u(t_j) = \tau\dot{u}(t_{j+\frac{1}{2}}) + \frac{\tau^3}{24}(\ddot{\ddot{u}}_+ + \ddot{\ddot{u}}_-).$$

$\ddot{\ddot{u}}$ ist stetig und nach dem Zwischenwertsatz existiert dann eine Zwischenstelle z , so daß $2\ddot{\ddot{u}}(z) = \ddot{\ddot{u}}(\xi_+) + \ddot{\ddot{u}}(\xi_-)$. Damit folgt

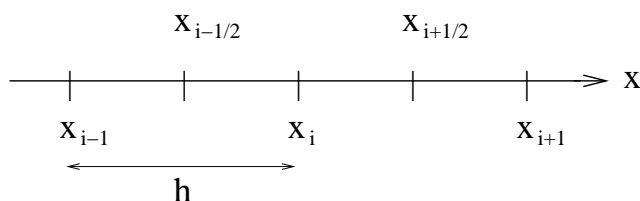
$$(2.3) \quad \frac{u(t_{j+1}) - u(t_j)}{\tau} = \dot{u}(t_{j+\frac{1}{2}}) + \frac{\tau^2}{24}(\ddot{\ddot{u}}(\xi_+) + \ddot{\ddot{u}}(\xi_-)) = \dot{u}(t_{j+\frac{1}{2}}) + \frac{\tau^2}{24}\ddot{\ddot{u}}(z)$$

also eine Approximation 2. Ordnung.

2. Ableitung

Zur Approximation der 2.ten Ableitung in Ortsrichtung verwenden wir den zentralen Differenzenquotienten in x_i , gebildet aus den zentralen Differenzenquotienten für die ersten Ableitungen in den Punkten $x_{i-\frac{1}{2}}$ und $x_{i+\frac{1}{2}}$.

Zentraler Differenzenquotient für u''



$$(2.4) \quad \frac{\partial^2 u(x_i, t_j)}{\partial x^2} \approx \frac{\frac{y_{i+1}^j - y_i^j}{h} - \frac{y_i^j - y_{i-1}^j}{h}}{h} = \frac{y_{i+1}^j - 2y_i^j + y_{i-1}^j}{h^2}$$

Um die Approximationsgüte abzuschätzen benutzen wir wieder die Taylorentwicklung: (die Indizes bezeichnen die x_i -Werte, das t -Argument wird unterdrückt.)

$$\begin{aligned}
 u_{i\pm 1} &= u_i \pm hu'_i + \frac{h^2}{2}u''_i \pm \frac{h^3}{6}u'''_i + \frac{h^4}{24}u^{(4)'}_i \pm \frac{h^5}{5!}u^{(5)'}_i + \frac{h^6}{6!}u^{(6)'}_{\pm} \\
 &\implies \\
 u_{i+1} + u_{i-1} &= 2u_i + h^2u''_i + \frac{h^4}{12}u^{(4)'}_i + 2\frac{h^6}{6!}u^{(6)'}_{\pm} \\
 &\quad \pm \text{ als Indizes bezeichnen Zwischenstellen.}
 \end{aligned}$$

Damit folgt

$$\begin{aligned}
 \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} &= \begin{cases} u''_i + \frac{h^2}{12}u^{(4)'}_i + \frac{h^4}{720}(u^{(6)'}_+ + u^{(6)'}_-) & \text{falls } u \in C^6 \\ u''_i + \frac{h^2}{24}(u^{(4)'}_+ + u^{(4)'}_-) & \text{falls } u \in C^4 \end{cases} \\
 (2.5) \quad &\text{Die Indizes „+“ und „-“ bezeichnen die Werte an} \\
 &\text{Zwischenstellen } x \pm \vartheta_{\pm}h, |\vartheta_{\pm}| \leq 1. \\
 &= \begin{cases} u''_i + \frac{h^2}{12}u^{(4)'}_i + \frac{h^4}{360}u^{(6)'}(z) & \text{falls } u \in C^6 \\ u''_i + \frac{h^2}{12}u^{(4)'}(z) & \text{falls } u \in C^4 \end{cases} \\
 &\text{jeweils an einer Zwischenstelle } z,
 \end{aligned}$$

also Approximationen 2. Ordnung.

Oft werden folgende **Bezeichnungen** benutzt.

$$(2.6) \quad y_i^j \begin{cases} \text{Indizes oben für die Zeitschicht} \\ \text{Indizes unten für die Ortsschicht} \end{cases}$$

oder indexlos

$$(2.7) \quad \begin{aligned}
 \mathbf{y} &\text{ für eine Zeitschicht (z.B. } t_j) \\
 \hat{\mathbf{y}} &\text{ für die folgende Zeitschicht (z.B. } t_{j+1}) \\
 \bar{\mathbf{y}} &\text{ für die vorhergehende Zeitschicht (also } t_{j-1})
 \end{aligned}$$

Die entsprechenden Differenzenquotienten für eine beliebige, aber feste Zeitschicht t_j , werden wie folgt abgekürzt:

vorwärtsgenommener Differenzenquotient:

$$\dot{u}(x_i, t_j) \approx \frac{y_i^{j+1} - y_i^j}{\tau} =: y_{t,i}^j \qquad \frac{\hat{\mathbf{y}} - \mathbf{y}}{\tau} \approx \frac{\partial u}{\partial t}$$

rückwärtsgenommener Differenzenquotient:

$$\dot{u}(x_i, t_j) \approx \frac{y_i^j - y_i^{j-1}}{\tau} =: y_{t,i}^j \qquad \frac{\mathbf{y} - \bar{\mathbf{y}}}{\tau} \approx \frac{\partial u}{\partial t}$$

Entsprechend wird bezeichnet auf der Zeitschicht t_j mit

$$\begin{aligned} y_{x,i}^j &\hat{=} \text{ vorwärts genommener Differenzenquotient in } x\text{-Richtung im Punkt } x_i \\ y_{\bar{x},i}^j &\hat{=} \text{ rückwärts genommener Differenzenquotient in } x\text{-Richtung im Punkt } x_i \end{aligned}$$

und entsprechend der zentrale Differenzenquotient 2. Ordnung

$$\frac{\partial^2 u(x_i, t_j)}{\partial x^2} \approx \frac{y_{i+1}^j - 2y_i^j + y_{i-1}^j}{h^2} =: y_{\bar{x}x,i}^j$$

bzw. ohne Auszeichnung der speziellen Schichten

$$\frac{\partial^2 u}{\partial x^2} \approx \mathbf{y}_{\bar{x}x}$$

Man beachte, daß $\mathbf{y}_{\bar{x}x}$ in \bar{x} und x symmetrisch ist: $\mathbf{y}_{\bar{x}x} = \mathbf{y}_{x\bar{x}}$. Im Sinne der Hintereinanderausführung gilt

$$\begin{aligned} (y_{x,i})_{\bar{x},i} &= (y_{\bar{x},i})_{x,i}, \quad \text{denn} \\ (y_{x,i})_{\bar{x},i} &= \left(\frac{y_{i+1} - y_i}{h} \right)_{\bar{x},i} = \left(\frac{\frac{y_{i+1} - y_i}{h} - \frac{y_i - y_{i-1}}{h}}{h} \right) = \left(\frac{y_i - y_{i-1}}{h} \right)_{x,i} = (y_{\bar{x},i})_{x,i}. \end{aligned}$$

Mit diesen Bezeichnungen erhalten wir für die näherungsweise Berechnung der Lösung von

$$\begin{aligned} \dot{u} &= u'' + f, \quad 0 < x < 1, \quad 0 < t < T \\ u(x, 0) &= u_0(x) \\ u(0, t) &= g_0(t), \\ u(1, t) &= g_1(t) \end{aligned}$$

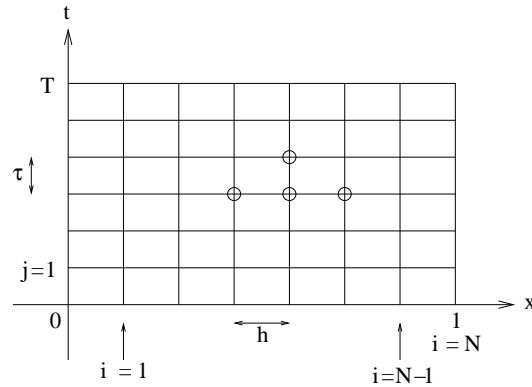
das **explizite Differenzschema**

$$(2.8) \quad \begin{aligned} y_{t,i}^j &= y_{\bar{x}x,i}^j + f_i^j, \quad i = 1, \dots, N-1, \quad j = 0, 1, \dots, \\ y_i^0 &= u_0(x_i), \quad i = 0, 1, \dots, N \\ \left. \begin{aligned} y_0^{j+1} &= g_0(t_{j+1}) \\ y_N^{j+1} &= g_1(t_{j+1}) \end{aligned} \right\} j \geq 0 \end{aligned}$$

unter den Minimalvoraussetzungen

$$\begin{aligned} u &\in C^4 \quad \text{bzgl. Ort} \quad \text{vgl. (2.5)} \\ &\in C^2 \quad \text{bzgl. Zeit} \quad \text{vgl. (2.2)}. \end{aligned}$$

Aus jeweils 3 Punkten einer Zeitschicht wird ein Punkt der nächst höheren Zeitschicht berechnet (vgl. Abb.).



Dies ist ein Verfahren, das im allgemeinen nicht viel taugt. Daß das so ist, wird durch Übungen belegt. Warum das so ist und wie man das verbessern kann, soll im folgenden untersucht werden.

Schreibt man die Differenzgleichung in Matrixgestalt, so erhält man für die Zeitschichten $j \geq 0$ (beachte: dies sind $N - 1$ Gleichungen für jede Zeitschicht)

$$(2.9) \quad \begin{pmatrix} \frac{y_1^{j+1} - y_1^j}{\tau} \\ \vdots \\ \frac{y_{N-1}^{j+1} - y_{N-1}^j}{\tau} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{1}{h^2} & -\frac{2}{h^2} & \frac{1}{h^2} & & & \\ & \ddots & \ddots & \ddots & & \\ & & \ddots & \ddots & \ddots & \\ & & & \ddots & \ddots & \\ & & & & \frac{1}{h^2} & -\frac{2}{h^2} & \frac{1}{h^2} \end{pmatrix}}_{N+1} \begin{pmatrix} y_0^j \\ \vdots \\ y_N^j \end{pmatrix} + \begin{pmatrix} f_1^j \\ \vdots \\ f_{N-1}^j \end{pmatrix}$$

In diesen Gleichungen werden die zu den Randwerten gehörigen Terme (die mit y_0^j, y_N^j) mit den Werten f_i^j zu einem Vektor φ^j zusammengefaßt.

Mit den Bezeichnungen

$$(2.10) \quad \mathbf{y}^j = (y_1^j, \dots, y_{N-1}^j)^T \quad (\text{nur die unbekanntenen Werte})$$

und der (symmetrischen) $(N - 1) \times (N - 1)$ Matrix (man beachte die Vorzeichen)

$$(2.11) \quad \mathbf{A}_h^0 := \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & \ddots & \ddots & \ddots & & \\ & & \ddots & \ddots & \ddots & \\ & & & \ddots & \ddots & -1 \\ & & & & -1 & 2 \end{pmatrix} =: \frac{1}{h^2} \text{tridiag}(-1, 2, -1)$$

erhält man für (2.9) die Darstellung

$$(2.12) \quad \left\{ \begin{array}{l} \mathbf{y}_t^j := \frac{\mathbf{y}^{j+1} - \mathbf{y}^j}{\tau} = -\mathbf{A}_h^0 \mathbf{y}^j + \boldsymbol{\varphi}^j, \quad \boldsymbol{\varphi}^j = \begin{pmatrix} \frac{1}{h^2} y_0^j + f_1^j \\ f_2^j \\ \vdots \\ f_{N-2}^j \\ \frac{1}{h^2} y_N^j + f_{N-1}^j \end{pmatrix}, \quad j \geq 0 \\ \text{oder (indexlos auch bzgl. der Zeitschicht)} \\ \mathbf{y}_t := \frac{1}{\tau}(\hat{\mathbf{y}} - \mathbf{y}) = -\mathbf{A}_h^0 \mathbf{y} + \boldsymbol{\varphi} \end{array} \right.$$

Es wird sich zeigen, daß \mathbf{A}_h^0 von ausschlaggebender Bedeutung für dieses, und später zu betrachtende Verfahren ist, weshalb wie zunächst Eigenschaften dieser Matrix studieren. Dies bedingt auch einen Exkurs in die Numerische Lineare Algebra.

§ 3 Hilfsmittel aus der linearen Algebra

Wir beginnen mit letzterem.

Satz 3.1 Neumann'sche Reihe

Sei $\mathbf{C} \in \mathbb{C}^{n \times n}$, $\|\mathbf{C}\| < 1$ (bezüglich einer Matrixnorm)

\Rightarrow

$$\exists (\mathbf{I} - \mathbf{C})^{-1} = \sum_{\nu=0}^{\infty} \mathbf{C}^{\nu}$$

Beweis:

1) Die Reihe konvergiert (absolut): $\|\sum \mathbf{C}^{\nu}\| \leq \sum \|\mathbf{C}^{\nu}\| \leq \sum \|\mathbf{C}\|^{\nu}$, $\|\mathbf{C}\| < 1$.

2) $(\mathbf{I} - \mathbf{C}) \sum_{\nu=0}^{\infty} \mathbf{C}^{\nu} = \sum_{\nu=0}^{\infty} (\mathbf{I} - \mathbf{C}) \mathbf{C}^{\nu} \stackrel{\text{abs. Kvg.}}{=} \sum_{\nu=0}^{\infty} \mathbf{C}^{\nu} - \sum_{\nu=1}^{\infty} \mathbf{C}^{\nu} = \mathbf{C}^0 = \mathbf{I}$.

Entsprechend für die rechtsseitige Inverse.

Satz 3.2

Sei $\mathbf{A} \in \mathbb{C}^{n \times n}$ hermite'sch (d.h. $\mathbf{A} = \bar{\mathbf{A}}^T =: \mathbf{A}^*$)

\Rightarrow

(3.1) Alle Eigenwerte $\lambda_i(\mathbf{A})$ von \mathbf{A} sind reell.

(3.2) \mathbf{A} besitzt ein orthonormales System von n Eigenvektoren (Basis).

(3.3) Extremaleigenschaft des *Rayleigh-Quotienten*.

Für alle $\mathbf{x} \in \mathbb{C}^n$, $\mathbf{x} \neq 0$ bzw. $\mathbf{y} \in \mathbb{C}^n$, $\|\mathbf{y}\|_2 = 1$ gilt

$$\lambda_{\min}(\mathbf{A}) \leq \frac{\bar{\mathbf{x}}^T \mathbf{A} \mathbf{x}}{\bar{\mathbf{x}}^T \mathbf{x}} = \bar{\mathbf{y}}^T \mathbf{A} \mathbf{y} \leq \lambda_{\max}(\mathbf{A}).$$

Die Grenzen werden für die Eigenvektoren zu λ_{\min} bzw. λ_{\max} angenommen.

Beweis (3.1): Die quadratische Form $\bar{\mathbf{x}}^T \mathbf{A} \mathbf{x}$ ist reell $\forall \mathbf{x} \in \mathbb{C}^n$:

$$\bar{\mathbf{x}}^T \mathbf{A} \mathbf{x} = \bar{\mathbf{x}}^T \bar{\mathbf{A}}^T \mathbf{x} = (\bar{\mathbf{A}} \bar{\mathbf{x}})^T \mathbf{x} = \mathbf{x}^T \bar{\mathbf{A}} \bar{\mathbf{x}} = \overline{\bar{\mathbf{x}}^T \mathbf{A} \mathbf{x}}, \quad \Rightarrow$$

$$(3.4) \quad \mathbf{A} \mathbf{x} = \lambda \mathbf{x} \quad \Rightarrow \quad \underbrace{\bar{\mathbf{x}}^T \mathbf{A} \mathbf{x}}_{\in \mathbb{R}} = \lambda \underbrace{\bar{\mathbf{x}}^T \mathbf{x}}_{> 0 \text{ für } \mathbf{x} \neq 0} \quad \Rightarrow \quad \lambda \in \mathbb{R}.$$

Beweis (3.2): Eine hermite'sche Matrix ist diagonalisierbar, besitzt also n linear unabhängige Eigenvektoren (Fischer: Lineare Algebra).

Wir zeigen zunächst: Eigenvektoren zu verschiedenen Eigenwerten sind orthogonal bezüglich des inneren Produkts $(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n x_i \bar{y}_i = \mathbf{y}^* \mathbf{x}$, $\mathbf{y}^* = \bar{\mathbf{y}}^T$.

$$\mathbf{A} \mathbf{x}^1 = \lambda_1 \mathbf{x}^1 \quad \Rightarrow \quad (\mathbf{A} \mathbf{x}^1, \mathbf{x}^2) = \lambda_1 (\mathbf{x}^1, \mathbf{x}^2)$$

$$\mathbf{A} \mathbf{x}^2 = \lambda_2 \mathbf{x}^2 \quad \stackrel{(3.1)}{\Rightarrow} \quad (\mathbf{x}^1, \mathbf{A} \mathbf{x}^2) = \lambda_2 (\mathbf{x}^1, \mathbf{x}^2).$$

Durch Subtraktion

$$(3.5) \quad (\mathbf{A}\mathbf{x}^1, \mathbf{x}^2) - (\mathbf{x}^1, \mathbf{A}\mathbf{x}^2) = (\lambda_1 - \lambda_2)(\mathbf{x}^1, \mathbf{x}^2)$$

Für jede quadratische Matrix gilt (vgl. vorseitige Definition des inneren Produkts)

$$(3.6) \quad (\mathbf{A}\mathbf{x}^1, \mathbf{x}^2) = \bar{\mathbf{x}}^{2T} \mathbf{A}\mathbf{x}^1 = (\mathbf{A}^T \bar{\mathbf{x}}^2)^T \mathbf{x}^1 = \overline{(\bar{\mathbf{A}}^T \mathbf{x}^2)^T} \mathbf{x}^1 = (\mathbf{x}^1, \bar{\mathbf{A}}^T \mathbf{x}^2).$$

Falls $\bar{\mathbf{A}}^T = \mathbf{A}$, folgt damit aus (3.5): $(\mathbf{x}^1, \mathbf{x}^2) = 0$ wegen $\lambda_1 \neq \lambda_2$.

Da Eigenvektoren zum gleichen Eigenwert orthogonalisiert werden können (Verfahren von Ehrhardt–Schmidt), existiert ein orthogonales System von n Eigenvektoren \mathbf{x}^μ zu den Eigenwerten λ_μ , das normiert werden kann: $((\mathbf{x}^\nu, \mathbf{x}^\mu) = \delta_{\nu\mu})$.

Beweis (3.3): $\forall \mathbf{x} \in \mathbb{C}^n, \|\mathbf{x}\|_2 = 1 \quad \exists!$ Darstellung $\mathbf{x} = \sum_{\nu=1}^n \alpha_\nu \mathbf{x}^\nu, \quad \sum_{\nu=1}^n |\alpha_\nu|^2 = 1$

Damit folgt:

$$(\mathbf{A}\mathbf{x}, \mathbf{x}) = \left(\sum \alpha_\nu \lambda_\nu \mathbf{x}^\nu, \sum \alpha_\nu \mathbf{x}^\nu \right) = \sum |\alpha_\nu|^2 \lambda_\nu \leq \lambda_{\max}(\mathbf{A}) \\ \geq \lambda_{\min}(\mathbf{A}).$$

■

Aus diesem Satz erhalten wir

Folgerung 3.3

(3.7) $\mathbf{A} \in \mathbb{C}^{n \times n}, \mathbf{A} \geq 0$ (positiv semidefinit: $\bar{\mathbf{x}}^T \mathbf{A}\mathbf{x} = (\mathbf{A}\mathbf{x}, \mathbf{x}) \geq 0 \quad \forall \mathbf{x} \in \mathbb{C}^n$)
 \implies alle Eigenwerte sind ≥ 0 .

(3.8) Ist \mathbf{A} hermite'sch, so gilt sogar

$$\mathbf{A} \text{ positiv} \quad \begin{cases} \text{definit} & (\mathbf{A} > 0) \\ \text{semidefinit} & (\mathbf{A} \geq 0) \end{cases}$$

$$\iff \text{Alle Eigenwerte von } \mathbf{A} \text{ sind} \quad \begin{cases} > \\ \geq \end{cases} 0.$$

(3.9) $\mathbf{A} \in \mathbb{C}^{n \times n}$: Die Spektralnorm $\|\mathbf{A}\|_2 = \|\mathbf{A}\|_S = \sqrt{\lambda_{\max}(\mathbf{A}^* \mathbf{A})}$ ist der euklidischen Vektornorm $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^* \mathbf{x}}$ zugeordnet.

(3.10) Ist $\mathbf{A} = \bar{\mathbf{A}}^T$, (hermitesch), so gilt

$$\|\mathbf{A}\|_S = \max_i |\lambda_i(\mathbf{A})|$$

Beweis (3.7) folgt aus (3.4).

Beweis (3.8) folgt aus (3.3). (Extremaleigenschaft des Rayleigh–Quotienten)

Beweis (3.9)

$$\|\mathbf{A}\|_2 = \sup_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} = \sup_{\mathbf{x} \neq 0} \sqrt{\frac{(\mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x})}{(\mathbf{x}, \mathbf{x})}} \stackrel{(3.6)}{=} \sup_{\mathbf{x} \neq 0} \sqrt{\frac{(\mathbf{x}, \bar{\mathbf{A}}^T \mathbf{A}\mathbf{x})}{(\mathbf{x}, \mathbf{x})}} \stackrel{(3.3)}{=} \sqrt{\lambda_{\max}(\bar{\mathbf{A}}^T \mathbf{A})},$$

denn $\bar{\mathbf{A}}^T \mathbf{A}$ ist hermite'sch.

Beweis (3.10)

Aus (3.9) folgt für $\mathbf{A} = \mathbf{A}^T$: $\mathbf{A}^T \mathbf{A} = \mathbf{A}^2$, also $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^2)}$ und $\mathbf{A}^2 \geq 0$, denn

$$\mathbf{A}^2 \mathbf{y} = \mu \mathbf{y} \implies \mu(\mathbf{y}, \mathbf{y}) = (\mathbf{A}^2 \mathbf{y}, \mathbf{y}) = (\mathbf{A} \mathbf{y}, \mathbf{A} \mathbf{y}) \geq 0 \implies \mu \geq 0 \implies \mathbf{A}^2 \geq 0.$$

$$\text{Weiterhin } \mathbf{A} \mathbf{x} = \lambda \mathbf{x} \implies \mathbf{A}^2 \mathbf{x} = \mathbf{A}(\mathbf{A} \mathbf{x}) = \mathbf{A}(\lambda \mathbf{x}) = \lambda^2 \mathbf{x},$$

also

(3.11) Ist λ Eigenwert von \mathbf{A} , so ist λ^2 Eigenwert von \mathbf{A}^2 zum gleichen EV.

Wir zeigen weiter

(3.12) Hat \mathbf{A}^2 nur nichtnegative Eigenwerte μ_i ,
so ist $\sqrt{\mu_i}$ oder $-\sqrt{\mu_i}$ Eigenwert von \mathbf{A} ,

denn (charakt. Polynom)

$$\begin{aligned} 0 = \det(\mathbf{A}^2 - \mu \mathbf{I}) &= \det\{(\mathbf{A} + \sqrt{\mu} \mathbf{I})(\mathbf{A} - \sqrt{\mu} \mathbf{I})\} \\ &= \det(\mathbf{A} + \sqrt{\mu} \mathbf{I}) \det(\mathbf{A} - \sqrt{\mu} \mathbf{I}), \end{aligned}$$

und einer der Faktoren muß verschwinden.

Damit gilt $\lambda_{\max}(\mathbf{A}^2) = \left(\max_i |\lambda_i(\mathbf{A})|\right)^2 \stackrel{(3.9)}{\implies} (3.10)$. ■

Eigenschaften von \mathbf{A}_h^0

Satz 3.4

$\mathbf{A}_h^0 > 0$ (d.h. *positiv definit*: $(\mathbf{A}_h^0 \mathbf{x}, \mathbf{x}) > 0 \quad \forall \mathbf{x} \neq 0$). (vgl. (2.11))

Zum Beweis benutzen wir ein gewichtetes Skalarprodukt

Definition 3.5 gewichtetes Skalarprodukt $(\cdot, \cdot)_{(0,h)}$

$$(\mathbf{y}, \mathbf{v})_{(0,h)} := \sum_{i=1}^{N-1} y_i \bar{v}_i h, \quad h > 0.$$

Bedeutung der Indizierung:

$0 \hat{=}$ es werden keine Ableitungen benutzt,
 $h \hat{=}$ Ortsschrittweite.

Beweis Satz 3.4

Wir zeigen $(\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)} > 0 \quad \forall \mathbf{y} \neq 0$. (Beachte dazu: $(\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)} = (\mathbf{A}_h^0 \mathbf{y}, \mathbf{y}) h$.)

Mit $\mathbf{y} = (y_1, \dots, y_{N-1})^T$, entsprechend den Dimensionen von \mathbf{A}_h^0 , ist ("auf jeder Zeitschicht")

$$y_{\bar{x},i} = \frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} = (y_{x,i} - y_{\bar{x},i}) \frac{1}{h}.$$

Wir setzen ergänzend fest: $y_0 = y_N = 0$, (entspricht Nullrandbedingung)

Dann gilt $(\mathbf{A}_h^0 \mathbf{y})_i = -\mathbf{y}_{\bar{x},i}$, $i = 1, \dots, N-1$ (vgl. (2.9)-(2.12)) und somit

$$\begin{aligned} (\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)} &= - \sum_{i=1}^{N-1} y_{\bar{x},i} \bar{y}_i h, \quad \text{und mit } y_{\bar{x},i} = (y_{x,i} - y_{\bar{x},i}) \frac{1}{h} \\ &= - \sum_{i=1}^{N-1} y_{x,i} \bar{y}_i + \sum_{i=1}^{N-1} y_{\bar{x},i} \bar{y}_i \\ &\quad \text{(das ist eine Art diskreter partieller Integration)} \\ &= - \sum_{i=1}^{N-1} y_{x,i} \bar{y}_i + \sum_{i=0}^{N-2} y_{x,i} \bar{y}_{i+1}, \quad \text{da } y_{\bar{x},i} = \frac{y_i - y_{i-1}}{h} = y_{x,i-1} \\ &= - \sum_{i=0}^{N-1} y_{x,i} \bar{y}_i + \sum_{i=0}^{N-1} y_{x,i} \bar{y}_{i+1}, \quad \text{da } y_0 = y_N = 0 \\ &= \sum_{i=0}^{N-1} (\bar{y}_{i+1} - \bar{y}_i) y_{x,i} = \sum_{i=0}^{N-1} \frac{(\bar{y}_{i+1} - \bar{y}_i)}{h} h y_{x,i} \\ (3.13) \quad (\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)} &= \sum_{i=0}^{N-1} (y_{x,i})^2 h \implies \mathbf{A}_h^0 \geq 0 \quad \text{(positiv semidefinit)} \end{aligned}$$

Wegen

$$\begin{aligned} (\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)} = 0 &\iff y_{x,i} = 0, \quad i = 0, \dots, N-1 \\ &\iff y_{i+1} - y_i = 0, \quad i = 0, \dots, N-1 \iff y_i = 0 \quad \forall i \quad \text{wegen } y_0 = y_N = 0 \\ &\iff \mathbf{y} = 0 \\ &\implies \mathbf{A}_h^0 > 0. \blacksquare \end{aligned}$$

Folgerung 3.6

$$\mathbf{A}_h^0 > 0 \implies \mathbf{A}_h^0 \text{ ist invertierbar}$$

denn alle Eigenwerte von \mathbf{A}_h^0 sind $> 0 \implies$ die Diagonale der Jordannormalform $J(\mathbf{A})$ enthält nur positive Elemente $\implies \exists J(\mathbf{A})^{-1} \implies \exists \mathbf{A}^{-1}$.

Beweis:

Mit $x_0 = 0 \cdot h$, $x_N = Nh = 1$ ist $y_0^k := \sin(k\pi x_0) = 0$, $y_N^k := \sin(k\pi x_N) = 0$.

Wir rechnen (komponentenweise) nach, daß die Vektoren

$$\tilde{\mathbf{y}}^k := (\tilde{y}_0^k, \dots, \tilde{y}_N^k)^T = (0, \mathbf{y}^k, 0)^T, k = 1, \dots, N-1$$

Eigenvektoren von $\tilde{\mathbf{A}}_h^0$ sind zu den Eigenwerten λ_k^h , $k = 1, 2, \dots, N-1$.

Die erste und letzte Zeile von (3.15) lauten $0 = \lambda \cdot 0$.

Für $i = 1, \dots, N-1$ gilt (wir unterdrücken die oberen Indizes $k = 1, \dots, N-1$)

$$\begin{aligned} y_{i+1} + y_{i-1} &= \sin(k\pi x_{i+1}) + \sin(k\pi x_{i-1}), \quad x_{i\pm 1} = x_i \pm h, \quad h = \frac{1}{N} \\ &= \sin(k\pi x_i + k\pi h) + \sin(k\pi x_i - k\pi h) \quad (\text{Additionstheoreme}) \\ &= 2 \underbrace{\sin(k\pi x_i)}_{y_i} \cos(k\pi h) = 2y_i \cos(k\pi h). \end{aligned}$$

\implies

$$\begin{aligned} y_{i+1} - 2y_i + y_{i-1} &= 2(\cos(k\pi h) - 1)y_i, \quad \text{mit} \quad \cos 2\varphi = 1 - 2\sin^2 \varphi \\ &= -4 \left(\sin^2 \frac{k\pi h}{2} \right) y_i \quad \text{Division durch } -h^2 \text{ liefert} \\ \frac{-y_{i+1} + 2y_i - y_{i-1}}{h^2} &= -y_{\bar{x},i} = \frac{4}{h^2} \left(\sin^2 \frac{k\pi h}{2} \right) y_i, \quad i = 1, \dots, N-1. \end{aligned}$$

Wegen $y_0 = y_N = 0$ ist dieses Gleichungssystem nichts anderes als

$$\mathbf{A}_h^0 \mathbf{y}^k = \lambda_k^h \mathbf{y}^k.$$

Wir haben also Eigenvektoren und Eigenwerte von \mathbf{A}_h^0 gefunden.

Wegen

$$k = 1, \dots, N-1, \quad h = \frac{1}{N} \quad \implies \quad kh \frac{\pi}{2} \leq (N-1) \frac{1}{N} \cdot \frac{\pi}{2} < \frac{\pi}{2}$$

sind die Eigenwerte mit k monoton wachsend (\sin wächst monoton in $[0, \frac{\pi}{2}]$). ■

Wegen $\lambda_{\min} = \frac{4}{h^2} \sin^2(\frac{\pi}{2N}) > 0$ folgt aus (3.8) nochmals die Invertierbarkeit von \mathbf{A}_h^0 .

Vergleich der Eigenwerte von kontinuierlicher und diskreter Aufgabe

Für kleine x ist $\sin x \approx x$ eine gute Approximation. Deshalb gilt für kleine Werte von $kh = k \cdot \frac{1}{N}$

$$\lambda_k^h = \frac{4}{h^2} \sin^2\left(\frac{k\pi h}{2}\right) \approx \frac{4}{h^2} \left(\frac{k\pi h}{2}\right)^2 = k^2 \pi^2 \quad (\text{gute Näherung})$$

Für große k , z.B. $k = N-1$, kann von Approximation, selbst für sehr kleines $h = 1/N$, keine Rede sein, wie man unmittelbar ersieht:

$$\lambda_{N-1}^h = \frac{4}{h^2} \left(\sin \frac{(N-1)\pi h}{2} \right)^2 \stackrel{h=\frac{1}{N}}{=} 4N^2 \left(\sin \left(\frac{(N-1)}{N} \cdot \frac{\pi}{2} \right) \right)^2 \stackrel{N \text{ groß}}{\approx} 4N^2,$$

aber

$$\lambda_{N-1} = (N-1)^2 \pi^2.$$

Eigenwertschranken

Unmittelbar folgt aus Lemma 3.7 mit $h = \frac{1}{N}$

$$(3.16) \quad \lambda_{\max}(\mathbf{A}_h^0) = \lambda_{N-1}^h = \frac{4}{h^2} \sin^2 \left(\frac{N-1}{N} \frac{\pi}{2} \right) \leq \frac{4}{h^2}.$$

Zur Abschätzung der Eigenwerte nach unten benutzen wir $N \geq 2$, also $h \leq \frac{1}{2}$ (es gibt mindestens einen inneren Punkt in $[0, 1]$).

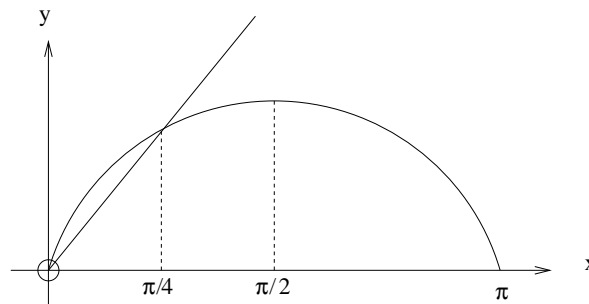
Nun gilt für $0 \leq x \leq \frac{\pi}{4}$

$$y = \frac{\sin \pi/4}{\pi/4} \cdot x \leq \sin x.$$

Wegen $\sin \frac{\pi}{4} = \sqrt{2}/2$ folgt

$$\frac{\sqrt{2}/2}{\pi/4} x \leq \sin x, \quad \text{bzw.}$$

$$\frac{2\sqrt{2}}{\pi} x \leq \sin x$$



Damit folgt für $x = \frac{\pi h}{2}$ (wegen $h \leq \frac{1}{2}$)

$$(3.17) \quad \lambda_{\min}(\mathbf{A}_h^0) = \lambda_1^h = \frac{4}{h^2} \sin^2 \frac{\pi h}{2} \geq \frac{4}{h^2} \frac{8}{\pi^2} \underbrace{\frac{\pi^2 h^2}{4}}_{x^2} = 8,$$

eine von h unabhängige untere Schranke.

Skalarprodukte, Normen und Abschätzungen

Definition 3.8 Skalarprodukte und Normen

$$(3.18) \quad (\mathbf{y}, \mathbf{v})_{(0,h)} := \sum_{i=1}^{N-1} y_i \bar{v}_i h \text{ ist ein Skalarprodukt (vgl. 3.5), } h = \frac{1}{N}.$$

Die zugehörige Vektornorm ergibt sich aus

$$(3.19) \quad \|\mathbf{y}\|_{(0,h)}^2 = (\mathbf{y}, \mathbf{y})_{(0,h)}.$$

Wegen $\mathbf{A}_h^0 > 0$ (pos. definit, (Satz 3.4)), ist $(\mathbf{A}_h^0 \mathbf{y}, \mathbf{x})$ ein Skalarprodukt mit

$$(3.20) \quad (\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)} \stackrel{(3.13)}{=} \sum_{i=0}^{N-1} (\mathbf{y}_{x,i})^2 h, \quad (y_0^j = y_N^j = 0 \ \forall j).$$

Es erzeugt die Norm

$$(3.21) \quad \|\mathbf{y}\|_{(1,h)}^2 = \sum_{i=0}^{N-1} (\mathbf{y}_{x,i})^2 h = (\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)} =: \|\mathbf{y}\|_{\mathbf{A}_h^0}^2, \quad (y_0^j = y_N^j = 0 \ \forall j)$$

liefert eine Vektornorm (energetische Norm).

Der Index $(1, h)$ in (3.21) weist auf eine Ableitung in der Definition hin.

Beachte: Die Darstellung (3.20) (bzw. (3.13)) gilt nur, wenn man $y_0 = y_N = 0$ setzt.

Mit (3.18), (3.20) erhält man

$$\frac{(\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)}}{(\mathbf{y}, \mathbf{y})_{(0,h)}} = \frac{(\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})}{(\mathbf{y}, \mathbf{y})}$$

und mit Satz 3.2,(3.3) (Eigenschaft des Rayleigh-Quotienten) und (3.16), (3.17) die Abschätzungen

$$(3.22) \quad 8 \leq \lambda_{\min}(\mathbf{A}_h^0) \leq \frac{(\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)}}{(\mathbf{y}, \mathbf{y})_{(0,h)}} \leq \lambda_{\max}(\mathbf{A}_h^0) \leq \frac{4}{h^2}.$$

Hieraus folgen wegen

$$\|\mathbf{y}\|_{(1,h)}^2 = \frac{(\mathbf{A}_h^0 \mathbf{y}, \mathbf{y})_{(0,h)}}{\|\mathbf{y}\|_{(0,h)}^2} \cdot \|\mathbf{y}\|_{(0,h)}^2$$

die Normvergleiche

$$(3.23) \quad 8 \|\mathbf{y}\|_{(0,h)}^2 \leq \|\mathbf{y}\|_{(1,h)}^2 \leq \frac{4}{h^2} \|\mathbf{y}\|_{(0,h)}^2$$

und mit $\|\mathbf{A}_h^0\|_S = \lambda_{\max}(\mathbf{A}_h^0)$ weiterhin

$$8 \leq \|\mathbf{A}_h^0\|_S \leq \frac{4}{h^2}.$$

Bemerkung: Solche Normabschätzungen kann man auch ohne Rückgriff auf die Eigenwertaussagen erhalten (vgl. (3.28)). Sie fallen dann etwas schlechter aus, sind aber auch bei Matrizen anwendbar, deren Eigenwerte man nicht kennt. (vgl. den § über die Behandlung der Differentialgleichung $u_\tau = \frac{\partial}{\partial t} (k(x) \frac{\partial u}{\partial t})$). Wir werden (später) Fehlerabschätzungen bzgl. dieser Skalarproduktnormen erhalten.

In praktischen Anwendungen möchte man gerne Abschätzungen bzgl. der Maximumnorm. Solche Abschätzungen erhält man oft nur mit Hilfe der Abschätzungen über Skalarprodukte. Wir zeigen, ohne Benutzung von Eigenwerten, die für die Anwendung wichtige Abschätzung

$$(3.24) \quad \|\mathbf{y}\|_\infty \leq \frac{1}{2} \|\mathbf{y}\|_{(1,h)} \quad (\mathbf{y} \in \mathbb{R}^{N-1})$$

Beweis: Unter Benutzung von $y_0 = y_N = 0$ gilt

$$\begin{aligned} y_i &= h \sum_{j=0}^{i-1} \frac{(y_{j+1} - y_j)}{h} = \sum_{j=0}^{i-1} y_{x,j} \cdot h \quad (\text{diskretes Analogon zu } y(x) = y(0) + \int_0^x y(\xi) d\xi) \\ &= \sum_{j=0}^{i-1} y_{x,j} \sqrt{h} \sqrt{h} \stackrel{\text{CSU}}{\leq} \sqrt{\sum_{j=0}^{i-1} (y_{x,j})^2 h \sum_{j=0}^{i-1} h} \quad \text{und mit } \sum_{j=0}^{i-1} h = x_i \end{aligned}$$

$$(3.25) \quad y_i^2 \leq x_i \sum_{j=0}^{i-1} (y_{x,j})^2 h.$$

Analog erhält man durch „rückwärtsintegrieren“: $y_i = \sum_{j=i}^{N-1} \frac{(y_{j+1} - y_j)}{h}$, die Identität

$$y_i = - \sum_{j=i}^{N-1} y_{x,j} h = \sum_{j=i}^{N-1} (-y_{x,j}) \sqrt{h} \sqrt{h}$$

woraus auf gleiche Weise wegen $\sum_{j=i}^{N-1} h = (N-i)h = (1-x_i)$ mit der CSU folgt

$$(3.26) \quad y_i^2 \leq (1-x_i) \sum_{j=i}^{N-1} (y_{x,j})^2 h.$$

Multipliziere (3.25) mit $(1-x_i)$ und (3.26) mit x_i und addiere dann (= Konvexkombination: $y_i^2(1-x_i) + y_i^2 \cdot x_i$) \implies

$$(3.27) \quad y_i^2 \leq x_i(1-x_i) \sum_{j=0}^{N-1} (y_{x,j})^2 h = x_i(1-x_i) \|\mathbf{y}\|_{(1,h)}^2.$$

Unter Beachtung von $\max_{x \in [0,1]} x(1-x) = \frac{1}{4}$ folgt hieraus

$$\|\mathbf{y}\|_\infty \leq \frac{1}{2} \|\mathbf{y}\|_{(1,h)}.$$

■

Zur späteren Verwendung leiten wir, ohne Rückgriff auf Eigenwerte, einen weiteren Vergleich her

$$(3.28) \quad \|\mathbf{y}\|_{(0,h)}^2 \stackrel{!}{\leq} \frac{1}{6} \sum_{j=0}^{N-1} (\mathbf{y}_{x_i})^2 h = \frac{1}{6} \|\mathbf{y}\|_{(1,h)}^2 \stackrel{(3.21)}{=} \frac{1}{6} \|\mathbf{y}\|_{\mathcal{A}_h^0}^2.$$

Beweis: Multipliziere (3.27) mit h und summiere: $\sum_1^{N-1} \implies$

$$\sum_{i=1}^{N-1} y_i^2 h = \|\mathbf{y}\|_{(0,h)}^2 \leq \|\mathbf{y}\|_{(1,h)}^2 \sum_{i=1}^{N-1} x_i(1-x_i)h.$$

Vermutung: $\sum_{i=1}^{N-1} x_i(1-x_i)h \approx \frac{1}{6}$ wegen $\int_0^1 x(1-x)dx = \frac{1}{6}$.

Nun ist

$$\begin{aligned} \sum_{i=1}^{N-1} x_i(1-x_i)h &= h^3 \sum_{i=1}^{N-1} i(N-i) = h^3 \left[N \frac{N(N-1)}{2} - \frac{N(N-1)}{2} \cdot \frac{(2N-1)}{3} \right] \\ &= h^3 \frac{N(N-1)}{2} \left[N - \frac{2N-1}{3} \right] \quad \text{und mit } hN = 1 \\ &= h^2 \frac{(N-1)}{6} [N+1] \\ &= h^2 \frac{(N^2-1)}{6} = \frac{1-h^2}{6} < \frac{1}{6}. \end{aligned}$$

■

Bemerkung: Die 8 aus (3.23) hat sich zu einer 6 verschlechtert.

§ 4 Stabilität (und „bessere“ Verfahren)

Stabilität eines Verfahrens bedeutet, daß bei beschränkten Anfangs- und Randwerten die fortlaufend errechneten Werte beschränkt bleiben. Eine genaue Definition folgt noch. Wir beweisen zunächst ein einfaches Stabilitätsergebnis, das zeigt, daß das beschriebene explizite Verfahren verbesserungsbedürftig ist.

Wir untersuchen das explizite Verfahren (2.8) für den Spezialfall $f = 0$ und Nullrandwerte $y_0^j = y_N^j = 0 \forall j$. Es läßt sich schreiben als

$$(4.1) \quad \mathbf{y}_t^j = -A_h^0 \mathbf{y}^j, \quad j \geq 0, \quad \mathbf{y}^0 = u_0,$$

und zeigen:

Notwendig und hinreichend dafür, daß alle \mathbf{y}^j gemäß (4.1) beschränkt sind, ist

$$(4.2) \quad \frac{\tau}{h^2} \leq \frac{1}{2} \left(\cos^2 \frac{\pi h}{2} \right)^{-1}.$$

Bemerkungen:

1. Für die zu (4.1) gehörige kontinuierliche Aufgabe

$$\dot{u} = u'', \quad u(x, 0) = u_0(x), \quad 0 < x < 1, \quad u(0, t) = u(1, t) = 0$$

gilt das Randmaximumprinzip, aus dem folgt, daß die Lösung $u(x, t)$ für $0 \leq x \leq 1$ und für alle t beschränkt ist durch $\max_{x \in [0, 1]} |u_0(x)|$. Deshalb ist die Beschränktheitsforderung notwendig für ein vernünftiges Verfahren.

2. Natürlich werden wir (4.2) verallgemeinern.
3. Für kleine h wächst $\cos^2 \frac{\pi h}{2} \xrightarrow{h \rightarrow 0} 1$, weshalb als hinreichende Bedingung üblicherweise $\frac{\tau}{h^2} \leq \frac{1}{2}$ genannt wird.
4. Diese Bedingung verlangt auf Grund der hohen Anzahl von Zeitschritten einen großen Rechenaufwand.

Beweis: (4.2)

A_h^0 hat ein Orthonormalsystem (ONS) von $N - 1$ Eigenvektoren (vgl. (3.2)).

$$A_h^0 \mathbf{v}^k = \lambda_k^h \mathbf{v}^k, \quad k = 1, \dots, N - 1, \quad \lambda_k^h > 0$$

Man kann also auf jeder Zeitschicht die durch (4.1) berechneten Vektoren darstellen durch (Fourier-Zerlegung)

$$(4.3) \quad \mathbf{y}^j = \sum_{k=1}^{N-1} c_k^j \mathbf{v}^k$$

Einsetzen dieser Darstellung in das Differenzschema ergibt

$$-\mathbf{A}_h^0 \sum_{k=1}^{N-1} c_k^j \mathbf{v}^k = \frac{\mathbf{y}^{j+1} - \mathbf{y}^j}{\tau} = \sum_{k=1}^{N-1} \frac{c_k^{j+1} - c_k^j}{\tau} \mathbf{v}^k,$$

Koeffizientenvergleich liefert

$$\begin{aligned} \frac{c_k^{j+1} - c_k^j}{\tau} &= -\lambda_k^h c_k^j \\ c_k^{j+1} &= (1 - \tau \lambda_k^h) c_k^j \quad (\text{Rekursionsformel}) \\ &= (1 - \tau \lambda_k^h)^2 c_k^{j-1} \\ &\vdots \\ (4.4) \quad c_k^{j+1} &= (1 - \tau \lambda_k^h)^{j+1} c_k^0. \end{aligned}$$

Gemäß (4.3) sind alle \mathbf{y}^j genau dann beschränkt, wenn dies auch für alle Koeffizienten c_k^j , $k = 1, \dots, N-1$, $j \geq 0$ gilt. Notwendig und hinreichend dafür ist (vgl. (4.4)):

$$|1 - \tau \lambda_k^h| \leq 1,$$

bzw.

$$-1 \leq 1 - \tau \lambda_k^h \leq 1 \quad \forall k.$$

Die linke Ungleichung besagt $\tau \lambda_k^h \leq 2 \quad \forall k$. Wegen $\tau > 0$, $\lambda_k^h > 0$ ist die rechte Ungleichung immer erfüllt.

Wir müssen, da die λ_k^h monoton geordnet sind, diese Ungleichung also für den maximalen Eigenwert von \mathbf{A}_h^0 fordern, also (vgl. Lemma 3.7)

$$(4.5) \quad \tau \frac{4}{h^2} \sin^2 \frac{(N-1)\pi h}{2} \leq 2.$$

Nun ist

$$\frac{(N-1)\pi h}{2} = \frac{N\pi h}{2} - \frac{\pi h}{2} = \frac{\pi}{2} - \frac{\pi h}{2}$$

und

$$\sin \frac{(N-1)\pi h}{2} = \sin \left(\frac{\pi}{2} - \frac{\pi h}{2} \right) = -\sin \left(\frac{\pi h}{2} - \frac{\pi}{2} \right) = \cos \frac{\pi h}{2}. \quad (\text{Phasenverschiebung})$$

Insgesamt liefert (4.5)

$$\frac{\tau}{h^2} \leq \frac{1}{2} \left(\cos^2 \frac{\pi h}{2} \right)^{-1}.$$

■

Ist also τ in der Größenordnung h^2 , so gibt es Ärger. Daß dieses Resultat nicht nur theoretisch sondern auch numerisch Ärger bereitet, zeigen die Übungen. Die Forderung $\frac{\tau}{h^2} \leq \frac{1}{2}$ ist sehr einschneidend. Man braucht „bessere“ Bedingungen (d.h. bessere Verfahren).

Motivation für neue Verfahren:

Der vorwärtsgenommene Differenzenquotient beim expliziten Verfahren hat nur eine lineare Konsistenzordnung (vgl. (2.2)). Derselbe Differenzenquotient als Näherung für die Ableitung auf der Zeitschicht $t_{j+\frac{1}{2}}$ hat quadratische Konsistenzordnung. Daher rührt der Vorschlag, das ganze Verfahren für die Zeitschicht $t_{j+\frac{1}{2}}$ zu formulieren, d.h.

Mittelung der Werte auf alter und neuer Zeitschicht

Wir untersuchen das Verfahren

$$(4.6) \quad \begin{aligned} \mathbf{y}_t &= -\mathbf{A}_h^0 \mathbf{y}^\sigma + \varphi \\ \mathbf{y}^\sigma &:= \sigma \mathbf{y}^{j+1} + (1 - \sigma) \mathbf{y}^j \quad \forall j, \quad \sigma \in [0, 1]. \end{aligned}$$

Bezeichnung: σ als Exponent bedeutet immer nur eine Mittelung, nie eine Potenz.

Bemerkung: φ enthält nur additive, bekannte Werte auf jeder Zeitschicht (Randwerte und f -Werte, vgl. (2.12)). Wir entscheiden später auf welcher Zeitschicht wir φ betrachten. Naheliegend ist $t_{j+\sigma}$.

Durch Taylorentwicklung untersuchen wir zunächst die Auswirkung der Ersetzung von \mathbf{y} durch \mathbf{y}^σ . Wir entwickeln u^σ (für eine differenzierbare Funktion u) an der Stelle $t_{j+\frac{1}{2}}$ (die Indizes $+$ und $-$ bezeichnen wieder Funktionswerte an einer Zwischenstelle)

$$\begin{aligned} u^\sigma &= \sigma \left(u(t_{j+\frac{1}{2}}) + \frac{\tau}{2} \dot{u} + \frac{\tau^2}{8} \ddot{u} + \frac{\tau^3}{8 \cdot 6} \ddot{u}_+ \right) + (1 - \sigma) \left(u(t_{j+\frac{1}{2}}) - \frac{\tau}{2} \dot{u} + \frac{\tau^2}{8} \ddot{u} - \frac{\tau^3}{8 \cdot 6} \ddot{u}_- \right) \\ &= u(t_{j+\frac{1}{2}}) + \frac{\tau}{2} (2\sigma - 1) \dot{u} + \frac{\tau^2}{8} \ddot{u} + O(\tau^3) \end{aligned}$$

\implies

Für $\sigma = \frac{1}{2}$ wird $u(t_{j+\frac{1}{2}})$ durch u^σ von 2. Ordnung approximiert. Da die Approximation der Ableitung 2. Ordnung in Zeitrichtung (durch den zentralen Differenzenquotienten) auch von 2. Ordnung war (vgl. (2.3)), ist keine Einbuße der Approximationsordnung zu befürchten.

In Abhängigkeit von der Wahl von σ erhält man aus (4.6) folgende Verfahren:

$$\sigma = \begin{cases} 0 & : \text{explizites Verfahren: IndexEuler-Verfahren} \\ \frac{1}{2} & : \text{Crank-Nicolson Schema (implizit)} \\ 1 & : \text{implizites Diff.-Schema: implizites Euler Verfahren.} \end{cases}$$

Umformung von (4.6):

$$\begin{aligned} \mathbf{y}^\sigma &= \sigma \hat{\mathbf{y}} + (1 - \sigma) \mathbf{y} \\ &= \sigma (\hat{\mathbf{y}} - \mathbf{y}) + \mathbf{y} \\ &= \sigma \tau \mathbf{y}_t + \mathbf{y} \end{aligned}$$

Einsetzen im (4.6) :

$$(4.7) \quad \begin{aligned} \mathbf{y}_t + \mathbf{A}_h^0(\sigma\tau\mathbf{y}_t + \mathbf{y}) &= \boldsymbol{\varphi} \quad \text{bzw.} \\ \underbrace{(\mathbf{I} + \sigma\tau\mathbf{A}_h^0)}_{=: \mathbf{B}} \mathbf{y}_t + \mathbf{A}_h^0 \mathbf{y} &= \boldsymbol{\varphi} \end{aligned}$$

Mit $\mathbf{B} = (\mathbf{I} + \sigma\tau\mathbf{A}_h^0)$ ordnet sich dieses Verfahren ein in die Klasse der Verfahren

$$(4.8) \quad \boxed{\mathbf{B}\mathbf{y}_t + \mathbf{A}\mathbf{y} = \boldsymbol{\varphi} \quad \text{mit Matrizen } \mathbf{A}, \mathbf{B}} \quad \text{Normalform des } 2\text{-Schicht-Verfahren}$$

die wir im folgenden untersuchen, zunächst unter der Voraussetzungen

$$\begin{aligned} \mathbf{B} &> 0 \quad (\text{d.h. pos. def.}), \quad \mathbf{B} = \mathbf{B}^T \\ \mathbf{A} &= \mathbf{A}^T, \quad \mathbf{A} > 0. \end{aligned}$$

Aus $\mathbf{B} > 0$ folgt \mathbf{B} ist invertierbar. Wünschenswert ist natürlich : \mathbf{B} leicht invertierbar.

Beachte: (4.8) ist i. allg. implizit.

Wir untersuchen nun getrennt

- die Stabilität der Verfahren (4.8) bzgl. der Anfangswerte (AWe), wobei $\boldsymbol{\varphi} \equiv 0$ gesetzt wird (d.h: rechte Seite der Differentialgleichung und Randwerte =0), und
- bezüglich der „rechten Seite $\boldsymbol{\varphi}$ “, wobei $y_0 = 0$, $y_0^j = y_N^j = 0$ (Nullanfangs- und Randwerte). Die Stabilität bzgl. Anfangswerten und rechter Seite erhält man dann durch Superposition.

Definition 4.1 Stabilität

Das Differenzenschema $\mathbf{B}\mathbf{y}_t + \mathbf{A}\mathbf{y} = \boldsymbol{\varphi}$ heißt

stabil bzgl. der $\left\{ \begin{array}{l} \text{Anfangswerte} \\ \text{rechten Seite} \end{array} \right.$

falls Abschätzungen der folgenden Art gelten

$$\|\mathbf{y}^j\|_{(a)} \leq \begin{cases} M_1 \|\mathbf{y}^0\|_{(b)} & \text{wobei } \boldsymbol{\varphi} \equiv 0 \quad (\text{insbesondere } y_0^j = y_N^j = 0) \\ M_2 \|\boldsymbol{\varphi}\|_{(c)} & \text{wobei } \mathbf{y}^0 = 0 \quad (\text{üblicherweise } y_0^j = y_N^j = 0) \end{cases}$$

mit **vernünftigen** Normen $\|\cdot\|_{(a)}, \|\cdot\|_{(b)}, \|\cdot\|_{(c)}$ und Konstanten $M_1, M_2 > 0$ die unabhängig von der Diskretisierung (d.h. von τ, h) sind.

Bemerkungen:

- Die Normen sollen in dem Sinn **vernünftig** sein, daß sie für $h \rightarrow 0$ und/oder $\tau \rightarrow 0$ weder gegen Null noch gegen ∞ gehen. Die L_2 -Norm (für Zeilenvektoren) ist deshalb nicht zulässig, wohl aber z.B.

$$\begin{aligned} \|\mathbf{y}\|_{(0,h)}^2 &= \sum_{i=0}^{N-1} y_i^2 h \xrightarrow{h \rightarrow 0} \int_0^1 u^2(x) dx \\ \|\mathbf{y}\|_{(1,h)}^2 &= \sum_{i=0}^{N-1} (y_{x,i})^2 h \xrightarrow[\substack{\uparrow \\ \text{beachte: } y_0=0}]{h \rightarrow 0} \int_0^1 u'^2(x) dx \quad (y_0 = 0, \text{ vgl. (3.21)}) \end{aligned}$$

$\sqrt{\int_0^1 u'^2(x) dx}$ ist, zusammen mit der Forderung $u(0) = 0$ auch eine Norm für differenzierbare Funktionen.

Die Stabilitätsdefinition macht klar, warum die Normen aus Definition 3.8 eingeführt werden mußten.

2. Stabilität sichert, daß die Lösung bei beschränkten Anfangswerten, Randwerten und rechter Seite beschränkt bleibt (Randmaximumprinzip) und liefert die stetige Abhängigkeit der Lösung von den Anfangswerten und der rechten Seite.
3. Verwunderlich mag erscheinen, daß man immer nur Nullrandwerte betrachtet. Der Grund dafür ist, wie wir sehen werden, daß die Randwerte für Konvergenzbetrachtungen keine Rolle spielen (vgl. § 5). Bei Konvergenzaussagen zeigt man, daß der Fehler *exakte Lösung - approximierte Lösung* gegen Null geht. Auf dem Rand werden immer die exakten Daten vorgegeben, weshalb dort der Fehler immer gleich Null ist.

Wir zeigen zunächst

Satz 4.2 Stabilität bzgl. Anfangswerten und rechter Seite

Für das Differenzensschema

$$(4.9) \quad \mathbf{B}\mathbf{y}_t + \mathbf{A}\mathbf{y} = \boldsymbol{\varphi}, \quad \mathbf{y} \in \mathbb{R}^{n-1} \text{ sei } \mathbf{A}^T = \mathbf{A} > 0, \quad \mathbf{B} > 0 \text{ und } y_0^j = y_N^j = 0 \quad \forall j.$$

Dann gilt

a) $\mathbf{B} - \frac{\tau}{2}\mathbf{A} \geq 0 \iff (4.9) \text{ ist stabil bzgl. der Anfangswerte mit } M_1 = 1 \text{ und (pos. semidef.)}$

$$\|\mathbf{y}^j\|_{(1,h)} \leq \|\mathbf{y}^{j-1}\|_{(1,h)} \leq \dots \leq \|\mathbf{y}^0\|_{(1,h)}.$$

b) $\mathbf{B} - \frac{\tau}{2}\mathbf{A} - \frac{\varepsilon}{2}\mathbf{I} \geq 0$ für ein $\varepsilon > 0 \implies$
 (4.9) ist stabil bzgl. Anfangswerten und rechter Seite und es gilt

$$\begin{aligned} \|\mathbf{y}^j\|_{(1,h)} &\leq \|\mathbf{y}^0\|_{(1,h)} + \frac{1}{\sqrt{\varepsilon}} \|\boldsymbol{\varphi}\|_{(b)}, \quad \text{wobei} \\ \|\boldsymbol{\varphi}\|_{(b)} &= \left(\sum_{k=0}^{j-1} \tau \|\boldsymbol{\varphi}^k\|_{(0,h)}^2 \right)^{1/2} \end{aligned}$$

Bemerkung: $\sqrt{\tau \sum_{r=0}^j \|\boldsymbol{\varphi}^r\|_{(0,h)}^2} =: \|\boldsymbol{\varphi}\|_{(b)}$ ist eine vernünftige Norm, denn

$$\|\varphi^\nu\|_{(0,h)}^2 = \sum_{i=1}^{N-1} |\varphi_i^\nu|^2 h \xrightarrow{h \rightarrow 0} \int_0^1 \varphi^2(x, t_\nu) dx \quad \text{und mit } j \cdot \tau = T$$

$$\tau \sum_{\nu=0}^j \|\varphi^\nu\|_{(0,h)}^2 \xrightarrow[\tau \rightarrow 0]{h \rightarrow 0} \int_0^T \int_0^1 \varphi^2(x, t) dx dt.$$

Also gilt $\|\varphi\|_{(b)} \approx \sqrt{\int_0^T \int_0^1 \varphi(x, t)^2 dx dt}$. j zählt die Zeitschichten, deshalb kann $\|\varphi\|_{(b)}$ mit j wachsen.

In dem φ^j aus (2.12) sind, je nach Komponente, f_i^j – Werte enthalten und Randwerte. Da wir nur Nullrandwerte betrachten, gilt für die Norm genauer

$$\|\varphi^\nu\|_{(0,h)}^2 \xrightarrow[\tau \rightarrow 0]{h \rightarrow 0} \int_0^T \int_0^1 f^2(x, t) dx dt.$$

Beweis (Ideen):

(4.9) wird zu einer Gleichung umgeformt in der $\|\hat{\mathbf{y}}\|$ und $\|\mathbf{y}\|$ auftreten und Ausdrücke, die vorzeichenmäßig beherrschbar sind (Skalarprodukt-Multiplikation). Danach werden Stabilität in Bezug auf Anfangswerte und rechten Seiten getrennt untersucht. Das allgemeine Ergebnis folgt durch Superposition.

Multipliziere (4.9) bzgl. $(\cdot)_{(0,h)}$ mit $2\tau \mathbf{y}_t \implies$

$$2\tau (\mathbf{B} \mathbf{y}_t, \mathbf{y}_t)_{(0,h)} + 2\tau \left(\mathbf{A} \mathbf{y}, \frac{\hat{\mathbf{y}} - \mathbf{y}}{\tau} \right)_{(0,h)} = 2\tau (\varphi, \mathbf{y}_t)_{(0,h)} \quad \text{bzw.}$$

$$2\tau \underbrace{(\mathbf{B} \mathbf{y}_t, \mathbf{y}_t)_{(0,h)}}_{>0} + 2(\mathbf{A} \mathbf{y}, \hat{\mathbf{y}} - \mathbf{y})_{(0,h)} = 2\tau (\varphi, \mathbf{y}_t)_{(0,h)}$$

Idee: Von \mathbf{B} kann man noch etwas subtrahieren, ohne daß die Vorzeichenbedingung verloren geht (Gershgorin), und damit den 2. Summanden umformen. In (4.7) ist $\mathbf{B} = \mathbf{I} + \sigma \tau \mathbf{A}_h^0$. Subtrahiert man $\frac{\tau}{2} \mathbf{A}_h^0$, erhält man

$$\mathbf{B} - \frac{\tau}{2} \mathbf{A}_h^0 = \mathbf{I} + \underbrace{\left(\sigma - \frac{1}{2} \right)}_{=0 \text{ für } \sigma = \frac{1}{2}} \tau \mathbf{A}_h^0$$

Im allgemeinen Fall von Satz 4.2 subtrahieren wir $\frac{\tau}{2} \mathbf{A}$ und erhalten

$$(4.10) \quad 2\tau \left(\left(\mathbf{B} - \frac{\tau}{2} \mathbf{A} \right) \mathbf{y}_t, \mathbf{y}_t \right)_{(0,h)} + \underbrace{\tau^2 (\mathbf{A} \mathbf{y}_t, \mathbf{y}_t)_{(0,h)} + 2(\mathbf{A} \mathbf{y}, \hat{\mathbf{y}} - \mathbf{y})_{(0,h)}}_{=: d} = 2\tau (\varphi, \mathbf{y}_t)$$

Umformung:

$$\begin{aligned}
 d &= (\mathbf{A}(\hat{\mathbf{y}} - \mathbf{y}), \hat{\mathbf{y}} - \mathbf{y})_{(0,h)} + 2(\mathbf{A}\mathbf{y}, \hat{\mathbf{y}} - \mathbf{y})_{(0,h)} \\
 &= (\mathbf{A}(\hat{\mathbf{y}} + \mathbf{y}), \hat{\mathbf{y}} - \mathbf{y})_{(0,h)} \\
 &= (\mathbf{A}\hat{\mathbf{y}}, \hat{\mathbf{y}})_{(0,h)} - (\mathbf{A}\mathbf{y}, \mathbf{y})_{(0,h)} + \underbrace{(\mathbf{A}\mathbf{y}, \hat{\mathbf{y}})_{(0,h)} - (\mathbf{A}\hat{\mathbf{y}}, \mathbf{y})_{(0,h)}}_{=0 \text{ wegen } \mathbf{A}=\mathbf{A}^T, \hat{\mathbf{y}}, \mathbf{y} \text{ reell}} \\
 &= \|\hat{\mathbf{y}}\|_{(1,h)}^2 - \|\mathbf{y}\|_{(1,h)}^2
 \end{aligned}$$

Damit folgt aus (4.10) die sog. *energetische Identität*

$$(4.11) \quad 2\tau \left((\mathbf{B} - \frac{\tau}{2}\mathbf{A})\mathbf{y}_t, \mathbf{y}_t \right)_{(0,h)} + \|\hat{\mathbf{y}}\|_{(1,h)}^2 - \|\mathbf{y}\|_{(1,h)}^2 = 2\tau(\boldsymbol{\varphi}, \mathbf{y}_t)_{(0,h)}.$$

Aus dieser Gleichung werden wir den Beweis ableiten.

Beweis a): Stabilität bzgl. der Anfangswerte ($\boldsymbol{\varphi} \equiv 0$)

„ \implies “

Aus (4.11) folgt sofort: Ist $\mathbf{B} - \frac{\tau}{2}\mathbf{A} \geq 0$ (positive Semidefinitheit), so gilt

$$(4.12) \quad \|\hat{\mathbf{y}}\|_{(1,h)}^2 =: \|\mathbf{y}^{j+1}\|_{(1,h)}^2 \leq \|\mathbf{y}^j\|_{(1,h)}^2 \leq \dots \leq \|\mathbf{y}^0\|_{(1,h)}^2$$

also Stabilität bzgl. der Anfangswerte mit $M_1 = 1$.

In Anbetracht des Randmaximumprinzips ist $M_1 = 1$ auch die richtige Konstante, wenn $\|\cdot\|_{(a)} = \|\cdot\|_{(b)}$ (vgl. Definition 4.1).

„ \longleftarrow “

Mit $\boldsymbol{\varphi} \equiv 0$ und der Stabilität (mit $M = 1$): $\|\hat{\mathbf{y}}\|_{(1,h)}^2 \leq \|\mathbf{y}\|_{(1,h)}^2$ folgt aus (4.11)

$$\left((\mathbf{B} - \frac{\tau}{2}\mathbf{A})\mathbf{y}_t, \mathbf{y}_t \right) \geq 0.$$

\mathbf{A} und \mathbf{B} sind invertierbar. Verfahren (4.9) und $\boldsymbol{\varphi} \equiv 0$ liefern $\mathbf{y}_t = -\mathbf{B}^{-1}\mathbf{A}\mathbf{y}$. Mit beliebigen \mathbf{y} durchläuft auch \mathbf{y}_t den ganzen Raum. Also gilt

$$(4.13) \quad \left((\mathbf{B} - \frac{\tau}{2}\mathbf{A})\mathbf{y}, \mathbf{y} \right) \geq 0.$$

Das bedeutet $\mathbf{B} - \frac{\tau}{2}\mathbf{A} \geq 0$.

Beweis b) Die Stabilität bzgl. der Anfangswerte folgt aus a).

Stabilität bzgl. der rechten Seite: ($\mathbf{y}^0 = 0$, $y_0^j = y_N^j = 0$)

Wir schätzen die rechte Seite von (4.11) ab mit der CSU

$$(u, v) \leq \|u\| \|v\|, \quad \|\cdot\| = \sqrt{(\cdot, \cdot)}$$

und der ε -Ungleichung:

$$2ab \leq \varepsilon a^2 + \frac{1}{\varepsilon} b^2 \text{ für } a, b, \varepsilon > 0.$$

Sie folgt aus $(\varepsilon a - b)^2 = \varepsilon^2 a^2 - 2\varepsilon ab + b^2 \geq 0$.

Somit folgt

$$2\tau(\boldsymbol{\varphi}^j, \mathbf{y}_t^j)_{(0,h)} \stackrel{\text{CSU}}{\leq} 2\tau \|\boldsymbol{\varphi}^j\|_{(0,h)} \|\mathbf{y}_t^j\|_{(0,h)} \stackrel{\varepsilon\text{-Ungleichung}}{\leq} \frac{\tau}{\varepsilon} \|\boldsymbol{\varphi}^j\|_{(0,h)}^2 + \varepsilon \tau \|\mathbf{y}_t\|_{(0,h)}^2.$$

Dies wird in (4.11) eingetragen, der letzte Summand wird unter Beachtung von $\|\mathbf{y}_t\|_{(0,h)}^2 = (\mathbf{I}\mathbf{y}_t, \mathbf{y}_t)_{(0,h)}$ auf die linke Seite gebracht. So folgt

$$2\tau \left((\mathbf{B} - \frac{\tau}{2}\mathbf{A} - \frac{\varepsilon}{2}\mathbf{I})\mathbf{y}_t, \mathbf{y}_t \right)_{(0,h)} + \|\hat{\mathbf{y}}\|_{(1,h)}^2 - \|\mathbf{y}\|_{(1,h)}^2 \leq \frac{\tau}{\varepsilon} \|\boldsymbol{\varphi}^j\|_{(0,h)}^2$$

Die Voraussetzung $\mathbf{B} - \frac{\tau}{2}\mathbf{A} - \frac{\varepsilon}{2}\mathbf{I} \geq 0$ liefert somit

$$\|\mathbf{y}^{j+1}\|_{(1,h)}^2 \leq \|\mathbf{y}^j\|_{(1,h)}^2 + \frac{\tau}{\varepsilon} \|\boldsymbol{\varphi}^j\|_{(0,h)}^2.$$

Wiederholte Anwendung dieses Schrittes ergibt (wegen $\mathbf{y}^0 = 0$)

$$(4.14) \quad \|\mathbf{y}^{j+1}\|_{(1,h)}^2 \leq \underbrace{\|\mathbf{y}^0\|_{(1,h)}^2}_{=0} + \frac{\tau}{\varepsilon} \sum_{\nu=0}^j \|\boldsymbol{\varphi}^\nu\|_{(0,h)}^2 = \frac{\tau}{\varepsilon} \sum_{\nu=0}^j \|\boldsymbol{\varphi}^\nu\|_{(0,h)}^2.$$

Damit liefert (4.14) die Behauptung b) im Fall $\mathbf{y}^0 \equiv 0$.

Nun kann man die Lösung \mathbf{y} der Aufgabe mit Nullrandwerten, Anfangswerten und rechter Seite durch Superposition erhalten $\mathbf{y} = \mathbf{y}_A + \mathbf{y}_R$, wobei

- $\mathbf{y}_A \hat{=}$ Lösung der Aufgabe mit Anfangswerten, Nullrandwerten und rechter Seite = 0
- $\mathbf{y}_R \hat{=}$ Lösung der Aufgabe mit Null-Anfangs- und - Randwerten, aber mit rechter Seite

Mit Hilfe der Dreieckungsgleichung. $\|\mathbf{y}\|_{(1,h)} \leq \|\mathbf{y}_A\|_{(1,h)} + \|\mathbf{y}_R\|_{(1,h)}$ folgt die Behauptung aus (4.13) und (4.14). ■

Bemerkungen

1. Dies ist die Theorie von Samarskij (in Birkhäuser gibt es 2 dicke Bücher von ihm).
2. Die Voraussetzung $\mathbf{B} = \mathbf{B}^T$ (vgl. nach (4.8)) wurde bisher nicht benötigt.
3. Der Satz zeigt, warum die Normvergleiche (3.22), (3.24), (3.28) wichtig sind. Dadurch ist es möglich die Aussagen des Satzes auch auf andere Normen zu übertragen.

Wir wenden Satz 4.2 an auf das Verfahren

$$(4.15) \quad \underbrace{(\mathbf{I} + \sigma\tau\mathbf{A}_h^0)}_B \mathbf{y}_t + \mathbf{A}_h^0 \mathbf{y} = \boldsymbol{\varphi}, \quad (\mathbf{A}_h^0)^T = \mathbf{A}_h^0, \quad y_0^j = y_N^j = 0 \quad \forall j$$

und zeigen, daß man σ , τ und ε -Werte angeben kann, welche die Stabilität sichern. Wir beginnen mit der **Stabilität bzgl. der Anfangswerte**.

Für Skalarproduktnormen und zugeordnete Matrixnormen gilt für beliebige $\mathbf{A} \in \mathbb{R}^n$ mit der CSU

$$(4.16) \quad (\mathbf{A}\mathbf{y}, \mathbf{y}) \leq \|\mathbf{A}\mathbf{y}\| \|\mathbf{y}\| \leq \|\mathbf{A}\| \|\mathbf{y}\|^2 = \|\mathbf{A}\| (\mathbf{y}, \mathbf{y}) \quad \text{bzw.} \quad \frac{\mathbf{A}}{\|\mathbf{A}\|} \leq \mathbf{I}.$$

Damit verschärfen wir die Bedingung $\mathbf{I} + \sigma\tau\mathbf{A}_h^0 - \frac{\tau}{2}\mathbf{A}_h^0 \geq 0$ aus Satz 4.2 a) zu

$$\frac{\mathbf{A}_h^0}{\|\mathbf{A}_h^0\|} + \sigma\tau\mathbf{A}_h^0 - \frac{\tau}{2}\mathbf{A}_h^0 \geq 0$$

und mit der Spektralnorm wegen $\|\mathbf{A}_h^0\| \leq \frac{4}{h^2}$ nochmals zu

$$\frac{h^2}{4}\mathbf{A}_h^0 + \sigma\tau\mathbf{A}_h^0 - \frac{\tau}{2}\mathbf{A}_h^0 \geq 0 \quad \text{bzw.} \quad \left(\frac{h^2}{4} + \tau\left(\sigma - \frac{1}{2}\right)\right)\mathbf{A}_h^0 \geq 0.$$

Dies ist wegen $\mathbf{A}_h^0 > 0$ richtig, falls $\frac{h^2}{4} + \tau\left(\sigma - \frac{1}{2}\right) \geq 0$ bzw.

$$(4.17) \quad \sigma \geq \frac{1}{2} - \frac{h^2}{4\tau}$$

Bedeutung

1. Das explizite Verfahren ($\sigma = 0$) ist stabil bzgl. der Anfangswerte, falls

$$0 \geq \frac{1}{2} - \frac{h^2}{4\tau}, \quad \text{bzw.} \quad \frac{\tau}{h^2} \leq \frac{1}{2}.$$

Dies war für den Spezialfall $\boldsymbol{\varphi} = 0$ schon bewiesen.

2. Das implizite Verfahren ist für $\sigma \geq \frac{1}{2}$ (insbesondere also Crank-Nicolson) *unbedingt stabil* bzgl. der Anfangswerte, d.h. stabil ohne Bedingungen an σ und τ .

Stabilität bzgl. der rechten Seite verlangt für ein $\varepsilon > 0$

$$\underbrace{\mathbf{I} + \sigma\tau\mathbf{A}_h^0}_B - \frac{\tau}{2}\mathbf{A}_h^0 - \frac{\varepsilon}{2}\mathbf{I} \geq 0.$$

Mit (4.16): $\frac{\mathbf{A}}{\|\mathbf{A}\|} \leq \mathbf{I}$, und damit

$$\mathbf{I} = \frac{\varepsilon}{2}\mathbf{I} + \left(1 - \frac{\varepsilon}{2}\right)\mathbf{I} \geq \frac{\varepsilon}{2}\mathbf{I} + \left(1 - \frac{\varepsilon}{2}\right)\frac{\mathbf{A}_h^0}{\|\mathbf{A}_h^0\|}$$

verschärfen wir die Bedingung unter der Voraussetzung $1 - \frac{\varepsilon}{2} \geq 0$ zu

$$\frac{\varepsilon}{2} \mathbf{I} + (1 - \frac{\varepsilon}{2}) \frac{\mathbf{A}_h^0}{\|\mathbf{A}_h^0\|} + \sigma \tau \mathbf{A}_h^0 - \frac{\tau}{2} \mathbf{A}_h^0 - \frac{\varepsilon}{2} \mathbf{I} \stackrel{!}{\geq} 0 \quad \text{bzw.}$$

$$\left(\frac{1 - \frac{\varepsilon}{2}}{\|\mathbf{A}_h^0\|} + \tau(\sigma - \frac{1}{2}) \right) \mathbf{A}_h^0 \geq 0.$$

Wegen $\|\mathbf{A}_h^0\| \leq \frac{4}{h^2}$ und $\mathbf{A}_h^0 > 0$ ist diese Bedingung sicher erfüllt, falls

$$(1 - \frac{\varepsilon}{2}) \frac{h^2}{4} + \tau(\sigma - \frac{1}{2}) \geq 0 \quad \text{bzw.}$$

$$(4.18) \quad \sigma \geq \frac{1}{2} - \frac{h^2}{4\tau} (1 - \frac{\varepsilon}{2}).$$

Insgesamt erhalten wir also

Folgerung 4.3

Für das Verfahren

$$(\mathbf{I} + \sigma \tau \mathbf{A}_h^0) \mathbf{y}_t + \mathbf{A}_h^0 \mathbf{y} = \boldsymbol{\varphi}, \quad (\mathbf{A}_h^0)^T = \mathbf{A}_h^0, \quad y_0^j = y_N^j = 0 \quad \forall j$$

gilt mit den Normen, die in Satz (4.2) verwandt wurden

a) Ist $\sigma \geq \frac{1}{2} - \frac{h^2}{4\tau} \left(1 - \frac{\varepsilon}{2}\right), \quad 0 < \varepsilon \leq 2,$

so ist das (implizite) Verfahren *stabil* bzgl. Anfangswerten und rechter Seite (vgl. (4.18)), bzw. *bedingt stabil*, d.h. in Abhängigkeit von der Wahl der Schrittweiten.

b) Ist $\sigma \geq \frac{1}{2}, \quad (\varepsilon \leq 2 \text{ ist keine echte Einschränkung})$

so ist das (implizite) Verfahren *unbedingt stabil* bzgl. Anfangswerten und rechter Seite (d.h. ohne Bedingungen an die Schrittweiten τ, h).

c) Ist $\sigma \geq \frac{1}{2} - \frac{h^2}{4\tau} \quad (\text{vgl. (4.17)})$

(also $\varepsilon = 0$ in a)), so ist das Verfahren *bedingt stabil* bzgl. der Anfangswerte.

d) Durch Verschärfung von c) zu $\sigma = 0$ folgt insbesondere die schon bekannte Bedingung

Ist $\frac{\tau}{h^2} \leq \frac{1}{2},$

(also $\sigma = 0, \varepsilon = 0$ in a)), so ist das (explizite) Verfahren *bedingt stabil* bzgl. der Anfangswerte.

Bemerkung:

$\varepsilon > 0$ kann in Voraussetzung a) beliebig klein sein. Daher liegt die Vermutung nahe, daß die Voraussetzung aus c) nicht nur die Stabilität bzgl. der Anfangswerte, sondern auch bzgl. der rechten Seite liefern könnte. Dies wird sich in der Tat bestätigen.

§ 5 Approximations- und Verfahrensfehler

Wir untersuchen das Verfahren (4.6), (4.7): $\mathbf{y}_t + \mathbf{A}_h^0 \mathbf{y}^{(\sigma)} = \varphi$ in der Gestalt

$$(5.1) \quad y_{t,i}^j - y_{\bar{x},i}^{j(\sigma)} - \varphi = 0, \quad i = 1, \dots, N-1, \quad \mathbf{y}^{(\sigma)} = \sigma \hat{\mathbf{y}} + (1-\sigma)\mathbf{y}, \quad \sigma \in [0, 1].$$

Wir behalten uns eine genaue Definition von φ noch vor. φ ergab sich aus einer Diskretisierung des Differentialgleichungsproblems (vgl. etwa (2.12)). Auch andere Diskretisierungen sind denkbar. Es ist nur darauf zu achten, daß die Diskretisierung für $h, \tau \rightarrow 0$ gegen die Differentialgleichung konvergiert, und daß φ so gewählt wird, daß der Stabilitätsbegriff (vgl. dazu die Norm $\|\varphi\|_{(b)}$) nicht darunter leidet.

Man beachte, daß gemäß der Definition von $\mathbf{y}^{(\sigma)}$ für jede Ableitung ∂_x^α ($\alpha \hat{=}$ Multiindex), und auch für jede entsprechende Diskretisierung d_x^α von ∂_x^α gilt

$$(\partial^\alpha u)^{(\sigma)} = \partial^\alpha (u^{(\sigma)}), \quad (d^\alpha \mathbf{u})^{(\sigma)} = d^\alpha (\mathbf{u}^{(\sigma)}).$$

Sei u die Lösung der kontinuierlichen Aufgabe ($\hat{=}$ exakte Lösung) und $\mathbf{u} = (u_i^j)$ die durch ihre Restriktion auf das Gitter entstandene Gitterfunktion.

Wir bezeichnen mit ψ (vgl. (5.2)) den Approximationsfehler des Differenzenschemas.

$$(5.2) \quad \psi_i^j := u_{t,i}^j - u_{\bar{x},i}^{j(\sigma)} - \varphi_i - \underbrace{\left(y_{t,i}^j - y_{\bar{x},i}^{j(\sigma)} - \varphi_i \right)}_{=0 \text{ laut Verfahren}}.$$

Da \mathbf{u} und \mathbf{y} die gleichen Randwerte haben, kann man ohne Einschränkung von **Nullrandwerten** ausgehen.

Definition 5.1

Sei u die exakte Lösung, so heißt die Gitterfunktion

$$\psi := \mathbf{u}_t - \mathbf{u}_{\bar{x}}^{(\sigma)} - \varphi \quad (\text{E Nullrandwerte})$$

Approximationsfehler (*Diskretisierungsfehler, truncation error*) des Differenzenschemas

$$\mathbf{y}_t - \mathbf{y}_{\bar{x}}^{(\sigma)} - \varphi = \mathbf{0}.$$

Zur Abschätzung des Approximationsfehlers benutzen wir (im Punkt x_i) die Taylorentwicklungen (2.3) und (2.5) bzgl. der Zeit an der Stelle $t^{j+\frac{1}{2}}$ (die Indizes z bezeichnen wieder geeignete Zwischenstellen) und erhalten (Raumindizes unterdrücken, φ wird

zunächst nur mitgeführt):

$$\begin{aligned}
\psi^j &= \dot{u}^{j+\frac{1}{2}} + \frac{\tau^2}{24} \ddot{u}_z - \left(u'' + \frac{h^2}{12} u^{(4)} + \frac{h^4}{360} u_z^{(6')} \right)^{(\sigma)} - \varphi \\
&= \dot{u}^{j+\frac{1}{2}} + \frac{\tau^2}{24} \ddot{u}_z - \left(\left(u'' + \frac{h^2}{12} u^{(4)} \right)^{(\sigma)} + \left(\frac{h^4}{360} u_z^{(6')} \right)^{(\sigma)} \right) - \varphi \\
&= \dot{u}^{j+\frac{1}{2}} + \frac{\tau^2}{24} \ddot{u}_z \\
&\quad - \left(\sigma \left(u''^{j+1} + \frac{h^2}{12} u^{(4),j+1} \right) + (1-\sigma) \left(u''^j + \frac{h^2}{12} u^{(4),j} \right) \right) - \varphi + O(h^4)
\end{aligned}$$

Zusammenfassend folgt unter Beachtung von $\frac{\tau^2}{2} \ddot{u}_z = O(\tau^2)$

$$\begin{aligned}
\psi^j &= \dot{u}^{j+\frac{1}{2}} - \left[\underbrace{\sigma u''^{j+1} + (1-\sigma) u''^j}_{u''^{(\sigma)}} + \frac{h^2}{12} \underbrace{\left(\sigma u^{(4),j+1} + (1-\sigma) u^{(4),j} \right)}_{u^{(4)'(\sigma)} \right) \right] - \varphi + O(\tau^2) + O(h^4) \\
(5.3) \quad \psi^j &= \dot{u}^{j+\frac{1}{2}} - u''^{(\sigma)} - \frac{h^2}{12} u^{(4)'(\sigma)} - \varphi + O(\tau^2) + O(h^4).
\end{aligned}$$

Wir berechnen $u''^{(\sigma)}$ und $u^{(4)'(\sigma)}$ durch Taylorentwicklung an $t^{j+\frac{1}{2}}$ bis auf Größenordnungen von $O(\tau^2)$ und $O(h^4)$. Der besseren Lesbarkeit halber schreiben wir - ausnahmsweise - die Zeitindizes unten:

$$\begin{aligned}
u''_{j+1} &= u''_{j+\frac{1}{2}} + \frac{\tau}{2} \dot{u}''_{j+\frac{1}{2}} + \frac{1}{2} \left(\frac{\tau}{2} \right)^2 \ddot{u}''_z \\
u''_j &= u''_{j+\frac{1}{2}} - \frac{\tau}{2} \dot{u}''_{j+\frac{1}{2}} + \frac{1}{2} \left(\frac{\tau}{2} \right)^2 \ddot{u}''_z \\
\implies u''^{(\sigma)} &= u''_{j+\frac{1}{2}} + \frac{\tau}{2} (2\sigma - 1) \dot{u}''_{j+\frac{1}{2}} + O(\tau^2) \\
u^{(4)}_{j+1} &= u^{(4)}_{j+\frac{1}{2}} + \frac{\tau}{2} \dot{u}^{(4)}_{j+\frac{1}{2}} + O(\tau^2) \\
u^{(4)}_j &= u^{(4)}_{j+\frac{1}{2}} - \frac{\tau}{2} \dot{u}^{(4)}_{j+\frac{1}{2}} + O(\tau^2) \\
\implies u^{(4)'(\sigma)} &= u^{(4)}_{j+\frac{1}{2}} + \frac{\tau}{2} (2\sigma - 1) \dot{u}^{(4)}_{j+\frac{1}{2}} + O(\tau^2)
\end{aligned}$$

In (5.3) ist einzutragen $\frac{h^2}{12} u^{(4)'(\sigma)}$. Für den 2. Summanden S_2 von $u^{(4)'(\sigma)}$ folgt somit

$$S_2 := \frac{h^2}{12} \cdot \frac{\tau}{2} (2\sigma - 1) \dot{u}^{(4)}_{j+\frac{1}{2}}.$$

Wegen $(\tau - h^2)^2 = \tau^2 - 2\tau h^2 + h^4 \geq 0$ ist $\tau h^2 \leq \frac{1}{2}(\tau^2 + h^4) \implies S_2 \leq O(\tau^2 + h^4)$.

Eintragen dieser Ergebnisse in (5.3) liefert

$$\begin{aligned}\psi &= \dot{u}_{j+\frac{1}{2}} - \left(u''_{j+\frac{1}{2}} + \frac{\tau}{2}(2\sigma - 1)\dot{u}''_{j+\frac{1}{2}} \right) - \frac{h^2}{12}u^{(4)}_{j+\frac{1}{2}} - \varphi + O(\tau^2 + h^4) \\ &= \underbrace{\dot{u}_{j+\frac{1}{2}} - u''_{j+\frac{1}{2}} - f_{j+\frac{1}{2}} - \varphi + f_{j+\frac{1}{2}}}_{=0 \text{ laut Dgl.}} - \frac{\tau}{2}(2\sigma - 1)\dot{u}''_{j+\frac{1}{2}} - \frac{h^2}{12}u^{(4)}_{j+\frac{1}{2}} + O(\tau^2 + h^4)\end{aligned}$$

$$(5.4) \quad \psi = -\varphi + f_{j+\frac{1}{2}} - \frac{\tau}{2}(2\sigma - 1)\dot{u}''_{j+\frac{1}{2}} - \frac{h^2}{12}u^{(4)}_{j+\frac{1}{2}} + O(\tau^2 + h^4)$$

Bemerkung zur Schreibweise: Im Zusammenhang mit Differenzenverfahren schreibt man oft zur Vereinfachung der Schreibweise z.B. (wie oben geschehen)

$$O(\tau^2) + O(h^4) + O(\tau^2 + h^4) = O(\tau^2 + h^4),$$

was auf Grund der Definition des Landausymbols nicht korrekt ist. Die linke Schreibweise bedeutet: Es gibt Terme mit $O(\tau^2)$, $O(h^4)$ und solche mit $O(\tau^2 + h^4)$. Da man bei der Konvergenzbehandlung üblicherweise τ und h gleichzeitig gegen Null gehen läßt, haben beide Schreibweisen den gleichen Effekt.

Setzt man in (5.4) $\varphi := f_{j+\frac{1}{2}}$, so folgt hieraus, falls $u \in C^4$ bzgl. des Ortes und $u \in C^3$ bzgl. der Zeit:

$$\left\{ \begin{array}{l} \psi = O(\tau + h^2), \text{ für beliebiges } \sigma \\ \psi = O(\tau^2 + h^2), \text{ für } \sigma = \frac{1}{2} \end{array} \right\} \quad \text{bzw. } \psi = O\left(\tau\left(\sigma - \frac{1}{2}\right) + \tau^2 + h^2\right)$$

Dieses Ergebnis läßt sich jedoch noch verbessern.

Aus der Differentialgleichung $\dot{u} - u'' - f = 0$ folgt $\dot{u}'' - u^{(4)} - f'' = 0$ oder

$$u^{(4)}_{j+\frac{1}{2}} = \dot{u}''_{j+\frac{1}{2}} - f''_{j+\frac{1}{2}}$$

Eintragen in (5.4) ergibt

$$\psi = \underbrace{-\varphi + f_{j+\frac{1}{2}} + \frac{h^2}{12} f''_{j+\frac{1}{2}}}_{\stackrel{!}{=} 0 \text{ durch Wahl von } \varphi} - \underbrace{\left[\frac{\tau}{2}(2\sigma - 1) + \frac{h^2}{12} \right]}_{\stackrel{!}{=} 0 \text{ durch Wahl von } \sigma} \dot{u}''_{j+\frac{1}{2}} + O(\tau^2 + h^4)$$

Insgesamt erhalten wir damit

Satz 5.2

Für den Approximationsfehler ψ des Verfahrens $\mathbf{y}_t - \mathbf{y}_{\bar{x}\bar{x}}^{(\sigma)} - \varphi = \mathbf{0}$ gilt

a) falls $\varphi = \mathbf{f}^{j+\frac{1}{2}}$, $u \in C^4$ bzgl. des Ortes und $u \in C^3$ bzgl. der Zeit

$$\left\{ \begin{array}{l} \psi = O(\tau + h^2), \text{ für beliebiges } \sigma \\ \psi = O(\tau^2 + h^2), \text{ für } \sigma = \frac{1}{2} \end{array} \right\} \text{ bzw. } \psi = O\left(\tau\left(\sigma - \frac{1}{2}\right) + \tau^2 + h^2\right)$$

b) falls $\varphi = \mathbf{f}^{j+\frac{1}{2}} + \frac{h^2}{12} \mathbf{f}''^{j+\frac{1}{2}}$, $\sigma = \frac{1}{2} - \frac{h^2}{12\tau}$, und

$u \in C^6$ bzgl. des Ortes und $u \in C^4$ bzgl. der Zeit gilt $\psi = O(\tau^2 + h^4)$

Bemerkungen:

1. Die Numerik zeigt im Fall b) in der Tat eine deutlich schnellere Konvergenz.
2. In der Praxis will man f'' nicht gerne berechnen. Das ist oft zu fehleranfällig, wennes denn überhaupt möglich ist. Man approximiert deshalb φ wie folgt durch:

$$(5.5) \quad \begin{aligned} \left(f + \frac{h^2}{12} f''\right)_i^{j+\frac{1}{2}} &\approx f_i^{j+\frac{1}{2}} + \frac{1}{12} (f_{i+1} - 2f_i + f_{i-1})^{j+\frac{1}{2}} \\ &= \frac{1}{12} (f_{i+1} + 10f_i + f_{i-1})^{j+\frac{1}{2}} \approx \varphi_i^j \end{aligned}$$

und erhält als Geschenk im allgemeinen die Folgerung

$$\mathbf{f} \geq \mathbf{0} \implies \varphi \geq \mathbf{0}.$$

3. Beachte:

Die in der Stabilitätsdefinition benutzte Norm lautet nun

$$\|\varphi\|_b^2 = \sum \sum \tau h |f_i^{j+\frac{1}{2}} + \frac{h^2}{12} f''^{j+\frac{1}{2}}|^2.$$

Mit $\tau, h \rightarrow 0$ strebt auch diese Norm gegen $\int_0^T \int_0^1 f^2(x, t) dx dt$, sodaß die Stabilität ungefährdet ist.

Nach Folgerung (4.3) war $\sigma \geq \frac{1}{2} - \frac{h^2}{4\tau} \left(1 - \frac{\varepsilon}{2}\right)$ hinreichend für Stabilität bzgl. der Anfangswerte und der rechten Seite. Die Forderung in Voraussetzung b) des vorigen Satzes lautet $\sigma = \frac{1}{2} - \frac{h^2}{12\tau}$. Für $\varepsilon \leq \frac{4}{3}$ gilt nun

$$\frac{1}{2} - \frac{h^2}{12\tau} \geq \frac{1}{2} - \frac{h^2}{4\tau} \left(1 - \frac{\varepsilon}{2}\right),$$

d.h. unter der Voraussetzung $\varepsilon \leq \frac{4}{3}$ ist die Stabilität bzgl. der Anfangswerte und der rechten Seite (gemeint ist die Abhängigkeit von f) gewährleistet. Wenn wir später zeigen können, daß in Voraussetzung a) aus Folgerung 4.3 auch $\varepsilon = 0$ zugelassen ist, so ist dies nochmals eine Verschärfung des obigen Ergebnisses.

Definition 5.3 Verfahrensfehler des Schemas $\mathbf{y}_t - \mathbf{y}^{(\sigma)} - \varphi = \mathbf{0}$

Sei $u(x, t)$ die exakte Lösung des kontinuierlichen Problems, (beliebige Anfangswerte, Randwerte und rechte Seite).

\mathbf{u} die zugehörige Gitterfunktion

\mathbf{y} die Gitterfunktion der Lösungsvektoren des Differenzschemas

$$\mathbf{y}_t - \mathbf{y}_{\bar{x}\bar{x}}^{(\sigma)} - \varphi = \mathbf{0}.$$

Dann definieren wir als *Verfahrensfehler* die Gitterfunktion $\mathbf{z} := \mathbf{u} - \mathbf{y}$.

Bemerkung: Da für die Approximation als Rand- und Anfangswerte die exakten Werte vorgegeben werden, hat \mathbf{z} Nullanfangs- und Randwerte.

Zur Abschätzung des Verfahrensfehlers (und damit zu Konvergenzaussagen) konstruieren wir ein Differenzschema für \mathbf{z} .

$$\begin{aligned} \mathbf{z}_t - \mathbf{z}_{\bar{x}\bar{x}}^{(\sigma)} &= \mathbf{u}_t - \mathbf{u}_{\bar{x}\bar{x}}^{(\sigma)} - \underbrace{[\mathbf{y}_t - (\mathbf{y}_{\bar{x}\bar{x}})^{(\sigma)}]}_{= \varphi \text{ laut Differenzschema}} \\ &= \mathbf{u}_t - \mathbf{u}_{\bar{x}\bar{x}}^{(\sigma)} - \varphi = \psi \quad (\text{vgl. Definition 5.1}) \end{aligned}$$

\mathbf{z} erfüllt also ein **Differenzschema mit Nullrandwerten und Nullanfangswerten** und dem Approximationsfehler als rechte Seite.

$$\mathbf{z}_t - \mathbf{z}_{\bar{x}\bar{x}}^{(\sigma)} - \psi = \mathbf{0}.$$

Mit Hilfe des Stabilitätssatzes 4.2 b)) (Stabilität bzgl. der rechten Seite), der Folgerung 4.3 und des Satzes 5.2 erhält man also die Konvergenzaussage

Satz 5.4 Konvergenz

Gegeben sei das Differenzenschema $\mathbf{y}_t - \mathbf{y}_{\bar{x}\bar{x}}^{(\sigma)} - \boldsymbol{\varphi} = \mathbf{0}$ mit beliebigen Anfangswerten und beliebigen Dirichlet-Randbedingungen.

Sei $\mathbf{z} = \mathbf{u} - \mathbf{y}$ der Verfahrensfehler und $\boldsymbol{\psi} = \mathbf{u}_t - \mathbf{u}_{\bar{x}\bar{x}}^{(\sigma)} - \boldsymbol{\varphi}$ der Approximationsfehler, so gilt für $\sigma \geq \frac{1}{2} - \frac{h^2}{4\tau} \left(1 - \frac{\varepsilon}{2}\right)$:

$$\|\mathbf{z}^j\|_{A_h^0} \leq \left(\frac{1}{\varepsilon} \sum_{\nu=0}^{j-1} \tau \|\boldsymbol{\psi}^\nu\|_{(0,h)}^2 \right)^{\frac{1}{2}} \leq \sqrt{\frac{j\tau}{\varepsilon}} \max_{\nu \leq j-1} \|\boldsymbol{\psi}^\nu\|_{(0,h)}$$

und für $T \geq j \cdot \tau$ folgt also:

Falls $u \in C^3$ bzgl. Zeit, $u \in C^4$ bzgl. Ort gilt

$$\|\mathbf{z}^j\|_{A_h^0} \leq \begin{cases} O(\tau + h^2) & \text{für } \sigma \neq \frac{1}{2} \\ O(\tau^2 + h^2) & \text{für } \sigma = \frac{1}{2}, \boldsymbol{\varphi} = \mathbf{f}^{j+\frac{1}{2}}, \\ & \text{(stabil für } \varepsilon \leq 2) \end{cases}$$

Falls $u \in C^4$ bzgl. Zeit, $u \in C^6$ bzgl. Ort gilt

$$\|\mathbf{z}^j\|_{A_h^0} \leq \begin{cases} O(\tau^2 + h^4) & \text{für } \sigma = \frac{1}{2} - \frac{h^2}{12}, \boldsymbol{\varphi} = \left(\mathbf{f} + \frac{h^2}{12}\mathbf{f}''\right)^{j+\frac{1}{2}}, \\ & \text{(stabil für } \varepsilon \leq \frac{4}{3}). \end{cases}$$

Bemerkungen:

1. Wesentlich für Konvergenz ist Stabilität bzgl. der rechten Seite
2. Beliebt (insbesondere bei Ingenieuren) sind Fehlerabschätzungen in der Maximumnorm (vgl. (3.24): $\|\mathbf{y}\|_\infty \leq \frac{1}{2}\|\mathbf{y}\|_{(1,h)}$). Eine Abschätzung des Approximationsfehlers $\|\boldsymbol{\psi}^k\|_{(0,h)}$ nach oben durch die Maximumnorm ist zwar herleitbar, doch ist der Approximationsfehler $\|\boldsymbol{\psi}^k\|_{(0,h)}$ nur bei speziellen Beispielen auswertbar.

In der Praxis ist man deshalb mit Konvergenz zufrieden und prüft die Genauigkeit durch Schrittweithalbierungen oder man testet das Verfahren (und natürlich damit auch die Programmierung) an geeigneten Testbeispielen.

Bemerkungen zur Konstruktion von Testbeispielen mit exakten Lösungen

Betrachte $r := \frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2}$.

Ersetze u durch ein Polynom $p(x, t)$, berechne $g := \frac{\partial p}{\partial t} - \frac{\partial^2 p}{\partial x^2}$.

Dann ist p die exakte Lösung der Aufgabe $\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = g$ mit den Rand- und Anfangswerten $p|_{\Gamma}$ und $p|_{t=0}$.

Hinweis: Solche Konstruktionen, ebenfalls unter Verwendung von Polynomen, funktionieren auch bei nicht linearen Aufgaben der Art

$$c(u) \frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left(k(u) \frac{\partial u}{\partial x} \right) + f(x) = 0.$$

Bevor wir auf solche Gleichungen eingehen, betrachten wir ein schnelles numerisches Verfahren zur Lösung der bisher behandelten Aufgaben.

§ 6 Spezielle Lösungsverfahren für lineare Gleichungssysteme

Der, das $\left\{ \begin{array}{l} \text{verkürzte Gauß - Algorithmus} \\ \text{Tridiagonalalgorithmus} \\ \text{Thomas - Algorithmus} \\ \text{Double sweep - Verfahren} \end{array} \right.$

zur rechenzeitsparenden Auflösung von tridiagonalen Gleichungssystemen.

Wir leiten zuerst das Verfahren her und geben dann leicht nachprüfbar hinreichende Bedingungen für seine Durchführbarkeit an. Vorgelegt sei das tridiagonale Gleichungssystem:

$$\begin{aligned}
 (6.1) \quad & b_0 y_0 - c_0 y_1 = f_0 \quad i = 0 \\
 & -a_1 y_0 + b_1 y_1 \quad \ddots = f_1 \\
 & \quad \quad \quad \ddots \quad \ddots \quad \ddots \\
 & \quad \quad \quad \quad -a_i y_{i-1} + b_i y_i - c_i y_{i+1} = f_i \quad i = 1, \dots, n-1 \\
 & \quad \quad \quad \quad \quad \quad \quad \ddots \quad \ddots \quad \ddots \\
 & \quad \quad \quad \quad \quad \quad \quad \quad -a_n y_{n-1} + b_n y_n = f_n \quad i = n.
 \end{aligned}$$

Man zeigt schnell, daß das GEV eine LU-Zerlegung liefert (Lower, Upper): $\mathbf{A} = \mathbf{LU}$ mit Matrizen (Diagonale + untere oder obere Nebendiagonale)

$$\mathbf{L} = \begin{bmatrix} \ddots & & & & \mathbf{0} \\ \ddots & \ddots & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ \mathbf{0} & & & \ddots & \ddots \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} \ddots & & & & \mathbf{0} \\ \ddots & \ddots & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ \mathbf{0} & & & \ddots & \ddots \end{bmatrix}$$

Die Gleichung $\mathbf{LUy} = \mathbf{f}$ wird dann gelöst durch

- 1) $\mathbf{Lv} = \mathbf{f}$
- 2) $\mathbf{Uy} = \mathbf{v}$

Das System 2) hat dann – mit $\mathbf{U} = (u_{ij})$ – (abgesehen von der letzten Gleichung) die Gestalt

$$\begin{aligned}
 u_{ii}y_i + u_{i,i+1}y_{i+1} &= v_i & i = 0, \dots, n-1 \\
 (6.2) \quad y_i &= \underbrace{-\frac{u_{i,i+1}}{u_{ii}}}_{\alpha_{i+1}} y_{i+1} + \underbrace{\frac{v_i}{u_{ii}}}_{\beta_{i+1}} \\
 &= \alpha_{i+1} y_{i+1} + \beta_{i+1}
 \end{aligned}$$

Wir leiten Rekursionsformeln zur Berechnung der $\alpha_{i+1}, \beta_{i+1}$ her und berechnen y_n . Dann können aus (6.2) die Werte y_i berechnet werden.

Wir setzen y_{i-1} gemäß (6.2) in die Zeilen $i = 1, \dots, n-1$ von (6.1) ein. \implies

$$-a_i(\alpha_i y_i + \beta_i) + b_i y_i - c_i y_{i+1} = f_i \quad i = 1, \dots, n-1$$

und erhalten durch Auflösen nach y_i die Rekursionsformel

$$(6.3) \quad y_i = \frac{c_i}{b_i - \alpha_i a_i} y_{i+1} + \frac{f_i + \alpha_i \beta_i}{b_i - \alpha_i a_i}, \quad i = 1, \dots, n-1, \quad \underline{b_i - \alpha_i a_i \neq 0}$$

Aus dem Vergleich von (6.2) und (6.3) erhalten wir Rekursionsformeln für α_i, β_i :

$$(6.4) \quad \alpha_{i+1} = \frac{c_i}{b_i - \alpha_i a_i}, \quad \beta_{i+1} = \frac{f_i + \beta_i a_i}{b_i - \alpha_i a_i}, \quad i = 1, \dots, n-1, \quad b_i - \alpha_i a_i \neq 0$$

Aus der 0.ten Zeile von (6.1) folgt

$$b_0 y_0 - c_0 y_1 = f_0 \implies y_0 = \frac{c_0}{b_0} y_1 + \frac{f_0}{b_0}, \quad \underline{b_0 \neq 0}$$

Aus dem Vergleich mit (6.2) also

$$(6.5) \quad \alpha_1 = \frac{c_0}{b_0}, \quad \beta_1 = \frac{f_0}{b_0}, \quad b_0 \neq 0.$$

Definieren wir $a_0 := 0$, so ist (6.5) in (6.4) enthalten für $i = 0$, und die Gleichung (6.3) gilt ebenfalls für $i = 0$.

Den Anfangswert y_n für (6.3) erhält man aus der letzten Zeile von (6.1), wenn man y_{n-1} aus (6.3), bzw. (6.2) in die letzte Zeile von (6.1) einsetzt.

$$\begin{aligned}
 -a_n(\alpha_n y_n + \beta_n) + b_n y_n &= f_n \implies \\
 (6.6) \quad y_n &= \frac{f_n + \beta_n a_n}{b_n - \alpha_n a_n} =: \beta_{n+1}, \quad \underline{b_n - \alpha_n a_n \neq 0}
 \end{aligned}$$

d.h. (6.6) ist auch der 2. Formel von (6.4) für β_{i+1} enthalten, wenn man dort $i = n$ zuläßt.

Insgesamt lautet das *Tridiagonalverfahren* also:

$$(I) \quad \alpha_1 = \frac{c_0}{b_0}, \quad \beta = \frac{f_0}{b_0}, \quad b_0 \neq 0 \quad (\text{vgl. (6.5)})$$

Berechne gemäß (6.4), (6.6)

$$(II) \quad \alpha_{i+1} = \frac{c_i}{b_i - \alpha_i a_i}, \quad i = 1, \dots, n-1, \quad b_i - \alpha_i a_i \neq 0 \quad i = 1, \dots, n$$

$$(III) \quad \beta_{i+1} = \frac{f_i + \beta_i a_i}{b_i - \alpha_i a_i}, \quad i = 1, \dots, n, \quad b_i - \alpha_i a_i \neq 0 \quad i = 1, \dots, n$$

und gemäß (6.3)

$$(IV) \quad y_i = \alpha_{i+1} y_{i+1} + \beta_{i+1}, \quad i = 0, 1, \dots, n-1 \quad y_n = \beta_{n+1} = \frac{f_n + \beta_n a_n}{b_n - \alpha_n a_n}$$

Bevor wir hinreichende Bedingungen für die Durchführbarkeit angeben, machen wir eine kurze Aufwandsbetrachtung.

Bezeichnet man als eine Operation den Aufwand für 1 Addition + Multiplikation oder für 1 Division, so benötigt man für

(I) 2 Operationen

(II) $2(n-1)$ Operationen (je 1 für den Nenner und eine für die Division)

(III) $2n$ Operationen (je 1 für den Zähler und eine für den Bruch,
(der Nenner ist aus (II) bekannt),

(IV) n Operationen,

insgesamt also ca. 5 mal (Dimension des Systems-1). Beachte: Die Dimension des Systems ist $(n+1) \times (n+1)$.

Satz 6.1

Das Tridiagonalverfahren für $\mathbf{A} = \text{tridiag}(a_i, b_i, c_i)$, $a_0 := 0 =: c_n$, ist durchführbar (d.h. $\exists A^{-1}$), falls

$$(1) \quad b_0 \neq 0$$

$$(2) \quad |b_i| \geq |a_i| + |c_i|, \quad i = 0, \dots, n$$

$$(3) \quad |b_i| > |a_i|, \quad i = 1, \dots, n$$

Die Forderung (3) kann ersetzt werden durch

$$(3') \quad |b_0| > |c_0| > 0 \text{ und } |a_i| > 0 \text{ für } i = 1, 2, \dots, n$$

oder

$$(3'') \quad \exists j \text{ sodaß: (2) mit "}" gilt, } |c_i| > 0, i = 1, \dots, j-1, \quad |a_i| > 0 \forall i > j.$$

Bemerkung: Das Beispiel $A = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix}$ zeigt, daß die Bedingungen für die Invertierbarkeit von A nur hinreichend sind. (3), (3'), (3'') versagen.

Beweis zunächst unter Benutzung von (3):

Zeige: Alle Nenner sind $\neq 0$ im Rahmen eines Induktionsbeweises für die Aussage $|\alpha_i| \leq 1 \quad \forall i$

$$\alpha_1 = \frac{c_0}{b_0} \stackrel{(1)(2)}{\implies} |\alpha_1| \leq 1.$$

Sei also $|\alpha_i| \leq 1$, dann gilt

$$(6.7) \quad |b_i - \alpha_i a_i| \geq |b_i| - |\alpha_i| |a_i| \geq |b_i| - |a_i| \stackrel{(3)}{\underset{\uparrow}{>}} 0, \quad i \geq 1$$

Damit sind alle $\alpha_{i+1}, \beta_{i+1}$ aus (II) und (III) erklärt und es gilt

$$|\alpha_{i+1}| = \frac{|c_i|}{|b_i - \alpha_i a_i|} \leq \frac{|c_i|}{|b_i| - |a_i|} \stackrel{(2)}{\leq} 1.$$

Das Verfahren ist durchführbar.

Ersetzt man (3) durch (3'), so ist $|\alpha_1| = \left| \frac{c_0}{b_0} \right| < 1$ und man erhält induktiv $|\alpha_i| < 1 \forall i$, denn mit $|\alpha_i| < 1$, für ein $i \leq n$, $|a_j| > 0$ für $j = 1, \dots, n$, erhält man im Induktionsschritt

$$|b_i - \alpha_i a_i| \geq |b_i| - |\alpha_i| |a_i| \stackrel{(3')}{\underset{\uparrow}{>}} |b_i| - |a_i| \stackrel{(2)}{\geq} |c_i| \geq 0 \implies |\alpha_{i+1}| < 1 \quad \forall i \geq 1$$

und alle Nenner sind $\neq 0$.

Benutzt man (3''), so zeigt man $|\alpha_i| \leq 1$ (induktiv) für $i = 1, \dots, n$ gemäß

$$\alpha_1 = \frac{c_0}{b_0}, \quad 0 \leq |\alpha_1| \leq 1,$$

und im Induktionsschritt

$$\begin{aligned} |b_i - \alpha_i a_i| &\geq |b_i| - |\alpha_i| |a_i| \geq |b_i| - |a_i| \geq |c_i| \underset{\uparrow}{>} 0 \text{ für } i = 1, \dots, j-1 \\ &\implies |\alpha_{i+1}| \leq 1, \quad i \leq j. \end{aligned}$$

Für das j gilt:

$$|b_j - \alpha_j a_j| \geq |b_j| - |\alpha_j| |a_j| \geq |b_j| - |a_j| \underset{\uparrow}{>} |c_j| \geq 0 \implies |\alpha_{j+1}| < 1$$

und mit $|\alpha_{j+1}| < 1, |a_{j+1}| > 0$ (vgl. 3'))

$$|b_{j+1} - \alpha_{j+1} a_j| \geq |b_{j+1}| - |\alpha_{j+1}| |a_{j+1}| \underset{\uparrow}{>} |b_{j+1}| - |a_{j+1}| \stackrel{(2)}{\geq} |c_j| \geq 0 \implies |\alpha_{j+2}| < 1$$

usw. ■

Bemerkung zu Satz 6.1

In der Praxis tritt üblicherweise folgender Fall auf:

$$b_i > 0 \quad i = 0, \dots, n$$

$$a_i < 0, \quad i = 1, \dots, n$$

$$c_i < 0, \quad i = 0, \dots, n-1$$

$$b_i + a_i + c_i \geq 0 \quad \forall i \quad (b_0 = c_n = 0)$$

> 0 für mindestens ein i .

Hierfür ist das Tridiagonalverfahren durchführbar (vgl. (3')).

Wir zeigen nun, daß man für spezielle Matrizen (M-Matrizen), die bei der Diskretisierung parabolischer Probleme entstehen, Invertierbarkeit und Normabschätzungen (in der Maximumnorm!) ohne die Kenntnis der Eigenwerte erhalten kann.

Definition 6.2

$A \in \mathbb{R}^{n \times n}$ heißt *M-Matrix* wenn

$$a_{ij} \leq 0, \text{ für } i \neq j \text{ und}$$

$$\exists \mathbf{p} \in \mathbb{R}^n, \mathbf{p} > 0 \text{ (komponentenweise), sodaß } \mathbf{Ap} > \mathbf{0} \text{ (komponentenweise).}$$

Satz 6.3

Ist $\mathbf{A} \in \mathbb{R}^{n \times n}$ eine *M-Matrix*, so gilt

a) $\exists \mathbf{A}^{-1}$ und $\mathbf{A}^{-1} \geq \mathbf{0}$ (elementweise)

b) $\|\mathbf{A}^{-1}\|_\infty \leq \frac{\|\mathbf{p}\|_\infty}{\min_i (\mathbf{Ap})_i}$, \mathbf{p} gemäß Definition 6.2 $\left(\|\mathbf{A}\|_\infty := \max_{i=1, \dots, n} \sum_{k=1}^n |a_{ik}| \right)$.

c) Ist $\mathbf{Ax} = \mathbf{b}$, so gilt schärfer $\|\mathbf{x}\|_\infty \leq \frac{\|\mathbf{p}\|_\infty}{\min_{i: b_i \neq 0} (\mathbf{Ap})_i} \|\mathbf{b}\|_\infty$

Beweis a): Wir zeigen zeigen $\exists \mathbf{D}^{-1}$, zerlegen

$$(6.8) \quad \mathbf{A} = \mathbf{D} - \mathbf{B} = \mathbf{D}(\mathbf{I} - \underbrace{\mathbf{D}^{-1}\mathbf{B}}_{\mathbf{C}}) =: \mathbf{D}(\mathbf{I} - \mathbf{C})$$

und zeigen: $\exists (\mathbf{I} - \mathbf{C})^{-1}$. Hieraus folgt die Behauptung a).

Wegen $\mathbf{p} > 0$, $\mathbf{Ap} > 0$, $a_{ij} \leq 0, i \neq j$ folgt $a_{ii} > 0 \forall i$, sonst $W!$, denn

$$(\mathbf{Ap})_i = \sum_{j=1}^{i-1} p_j a_{i,j} + p_{ii} a_{i,i} + \sum_{j=i+1}^n p_j p_j a_{i,j} > 0, \quad a_0 = a_n = 0 \xrightarrow{a_{ij} \leq 0, i \neq j} a_{ii} > 0.$$

$\implies \mathbf{D} = \text{diag}(\mathbf{A}) > 0$ (elementweise) ist invertierbar.

Wegen $\mathbf{B} \geq \mathbf{0}$, $\mathbf{D}^{-1} > \mathbf{0}$ (jeweils elementweise), folgt $\mathbf{C} := \mathbf{D}^{-1}\mathbf{B} \geq \mathbf{0}$ (elementweise).

Multiplikation von $\mathbf{0} < \mathbf{A}\mathbf{p} = \mathbf{D}(\mathbf{I} - \mathbf{C})\mathbf{p}$, (vgl. (6.8)), mit \mathbf{D}^{-1} von links liefert

$$\mathbf{0} < \mathbf{D}^{-1}\mathbf{A}\mathbf{p} = \mathbf{p} - \mathbf{C}\mathbf{p} \quad \text{und mit} \quad \mathbf{C}\mathbf{p} \geq \mathbf{0}$$

$$(6.9) \quad \mathbf{0} \leq \mathbf{C}\mathbf{p} < \mathbf{p}.$$

Wegen $\mathbf{p} > \mathbf{0}$ ist $\mathbf{P} := \text{diag}(p_i)$ regulär und $\exists \mathbf{P}^{-1} \geq \mathbf{0}$ (elementweise).

$$\begin{aligned} \implies \quad \|\mathbf{x}\|_p &:= \|\mathbf{P}^{-1}\mathbf{x}\|_\infty && \text{ist eine Vektornorm,} \\ \|\mathbf{A}\|_p &:= \|\mathbf{P}^{-1}\mathbf{A}\mathbf{P}\|_\infty && \text{ist die zugeordnete Matrixnorm.} \end{aligned}$$

Wir zeigen

$$\|\mathbf{C}\|_p < 1.$$

Mit $\mathbf{e} = (1, 1, \dots, 1)^T > \mathbf{0}$ gilt $\mathbf{P}\mathbf{e} = \mathbf{p}$

$$\stackrel{(6.9)}{\implies} \quad \mathbf{0} \leq \mathbf{C}\mathbf{P}\mathbf{e} < \mathbf{P}\mathbf{e}.$$

Multiplikation mit \mathbf{P}^{-1} von links ergibt

$$\mathbf{0} \leq (\mathbf{P}^{-1}\mathbf{C}\mathbf{P})\mathbf{e} < \mathbf{e} \quad (\text{komponentenweise}),$$

d.h. in jeder Zeile der Matrix $\mathbf{P}^{-1}\mathbf{C}\mathbf{P}$ ist die Betragssumme der Elemente < 1 , d.h.

$$\|\mathbf{C}\|_p = \|\mathbf{P}^{-1}\mathbf{C}\mathbf{P}\|_\infty < 1 \quad \stackrel{\text{Satz 3.1}}{\implies} \quad \exists (\mathbf{I} - \mathbf{C})^{-1} = \sum_{\nu=0}^{\infty} \mathbf{C}^\nu \geq \mathbf{0} \quad (\text{elementweise})$$

$$\stackrel{(6.8)}{\implies} \quad \exists \mathbf{A}^{-1} = \underbrace{(\mathbf{I} - \mathbf{C})^{-1}}_{\geq \mathbf{0}} \underbrace{\mathbf{D}^{-1}}_{\geq \mathbf{0}} \geq \mathbf{0} \quad (\text{elementweise}).$$

Beweis c)

Für $\mathbf{A} \in \mathbb{R}^{n \times n}$ und $\mathbf{A}^{-1} \geq \mathbf{0}$ (elementweise) suchen wir für $\mathbf{A}\mathbf{x} = \mathbf{b}$ eine Abschätzung

$$\|\mathbf{x}\|_\infty \leq k \|\mathbf{b}\|_\infty$$

Bemerkung: Ist k unabhängig von \mathbf{b} , so gilt laut Definition der Matrixnorm

$$\|\mathbf{A}^{-1}\|_\infty \leq k,$$

denn

$$\|\mathbf{x}\|_\infty = \|\mathbf{A}^{-1}\mathbf{b}\|_\infty \leq \|\mathbf{A}^{-1}\|_\infty \|\mathbf{b}\|_\infty \quad \text{und} \quad \|\mathbf{A}^{-1}\|_\infty = \inf_{\mathbf{b} \neq \mathbf{0}} \{k; \|\mathbf{A}^{-1}\mathbf{b}\|_\infty \leq k \|\mathbf{b}\|_\infty \forall \mathbf{b} \in \mathbb{R}^n\}.$$

Ziel: Für einen Vektor $\mathbf{p} > \mathbf{0}$ mit $(\mathbf{A}\mathbf{p})_i > 0 \forall_i$ bestimmen wir ein $m > 0$ so, daß

$$(6.10) \quad \mathbf{A}(m\mathbf{p} \pm \mathbf{x}) = m\mathbf{A}\mathbf{p} \pm \mathbf{A}\mathbf{x} = m\mathbf{A}\mathbf{p} \pm \mathbf{b} \stackrel{!}{\geq} \mathbf{0}$$

Es genügt $m > 0$ so zu bestimmen, daß (6.10) für die Komponenten i mit $b_i \neq 0$ erfüllt ist, denn für die anderen Komponenten gilt (6.10) ohnehin. Fordere also

$$m(\mathbf{A}\mathbf{p})_i \geq \pm b_i \quad \forall_i \quad \text{mit } b_i \neq 0.$$

Der „schlimmste“ Fall liegt vor, wenn links das Minimum, rechts das Maximum angenommen wird. Wir definieren deshalb m durch die Forderung

$$m \min_{i: b_i \neq 0} (\mathbf{A}\mathbf{p})_i = \|\mathbf{b}\|_\infty \quad \text{bzw.} \quad m := \frac{\|\mathbf{b}\|_\infty}{\min_{i: b_i \neq 0} (\mathbf{A}\mathbf{p})_i}.$$

Dann gilt für alle Komponenten $m\mathbf{A}\mathbf{p} \geq \pm \mathbf{b}$ und wegen $\mathbf{A}^{-1} \geq \mathbf{0}$, $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$, folgt

$$m\mathbf{p} \geq \pm \mathbf{x} \quad \text{bzw. wegen } \mathbf{p} > \mathbf{0} \quad |x_i| \leq m p_i$$

$$\Rightarrow \|\mathbf{x}\|_\infty \leq m \|\mathbf{p}\|_\infty = \frac{\|\mathbf{p}\|_\infty}{\min_{i: b_i \neq 0} (\mathbf{A}\mathbf{p})_i} \|\mathbf{b}\|_\infty, \quad \text{also Behauptung c).}$$

Beweis b)

Wird insbesondere m unabhängig von \mathbf{b} bestimmt, d.h. unabhängig von den Nullkomponenten von \mathbf{b} , so folgt aus der letzten Abschätzung

$$\|\mathbf{x}\|_\infty \leq \frac{\|\mathbf{p}\|_\infty}{\min_{i=1, \dots, n} (\mathbf{A}\mathbf{p})_i} \|\mathbf{b}\| \quad \forall \mathbf{b}, \quad \text{also Behauptung b).}$$



Im Anschluß an Satz 6.1 und der Bemerkung von Satz 6.1 zeigen wir

Satz 6.4

Für $\mathbf{A} = \text{tridiag}(-a_i, b_i, -c_i)$ sei

$$a_i, b_i, c_i > 0 \quad \forall_i \quad \text{außer } a_0 = c_n = 0$$

$$b_i \geq a_i + c_i \quad \forall_i$$

$$> \quad \text{für mindestens ein } i$$

$$\Rightarrow \quad \mathbf{A} \text{ ist eine } M\text{-Matrix.}$$

Bemerkungen:

1. Solche Matrizen entstehen bei der Diskretisierung eines parabolischen Problems.
2. Satz 6.1 zeigt, daß für diese Matrizen das Tridiagonalverfahren durchführbar ist, daß also die Werte auf der neuen Zeitschicht berechenbar sind.

3. Satz 6.4 liefert dann Abschätzungen für die Lösung von $\mathbf{Ax} = \mathbf{b}$.
4. „ \Leftarrow “ gilt nicht. Beispiel: $\mathbf{A} = \mathbf{I}$ ist M -Matrix ($\mathbf{p} = \mathbf{e} = (1, \dots, 1)^T$) und $\mathbf{A} = \begin{pmatrix} 1 & -2 \\ 0 & 1 \end{pmatrix}$ mit $\mathbf{p} = (1, 0.25)^T$ ebenfalls.

Beweis:

1. Gemäß der Definition der M -Matrizen (vgl. Definition 6.2) ist ein Vektor $\mathbf{p} > 0$ (komponentenweise) zu finden mit $\mathbf{Ap} > 0$.

Mit $\mathbf{e} = (1, \dots, 1)^T$ gilt $\mathbf{Ae} \geq 0$ (komponentenweise, da nach Voraussetzung $b_i \geq a_i + c_i$).

Für alle $\varepsilon > 0$ gilt deshalb

$$(6.11) \quad (\varepsilon \mathbf{I} + \mathbf{A})\mathbf{e} \geq \varepsilon \mathbf{e} > \mathbf{0}, \quad \varepsilon \mathbf{I} + \mathbf{A} \text{ ist eine } M\text{-Matrix mit } \mathbf{p} = \mathbf{e} \xrightarrow{\text{Satz 6.3}} \\ \exists (\varepsilon \mathbf{I} + \mathbf{A})^{-1} \geq \mathbf{0} \quad (\text{elementweise})$$

Satz 6.1 mit (3') zeigt $\exists \mathbf{A}^{-1}$.

Deshalb kann man im (6.11) den Grenzübergang $\varepsilon \rightarrow 0$ machen. (Die Existenz von \mathbf{A}^{-1} bleibt ja erhalten) \implies

$$(6.12) \quad \mathbf{A}^{-1} \geq \mathbf{0} \quad (\text{elementweise})$$

2. Wir konstruieren nun das gesuchte $\mathbf{p} > 0$. Für $\mathbf{e} = (1, \dots, 1)^T$ existiert (vgl. (6.12)), die Lösung von

$$\mathbf{Ay} = \mathbf{e} \xrightarrow{(6.12)} \mathbf{y} = \mathbf{A}^{-1}\mathbf{e} \geq \mathbf{0} \quad (\text{elementweise}).$$

Sei $\mathbf{A}^{-1} = (\alpha_{ij})$. Annahme: $\exists i : y_i = 0$ d.h. $\sum_{j=1}^n \alpha_{ij} \cdot 1 = 0$

$\xrightarrow{(6.12)} \alpha_{ij} = 0 \quad \forall j$, also existiert \mathbf{A}^{-1} nicht W!

Also kann $\mathbf{p} = \mathbf{y}$ gewählt werden. ■

Damit erhalten wir, insbesondere aus Beweisteil 2) die Charakterisierung:

Satz 6.5

Für $\mathbf{A} \in \mathbb{R}^{n \times n}$ gilt

$$\mathbf{A} \text{ ist } M\text{-Matrix} \iff \exists \mathbf{A}^{-1} \geq \mathbf{0} \quad (\text{elementweise}).$$

Beweis: „ \implies “ liefert sofort Satz 6.3 a).

„ \Leftarrow “ liefert Beweisteil 2) des vorigen Satzes, der keinen Gebrauch von der Tridiagonalform machte. ■

§ 7 Die Gleichung $u_t = \frac{\partial}{\partial x} (k(x) \frac{\partial u}{\partial x}) + f$

In der Anwendung ist k oft stückweise stetig. Wie kann $\frac{\partial}{\partial x} (k(x) \frac{\partial u}{\partial x}) =: R$ geeignet diskretisiert werden? Physikalisch stellt dieser Ausdruck eines Wärmestrom dar, wenn u die Temperatur ist. Deshalb sollte er bei der Diskretisierung nicht durch Ausdifferenzieren ($R = ku_{\bar{x}x} + k_{x^0}u_{x^0}$, wobei etwa $y_{x^0} = \frac{y_{i+1} - y_{i-1}}{2} = \frac{1}{2}(y_x + y_{\bar{x}})$ – zentral wegen besserer Approximationsordnung –) in nicht interpretierbare Summanden zerlegt werden. Würde man dies trotzdem tun, so ergäbe sich für R eine nicht symmetrische Matrix. Frau(Man) kann das leicht nachrechnen. Dies ist auch aus mathematischen Gründen ungünstig, da der Operator $Lu = \frac{\partial}{\partial x} (k(x) \frac{\partial u}{\partial x})$ bei geeigneten Randbedingungen selbstadjungiert ist, wie eine kurzer Rechnung zeigt.

Selbstadjungiert (das ist für Operatoren die Verallgemeinerung des Begriffs symmetrisch, bzw. hermite'sch bei Matrizen) bedeutet (z.B. für $x \in [0, 1]$)

$$(Lu, v) = (u, Lv), \quad (u, v) = \int_0^1 u(x)v(x)dx;$$

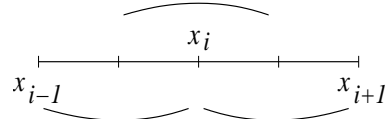
damit

$$\begin{aligned} (Lu, v) &= \int_0^1 (ku')'v \, dx = [ku'v]_0^1 - \int_0^1 ku'v' \, dx \\ &= \underbrace{[ku'v]_0^1 - [kv'u]_0^1}_{=0 \text{ bei entsprechenden Randwerten}} + \int_0^1 (kv')'v \, dx = \int_0^1 u(x)u(x)dx \end{aligned}$$

In der Tat zeigen auch numerische Beispiele, daß eine nichtsymmetrische Diskretisierung Konvergenz des Verfahrens gegen falsche Werte liefern kann. Deshalb wird der Wärmestrom als Ganzes diskretisiert. Da wir beim Wärmestrom nur mit Ortsableitungen zu tun haben, unterdrücken wir in der Bezeichnung die Ortsabhängigkeit. Dabei benutzen wir, so weit als möglich, zentrale Differenzenquotienten der besseren Approximationseigenschaften wegen.

Mit $y_i \approx u(x_i)$ approximieren wir

$$\begin{aligned} \left(k \frac{\partial u}{\partial x}\right)_{i+\frac{1}{2}} &\approx k_{i+\frac{1}{2}} \frac{y_{i+1} - y_i}{h} = k_{i+\frac{1}{2}} y_{x,i}, \\ (ku')'_i &\approx \frac{1}{h} (k_{i+\frac{1}{2}} y_{x,i} - k_{i-\frac{1}{2}} y_{\bar{x},i}) =: (k \cdot y_x)_{\bar{x},i} \quad (\text{abkürzende Bezeichnung}) \end{aligned}$$



$$(7.1) \quad (k \cdot y_x)_{\bar{x},i} = \frac{1}{h} \left(k_{i+\frac{1}{2}} \frac{y_{i+1} - y_i}{h} - k_{i-\frac{1}{2}} \frac{y_i - y_{i-1}}{2} \right)$$

$$(7.2) \quad = \frac{1}{h^2} \left(k_{i+\frac{1}{2}} y_{i+1} - \frac{1}{h^2} \left(k_{i+\frac{1}{2}} + k_{i-\frac{1}{2}} \right) y_i + \frac{1}{h^2} k_{i-\frac{1}{2}} y_{i-1} \right) =: -(\mathbf{A}_h(k) \mathbf{y})_i$$

Als Diskretisierungsmatrix für $A_h(k)$ erhalten wir – nachdem die Randwerte aus der Matrix eliminiert wurden (d.h. “der rechten Seite zugeschlagen”) – eine tridiagonale, symmetrische Matrix

$$A_h(k) = \frac{1}{h^2} \begin{pmatrix} k_{\frac{1}{2}} + k_{\frac{3}{2}} & -k_{\frac{3}{2}} & & & \\ -k_{\frac{3}{2}} & k_{\frac{3}{2}} + k_{\frac{5}{2}} & -k_{\frac{5}{2}} & & \\ & & \ddots & \ddots & \\ & & & \ddots & \ddots \\ & & & & \ddots & \ddots \end{pmatrix}$$

Wir untersuchen

1. die Approximationseigenschaften, d.h. den Diskretisierungsfehler der Diskretisierung,
2. die Stabilität des resultierenden Verfahrens,
3. die Konvergenzeigenschaften.

Diskretisierungsfehler

Für $u \in C^4$, $k \in C^3$ (bzgl. Ort) untersuchen wir die Diskretisierung von (7.1).

$$u_{i+\frac{1}{2} \pm \frac{1}{2}} = u_{i+\frac{1}{2}} \pm \frac{h}{2} u'_{i+\frac{1}{2}} + \frac{h^2}{8} u''_{i+\frac{1}{2}} \pm \frac{h^3}{48} u'''_{i+\frac{1}{2}} + \frac{1}{4!} \left(\frac{h}{2}\right)^4 u''''_{\pm}$$

liefert

$$u_{x,i} = \frac{u_{i+1} - u_i}{h} = u'_{i+\frac{1}{2}} + \frac{h^2}{24} u'''_{i+\frac{1}{2}} + O(h^3).$$

Damit erhält man

$$\begin{aligned} k_{i+\frac{1}{2}} u_{x,i} &= (ku')_{i+\frac{1}{2}} + \frac{h^2}{24} (ku''')_{i+\frac{1}{2}} + O(h^3), \\ k_{i-\frac{1}{2}} u_{\bar{x},i} &= (ku')_{i-\frac{1}{2}} + \frac{h^2}{24} (ku''')_{i-\frac{1}{2}} + O(h^3). \end{aligned}$$

Damit folgt (vgl. die Bezeichnung in (7.1))

$$\begin{aligned} (ku_x)_{\bar{x},i} &= \frac{1}{h} \left(k_{i+\frac{1}{2}} u_{x,i} - k_{i-\frac{1}{2}} u_{\bar{x},i} \right) \\ &= \frac{1}{h} \left[(ku')_{i+\frac{1}{2}} - (ku')_{i-\frac{1}{2}} + \underbrace{\frac{h}{24} \left\{ (ku''')_{i+\frac{1}{2}} - (ku''')_{i-\frac{1}{2}} \right\}}_{\frac{h^2}{24} (ku''')'_z = O(h^2)} \right] + O(h^2). \end{aligned}$$

Die ersten beiden Summanden der rechten Seite müssen an der Stelle x_i entwickelt werden.

$$(ku')_{i \pm \frac{1}{2}} = (ku')_i \pm \frac{h}{2} (ku')'_i + \frac{1}{2!} \left(\frac{h}{2}\right)^2 (ku')''_i \pm \frac{h^3}{48} (ku')'''_{z \pm}.$$

Damit folgt

$$(ku_x)_{\bar{x},i} = (ku')'_i + \underbrace{\frac{h^2}{48} [(ku')'''_{z_+} - (ku')'''_{z_-}]}_{O(h^2)} + O(h^2) \quad \text{bzw.}$$

$$(7.3) \quad (ku_x)_{\bar{x},i} = (ku')'_i + O(h^2) \quad \text{für } u \in C^2, k \in C^3 \text{ (bzgl. Ort)}$$

Wir untersuchen nun gleich **das gewichtete Verfahren**

$$(7.4) \quad \mathbf{y}_t = (k\mathbf{y}_x)_{\bar{x}}^{(\sigma)} + \boldsymbol{\varphi}, \quad \mathbf{y}^{(\sigma)} = \sigma \hat{\mathbf{y}} + (1 - \sigma)\mathbf{y} = \sigma\tau\mathbf{y}_t + \mathbf{y}$$

Beachte: $k = k(x)$ ist zeitlich konstant, deshalb kann man numerische Differentiation bzgl. x und Mittelwertbildung mittels σ vertauschen. Auf Grund der Matrixdarstellung (7.2) folgt aus der Linearität in \mathbf{y} aus (7.4)

$$\begin{aligned} \mathbf{y}_t &= -\mathbf{A}_h(k)[\sigma \hat{\mathbf{y}} + (1 - \sigma)\mathbf{y}] + \boldsymbol{\varphi} \\ &= -\mathbf{A}_h(k)[\sigma\tau\mathbf{y}_t + \mathbf{y}] + \boldsymbol{\varphi} \quad \text{oder} \end{aligned}$$

$$(7.5) \quad (\mathbf{I} + \sigma\tau\mathbf{A}_h(k))\mathbf{y}_t + \mathbf{A}_h(k)\mathbf{y} = \boldsymbol{\varphi}.$$

Dies ist eine Gestalt, die unter die Verfahrensklasse

$$\mathbf{B}\mathbf{y}_t + \mathbf{A}\mathbf{y} = \boldsymbol{\varphi}, \quad \mathbf{B} = \mathbf{I} + \sigma\tau\mathbf{A}, \quad \text{mit } \mathbf{A} = \mathbf{A}_h(k)$$

fällt, für die Stabilitätssätze gelten. Wir müssen die Voraussetzungen prüfen. \mathbf{A} symmetrisch ist klar (vgl. (7.2)). Wir zeigen zunächst

$\mathbf{A}_h(k)$ ist positiv definit. (Beweis analog zu (3.13))

$$\begin{aligned} (\mathbf{A}_h(k)\mathbf{y}, \mathbf{y})_{(0,h)} &= \sum_{i=1}^{N-1} (\mathbf{A}_h(k)\mathbf{y})_i y_i h \stackrel{(7.1)}{=} - \sum_{i=1}^{N-1} (ky_x)_{\bar{x},i} y_i h \\ &\stackrel{(7.1)}{=} - \sum_{i=1}^{N-1} k_{i+\frac{1}{2}} y_{x,i} y_i + \sum_{i=1}^{N-1} k_{i-\frac{1}{2}} y_{x,i-1} y_i \quad \text{und mit } y_0 = y_N = 0 \\ &= - \sum_{i=0}^{N-1} k_{i+\frac{1}{2}} y_{x,i} y_i + \sum_{i=1}^N k_{i-\frac{1}{2}} y_{x,i-1} y_i \\ &= - \sum_{i=1}^N k_{i-\frac{1}{2}} y_{x,i-1} y_{i-1} + \sum_{i=1}^N k_{i-\frac{1}{2}} y_{x,i-1} y_i \\ &= \sum_{i=1}^N k_{i-\frac{1}{2}} y_{x,i-1} (y_i - y_{i-1}) \\ (7.6) \quad &= \sum_{i=1}^N k_{i-\frac{1}{2}} (y_{x,i-1})^2 h = \sum_{i=0}^{N-1} k_{i+\frac{1}{2}} (y_{x,i})^2 h > 0 \\ &= 0 \text{ nur für } \mathbf{y} \equiv 0, \text{ da } y_0 = y_N = 0 \end{aligned}$$

Unter der Voraussetzung $0 < c_0 \leq k(x) \leq c_1 \quad \forall x \in [0, 1]$ folgt somit aus (7.6) (vgl. (3.8))

$$(7.7) \quad \mathbf{A}_h(k) > 0 \text{ und } (\mathbf{A}_h(k)\mathbf{y}, \mathbf{y})_{(0,h)} \geq c_0 \|\mathbf{y}\|_{(1,h)}^2 =: c_0 \|\mathbf{y}\|_{\mathbf{A}_h^0}^2$$

Wir können also eine Vektornorm definieren durch

$$(7.8) \quad \|\mathbf{y}\|_{\mathbf{A}_h(k)}^2 = \|\mathbf{y}\|_{(1,h,k)}^2 = (\mathbf{A}_h(k)\mathbf{y}, \mathbf{y})_{(0,h)} := \sum_{i=0}^{N-1} k_{j+\frac{1}{2}}(y_{x,i})^2 h, \quad y_0 = y_N = 0.$$

Abschätzungen von $\mathbf{A}_h(k)$ nach oben und unten (im Sinne $(\cdot, \cdot)_{(0,h)}$).

Unter Benutzung von (3.23): $\|\mathbf{y}\|_{(1,h)}^2 \geq 8\|\mathbf{y}\|_{(0,h)}^2$ folgt aus (7.7)

$$(\mathbf{A}_h(k)\mathbf{y}, \mathbf{y})_{(0,h)} \geq c_0 \|\mathbf{y}\|_{(1,h)}^2 \geq 8c_0 \|\mathbf{y}\|_{(0,h)}^2,$$

also

$$(7.9) \quad \mathbf{A}_h(k) \geq 8c_0 \mathbf{I} \quad (\text{im Sinne positiv semidefinit}).$$

Abschätzung nach oben

$$(\mathbf{A}_h(k)\mathbf{y}, \mathbf{y})_{(0,h)} \stackrel{(7.6)}{=} \sum_{i=0}^{N-1} k_{i+\frac{1}{2}}(y_{x,i})^2 h \leq c_1 \|\mathbf{y}\|_{(1,h)}^2 \stackrel{(3.23)}{\leq} \frac{4c_1}{h^2} \|\mathbf{y}\|_{(0,h)}^2 = \frac{4c_1}{h^2} (\mathbf{y}, \mathbf{y})_{(0,h)},$$

also insgesamt (im Sinne positiv semidefinit)

$$(7.10) \quad 8c_0 \mathbf{I} \leq \mathbf{A}_h(k) \leq \frac{4c_1}{h^2} \mathbf{I}$$

Stabilität:

Die allgemeine Stabilitätsbedingung für das Verfahren 7.5 lautet gemäß Satz 4.2 b)

$$\begin{aligned} \mathbf{B} = \mathbf{I} + \sigma\tau \mathbf{A}_h(k) &\geq \frac{\tau}{2} \mathbf{A}_h(k) + \frac{\varepsilon}{2} \mathbf{I} \quad \text{bzw.} \\ (1 - \frac{\varepsilon}{2}) \mathbf{I} + \sigma\tau \mathbf{A}_h(k) - \frac{\tau}{2} \mathbf{A}_h(k) &\geq \mathbf{0} \end{aligned}$$

Mit Hilfe von 7.10 verschärfen wir sie (mit $\varepsilon \leq 2$) zu

$$(1 - \frac{\varepsilon}{2}) \frac{h^2}{4c_1} \mathbf{A}_h(k) + \sigma\tau \mathbf{A}_h(k) - \frac{\tau}{2} \mathbf{A}_h(k) \geq \mathbf{0}$$

bzw.

$$\left[(1 - \frac{\varepsilon}{2}) \frac{h^2}{4c_1} + \tau(\sigma - \frac{1}{2}) \right] \mathbf{A}_h(k) \geq \mathbf{0},$$

Wegen $\mathbf{A}_h(k) \geq \mathbf{0}$ ist dies erfüllt, wenn

$$(1 - \frac{\varepsilon}{2}) \frac{h^2}{4c_1} + \tau(\sigma - \frac{1}{2}) \geq 0 \quad \text{bzw.}$$

$$(7.11) \quad \sigma \geq \frac{1}{2} - \frac{h^2}{4c_1\tau} \left(1 - \frac{\varepsilon}{2}\right) \quad \text{vgl. Folgerung 4.3 b)}$$

Für das explizite Verfahren ($\sigma = 0$) bedeutet dies

$$\tau \leq \frac{h^2}{2c_1} \left(1 - \frac{\varepsilon}{2}\right).$$

In der Praxis kann c_1 sehr groß sein \implies fatale Auswirkungen auf $\tau \implies$.

Dringende Empfehlung: Kein explizites Verfahren bei variablen Koeffizienten $k(x)$.

Konvergenz:

Analog zum Vorgehen in Definition 5.3 sei

$$\mathbf{z} = \mathbf{u} - \mathbf{y} \quad (\mathbf{u} \text{ die Gitterfunktion zur exakten Lösung } u)$$

Anfangs- und Randwerte sind = Null : $\mathbf{z}_0 = 0$, $z_0^j = z_N^j = 0 \forall j \geq 0$.

Für den Approximationsfehler $\boldsymbol{\psi} = \mathbf{u}_t - (k\mathbf{u}_x)_{\bar{x}}^{(\sigma)} - \boldsymbol{\varphi}$ zeigt man durch Taylorentwicklung wie früher

$$\boldsymbol{\varphi} = O\left(\tau\left(\sigma - \frac{1}{2}\right) + \tau^2 + h^2\right)$$

Der Anteil $\tau\left(\sigma - \frac{1}{2}\right) + \tau^2$ kommt aus der Zeitdiskretisierung, ist also unverändert gegenüber dem einfachen Verfahren. Der „Anteil h^2 “ wurde in (7.3) hergeleitet. Unter Benutzung der Vektornorm (7.8): $\|\mathbf{y}\|_{\mathbf{A}_h(k)} = \|\mathbf{y}\|_{(1,h,k)} \leq \sqrt{c_1} \|\mathbf{y}\|_{(1,h)}$ gilt analog zu Satz 5.4 somit die

Konvergenzabschätzung

Für $T \geq j\tau$ gilt mit $\sigma \geq \frac{1}{2} - \frac{h^2}{4c_1\tau} \left(1 - \frac{\varepsilon}{2}\right)$

$$\frac{1}{\sqrt{c_1}} \|\mathbf{z}^j\|_{\mathbf{A}_h(k)} \leq \|\mathbf{z}^j\|_{(1,h)} \leq \|\mathbf{z}^j\|_{(1,h)} = \sqrt{\frac{T}{\varepsilon}} O\left(\tau\left(\sigma - \frac{1}{2}\right) + \tau^2 + h^2\right),$$

und damit auch

$$\|\mathbf{z}^j\|_{\mathbf{A}_h(k)} = \sqrt{\frac{T}{\varepsilon}} O\left(\tau\left(\sigma - \frac{1}{2}\right) + \tau^2 + h^2\right).$$

Beachte: $\theta < \varepsilon \leq 2$ und $\sigma = \frac{1}{2}$ sind möglich.

§ 8 Die allgemeine 1-dimensionale Wärmeleitungsgleichung

$$(8.1) \quad c\rho_p \frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left(k \frac{\partial u}{\partial x} \right) - v \frac{\partial u}{\partial x} - qu + f$$

mit variablen Koeffizienten $c, \rho_p, k, v, d > 0$

Vorbemerkung und Überblick:

Man kann Stabilitätsbedingungen und Konvergenz auch zeigen, wenn die Diskretisierungsmatrix \mathbf{A} der rechten Seite von (8.1), abgesehen von f , eine M -Matrix ist (vgl. dazu den nächsten Paragraphen). Vorteil: Mit der M -Matrizentheorie vermeidet man Oszillationen, die bei der L_2 -Theorie (Skalarprodukttheorie) auftauchen.

Nachteil: Man erhält oft schlechtere Konvergenzordnungen, muß gelegentlich empfindliche Einschränkungen bzgl. der Schrittweiten hinnehmen, und die M -Matrizentheorie ist nicht immer anwendbar.

Wir führen in diesem Paragraphen die L_2 -Theorie weiter, geben jedoch gelegentlich schon Hinweise auf die M -Matrizentheorie, die wir im nächsten Paragraphen untersuchen.

Wir betrachten zunächst nur eine x -Abhängigkeit der Koeffizienten. Bei der Diskretisierung des Konvektionsanteils $v \frac{\partial u}{\partial x}$, wird sich zeigen, dass wir bei $v \neq 0$ die Symmetrie der Diskretisierungsmatrix verlieren. Der Stabilitätssatz (4.2) muß abgeändert werden, liefert dann aber nur noch die Stabilität bzgl. der Anfangswerte.

Bei zeitabhängigen Koeffizienten beschränken wir uns auf die Untersuchung des Diffusionsanteils $\frac{\partial}{\partial x} \left(k \frac{\partial u}{\partial x} \right)$. Wir stellen fest, dass die Diskretisierungsmatrix nun zeitabhängig wird. Deshalb ist der Stabilitätssatz (4.2) nicht mehr anwendbar. Ein neuer Stabilitätssatz muß bewiesen werden, der allerdings auch nur die Stabilität bzgl. der Anfangswerte liefert. Man rufe sich ins Gedächtnis zurück, dass der Konvergenzssatz (5.4) wesentlich auf der Stabilität bzgl. der rechten Seite beruht. Wir beweisen deshalb zum Abschluß des Paragraphen, dass man die Stabilität bzgl. der rechten Seite aus der Stabilität bzgl. der Anfangswerte folgern kann.

Wir behandeln im Folgenden die einzelnen Summanden der Differentialgleichung (8.1) getrennt.

1) Die Abbaurrate q ($d > 0$) wird in \mathbf{A}_h integriert.

Sie liefert bei der Diskretisierung eine Diagonalmatrix $\mathbf{Q} = \text{diag}(q_i)$. Diese wird dem Diskretisierungsterm

$$\frac{\partial}{\partial x} \left(k(x) \frac{\partial u}{\partial x} \right) \approx -\mathbf{A}_h(k)\mathbf{y}$$

zugeschlagen. Man erhält dann die symmetrische Matrix

$$\mathbf{A} = \mathbf{A}_h(k) + \mathbf{Q}.$$

Dadurch wird die positive Definitheit von $\mathbf{A}_h(k)$ sogar gestärkt. Ist $\mathbf{A}_h(k)$ eine M -Matrix, so auch $\mathbf{A}_h(k) + \mathbf{Q}$. Kein Problem.

2) Der Term $c\rho_p =: \kappa \geq \delta_0 > 0 \quad \forall x, t$. (Wärmekapazität und Dichte)

Er liefert bei der Diskretisierung $\kappa u_t \approx \mathbf{K}\mathbf{I} \mathbf{y}_t$ mit einer Diagonalmatrix $\mathbf{K} = \text{diag}(\kappa_i)$. Im Stabilitätssatz 4.2 hat dann \mathbf{B} die Gestalt

$$\mathbf{B} = (\mathbf{K}\mathbf{I} + \sigma\tau\mathbf{A})$$

(der Abbauterm ist in \mathbf{A} enthalten, vgl. oben).

Aus $(\mathbf{A}\mathbf{x}, \mathbf{x}) > 0$ folgt also $(\mathbf{B}\mathbf{x}, \mathbf{x}) > 0$.

Die Stabilitätsbedingung (vgl. Satz 4.2) lautet (setze $\mathbf{A} := \mathbf{A}_h(k)$ bzw. $\mathbf{A}_h(k) + \mathbf{D}$)

$$\mathbf{B} := \mathbf{K}\mathbf{I} + \sigma\tau\mathbf{A} \stackrel{!}{\geq} \frac{\tau}{2}\mathbf{A} + \frac{\varepsilon}{2}\mathbf{I} \quad \text{für ein } \varepsilon > 0.$$

Mit der Bedingung $\kappa \geq \delta_0 > 0$ muß man fordern

$$\mathbf{B} \geq \delta_0\mathbf{I} + \sigma\tau\mathbf{A} \stackrel{!}{\geq} \frac{\tau}{2}\mathbf{A} + \frac{\varepsilon}{2}\mathbf{I}$$

und mit der Verschärfung $\mathbf{I} \geq \frac{\mathbf{A}}{\|\mathbf{A}\|}$ (vgl. (4.16) sogar

$$\left(\delta_0 - \frac{\varepsilon}{2}\right) \frac{\mathbf{A}}{\|\mathbf{A}\|} + \sigma\tau\mathbf{A} - \frac{\tau}{2}\mathbf{A} \geq 0,$$

was durch

$$\frac{\delta_0 - \frac{\varepsilon}{2}}{\|\mathbf{A}\|} + \sigma\tau - \frac{\tau}{2} \geq 0 \quad \text{bzw.}$$

$$(8.2) \quad \sigma \geq \frac{1}{2} - \frac{\delta_0 - \frac{\varepsilon}{2}}{\tau\|\mathbf{A}\|}$$

garantiert wird.

Beachte: $\sigma = \frac{1}{2}$ ist möglich.

Für das explizite Verfahren ($\sigma = 0$) bedeutet dies

$$\tau \leq \frac{2(\delta_0 - \frac{\varepsilon}{2})}{\|\mathbf{A}\|}$$

Dies ist fatal, wenn δ_0 klein ist, was realistisch ist.

⇒

**Kein explizites Verfahren für kleines $\kappa = c\rho_p$
Verfahren mit $\sigma \geq \frac{1}{2}$ (vgl. (8.2)) sind nicht berührt!**

setzt, so folgt für konstantes v (der Einfachheit halber)

$$\tilde{u}(x') = u(x), \quad \frac{\partial u}{\partial x} = \frac{\partial \tilde{u}}{\partial x'} \cdot \frac{1}{L}, \quad \frac{\partial^2 u}{\partial x^2} = \frac{\partial^2 \tilde{u}}{\partial x'^2} \cdot \frac{1}{L^2}$$

$$\frac{\partial \tilde{u}}{\partial t} = k \frac{\partial^2 \tilde{u}}{\partial x'^2} \cdot \frac{1}{L^2} + \frac{v}{L} \frac{\partial \tilde{u}}{\partial x'}$$

$$\text{d.h. es werden ersetzt } v \rightarrow \frac{v}{L} \quad \text{und} \quad k \rightarrow \frac{k}{L^2}$$

d.h. obige Bedingung geht über in

$$\frac{2k_0}{|v|} \rightarrow \frac{2k_0}{Lv} \quad \text{also} \quad \underline{\underline{h \leq \frac{2k_0}{Lv}}}$$

Bei Flußlängen von ca. 1000 km hat dies katastrophale Folgen, selbst bei großen und schnellen Rechnern, denn für Rechnungen mit so kleinen Schrittweiten ist das Modell zu ungenau (man bräuchte u.a. die Information (Anfangswerte) für die sehr kleinen Schrittweiten).

Noch problematischer wird es, falls Nebenflüsse eingeschlossen werden (Nebenrohre in der Industrie). Das h ist für die Praxis zu klein.

Abhilfe: Statt des zentralen Differenzenquotienten \mathbf{y}_x für $\frac{\partial u}{\partial x}$ könnte man einen einseitigen Differenzenquotienten $\mathbf{y}_{\bar{x}}$ (rückwärtsgenommen) wählen.

physikalische Begründung: Läuft die Strömung von links nach rechts ($v > 0$), so enthält die ankommende Strömung wertvollere Informationen, als die abfließende (stromaufwärts schauen).

mathematische Begründung: Bei einer einseitigen Approximation für die x -Ableitung gilt für die Matrix $(\mathbf{A}_h \mathbf{y})_i := -k y_{\bar{x},i} + (v y_{\bar{x}})_i$

$$\begin{aligned} (\mathbf{A}_h \mathbf{y})_i &= -\frac{k_0}{h^2}(y_{i+1} - 2y_i + y_{i-1}) + \frac{v_i}{2h}(y_i - y_{i-1}) \\ &= -\frac{k}{h^2} y_{i+1} + \left(\frac{2k}{h^2} + \frac{v}{h}\right) y_i - \left(\frac{k}{h^2} + \frac{v}{h}\right) y_{i-1} \end{aligned}$$

$$(8.4) \quad \mathbf{A}_h = \begin{pmatrix} \ddots & & & & \\ & -\frac{k}{h^2} - \frac{v}{h} & & & \\ & & \frac{2k}{h^2} + \frac{v}{h} & & \\ & & & \ddots & \\ & & & & -\frac{k}{h^2} & \ddots \end{pmatrix}$$

und man erkennt sofort:

Vorteil: Für $v \geq 0$ liegt eine M -Matrix vor (vgl. Satz (6.4)), auch für variable Koeffizienten, denn

1. Die Vorzeichenbedingung ist erfüllt
2. Die strenge Diagonaldominanz gilt in der ersten und letzten Zeile

Die Matrix ist also auch invertierbar.

Nachteil: $(\mathbf{A}_h \mathbf{y})_i \approx -L u + \underset{\substack{\uparrow \\ \text{von einseitiger} \\ \text{Approx.}}}{O(h)} \quad , \quad L u = -k \frac{\partial^2 u}{\partial x^2} + v \frac{\partial u}{\partial x}$

d.h., man hat Diskretisierungsordnung (Konvergenzgeschwindigkeit) für Oszillationsfreiheit geopfert, wie wir noch zeigen werden).

Beachte: Verfahren, die ihre Stabilität aus der L_2 -Theorie beziehen (d.h. aus Normen, die von Skalarprodukten herrühren), können Oszillationen nicht verhindern, wenn es keine Abschätzung in der Maximumnorm gibt. Sie werden in der Praxis daher oft (obwohl von besserer Papierform bzgl. der Konvergenz) gegenüber Verfahren niedrigerer Ordnung, die Oszillationen ausschließen, hinten an gestellt.

Trotzdem führen wir hier auch die zugehörige L_2 -Theorie vor, denn oft rechnet man in Randnähe bei Vorhandensein von Oszillationen mit oszillationsfreien Verfahren und steigt auf schneller konvergente Verfahren um in einem gewissen Abstand von der Anfangswertgeraden.

Erweiterung des Stabilitätssatzes (4.2)

Der Stabilitätssatz 4.2 für das Verfahren $\mathbf{B} y_t + \mathbf{A} y = \varphi$ wurde hergeleitet aus der energetischen Identität (4.10), zu deren Beweis die Symmetrie von \mathbf{A} benötigt wurde, die im allgemeinen Fall nicht notwendig vorliegt (vgl. (8.3) für $v \neq 0$, und (8.4)). Wir schreiben deshalb unser Verfahren um, können dann allerdings nur noch die Stabilität bzgl. der Anfangswerte zeigen aber nicht bzgl. der rechten Seite, was für die Konvergenz nötig war. Man (bzw. Samarskij) kann jedoch einen Satz beweisen, der die Stabilität der rechten Seite zurückführt auf die Stabilität bzgl. der Anfangswerte. Damit ist unser weiteres Vorgehen vorgezeichnet.

Im parabolischen Fall existiert unter der Voraussetzung $y_0 = y_N = 0$ üblicherweise \mathbf{A}^{-1} , entweder durch den Nachweis, daß \mathbf{A}^{-1} eine M -Matrix ist, vgl. (8.3), (8.4), oder mittels der reellen positiven Definitheit von \mathbf{A} (vgl. (8.7), (8.12)).

Dann kann man unser Verfahren mit einem geeigneten \mathbf{A} , (vgl. die vorigen Abschnitte, aber $\kappa = \text{const.}$, $v = \text{const.}$, vgl. dazu auch (8.11)-(8.13))

$$(\kappa \mathbf{I} + \sigma \tau \mathbf{A}) \mathbf{y}_t + \mathbf{A} \mathbf{y} = \varphi$$

von links mit \mathbf{A}^{-1} multiplizieren und erhält

$$(8.5) \quad \underbrace{(\kappa \mathbf{A}^{-1} + \sigma \tau \mathbf{I})}_{\tilde{\mathbf{B}}} \mathbf{y}_t + \underbrace{\mathbf{I}}_{\tilde{\mathbf{A}}} \mathbf{y} = \underbrace{\mathbf{A}^{-1} \varphi}_{\tilde{\varphi}}$$

also das Verfahren

$$(8.6) \quad \tilde{\mathbf{B}} \mathbf{y}_t + \tilde{\mathbf{A}} \mathbf{y} = \tilde{\varphi} \quad \text{mit} \quad \tilde{\mathbf{B}} = \kappa \mathbf{A}^{-1} + \sigma \tau \mathbf{I}, \quad \tilde{\mathbf{A}} = \mathbf{I}, \quad \tilde{\varphi} = \mathbf{A}^{-1} \varphi.$$

$\tilde{\mathbf{A}}$ ist nun wieder symmetrisch. Man kann also vorgehen wie bei Satz (4.2). Es zeigt sich allerdings, daß man für \mathbf{A} , bzw. \mathbf{A}^{-1} , und damit auch für \mathbf{B} die positive Definitheit nicht mehr immer nachweisen kann (z.B. für $v \neq 0$). Man muß sie abschwächen

Dies beweist

Satz 8.1

Das Verfahren

$$(\kappa \mathbf{I} + \sigma \tau \mathbf{A}) \mathbf{y}_t + \mathbf{A} \mathbf{y} = \boldsymbol{\varphi}$$

bzw. $(\kappa \mathbf{A}^{-1} + \sigma \tau \mathbf{I}) \mathbf{y}_t + \mathbf{I} \mathbf{y} = \mathbf{A}^{-1} \boldsymbol{\varphi}$

(mit $\kappa = \text{const.}$) ist stabil bzgl. der Anfangswerte (also $y_0^j = y_N^j = 0 \forall j$, $\boldsymbol{\varphi} \equiv \mathbf{0}$), falls

- (i) $(\mathbf{A} \mathbf{y}, \mathbf{y})_{(0,h)} > 0 \forall \mathbf{y} \in \mathbb{R}^n, \mathbf{y} \neq \mathbf{0}$ (\mathbf{A} reell positiv definit) und
- (ii) $\sigma \geq \frac{1}{2}$

Bemerkungen Fordert man in die stärkere Bedingung (Stabilität bzgl. der rechten Seite)

$$\tilde{\mathbf{B}} \geq \frac{\tau}{2} \tilde{\mathbf{A}} + \frac{\varepsilon}{2} \mathbf{I}$$

so führt dies auf obigem Wege zu der Bedingung $\sigma \geq \frac{1}{2} + \frac{\varepsilon}{2\tau}$. Wegen $\tau \rightarrow 0$ kann man hieraus keine Stabilitätsbedingung bzgl. der rechten Seite gewinnen. Man kann nicht ohne weiteres $\varepsilon \rightarrow 0$ gehen lassen, da sonst die Stabilität bezl. der rechten Seite ebenfalls verloren geht, vgl. dazu Satz (4.2) b), wo ε im Nenner steht.

Wir zeigen zunächst, dass (i) für die Matrizen (8.3) und (8.4) erfüllt ist.

Reelle pos. Definit. für $(\mathbf{A} \mathbf{y})_i := -(ky_x)_{\bar{x},i} + v y_{0,x}$, $v = \text{const}$, $y_0 = y_N = 0$

$$= \frac{1}{h} \left(k_{i+\frac{1}{2}} \frac{y_{i+1} - y_i}{h} - k_{i-\frac{1}{2}} \frac{y_i - y_{i-1}}{2} \right) + v \frac{y_{i+1} - y_i}{h}.$$

(vgl. (7.2), (8.3))

$$\begin{aligned} (\mathbf{A} \mathbf{y}, \mathbf{y})_{(0,h)} &= \sum_{i=1}^{N-1} (\mathbf{A} \mathbf{y})_i y_i h \\ (8.11) \quad &= - \sum_{i=1}^{N-1} (ky_x)_{\bar{x},i} y_i h + \sum_{i=1}^{N-1} v (y_{0,x,i}) y_i h \quad \text{und mit } y_0^j = y_N^j = 0 \forall j \end{aligned}$$

Bekannt ist (vgl. (8.11)):

$$\sum_{i=1}^N (ky_x)_{\bar{x},i} = \sum_{i=1}^N \underbrace{k_{i-\frac{1}{2}}}_{\geq c_0} (y_{x,i-1})^2 h \geq c_0 \|\mathbf{y}\|_{(1,h)}^2.$$

Wir formen die zweite Summe um für $v := \text{const.}$

$$\begin{aligned} \sum_{i=1}^{N-1} v(y_{\bar{x},i}) y_i h &= v \sum_{i=1}^{N-1} \frac{y_{i+1} - y_{i-1}}{2} y_i = \frac{v}{2} \left(\sum_{i=1}^{N-1} y_{i+1} y_i - \sum_{i=1}^{N-1} y_i y_{i-1} \right) \\ &= \frac{v}{2} \left(\sum_{i=2}^N y_i y_{i-1} - \sum_{i=1}^{N-1} y_i y_{i-1} \right) \\ &= \frac{v}{2} \left(\underset{=0}{y_N} y_{N-1} - y_1 \underset{=0}{y_0} \right), \quad \text{klappt nicht für variables } v! \end{aligned}$$

Es bleibt

$$(8.12) \quad (\mathbf{A}\mathbf{y}, \mathbf{y})_{(0,h)} = \sum_{i=1}^N k_{i-\frac{1}{2}} (y_{x,i-1})^2 h \stackrel{(7.7)}{\geq} c_0 \|\mathbf{y}\|_{(1,h)}^2 > 0 \quad \text{für } k \geq c_0.$$

Fazit: Der Konvektionsterm mit konstanter Geschwindigkeit berührt die reelle positive Definitheit nicht.

Reelle pos. Definitheit für $(\mathbf{A}_h \mathbf{y})_i := -(ky_x)_{\bar{x},i} + v y_{\bar{x}}, v = \text{const}, y_0 = y_N = 0$
(vgl. (8.4))

Analog zum vorigen brauchen wir nur den Konvektionsterm zu untersuchen, also den Term (vgl. (8.11))

$$v \sum_{i=1}^{N-1} (y_i - y_{i-1}) y_i$$

Nun ist

$$\begin{aligned} (y_i - y_{i-1}) y_i &= y_i^2 - y_i y_{i-1}, \quad \text{wegen } 2ab = -(a-b)^2 + a^2 + b^2 \\ &= y_i^2 - \frac{1}{2} [y_i^2 + y_{i-1}^2 - (y_i - y_{i-1})^2] \\ &= y_i^2 + \frac{1}{2} [-y_i^2 - y_{i-1}^2 + (y_{\bar{x},i})^2 \cdot h^2] \\ &= \frac{1}{2} [y_i^2 - y_{i-1}^2 + (y_{\bar{x},i})^2 h^2] \\ \sum_{i=1}^{N-1} (y_i - y_{i-1}) y_i &= \sum_{i=1}^{N-1} \frac{1}{2} [y_i^2 - y_{i-1}^2 + (y_{\bar{x},i})^2 h^2] = \frac{1}{2} \left(y_{N-1}^2 + \sum_{i=1}^{N-1} (y_{\bar{x},i})^2 h^2 \right), \\ &\text{und wegen } y_N = 0, y_{N-1}^2 = (y_{n-1} - y_N)^2 = h^2 y_{\bar{x},N} \\ &= \frac{1}{2} \sum_{i=1}^N (y_{\bar{x},i})^2 h^2 \stackrel{(3.21)}{=} \frac{h}{2} \|\mathbf{y}\|_{(1,h)}^2 \xrightarrow{\text{insgesamt}} \end{aligned}$$

$$(8.13) \quad (\mathbf{A}_h \mathbf{y}, \mathbf{y})_{(0,h)} \geq \|\mathbf{y}\|_{(1,h)}^2 + \frac{vh}{2} \|\mathbf{y}\|_{(1,h)}^2 = \left(c_0 + \frac{vh}{2} \right) \|\mathbf{y}\|_{(1,h)}^2 > 0 \quad \forall \mathbf{y} \neq \mathbf{0}$$

Also reell positiv definit.

Fazit: wie oben!

Zusammenfassend haben wir also folgendes (Zwischen-) Ergebnis

Satz 8.2

Für den Differentialoperator

$$Lu := \kappa \frac{\partial u}{\partial t} - \frac{\partial}{\partial x} \left(k(x) \frac{\partial u}{\partial x} \right) + v \frac{\partial u}{\partial x} + q(x) u$$

seien $0 < c_0 \leq k(x) \leq c_1$, $q(x) > 0$, $\kappa = \text{const}$, $v = \text{const}$, $\kappa, v > 0$.

Dann liefern die Diskretisierungen

$$(*) \quad (\mathbf{A}y)_i := -(ky_x)_{\bar{x},i} + v y_{x,i}^o$$

$$(**) \quad (\mathbf{A}y)_i := -(ky_x)_{\bar{x},i} + v y_{\bar{x},i}$$

reell positiv definite Matrizen.

Mit diesen Diskretisierungen ist das implizite Verfahren

$$(\kappa \mathbf{I} + \sigma \tau \mathbf{A})y_t + \mathbf{A}y = \varphi, \quad \sigma \geq \frac{1}{2}$$

für die Aufgabe $Lu = f$ mit Rand- und Anfangswerten stabil bzgl. der Anfangswerte, und von der Konvergenzordnung $O(\tau^2 + h^2)$ für $(*)$ und $O(\tau^2 + h)$ für $(**)$.

Die Wärmeleitung mit zeitabhängigem Diffusionskoeffizienten

Der Einfachheit halber verzichten wir (zunächst) auf Konvektionsterm, Abbaurate und rechte Seite und betrachten die Aufgabe

$$(8.14) \quad \begin{aligned} \frac{\partial u}{\partial t} &= \frac{\partial}{\partial x} \left(k(x, t) \frac{\partial u}{\partial x} \right), & 0 < x < 1, \quad t > 0 \\ u(x, 0) &= u_0(x), & 0 \leq x \leq 1 \\ u(0, t) &= g_0(t), \quad u(1, t) = g_1(t), & t > 0. \end{aligned}$$

Für die Aufgabe (8.14) betrachten wir für die Punkte x_i , $i = 1, \dots, N - 1$, die Approximation (vgl. (7.3))

$$y_{t,i}^j = \frac{1}{h} \left[\left(k_{i+\frac{1}{2}}^{j+\frac{1}{2}} y_x \right)_{i+\frac{1}{2}} - \left(k_{i-\frac{1}{2}}^{j+\frac{1}{2}} y_{\bar{x},i} \right)_{i-\frac{1}{2}} \right]^{(\sigma)}$$

(8.15) Hierbei wird für jede Zeitschicht j der Diffusionskoeffizient k auf die Schicht $j + \frac{1}{2}$ gesetzt (im Hinblick auf Crank- Nicolson) und unabhängig von der Wahl von σ festgehalten. Die Wirkung von σ beschränkt sich auf die y -Werte.

Wir erhalten somit (vgl. (7.2)) für jede Zeitschicht t^j eine **symmetrische Matrix** $\mathbf{A}_h^{(j+\frac{1}{2})}(k)$ (zum Aussehen vgl. (7.3) mit k auf der Zeitschicht $t^{j+\frac{1}{2}}$), deren **positive Definitheit** wie in (7.6) nachgerechnet wird.

Wir können wieder die zeitliche Mittelwertbildung durch σ mit der numerischen Differentiation bzgl. x vertauschen und unser Verfahren erhält die Gestalt

$$(8.16) \quad \begin{aligned} \mathbf{y}_t^j &= -\mathbf{A}_h^{(j+\frac{1}{2})}(k) \mathbf{y}^\sigma \quad \text{bzw.} \\ \left(\mathbf{I} + \sigma \tau \mathbf{A}_h^{(j+\frac{1}{2})}(k) \right) \mathbf{y}_t^j + \mathbf{A}_h^{(j+\frac{1}{2})}(k) \mathbf{y}^j &= 0. \\ \underbrace{\left(\mathbf{A}_h^{(j+\frac{1}{2})}(k) \right)^{-1} + \sigma \tau \mathbf{I}}_{\tilde{\mathbf{B}}^j} \mathbf{y}^j + \mathbf{I} \mathbf{y}^j &= 0. \end{aligned}$$

Natürlich ist $\mathbf{A}_h^{(j+\frac{1}{2})}(k)$, und damit auch die Inverse, reell positiv definit als Folge der positiven Definitheit. Wie zu Beginn des Paragraphen kann man die Abbauraten für beliebiges $q(x, t)$ in $\mathbf{A}_h^{(j+\frac{1}{2})}(k)$ integrieren, ein konstantes $\kappa > 0$ einschließen und einen Geschwindigkeitsterm $v \frac{\partial u}{\partial x}$ mit konstantem v berücksichtigen. Für die zu (8.3), (8.4) analogen Matrizen zeigt man wie in (8.11)-(8.13) die reell positive Definitheit. Dann kann man wieder Satz 4.2 analog zu (8.9) anwenden. Dabei spielt es keine Rolle, daß $\tilde{\mathbf{B}}$ nun zeitabhängig ist, wenn nur für jede Zeitschicht die zu (8.10) analogen Voraussetzungen für $\mathbf{A}_h^{(j+\frac{1}{2})}(k)$ erfüllt sind. Damit erhalten wir die zu Satz 8.1 analoge Aussage

Satz 8.3

Es seien

(i) $(\mathbf{A}_h^{(j+\frac{1}{2})}(k) \mathbf{y}, \mathbf{y})_{(0,h)} > 0 \quad \forall \mathbf{y} \in \mathbb{R}^n, \mathbf{y} \neq \mathbf{0} \quad (\mathbf{A}_h^{(j+\frac{1}{2})}(k) \text{ reell positiv definit})$
und

(ii) $\sigma \geq \frac{1}{2}$.

Dann ist das Verfahren

$$\begin{aligned} (\kappa \mathbf{I} + \sigma \tau \mathbf{A}_h^{(j+\frac{1}{2})}(k)) \mathbf{y}_t + \mathbf{A}_h^{(j+\frac{1}{2})}(k) \mathbf{y} &= \varphi \\ \text{bzw.} \quad \left(\kappa \left(\mathbf{A}_h^{(j+\frac{1}{2})}(k) \right)^{-1} + \sigma \tau \mathbf{I} \right) \mathbf{y}_t + \mathbf{I} \mathbf{y} &= \left(\mathbf{A}_h^{(j+\frac{1}{2})}(k) \right)^{-1} \varphi \end{aligned}$$

(mit $\kappa = \text{const.}$) ist stabil bzgl. der Anfangswerte (also $y_0^j = y_N^j = 0 \quad \forall j, \quad \varphi \equiv \mathbf{0}$).

Die Stabilität bzgl. der rechten Seite

wird zurückgeführt auf die Stabilität bzgl. der Anfangswerte. Stabilität bzgl. Anfangswerten und rechter Seite wird wieder durch Superposition gezeigt. Da die Stabilität bzgl. der Anfangswerte bekannt ist, genügt es die Stabilität bzgl. der rechten Seite für eine Aufgabe mit Nullanfangswerten zu betrachten.

Wir behandeln die Aufgabe

$$(8.17) \quad u_t - \frac{\partial}{\partial x} \left(k(x, t) \frac{\partial u}{\partial x} \right) + v(x, t) \frac{\partial u}{\partial x} + q(x, t) = f(x, t), \quad u(0, t) = u(1, t) = 0.$$

und zu ihrer Lösung ein Differenzenschema der Form (vgl. (8.5), (8.6))

$$(8.18) \quad \tilde{\mathbf{B}}^j \mathbf{y}_t^j + \tilde{\mathbf{A}} \mathbf{y}^j = \tilde{\boldsymbol{\varphi}}^j, \quad \text{mit}$$

$$\tilde{\mathbf{B}}^j = (\mathbf{A}^j)^{-1} + \sigma \tau \mathbf{I}, \quad \tilde{\mathbf{A}} = \mathbf{I}, \quad \tilde{\boldsymbol{\varphi}} = (\mathbf{A}^j)^{-1} \boldsymbol{\varphi}, \quad (\mathbf{A}^j \mathbf{y})^{(\sigma)} = \mathbf{A}^j \mathbf{y}^{(\sigma)}, \quad \sigma \geq \frac{1}{2}.$$

Die Diskretisierungen aus § 7, § 8, welche die Gestalt von \mathbf{A} bestimmen, sind zugelassen (also die Matrizen (7.2), (8.3), (8.4), (8.14)). Welche Zeitschichten für $\boldsymbol{\varphi}$ zugelassen werden ist ohne Belang. Wir treffen Voraussetzungen, welche die Stabilität bzgl. der Anfangswerte sichern (vgl. dazu die Sätze 4.2 und 8.1 und 8.4).

Der einfacheren Schreibweise wegen unterdrücken wir im folgenden Satz bei den Matrizen die Indizes und Argumente, die auf Zeit- bzw. Ortsabhängigkeit hinweisen.

Satz 8.4 (Samarskij)

Für das Differenzenschema

$$(8.19) \quad \frac{1}{\tau} \tilde{\mathbf{B}}(\mathbf{y}^j - \mathbf{y}^{j-1}) + \tilde{\mathbf{A}} \mathbf{y}^{j-1} = \tilde{\boldsymbol{\varphi}}^{j-1},$$

$$\tilde{\mathbf{B}} = \mathbf{A}^{-1} + \sigma \tau \mathbf{I}, \quad \tilde{\mathbf{A}} = \mathbf{I}, \quad \tilde{\boldsymbol{\varphi}} = \mathbf{A}^{-1} \boldsymbol{\varphi}, \quad \sigma \geq \frac{1}{2}, \quad y_0^j = y_N^j = 0 \quad \forall j$$

seien folgende Voraussetzungen erfüllt

$$\tilde{\mathbf{A}} > 0 \quad (\text{im Sinne reell positiv definit}), \quad \sigma \geq \frac{1}{2}.$$

Dann ist (8.19) stabil bzgl. Anfangswerten und rechter Seite:

$$(8.20) \quad \|\mathbf{y}^j\|_{(0,h)} \leq \|\mathbf{y}^0\|_{(0,h)} + \sum_{k^1}^j \tau \|\boldsymbol{\varphi}^{k-1}\|_{(0,h)}.$$

Beweis

Die Stabilität bzgl. Anfangswerten (bereits bekannt) und rechter Seite erhält man durch Superposition. Es genügt also die Stabilität bzgl. der rechten Seite zu beweisen für $\mathbf{y}^0 = \mathbf{0}$, $y_0^j = y_N^j = 0 \quad \forall j$.

Für die Lösung von (8.19) machen wir den Ansatz für $m > j \quad \forall j$ mit $j\tau \leq T$:

$$\mathbf{y}^j = \sum_{k=1}^m \mathbf{y}_{(k)}^j,$$

die $\mathbf{y}_{(k)}^j$ seien Lösungen von

$$(8.21) \quad \frac{1}{\tau} \tilde{\mathbf{B}} \left(\mathbf{y}_{(k)}^j - \mathbf{y}_{(k)}^{j-1} \right) + \tilde{\mathbf{A}} \mathbf{y}_{(k)}^{j-1} = \delta_{kj} \tilde{\varphi}^{k-1}, \quad k = 1, \dots, m, \quad \mathbf{y}_{(k)}^0 = \mathbf{0}$$

mit

$$(8.22) \quad \delta_{kj} \tilde{\varphi}^{k-1} = \begin{cases} \mathbf{0} & \text{für } k \neq j \\ \tilde{\varphi}^{j-1} & \text{für } k = j \end{cases}$$

Man erkennt: Aufsummieren von (8.21) über $k = 1, \dots, m$ liefert unter Berücksichtigung von (8.22) das Verfahren (8.19), wenn alle $\mathbf{y}_{(k)}^j$ berechnet sind.

Berechnung der $\mathbf{y}_{(k)}^j$: $k \in \{1, \dots, m\}$ fest, $j = 1, 2, \dots$. Aus (8.21), (8.22) folgt

- (i) $j < k$: $\mathbf{y}_{(k)}^j = \mathbf{0}$ (da Anfangswerte = 0, rechte Seite = 0),
- (ii) $j = k$: $\frac{1}{\tau} \tilde{\mathbf{B}} \mathbf{y}_{(k)}^k = \tilde{\varphi}^{k-1} \implies \mathbf{y}_{(k)}^k = \tau \tilde{\mathbf{B}}^{-1} \tilde{\varphi}^{k-1}$
- (iii) $j > k$: $\frac{1}{\tau} \tilde{\mathbf{B}} \left(\mathbf{y}_{(k)}^j - \mathbf{y}_{(k)}^{j-1} \right) + \tilde{\mathbf{A}} \mathbf{y}_{(k)}^{j-1} = \mathbf{0}$

Es können also für jedes k alle $\mathbf{y}_{(k)}^j$ berechnet werden und zwar, sofern sie $\neq \mathbf{0}$ sind als Lösungen von Differenzenverfahren, deren rechte Seite = 0 ist.

$\tilde{\mathbf{A}} > 0$, $\sigma \geq \frac{1}{2}$ sichern nach den Sätzen 8.1 und 4.2 die Ungleichungen

$$(8.23) \quad \|\mathbf{y}_{(k)}^j\|_{(0,h)} \leq \|\mathbf{y}_{(k)}^{j-1}\|_{(0,h)} \leq \dots \leq \|\mathbf{y}_{(k)}^k\|_{(0,h)} = \tau \left\| \tilde{\mathbf{B}}^{-1} \tilde{\varphi}^{k-1} \right\|_{(0,h)}$$

eigentlich mit dem Index $(1, h)$, vgl. (4.12). Beachte jedoch Für $\tilde{\mathbf{A}} = \mathbf{I}$ gilt

$$\|\mathbf{y}\|_{(0,h)}^2 = (\mathbf{y}, \mathbf{y})_{(0,h)}, \quad \|\mathbf{y}\|_{(1,h)}^2 = (\tilde{\mathbf{A}}\mathbf{y}, \mathbf{y})_{(0,h)} = (\mathbf{y}, \mathbf{y})_{(0,h)}$$

Einsetzen von (8.23) in die Darstellung von \mathbf{y}^j liefert (vgl. (i))

$$(8.24) \quad \|\mathbf{y}^j\|_{(0,h)} \leq \sum_{k=1}^j \|\mathbf{y}_{(k)}^j\|_{(0,h)} \leq \sum_{k=1}^j \tau \left\| \tilde{\mathbf{B}}^{-1} \tilde{\varphi}^{k-1} \right\|_{(0,h)} \quad (\mathbf{y}_{(k)}^j = \mathbf{0} \text{ für } k > j)$$

Nun ist wegen $\tilde{\mathbf{B}} = \mathbf{A}^{-1}(\mathbf{I} + \sigma\tau\mathbf{A})$

$$\underbrace{\tilde{\mathbf{B}}^{-1} \tilde{\varphi}^{k-1}}_{=: \mathbf{v}} = (\mathbf{I} + \sigma\tau\mathbf{A})^{-1} \mathbf{A} \underbrace{\mathbf{A}^{-1} \tilde{\varphi}^{k-1}}_{=: \tilde{\varphi}^{k-1}} = (\mathbf{I} + \sigma\tau\mathbf{A})^{-1} \underbrace{\tilde{\varphi}^{k-1}}_{=: \mathbf{w}}$$

Zur Abschätzung von $(\mathbf{I} + \sigma \tau \mathbf{A})^{-1}$ schätzen wir die Lösung \mathbf{v} ab von

$$(8.25) \quad (\mathbf{I} + \sigma \tau \mathbf{A})\mathbf{v} = \mathbf{w}, \quad (\mathbf{v} = \tilde{\mathbf{B}}^{-1} \tilde{\boldsymbol{\varphi}}^{k-1}, \mathbf{w} = \boldsymbol{\varphi}^{k-1})$$

Multiplizieren dieser Gleichung mit $(\cdot, \mathbf{v})_{(0,h)}$ liefert

$$\begin{aligned} \|\mathbf{v}\|_{(0,h)}^2 + \underbrace{\sigma \tau (\mathbf{A}\mathbf{v}, \mathbf{v})_{(0,h)}}_{>0} &= (\mathbf{w}, \mathbf{v})_{(0,h)} \stackrel{\text{CSU}}{\leq} \|\mathbf{w}\|_{(0,h)} \|\mathbf{v}\|_{(0,h)} \\ \Rightarrow \|\mathbf{v}\|_{(0,h)}^2 &\leq \|\mathbf{w}\|_{(0,h)} \|\mathbf{v}\|_{(0,h)}, \quad \text{bzw. nach Kürzen} \\ \left\| \tilde{\mathbf{B}}^{-1} \tilde{\boldsymbol{\varphi}}^{k-1} \right\|_{(0,h)} &= \|\mathbf{v}\|_{(0,h)} \leq \|\mathbf{w}\|_{(0,h)} = \|\boldsymbol{\varphi}^{k-1}\|_{(0,h)} \end{aligned}$$

Damit erhält man aus (8.24)

$$\|\mathbf{y}^j\|_{(0,h)} \leq \sum_{k=1}^j \tau \|\boldsymbol{\varphi}^{k-1}\|_{(0,h)}$$

□

Bemerkung Die Stabilitätsvoraussetzung $\mathbf{A} > 0$, $\sigma \geq \frac{1}{2}$ des Satzes kann durch jeden anderen Satz von Voraussetzungen ersetzt werden, welcher die Stabilität bzgl. der Anfangswerte garantiert.

§ 9 Gleichmäßige Stabilität und Konvergenz

Satz 9.1

Für eine parabolische Differentialgleichung mit Nullrandwerten erfülle das Verfahren

$$(9.1) \quad (\mathbf{I} + \sigma \tau \mathbf{A}^j) \mathbf{y}_t + \mathbf{A}^j \mathbf{y} = \boldsymbol{\varphi}^j, \quad (\mathbf{A}^j \mathbf{y})^{(\sigma)} = \mathbf{A}^j \mathbf{y}^{(\sigma)}, \quad j \geq 0, \quad 0 \leq \sigma \leq 1$$

folgende Voraussetzungen

$$(9.2) \quad \mathbf{A}^j = (a_{i,k}^j)_{i,k=1,\dots,N-1} \text{ sei } \forall j \text{ eine tridiagonale } M\text{-Matrix, } a_{ik}^j \leq 0 \quad \forall i \neq k$$

$$(9.3) \quad a_{ii}^j + a_{i,i-1}^j + a_{i,i+1}^j \geq 0, \quad a_{1,0} = a_{N-1,N} = 0, \quad \forall i = 1, \dots, N-1, \quad j \geq 0$$

$$(9.4) \quad (1 - \sigma)\tau a_{ii}^j \leq 1 \quad (\text{Schrittweitenbeschränkung})$$

Dann gelten

$$(9.5) \quad \|\hat{\mathbf{y}}\|_\infty \leq \|\mathbf{y}\|_\infty + \tau \|\boldsymbol{\varphi}\|_\infty \quad \text{und}$$

$$(9.6) \quad \|\mathbf{y}^{j+1}\|_\infty \leq \|\mathbf{y}^0\|_\infty + \tau \sum_{k=0}^j \|\boldsymbol{\varphi}^k\|_\infty$$

(Stabilität bzgl. Anfangswerten und rechter Seite)

Bemerkungen:

α) \mathbf{A}^j bedeutet, dass \mathbf{A} zeitabhängig sein darf (vgl. etwa (8.14)). Aufgaben mit variablen Koeffizienten sind erfaßbar, die Diskretisierung muß keine symmetrische Matrix \mathbf{A}^j liefern und es sind auch keine Definitheitsvoraussetzungen für \mathbf{A}^j erforderlich.

β) Für das einfache Verfahren (4.1): $\mathbf{y}_t^j = -\mathbf{A}_h^0 \mathbf{y}^j$ mit $\mathbf{y}^0 = u_0$, $y_0^j = y_N^j = 0$, $j \geq 0$ bedeutet (9.4) wegen $\sigma = 0$, $a_{ii} = \frac{2}{h^2}$, dass $\tau \leq \frac{h^2}{2}$ eine Voraussetzung ist, die sich – vgl. (4.2) – kaum abschwächen läßt.

γ) Auf Grund der Nullrandwerte hängt $\boldsymbol{\varphi}$ nur von der rechten Seite der Differentialgleichung ab, z.B. $\boldsymbol{\varphi}^j = \tilde{\mathbf{f}}^j$, wobei die „ \sim “ bedeutet, dass \mathbf{f} auch an einer Zeitschicht zwischen $j+1$ und j betrachtet werden kann.

δ) **Nur für rein implizite Verfahren keine Schrittweitenbeschränkung!**
(d.h. $\sigma = 1$)

Beweis

$\forall \mathbf{A}^j$ gilt: $\mathbf{B}^j = (\mathbf{I} + \sigma \tau \mathbf{A}^j)$ ist eine M -Matrix, denn (wir unterdrücken den Index j) die Vorzeichenregel ist erfüllt und mit $\mathbf{p} = (1, \dots, 1)^T$ folgt

$$(9.7) \quad (\mathbf{B}\mathbf{p})_i = ((\mathbf{I} + \sigma \tau \mathbf{A})\mathbf{p})_i = 1 + \sigma \tau (a_{i,i-1} + a_{ii} + a_{i,i+1}) \stackrel{(9.3)}{\geq} 1 > 0 \quad \forall i$$

Damit gilt nach Satz (6.4)

$$(9.8) \quad \|\mathbf{B}^{-1}\|_{\infty} \leq \frac{1}{\min_i (\mathbf{B}\mathbf{p})_i} = 1$$

↑
Zeilensummennorm

Auflösen von (9.1) nach $\hat{\mathbf{y}}$ liefert

$$(9.9) \quad \underbrace{(\mathbf{I} + \sigma \tau \mathbf{A})}_{\mathbf{B}} \hat{\mathbf{y}} = (\mathbf{I} + \tau(\sigma - 1)\mathbf{A}) \mathbf{y} + \tau \boldsymbol{\varphi}$$

und

$$\begin{aligned} |\{(\mathbf{I} + \tau(\sigma - 1)\mathbf{A}) \mathbf{y}\}_i| &= |(1 + \tau(\sigma - 1) a_{ii}) y_i + \tau(\sigma - 1) a_{i,i-1} y_{i-1} + \tau(\sigma - 1) a_{i,i+1} y_{i+1}| \\ &\leq \underbrace{|(1 + \tau(\sigma - 1) a_{ii})|}_{\geq 0 \text{ nach (9.4), } a_{ii} \geq 0} |y_i| + \underbrace{|\tau(\sigma - 1) a_{i,i-1}|}_{\geq 0, \text{ da } a_{i,i-1} \leq 0} |y_{i-1}| \\ &\quad + \underbrace{|\tau(\sigma - 1) a_{i,i+1}|}_{\geq 0, \text{ da } a_{i,i+1} \leq 0} |y_{i+1}| \\ &\leq (1 + \underbrace{\tau(\sigma - 1)}_{\leq 0} \underbrace{(a_{ii} + a_{i,i-1} + a_{i,i+1})}_{\geq 0 \text{ nach (9.3)}}) \|\mathbf{y}\|_{\infty} \\ &\leq \|\mathbf{y}\|_{\infty}, \end{aligned}$$

deshalb folgt aus (9.9) mit (9.8)

$$\begin{aligned} \|\hat{\mathbf{y}}\|_{\infty} &\leq \underbrace{\|\mathbf{B}^{-1}\|_{\infty}}_{\leq 1} (\|\mathbf{I} + \tau(\sigma - 1)\mathbf{A}\mathbf{y}\|_{\infty} + \tau\|\boldsymbol{\varphi}\|_{\infty}) \\ &\leq \|\mathbf{y}\|_{\infty} + \tau\|\boldsymbol{\varphi}\|_{\infty} \quad (\text{also (9.5)}) \end{aligned}$$

Wendet man diese Ungleichung von Zeitschicht zu Zeitschicht an, so folgt

$$\|\mathbf{y}^{j+1}\|_{\infty} \leq \|\mathbf{y}^0\|_{\infty} + \sum_{k=0}^j \tau \|\boldsymbol{\varphi}^k\|_{\infty}.$$

□

Konvergenzbetrachtung

Unter den Voraussetzungen von Satz 9.1 gilt mit den Bezeichnungen aus Definition 5.3:

Mit

$$\mathbf{z} = \mathbf{u} - \mathbf{y}$$

↑
Gitterfkt. der exakten Lösung

folgt

$$\begin{aligned} (\mathbf{I} + \sigma \tau \mathbf{A}^j) \mathbf{z}_t + \mathbf{A}^j \mathbf{z} &= (\mathbf{I} + \sigma \tau \mathbf{A}^j) \mathbf{u}_t + \mathbf{A}^j \mathbf{u} - \underbrace{[(\mathbf{I} + \sigma \tau \mathbf{A}^j) \mathbf{y}_t + \mathbf{A}^j \mathbf{y}]}_{=\boldsymbol{\varphi}^j} \\ &= \boldsymbol{\psi}^j \quad (\text{Diskretisationsfehler, vgl. Definition (5.1)}) \end{aligned}$$

Wegen $z^0 = 0, z_0^j = z_N^0 = 0$ liefert die Abschätzung (9.6)

$$(9.10) \quad \|z^{i+1}\|_\infty \leq \tau \sum_{k=0}^j \|\psi^k\|_\infty,$$

also gleichmäßige Konvergenz von der Ordnung des Diskretisierungsfehlers unter entsprechenden Differenzierbarkeitsvoraussetzungen an u .

Bemerkung:

Wählt man $\sigma = \frac{1}{2}$ (wie bei Crank-Nicolson), so erhält man aus (9.10) unter der Schrittweitenbeschränkung (9.4): $\tau \leq \frac{2}{a_{ii}^j}$ ($= h^2$ im einfachsten Fall, vgl. (2.11)) eine Verfahrensordnung $O(\tau^2 + h^2)$.

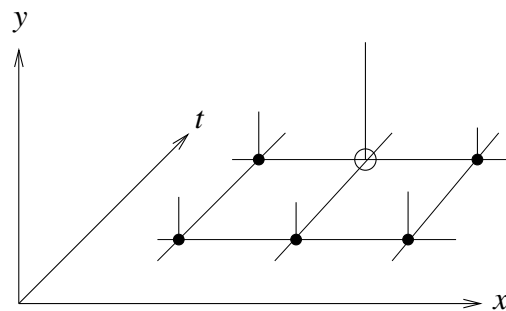
Fazit:

Crank-Nicolson ohne Schrittweitenbeschränkung erst anwenden, wenn das Verfahren schon geglättet ist – sofern man für kleine Zeiten Oszillationen vermeiden will.

Oszillationsfreiheit

Hinweise: Im Abschnitt über elliptische Differentialgleichungen werden wir ein diskretes Maximumprinzip beweisen, das im Wesentlichen für M -Matrizen gilt (vgl. Satz 12.5 bis 12.5. Es hängt nur ab von der Diskretisierungsmatrix (die allerdings auch die Randwerte mit einschließen muß) für die entsprechende Aufgabe. Im parabolischen Fall muß man dazu eine “Blockmatrix” aufstellen, die über alle Zeitschichten geht. Das Maximumprinzip besagt dann, daß das Maximum nur am parabolischen Rand angenommen werden kann. Dieses Maximumprinzip gilt auch lokal.

Die folgende 2-dimensionale Abbildung (die y -Werte werden in den Gitterpunkten abgetragen) entspricht einer Oszillation. Gilt das Maximumprinzip (hier lokal), so muß das Maximum im parabolischen Rand angenommen werden, d.h. hier in einem der “dicken” Gitterpunkte. Die Zeichnung steht also im Widerspruch zum Maximumprinzip.



§ 10 Die mehrdimensionale Wärmeleitungsgleichung

Wir beschränken uns (zunächst) auf die Aufgabe

$$\frac{\partial u}{\partial t} = \operatorname{div}(k(x) \operatorname{grad} u) + f = \sum_{i=1}^n \frac{\partial}{\partial x_i} \left(k(x) \frac{\partial u}{\partial x_i} \right) + f$$

$$(10.1) \quad x \in \Omega \subset \mathbb{R}^p, \quad 0 < t < T$$

AWe: $t = 0 : u(x, 0) = u_0(x), \quad x \in \overline{\Omega}$

RWe: $u|_{\Gamma \times [0, T]} = g, \quad \Gamma = \partial \Omega$

Insbesondere betrachten wir als Ω zunächst einen Quader (analog zum 1D-Fall). Wir werden zeigen, daß diese Aufgabe im Wesentlichen auf den 1D-Fall (bzgl. des Orts) reduzierbar ist. Danach kann man wie im 1D-Fall auch $k = k(x, t)$ zulassen, sowie Konvektion und Abbaurrate.

Im Weiteren beschränken wir uns zunächst (im Wesentlichen) auf den 2D-Fall, die Verallgemeinerung auf mehr Dimensionen ist dann offensichtlich.

Wir beschreiben zunächst die

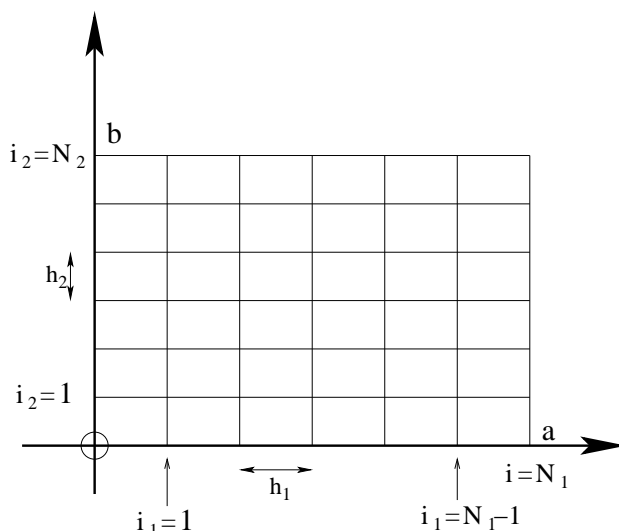
Diskretisierung des $\operatorname{div} \operatorname{grad}$ -Terms im Rechteck

Sei $\overline{\Omega} = [0, a] \times [0, b]$, $a = N_1 \cdot h_1$, $b = N_2 \cdot h_2$.

$\overline{\omega}$ sei die Menge aller Gitterpunkte aus $\overline{\Omega}$,

ω die Menge der inneren Gitterpunkte: $\omega \subset \Omega$,

γ die Menge der Randpunkte: $\gamma \in \partial \Omega$, also $\overline{\omega} = \omega + \gamma$.



Ist u eine auf Ω (bzw. $\overline{\Omega}$) definierte Funktion, so definiert sie eine Gitterfunktion \mathbf{u} , die durch die Restriktion von u auf die Gitterpunkte entsteht. Mit der Gitterfunktion \mathbf{y} bezeichnen wir die zugehörigen Näherungswerte.

Wir bezeichnen die inneren Gitterpunkte ($\in \omega$) auf verschiedene Weisen:

- (i) Wir können die $N = (N_1 - 1) \cdot (N_2 - 1)$ inneren Gitterpunkte durchnummerieren (z.B. zeilenweise von links nach rechts und von unten nach oben oder spaltenweise von unten nach oben und von links nach rechts). Der einzelne Gitterpunkt bekommt dann seinen Zählindex

$$x_\ell, \quad \ell \in [1, \dots, N], \quad N = (N_1 - 1)(N_2 - 1)$$

oder

- (ii) entsprechend den (i_1, i_2) Indizes, die die Stellung des Punktes im Gitter zeigen

$$x_\ell = x_{i_1, i_2} = (i_1 h_1, i_2 h_2), \quad 1 \leq i_\nu \leq N_\nu - 1, \quad \nu = 1, 2.$$

Damit definieren wir den Gittervektor

$$\mathbf{y} = (x_1, \dots, y_N)^T, \quad N = (N_1 - 1)(N_2 - 1).$$

Nun entspricht der Aufgabe (10.1) das Differenzenschema

$$\mathbf{y}_t = -\mathbf{A} \mathbf{y}^\sigma + \mathbf{f}^{j+\frac{1}{2}}, \quad \mathbf{A} = \sum_{i=1}^2 \mathbf{A}_i \quad \left(\sum_{i=1}^d \text{ im } d\text{-dim. Fall} \right)$$

Dabei beschreibt \mathbf{A}_i die Diskretisierung der Ableitungen in x_i -Richtung. Für \mathbf{A}_i werden die Gitterpunkte in Richtung der x_i -Achse durchnummeriert (also zeilenweise von links nach rechts und von unten nach oben für x_1 , und spaltenweise von unten nach oben und links nach rechts für x_2). Daß damit die Gitterpunkte für die Ableitungen in x_1 -Richtung und x_2 -Richtung unterschiedlich nummeriert sind, soll uns im Augenblick nicht stören. Wir wollen zunächst nur Abschätzungen für die Normen der \mathbf{A}_i herleiten und dafür spielt die Reihenfolge der Nummerierung keine Rolle. Dann ist

$$\begin{aligned} (\mathbf{A}_i \mathbf{y})_\ell &= -(k y_{\bar{x}_\ell})_{x_\ell} \quad (\text{vgl. (7.1)}) \\ &:= \frac{1}{h_i} \left((k(x_\ell + \frac{h_i}{2} \mathbf{e}^i) \frac{y_\ell(x_\ell + h_i \mathbf{e}^i) - y_\ell(x_\ell)}{h_i} - k(x_\ell - \frac{h_i}{2} \mathbf{e}^i) \frac{y_\ell(x_\ell) - y_\ell(x_\ell - h_i \mathbf{e}^i)}{h_i}) \right) \end{aligned}$$

mit den Einheitsvektoren \mathbf{e}^i , deren Koordinaten sich natürlich auch an der Zählweise orientieren. Abgesehen von den „Randpunkten“ unter den inneren Randpunkten gilt dann

$$y_\ell(x_\ell \pm h_i \mathbf{e}^i) = y_{\ell \pm 1}.$$

$y_{\ell \pm 1}$ bedeutet hier nur – ausgehend von x_ℓ – einen Schritt in $+x_i$ - bzw. $-x_i$ -Richtung entsprechend der jeweiligen Durchnummerierung.

Im 2D-Fall benutzen wir das Skalarprodukt $(\mathbf{v}, \mathbf{w})_{(0,H)}$, $H = h_1 h_2$:

$$(\mathbf{v}, \mathbf{w})_{(0,H)} = \sum_{\ell=1}^N v_\ell w_\ell h_1 h_2 = \sum_{k=1}^{N_2-1} \sum_{i=1}^{N_1-1} v_{i,k} w_{i,k} \cdot h_1 h_2 = \sum_{x \in \omega} (v(x), w(x)) H.$$

Man beachte, dass die Definition dieser Skalarprodukte und insbesondere ihr Wert nicht von der Nummerierung der Gitterpunkte abhängen. Wir setzen nun voraus

$$0 < c_0 \leq k(x) \leq c_1 \quad \forall x \in \Omega$$

und leiten Normabschätzungen her, basierend auf den bekannten Abschätzungen aus dem 1D-Fall.

Wir betrachten zunächst \mathbf{A}_1 bei zeilenweiser Durchnummerierung des Gittervektors. Die Wirkung von \mathbf{A}_1 auf den Teilgittervektor $\mathbf{v}^{(k)}$ der k -ten Zeile wird durch eine symmetrische, positiv definite, tridiagonale Matrix $\mathbf{A}_1^{(k)}$ (vgl. (7.3)) beschrieben. Die Gesamtmatrix \mathbf{A}_1 ergibt sich damit als Blockdiagonalmatrix, deren Blöcke durch die „Zeilenblöcke“ $\mathbf{A}_1^{(k)}$ gegeben werden. \mathbf{A}_1 ist damit eine symmetrische Tridiagonalmatrix.

Wir setzen nun voraus: $c_0 \leq k \leq c_1$ und erhalten unter Beachtung von $\mathbf{v}|_\Gamma = 0$ für den Teilgittervektor $\mathbf{v}^{(k)}$

$$\begin{aligned} \left(\mathbf{A}_1^{(k)} \mathbf{v}^{(k)}, \mathbf{v}^{(k)} \right)_{(0,H)} &\stackrel{(7.6)}{=} \sum_{i=1}^{N_1} k_{i-\frac{1}{2}}^{(k)} \left(v_{\bar{x}_1, i}^{(k)} \right)^2 h_1 h_2 \begin{cases} \leq & c_1 \sum_{i=1}^{N_1} \left(v_{\bar{x}_1, i}^{(k)} \right)^2 h_1 \cdot h_2 \\ \geq & c_0 \underbrace{\sum_{i=1}^{N_1} \left(v_{\bar{x}_1, i}^{(k)} \right)^2 h_1 \cdot h_2}_{\|\mathbf{v}^{(k)}\|_{(1,h_1)}^2} \end{cases} \\ &= \begin{cases} c_1 \\ c_0 \end{cases} \|\mathbf{v}^{(k)}\|_{(1,h_1)}^2 \cdot h_2 \begin{cases} \leq & \frac{4}{h_1^2} c_1 \|\mathbf{v}^{(k)}\|_{(0,h_1)}^2 \cdot h_2 \\ \geq & 8c_0 \end{cases} \end{aligned}$$

Für den ganzen Gittervektor \mathbf{v} und die Matrix \mathbf{A}_1 gilt also (Summation über die Zeilen)

$$\left(\mathbf{A}_1 \mathbf{v}, \mathbf{v} \right)_{(0,H)} \begin{cases} \leq & c_1 \sum_{k=1}^{N_2-1} h_2 \|\mathbf{v}^{(k)}\|_{(1,h_1)}^2 \\ \geq & c_0 \sum_{k=1}^{N_2-1} h_2 \|\mathbf{v}^{(k)}\|_{(0,h_1)}^2 \end{cases} \begin{cases} \leq & \frac{4c_1}{h_1^2} \sum_{k=1}^{N_2-1} h_2 \|\mathbf{v}^{(k)}\|_{(0,h_1)}^2 \\ \geq & 8c_0 \end{cases}$$

Wir definieren im 2D-Fall die Normen entsprechend dem Skalarprodukt durch

$$\|\mathbf{v}\|_{(1,H)}^2 = \sum_{k=1}^{N_2-1} h_2 \|\mathbf{v}^{(k)}\|_{(1,h_1)}^2, \quad \|v\|_{(0,H)}^2 = \sum_{k=1}^{N_2-1} h_2 \|\mathbf{v}^{(k)}\|_{(0,h_1)}^2$$

und erhalten damit

$$(10.2) \quad \left. \begin{aligned} \left(\mathbf{A}_1 \mathbf{v}, \mathbf{v} \right)_{(0,H)} &\leq \frac{4c_1}{h_1^2} \\ &\geq 8c_0 \end{aligned} \right\} \|\mathbf{v}\|_{(0,H)}^2$$

Beachte: Ist \mathbf{v} die Gitterfunktion einer Funktion $v(x)$, so gilt

$$\|\mathbf{v}\|_{(0,H)}^2 \xrightarrow{h_1, h_2 \rightarrow 0} \int_0^1 \int_0^1 v(x, y)^2 dy dx.$$

Durch analoges Vorgehen erhalten wir für \mathbf{A}_2 , indem wir nun den Gittervektor spaltenweise durchnummerieren:

$$(10.3) \quad (\mathbf{A}_2 \mathbf{v}, \mathbf{v})_{(0,H)} \begin{cases} \leq \frac{4c_1}{h_2^2} \|v\|_{(0,H)}^2 \\ \geq 8c_0 \end{cases}$$

Bei spaltenweiser Nummerierung ist \mathbf{A}_2 eine symmetrische (, tridiagonale) Matrix, d.h.

$$(*) \quad (\mathbf{A}_2 \mathbf{v}, \mathbf{v})_{(0,H)} = (\mathbf{v}, \mathbf{A}_2 \mathbf{v})_{(0,H)}$$

Wird nun \mathbf{v} anders nummeriert und die Indizierung der Matrixelemente entsprechend geändert, so geht die tridiagonale Form von \mathbf{A}_2 verloren, die Gleichung (*) bleibt jedoch richtig, da das Skalarprodukt unabhängig ist von der Komponentennummerierung der Vektoren, d.h. auch bei zeilenweiser Nummerierung bleibt \mathbf{A}_2 (eigentlich müßte man nach der Umnummerierung den Namen ändern) eine symmetrische Matrix. Dasselbe Argument zeigt, dass \mathbf{A}_2 positiv definit ist, denn die „Spaltenblöcke“ $\mathbf{A}_2^{(i)}$ der Matrix \mathbf{A}_2 (ebenso wie die „Zeilenblöcke“ von \mathbf{A}_1) sind positiv definit.

Aufgabe: Wie sieht \mathbf{A}_2 bei zeilenweiser Nummerierung aus?

Wir können nun wie früher zur Lösung der Aufgabe (10.1) das Verfahren

$$(10.4) \quad \mathbf{B} \mathbf{y}_t + \mathbf{A} \mathbf{y} = \boldsymbol{\varphi}, \quad \mathbf{B} = \mathbf{I} + \sigma \tau \mathbf{A}, \quad \mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2 \text{ symmetr., pos. def.}$$

anwenden und erhalten dafür die bekannten Stabilitäts- und Konvergenzaussagen

$\mathbf{B} \geq \frac{\tau}{2} \mathbf{A}$ sichert die Stabilität bzgl. der Anfangswerte (vgl. Satz 8.4), die wichtigste Eigenschaft. Diese Bedingung wird mit $\mathbf{I} \geq \frac{\mathbf{A}}{\|\mathbf{A}\|}$ verschärft zu

$$\mathbf{B} \geq \left(\frac{1}{\|\mathbf{A}\|} + \sigma \tau \right) \mathbf{A} \stackrel{!}{\geq} \frac{\tau}{2} \mathbf{A}$$

was wegen $\mathbf{A} > 0$ durch

$$\sigma \geq \frac{1}{2} - \frac{1}{\tau \|\mathbf{A}\|}$$

gesichert wird. Offensichtlich ist $\sigma \geq \frac{1}{2}$ problemlos.

Für $\sigma = 0$ (explizites Verfahren) folgt hieraus die Beschränkung

$$\tau \leq \frac{2}{\|\mathbf{A}\|}$$

Setzt man $h = h_1 = h_2$ so folgt aus (10.2), (10.3)

$$\|\mathbf{A}\|_S \leq \|\mathbf{A}_1\|_S + \|\mathbf{A}_2\|_S \leq 2 \cdot \frac{4c_1}{h^2} = \underset{p=\text{Raumdimension}}{p} \cdot \frac{4c_1}{h^2}$$

und damit

$$\tau \leq \frac{h^2}{2c_1 p}$$

Dies ist um den Faktor p schlimmer als im 1D-Fall

$\implies \sigma = 0$ indiskutabel im pD -Fall.

Lösungsmöglichkeiten für $\sigma > 0$ insbesondere $\sigma = \frac{1}{2}$

Das Verfahren (10.4) läßt sich in der Form schreiben

$$(\mathbf{I} + \sigma \tau \mathbf{A}) \mathbf{y}^{j+1} = (\mathbf{I} - (1 - \sigma)\tau \mathbf{A}) \mathbf{y}^j + \tau \boldsymbol{\varphi} =: \mathbf{F}^j$$

Die rechte Seite der Gleichung ist bekannt. Für jeden Zeitschritt kann man dies als stationäres parabolisches, d.h. elliptisches Problem betrachten und darauf das Mehrgitterverfahren (MGV) anwenden, das wir im Kapitel für elliptische Aufgaben beschreiben. Damit läßt sich auch der Zeitaufwand für die Lösung des Gleichungssystems abschätzen. Solche Verfahren werden in der Praxis gerechnet (beachte, dass \mathbf{A} nun keine Tridiagonalmatrix mehr ist).

Es gibt jedoch schnellere Verfahren, die auf tridiagonalen Gleichungssystemen beruhen, die Verfahren der alternativen Richtungen: (ADI-method, alternating direction implicit method) auf die wir nun eingehen.

Verfahren der Alternierenden Richtungen

(ADI-Verfahren: Alternating direction implicit method)

Das Verfahren wurde Anfang der 50er Jahre von Peachman-Rachford („Erdölleuten“) entwickelt für eine AWA mit der Differentialgleichung

$$(10.5) \quad \begin{aligned} \frac{\partial u}{\partial t} &= \sum_{i=1}^2 L_i u + f \\ L_i u &= \frac{\partial}{\partial x_i} \left(k_i \frac{\partial u}{\partial x_i} \right) - q_i u \quad (\text{nicht konstante Koeffizienten möglich}) \end{aligned}$$

Die diskretisierte Aufgabe führt zum Differenzenchema (vgl. (4.7) bzw. (10.4) und setze $\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2$). Wären Abbauterme $q = (q_1, q_2)$ in der Differentialgleichung enthalten, so könnten wir sie uns, wie früher, als bereits in \mathbf{A}_i integriert vorstellen, da sie nur zur Hauptdiagonale beitragen.

$$(10.6) \quad (\mathbf{I} + \sigma \tau (\mathbf{A}_1 + \mathbf{A}_2)) \mathbf{y}_t + (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{y} = \boldsymbol{\varphi}, \quad (\mathbf{A}_i \hat{=} \text{Diskretisierung von } L_i)$$

bzw. für $\sigma = \frac{1}{2}$ (q_i trägt nur zur Hauptdiagonalen bei)

$$(10.7) \quad \left(\mathbf{I} + \frac{\tau}{2} (\mathbf{A}_1 + \mathbf{A}_2) \right) \mathbf{y}_t + (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{y} = \boldsymbol{\varphi}.$$

Wegen $\mathbf{B} = \mathbf{I} + \frac{\tau}{2} \mathbf{A} \geq \frac{\tau}{2} \mathbf{A} + \frac{\varepsilon}{2} \mathbf{I}$ ist die Stabilität für $\varepsilon \leq 2$ problemlos. Das $\frac{\varepsilon}{2} \mathbf{I}$ wird durch den Beitrag von q geliefert. Allerdings ist $\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2$ keine Tridiagonalmatrix mehr, weshalb jeder Iterationsschritt ziemlichen Aufwand erfordert.

Deshalb schlugen Peachman-Rachford statt (10.6) folgendes 2-stufige Verfahren vor

$$(10.8) \quad \begin{aligned} \frac{\mathbf{y}^{j+\frac{1}{2}} - \mathbf{y}^j}{\tau/2} + \mathbf{A}_1 \mathbf{y}^{j+\frac{1}{2}} + \mathbf{A}_2 \mathbf{y}^j &= \boldsymbol{\varphi}^{j+\frac{1}{2}} \\ \frac{\mathbf{y}^{j+1} - \mathbf{y}^{j+\frac{1}{2}}}{\tau/2} + \mathbf{A}_1 \mathbf{y}^{j+\frac{1}{2}} + \mathbf{A}_2 \mathbf{y}^{j+1} &= \boldsymbol{\varphi}^{j+\frac{1}{2}}. \end{aligned}$$

Numerisch geht man bei der Rechnung wie folgt vor:

Die erste Gleichung schreibt man als

$$(*) \quad \left(\frac{2}{\tau} \mathbf{I} + \mathbf{A}_1 \right) \mathbf{y}^{j+\frac{1}{2}} = \mathbf{F} \quad \text{mit} \quad \mathbf{F} = \left(\frac{2}{\tau} - \mathbf{A}_2 \right) \mathbf{y}^j + \boldsymbol{\varphi}^{j+\frac{1}{2}}.$$

Dazu werden die \mathbf{A}_i und entsprechend \mathbf{y}^j zeilenweise von links nach rechts und unten nach oben geordnet, \mathbf{A}_1 wird dadurch zu einer Tridiagonalmatrix. Man berechnet \mathbf{F} und berechnet dann $\mathbf{y}^{j+\frac{1}{2}}$ relativ billig mit dem Tridiagonalverfahren.

Die Matrizen \mathbf{A}_i werden nun spaltenweise von links nach rechts und von unten nach oben umgeordnet. Ebenso wird $\mathbf{y}^{j+\frac{1}{2}}$ entsprechend umgeordnet. Die umgeordneten Größen werden (nicht ganz korrekt) wieder mit \mathbf{A}_i bzw. $\mathbf{y}^{j+\frac{1}{2}}$ bezeichnet.

Die zweite Gleichung aus (10.8) ist äquivalent zu

$$(**) \quad \left(\frac{2}{\tau} + \mathbf{A}_2 \right) \mathbf{y}^{j+1} = \bar{\mathbf{F}} \quad \text{mit} \quad \bar{\mathbf{F}} = \left(\frac{2}{\tau} - \mathbf{A}_1 \right) \mathbf{y}^{j+\frac{1}{2}} + \boldsymbol{\varphi}^{j+\frac{1}{2}}.$$

Man berechnet $\bar{\mathbf{F}}$ bzgl. der neuen Nummerierung der Komponenten von $\mathbf{y}^{j+\frac{1}{2}}$ und danach wieder mit dem Tridiagonalalgorithmus \mathbf{y}^{j+1} . Dieses Spiel wiederholt sich von Zeitschritt zu Zeitschritt.

Bemerkungen:

1. Der Schritt von $t_j \rightarrow t_{j+1}$ wird in 2 Teilschritte zerlegt, der erste in x_1 -Richtung, der zweite in x_2 -Richtung (**alternierende Richtungen**), wobei in jedem Schritt ein Gleichungssystem mit einer **Tridiagonalmatrix** zu lösen ist, was mit geringem Aufwand (vgl. den Abschnitt: Tridiagonalverfahren) möglich ist. Physikalisch bedeutet dieses Vorgehen:
 1. Schritt: Wärmeausbreitung (Diffusion) in x_1 -Richtung,
 2. Schritt: Wärmeausbreitung (Diffusion) in x_2 -Richtung.
2. Beide Gleichungen (10.8) sind Approximationen an die Ausgangsgleichung, jedoch nur von 1. Ordnung. Wir werden jedoch sehen, dass das Gesamtverfahren von 2. Ordnung ist und stabil.

Wir zeigen die Stabilität gleich für die etwas allgemeinere Verfahrensklasse der **Splitting-Verfahren**, wobei wir in (10.8) verschiedene rechte Seiten zulassen, also

$$(10.9) \quad \begin{aligned} (\alpha) \quad & \frac{\mathbf{y}^{j+\frac{1}{2}} - \mathbf{y}^j}{\tau/2} + \mathbf{A}_1 \mathbf{y}^{j+\frac{1}{2}} + \mathbf{A}_2 \mathbf{y}^j = \varphi_1 \\ (\beta) \quad & \frac{\mathbf{y}^{j+1} - \mathbf{y}^{j+\frac{1}{2}}}{\tau/2} + \mathbf{A}_1 \mathbf{y}^{j+\frac{1}{2}} + \mathbf{A}_2 \mathbf{y}^{j+1} = \varphi_2. \end{aligned}$$

Zur Stabilitätsberechnung wird dieses 2-stufige Verfahren auf ein 1-stufiges zurückgeführt durch Elimination von $\mathbf{y}^{j+\frac{1}{2}}$. Danach lassen sich die bekannten Stabilitätssätze anwenden.

Wir eliminieren $\mathbf{y}^{j+\frac{1}{2}}$. Subtraktion von $(\beta) - (\alpha)$ liefert

$$\begin{aligned} \frac{\mathbf{y}^{j+1} - 2\mathbf{y}^{j+\frac{1}{2}} + \mathbf{y}^j}{\tau/2} &= -\mathbf{A}_2 (\mathbf{y}^{j+1} - \mathbf{y}^j) + \varphi_2 - \varphi_1 \\ \mathbf{y}^{j+\frac{1}{2}} &= \frac{\tau}{4} \mathbf{A}_2 (\mathbf{y}^{j+1} - \mathbf{y}^j) - \frac{\tau}{4} (\varphi_2 - \varphi_1) + \frac{1}{2} (\mathbf{y}^{j+1} + \mathbf{y}^j) \end{aligned}$$

Einsetzen in (β) :

$$\frac{\mathbf{y}^{j+1} - \frac{1}{2} (\mathbf{y}^{j+1} + \mathbf{y}^j)}{\tau/2} + \mathbf{A}_1 \mathbf{y}^{j+\frac{1}{2}} + \mathbf{A}_2 \left(\mathbf{y}^{j+1} - \frac{1}{2} (\mathbf{y}^{j+1} - \mathbf{y}^j) \right) + \frac{1}{2} (\varphi_2 - \varphi_1) = \varphi_2$$

Auflösen nach \mathbf{y}_t^j :

$$\frac{\frac{1}{2}(\mathbf{y}^{j+1} - \mathbf{y}^j)}{\tau/2} + \mathbf{A}_1 \left(\frac{\tau}{4} \mathbf{A}_2 \mathbf{y}^{j+1} - \frac{\tau}{4} \mathbf{A}_2 \mathbf{y}^j - \frac{\tau}{4}(\varphi_2 - \varphi_1) + \frac{1}{2}(\mathbf{y}^{j+1} + \mathbf{y}^j) \right) + \mathbf{A}_2 \left(\frac{1}{2}(\mathbf{y}^{j+1} + \mathbf{y}^j) \right) = \frac{1}{2}(\varphi_1 + \varphi_2)$$

$$\frac{\mathbf{y}^{j+1} - \mathbf{y}^j}{\tau} + \frac{1}{2} \mathbf{A}_1 (\mathbf{y}^{j+1} + \mathbf{y}^j) + \frac{1}{2} \mathbf{A}_2 (\mathbf{y}^{j+1} + \mathbf{y}^j) + \frac{\tau}{4} \mathbf{A}_1 \mathbf{A}_2 (\mathbf{y}^{j+1} - \mathbf{y}^j) = \frac{1}{2}(\varphi_1 + \varphi_2) + \frac{\tau}{4} \mathbf{A}_1 (\varphi_1 - \varphi_2)$$

$$\frac{\mathbf{y}^{j+1} - \mathbf{y}^j}{\tau} + \frac{1}{2}(\mathbf{A}_1 + \mathbf{A}_2) (\mathbf{y}^{j+1} - \mathbf{y}^j) + \frac{\tau}{4} \mathbf{A}_1 \mathbf{A}_2 (\mathbf{y}^{j+1} - \mathbf{y}^j) + (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{y}^j = \frac{\tau}{4} \mathbf{A}_1 (\varphi_1 - \varphi_2) + \frac{1}{2}(\varphi_1 + \varphi_2)$$

Wir erhalten mit $\mathbf{y}^{j+1} - \mathbf{y}^j = \tau \mathbf{y}_t^j$

$$(10.10) \quad \underbrace{\left(\mathbf{I} + \frac{\tau}{2}(\mathbf{A}_1 + \mathbf{A}_2) + \frac{\tau^2}{4} \mathbf{A}_1 \mathbf{A}_2 \right)}_{\tilde{\mathbf{B}}} \mathbf{y}_t^j + \underbrace{(\mathbf{A}_1 + \mathbf{A}_2)}_{\mathbf{A}} \mathbf{y}^j = \underbrace{\left(\frac{\tau}{4} \mathbf{A}_1 (\varphi_1 - \varphi_2) + \frac{1}{2}(\varphi_1 + \varphi_2) \right)}_{\varphi}$$

bzw.

$$(10.11) \quad \left(\mathbf{I} + \frac{\tau}{2} \mathbf{A}_1 \right) \left(\mathbf{I} + \frac{\tau}{2} \mathbf{A}_2 \right) \mathbf{y}_t^j + (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{y}^j = \frac{\tau}{4} \mathbf{A}_1 (\varphi_1 - \varphi_2) + \frac{1}{2}(\varphi_1 + \varphi_2).$$

Bemerkung:

Setzt man nicht von vorneherein $\sigma = \frac{1}{2}$ wie beim Übergang von (10.6) zu (10.7), sondern betrachtet das leicht verallgemeinerte Verfahren, das man erhält, wenn man in (10.8), (10.9) $\tau/2$ durch $\tau \sigma$ ersetzt, so lautet die (10.10) entsprechende Form

$$(10.12) \quad (\mathbf{I} + \sigma \tau (\mathbf{A}_1 + \mathbf{A}_2) + \sigma^2 \tau^2 \mathbf{A}_1 \mathbf{A}_2) \mathbf{y}_t^j + (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{y}^j = \frac{\sigma \tau}{2} \mathbf{A}_1 (\varphi_1 - \varphi_2) + \frac{1}{2}(\varphi_1 + \varphi_2).$$

Diese Form liefert bei der Abschätzung des Diskretisierungsfehlers die Möglichkeit, analog zum Vorgehen in Satz 5.2, durch entsprechende Wahlen von σ und φ höhere Konvergenzgeschwindigkeiten zu erreichen.

Nun ist (10.12) wieder ein Verfahren der Art

$$(10.13) \quad \tilde{\mathbf{B}} \mathbf{y}_t + \mathbf{A} \mathbf{y} = \varphi, \quad \tilde{\mathbf{B}} = \mathbf{I} + \sigma \tau \mathbf{A} + \sigma^2 \tau^2 \mathbf{A}_1 \mathbf{A}_2$$

$$\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2$$

$$\varphi = \frac{\sigma \tau}{2} \mathbf{A}_1 (\varphi_1 - \varphi_2) + \frac{1}{2}(\varphi_1 + \varphi_2)$$

(wobei $\varphi_1 = \varphi_2 = \mathbf{f}^{j+\frac{1}{2}}$ üblicherweise)

Da $\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2$ symmetrisch und positiv definit ist, liegt Stabilität vor falls $\tilde{\mathbf{B}} \geq \frac{\tau}{2} \mathbf{A} + \frac{\varepsilon}{2} \mathbf{I}$ ist (vgl. Satz 4.2), bzw.

$$(10.14) \quad \mathbf{I} + \tau \left(\sigma - \frac{1}{2} \right) \mathbf{A} + \sigma^2 \tau^2 \mathbf{A}_1 \mathbf{A}_2 \geq \frac{\varepsilon}{2} \mathbf{I}.$$

Diese Bedingung ist für $\sigma \geq \frac{1}{2}$ sicher erfüllt für $\varepsilon \leq 2$ falls $\mathbf{A}_1 \mathbf{A}_2 \geq 0$ ist.

Letztere Eigenschaft werden wir gleich zeigen. Sie bringt jedoch einschneidende Voraussetzungen bzgl. der Wahl von Ω mit sich (vgl. Satz 10.1), weshalb wir darauf hinweisen, dass in der Wahl von $\varepsilon = 2$ noch „Luft“ liegt. Wir zeigen nun

Satz 10.1

In einem Rechteck $\bar{\Omega} = [0, a_1] \times [0, a_2] \subset \mathbb{R}^2$ seien \mathbf{A}_i die Diskretisierungsmatrizen von $L_i u = \frac{\partial}{\partial x_i} \left(k_i \frac{\partial}{\partial x_i} u \right) - q_i u$, $i = 1, 2$.

Es sei $k_i(x) = k_i(x_i)$, $q_i(x) = q_i(x_i)$ und $\mathbf{A}_i = \mathbf{A}_i^T > 0$, $i = 1, 2$.

Dann gilt:

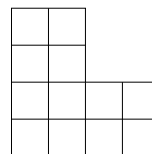
- a) $\mathbf{A}_1 \mathbf{A}_2 = \mathbf{A}_2 \mathbf{A}_1$ und
- b) $(\mathbf{A}_1 \mathbf{A}_2 \mathbf{y}, \mathbf{y}) \geq 0 \quad \forall \mathbf{y} \in \mathbb{R}^n$.

Bemerkungen:

1. Dass $k_i(x) = k_i(x_i)$, d.h. die Diffusionskonstante $k_1(x)$ nur von x_1 - nicht von x_2 - abhängt, entsprechend für k_2 , ist realistisch wenn z.B. $\bar{\Omega}$ ein Schnitt in einer geschichteten Fläche darstellt.
2. Entsprechende Diskretisierungsmatrizen wurden in § 7, (7.3) hergeleitet. Dazu beachte man, dass gemäß § 8 q_i nur einen Beitrag zur Hauptdiagonalen von \mathbf{A}_i liefert und somit weder Symmetrie noch Definitheit stört.
3. Die Voraussetzung eines Rechteckgebiets ist wesentlich, wie folgende Aufgabe zeigt.

Aufgabe:

1. Man beweise a) aus Satz 10.1
2. Für $L_1 u = u_{x_1 x_1}$, $L_2 u = u_{x_2 x_2}$ und das L -Gebiet, samt eingezeichnetem Gitter zeige man, dass a) für die entsprechenden Matrizen \mathbf{A}_i nicht gilt.



3. Für $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{A}^T = \mathbf{A} > 0$ gibt es genau eine Wurzel $\mathbf{A}^{1/2}$ (d.h. $\mathbf{A}^{1/2} \mathbf{A}^{1/2} = \mathbf{A}$) mit $\mathbf{A}^{1/2} = (\mathbf{A}^{1/2})^T > 0$.

4. Das Verfahren (Newton-Verfahren) für $\mathbf{A} = \mathbf{A}^T > 0$

$$(10.15) \quad \mathbf{X}_{n+1} = \frac{1}{2} (\mathbf{X}_n + \mathbf{X}_n^{-1} \mathbf{A}), \quad \mathbf{X}_1 = \mathbf{I}$$

ist durchführbar und konvergiert gegen das symmetrische $\mathbf{A}^{1/2} > 0$. Dazu gehört insbesondere der Nachweis der Existenz der \mathbf{X}_n^{-1} .

5. $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, $\det \mathbf{A} \neq 0 \implies \mathbf{AB}$ und \mathbf{BA} haben dieselben Eigenwerte.

Beweis Satz 10.1

a) laut Aufgabe

b) wird aus a) gefolgert.

Wir zeigen zunächst für beliebige Matrizen $\mathbf{A} = \mathbf{A}^T > 0$, $\mathbf{B} = \mathbf{B}^T > 0$:

$$(10.16) \quad \mathbf{AB} = \mathbf{BA} \implies \sqrt{\mathbf{A}}\mathbf{B} = \mathbf{B}\sqrt{\mathbf{A}} \implies \sqrt{\mathbf{A}}\sqrt{\mathbf{B}} = \sqrt{\mathbf{B}}\sqrt{\mathbf{A}}$$

Für das Verfahren (10.15)

$$(10.17) \quad \begin{aligned} \mathbf{X}_{n+1} &= \frac{1}{2} (\mathbf{X}_n + \mathbf{X}_n^{-1} \mathbf{A}), \quad \mathbf{X}_1 = \mathbf{I} \\ &\text{gilt} \end{aligned}$$

$$\mathbf{X}_n \mathbf{B} = \mathbf{B} \mathbf{X}_n \quad \forall n.$$

Die Behauptung ist richtig für $n = 1$, da $\mathbf{X}_1 = \mathbf{I}$.

Sie sei richtig für n , dann gilt

$$\mathbf{X}_{n+1} \mathbf{B} = \frac{1}{2} (\mathbf{X}_n \mathbf{B} + \mathbf{X}_n^{-1} \mathbf{A} \mathbf{B}) \stackrel{\text{a)}}{=} \frac{1}{2} (\mathbf{X}_n \mathbf{B} + \mathbf{X}_n^{-1} \mathbf{B} \mathbf{A})$$

und mit der Induktionsvoraussetzung folgt wegen

$$\mathbf{B} \mathbf{X}_n = \mathbf{X}_n \mathbf{B} \iff \mathbf{B} = \mathbf{X}_n \mathbf{B} \mathbf{X}_n^{-1} \iff \mathbf{X}_n^{-1} \mathbf{B} = \mathbf{B} \mathbf{X}_n^{-1}$$

also

$$\mathbf{X}_{n+1} \mathbf{B} = \frac{1}{2} (\mathbf{B} \mathbf{X}_n + \mathbf{B} \mathbf{X}_n^{-1} \mathbf{A}) = \mathbf{B} \frac{1}{2} (\mathbf{X}_n + \mathbf{X}_n^{-1} \mathbf{A}) = \mathbf{B} \mathbf{X}_{n+1}.$$

Wegen $\mathbf{X}_n \rightarrow \sqrt{\mathbf{A}}$ folgt aus (10.5) die erste Behauptung (10.16): $\sqrt{\mathbf{A}}\mathbf{B} = \mathbf{B}\sqrt{\mathbf{A}}$. Unter Benutzung dieser Behauptung zeigt man auf dieselbe Weise, in dem man das Verfahren (10.15) zur Berechnung von $\sqrt{\mathbf{B}}$ aufstellt, zuerst $\mathbf{X}_{n+1} \sqrt{\mathbf{A}} = \sqrt{\mathbf{A}} \mathbf{X}_{n+1}$ und durch Grenzübergang schließlich $\sqrt{\mathbf{A}}\sqrt{\mathbf{B}} = \sqrt{\mathbf{B}}\sqrt{\mathbf{A}}$, also Behauptung (10.16).

Wie wenden (10.16) an auf $\mathbf{A}_1, \mathbf{A}_2$ und erhalten

$$\begin{aligned} (\mathbf{A}_1 \mathbf{A}_2 \mathbf{y}, \mathbf{y}) &= (\mathbf{A}_1^{1/2} \mathbf{A}_1^{1/2} \mathbf{A}_2^{1/2} \mathbf{A}_2^{1/2} \mathbf{y}, \mathbf{y}) = (\mathbf{A}_1^{1/2} \mathbf{A}_2^{1/2} \mathbf{A}_1^{1/2} \mathbf{A}_2^{1/2} \mathbf{y}, \mathbf{y}) = \\ &= (\mathbf{A}_1^{1/2} \mathbf{A}_2^{1/2} \mathbf{y}, \mathbf{A}_2^{1/2} \mathbf{A}_1^{1/2} \mathbf{y}) = (\mathbf{A}_1^{1/2} \mathbf{A}_2^{1/2} \mathbf{y}, \mathbf{A}_1^{1/2} \mathbf{A}_2^{1/2} \mathbf{y}) = \\ &= \|\mathbf{A}_1^{1/2} \mathbf{A}_2^{1/2} \mathbf{y}\|^2 > 0 \quad \forall \mathbf{y} \neq 0. \end{aligned}$$

□

Zusammenfassend gilt

Satz 10.2 Stabilität des ADI-Verfahrens (10.9)

$\mathbf{A}_i = \mathbf{A}_i^T > 0$ seien die Diskretisierungsmatrizen von $L_i u = \frac{\partial}{\partial x_i} \left(k_i \frac{\partial u}{\partial x_i} \right) + q_i$, $i = 1, 2$ mit $k_i(\mathbf{x}) = k_i(x_i)$, $q_i(\mathbf{x}) = q_i(x_i)$ **in einem Rechteck.**

Dann ist das ADI-Verfahren (10.9) bzw. (10.10) stabil bzgl. Anfangswerten und rechter Seite (Satz 4.2 und $\varepsilon \leq 2$)

$$\|\mathbf{y}^j\|_{(1,h)} \leq \|\mathbf{y}^0\|_{(1,h)} + \frac{1}{\sqrt{\varepsilon}} \left(\sum_{k=0}^j \tau \|\varphi^k\|_{(0,h)}^2 \right)^{1/2} \quad \text{mit } \varphi = \frac{\tau}{4} \mathbf{A}_1 (\varphi_1 - \varphi_2) + \frac{1}{2} (\varphi_1 + \varphi_2)$$

Konvergenz der mehrdimensionalen Verfahren

a) das Verfahren (10.1)

Sind die \mathbf{A}_i die Diskretisierungsmatrizen 2. Ordnung der L_i ($i = 1, 2$), dann lautet die Diskretisierung der Aufgabe (10.1) (vgl. (4.6), (4.7) und (10.4), (10.6)) für ein geeignetes φ :

$$(10.18) \quad \underbrace{(\mathbf{I} + \sigma \tau (\mathbf{A}_1 + \mathbf{A}_2))}_{\mathbf{B}} \mathbf{y}_t + \underbrace{(\mathbf{A}_1 + \mathbf{A}_2)}_{\mathbf{A}} \mathbf{y} = \mathbf{I} \mathbf{y}_t + \mathbf{A} \mathbf{y}^{(\sigma)} = \varphi$$

Für den Verfahrensfehler $\mathbf{z} = \mathbf{u} - \mathbf{y}$ dieser Diskretisierung gilt dann analog zum Beweis Satz 5.2

$$\begin{aligned} \mathbf{z}_t + \mathbf{A} \mathbf{z}^{(\sigma)} &= \mathbf{u}_t + \mathbf{A} \mathbf{u}^{(\sigma)} - \underbrace{(\mathbf{y}_t + \mathbf{A} \mathbf{y}^{(\sigma)})}_{=\varphi} = \mathbf{u}_t + \mathbf{A} \mathbf{u}^{(\sigma)} - \varphi \\ &= \underbrace{\left(\frac{\partial u}{\partial t} - L u - f \right)}_{=0}^{j+\frac{1}{2}} + O \left(\tau \left(\sigma - \frac{1}{2} \right) + \tau^2 + h_1^2 + h_2^2 \right) \quad \text{für } \varphi = \mathbf{f}^{j+\frac{1}{2}}. \end{aligned}$$

vgl. die Abschätzungen zu Satz 5.2 wobei \mathbf{A}_i geeignete Diskretisierungsmatrizen 2. Ordnung sind. Hieraus folgt die Konvergenz von der Ordnung des Diskretisierungsfehlers (vgl. Satz 5.4.)

Wählt man $h = h_1 = h_2$, $\sigma = \frac{1-h^2}{12}$ und $\varphi = \left(\mathbf{I} + \sum_{i=1}^2 \frac{h^2}{12} L_i \right) \mathbf{f}$, so erhält man analog zu Satz 5.4 eine Konvergenz der Ordnung $O(\tau^2 + h^4)$.

b) das ADI-Verfahren

Für den Verfahrensfehler des Verfahrens (10.12) folgt mit

$$\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2, \quad \tilde{\mathbf{B}} = \mathbf{I} + \sigma, \quad \tau \mathbf{A} + \tau^2 \sigma^2 \mathbf{A}_1 \mathbf{A}_2, \quad \mathbf{B} = \mathbf{I} + \sigma \tau \mathbf{A}$$

unter Verwendung der Abschätzung aus Teil a):

$$\begin{aligned}\tilde{\mathbf{B}}\mathbf{z} + \mathbf{A}\mathbf{z} &= \mathbf{B}\mathbf{z} + \mathbf{A}\mathbf{z} + \tau^2\sigma^2 \mathbf{A}_1 \mathbf{A}_2 \mathbf{z} \\ &= \mathbf{B}\mathbf{u}_t + \mathbf{A}\mathbf{u} + \tau^2\sigma^2 \mathbf{A}_1 \mathbf{A}_2 \mathbf{u} - \underbrace{(\mathbf{B}\mathbf{y}_t + \tau^2\sigma^2 \mathbf{A}_1 \mathbf{A}_2 \mathbf{y}_t + \mathbf{A}\mathbf{y})}_{\varphi} \\ &= \underbrace{\left(\frac{\partial u}{\partial t} - L u - f\right)}_{=0}^{j+\frac{1}{2}} + \tau^2\sigma^2 \mathbf{A}_1 \mathbf{A}_2 \mathbf{u}_t + O\left(\tau\left(\sigma - \frac{1}{2}\right) + \tau^2 + h_1^2 + h_2^2\right)\end{aligned}$$

Beachte: Das φ ist dasselbe wie in (10.8).

Ist $\mathbf{A}_1 \mathbf{A}_2 \mathbf{u}_t$ beschränkt (das kann man zeigen, wenn $u \in C^6$, beachte: $\mathbf{A}_1 \mathbf{A}_2 \mathbf{u}_t$ ist Diskretisierung einer Ableitung 5. Ordnung), so erhält man über den Stabilitätssatz dieselbe Stabilitäts- und Konvergenzordnung wie für (10.8).

Diese Überlegungen müssen (und können) für $h_1 \neq h_2$ etwas modifiziert werden. Insgesamt liefert das Vorgehen von Satz 5.2 dann

Satz 10.3 Konvergenz der ADI-Methode (Peachman/Rachford)

Für die Aufgabe

$$\begin{aligned}\frac{\partial u}{\partial t} &= \sum_{i=1}^2 L_i u + f \quad \text{in } \Omega = (0, a_1) \times (0, a_2) \subset \mathbb{R}^2 \quad \text{mit} \\ L_i u &= \frac{\partial}{\partial x_i} \left(k_i \frac{\partial}{\partial x_i} u \right) + q_i u, \quad k_i(\mathbf{x}) = k_i(x_i), \quad q_i(\mathbf{x}) = q_i(x_i)\end{aligned}$$

seien \mathbf{A}_i geeignete symmetrische, positiv definite Diskretisierungsmatrizen 2. Ordnung für L_i , $i = 1, 2$ (vgl. dazu z.B. (7.2) und § 8, 1. Abschnitt). Dann gilt:

Das ADI-Verfahren

$$\begin{aligned}\frac{\mathbf{y}^{j+\frac{1}{2}} - \mathbf{y}^j}{\tau\sigma} + \mathbf{A}_1 \mathbf{y}^{j+\frac{1}{2}} + \mathbf{A}_2 \mathbf{y}^j &= \varphi \\ \frac{\mathbf{y}^{j+1} - \mathbf{y}^{j+\frac{1}{2}}}{\tau\sigma} + \mathbf{A}_1 \mathbf{y}^{j+\frac{1}{2}} + \mathbf{A}_2 \mathbf{y}^{j+\frac{1}{2}} &= \varphi\end{aligned}$$

konvergiert von der Ordnung

$O\left(\tau\left(\sigma - \frac{1}{2}\right) + \tau^2 + h_1^2 + h_2^2\right)$ falls $\varphi = \mathbf{f}^{j+\frac{1}{2}}$, $u \in C^{3,4}(\Omega)$ (3 bzgl. Zeit, 4 bzgl. Ort)

$O(\tau^2 + h^4)$ für $h = h_1 = h_2$, $\sigma = \frac{1}{12} - \frac{h^2}{12\tau}$, $\varphi = \left(\mathbf{I} + \sum_{i=1}^2 \frac{h^2}{12} L_i\right) \mathbf{f}$ und $u \in C^{3,6}(\Omega)$

$O(\tau^2 + h^6)$ für noch kompliziertere Wahlen von σ und φ und hinreichend hohe Differenzierbarkeitsordnung.

Bemerkungen

1. Damit man im 2. Schritt von (10.8) das Tridiagonalverfahren anwenden kann, muß zuerst das Gitter (und damit auch \mathbf{A}_1) spaltenweise unnummeriert werden.

2. Die Verallgemeinerung von (10.11) auf den nD -Fall lautet:

$$\prod_{i=1}^n \left(\mathbf{I} - \frac{\tau}{2} \mathbf{A}_i \right) \mathbf{y}_t + \sum_{i=1}^n \mathbf{A}_i \mathbf{y} = \boldsymbol{\varphi}.$$

Stabilität und Konvergenz müssen neu untersucht werden (Arbeiten von Janenko, Fairwather, Mitchel, Dyakonov - Mitte der 60er Jahre).

3. Unter mehreren Varianten des Verfahrens erwähnen wir eine von Janenko:

$$(\mathbf{I} + \sigma \tau \mathbf{A}_i) \mathbf{y}^{j+\frac{i}{n}} = (\mathbf{I} - (1 - \sigma) \tau \mathbf{A}_i) \mathbf{y}^{j+\frac{i-1}{n}} + \underset{\substack{\uparrow \\ \text{Kronecker}}}{\delta_{i1}} \boldsymbol{\varphi}, \quad i = 1, \dots, n,$$

die auch mit tridiagonalen Matrizen arbeitet (wie auch weitere Varianten von Stoyan und Dyakonov). Bemerkenswert ist, dass die Zwischenschritte dieser Verfahren die ursprüngliche Aufgabe nicht mehr approximieren. Die Approximation kommt erst nach Elimination der Zwischenschritte zustande.

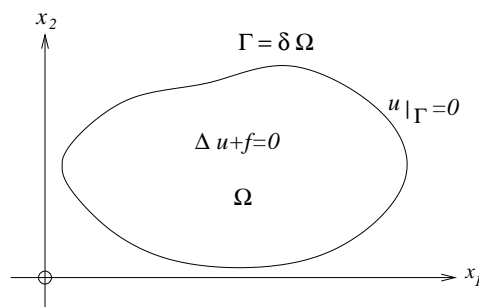
4. Man kann die Stabilität und Konvergenz der ADI-Verfahren auch für nicht vertauschbare Matrizen und andere als Rechtecksgebiete zeigen, muß dann aber Einbußen in der Ordnung in Kauf nehmen.
5. Der Stabilitätssatz 10.2 läßt sich auch mit Hilfe einer diskreten Fourier-Analyse beweisen (vgl. hierzu die einfache Demonstration des Verfahrens zum Beweis von (4.2)).
6. ADI-Verfahren lassen sich auch zur Behandlung elliptischer Probleme anwenden.
7. Literatur und einige Varianten der ADI-Verfahren (für Rechtecksgebiete) findet man in Morton/Mayers: Num. Sol. of PDEs, Cambridge Univ. Press, Abschnitt 3.2 f.

Kapitel II

Elliptische Gleichungen

§ 11 Die Poissongleichung - Einleitung

Wir stellen 2 einfache Motivationsbeispiele für die Poissonaufgabe vor:



Ein einfaches Beispiel liefert eine stationäre Flüssigkeitsströmung.

Ist $\mathbf{u} = (u, v)^T$ der Strömungsvektor, so bedeuten

$$-\frac{\partial u}{\partial y} + \frac{\partial v}{\partial x} = 0 \quad \text{die Wirbelfreiheit}$$

$$\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} = 0 \quad \text{das Nichtvorhandensein von} \\ \text{Massenquellen (Massenerhaltung)}$$

Beide Gleichungen zusammen bilden das Cauchy-Riemann'sche System (vgl. Ableitungsbedingungen für Real- und Imaginärteil komplexwertiger Funktionen).

Wird die 1. Gleichung nach x , die zweite nach y differenziert und danach die Differenz gebildet, so erhält man

$$\Delta v = 0$$

Vertauschung von x und y liefert $\Delta u = 0$.

Kompliziertere Anwendungen folgen aus den Navier-Stokes-Gleichungen (wir beschränken uns auf den 2-dimensionalen Fall). Anwendungen z.B. im Umweltschutz, z.B. Chemieunfall.

Seien $\mathbf{u} = (u, v)^T$ der Strömungsvektor, p der Druck, ν die kinematische Viskosität der Flüssigkeit (sie ist bei Gasen klein und kann bei Flüssigkeiten sehr groß sein) und f_i einwirkende Kräfte, so stellen die ersten beiden der folgenden Gleichungen die Impulserhaltung für die 1. bzw. 2. Komponente des Geschwindigkeitsvektors dar, die 3. Gleichung beschreibt die Massenerhaltung.

$$\begin{aligned}\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} + \frac{\partial p}{\partial x} &= \nu \Delta u + f_1 \\ \frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} + \frac{\partial p}{\partial y} &= \nu \Delta v + f_2 \\ \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} &= 0.\end{aligned}$$

Vereinfachungsmöglichkeiten:

1. Bei kleinen Geschwindigkeitsänderungen ($\frac{\partial u}{\partial t}, \frac{\partial v}{\partial t} = 0$) und bei kleinen Geschwindigkeiten (\Rightarrow Vernachlässigung der nicht linearen Terme) erhält man eine Poissongleichung für die Geschwindigkeitskomponenten.
2. Daß auch die einfache Poissongleichung für die Anwendung interessant sein kann, zeigt folgende mathematische Umrechnung der Navier-Stokes-Gleichungen: Differenziert man (bei konstantem ν) die 1. Gleichung partiell nach x , die zweite nach y , vertauscht die Differentiationsreihenfolge und addiert die ersten beiden Gleichungen, so erhält man

$$\begin{aligned}\frac{\partial}{\partial t} \underbrace{\left(\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right)}_{=0} + \left(\frac{\partial u}{\partial x} \right)^2 &+ \frac{\partial v}{\partial x} \frac{\partial u}{\partial y} \left[\begin{array}{c} + u \frac{\partial^2 u}{\partial x^2} \\ + u \frac{\partial^2 v}{\partial x \partial y} \end{array} \right] + \left(\frac{\partial v}{\partial y} \right)^2 &+ \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} \left[\begin{array}{c} + v \frac{\partial^2 u}{\partial x \partial y} \\ + v \frac{\partial^2 v}{\partial y^2} \end{array} \right] + \Delta p &= \nu \Delta \underbrace{\left(\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right)}_{=0} \\ &+ \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} &+ \frac{\partial}{\partial x} f_1 + \frac{\partial}{\partial y} f_2 \\ &+ u \frac{\partial}{\partial x} \underbrace{\left(\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right)}_{=0} &+ v \frac{\partial}{\partial y} \underbrace{\left(\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right)}_{=0}\end{aligned}$$

also folgt

$$\left(\frac{\partial u}{\partial x} \right)^2 + \frac{\partial v}{\partial x} \frac{\partial u}{\partial y} + \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} + \left(\frac{\partial v}{\partial y} \right)^2 + \Delta p = \frac{\partial f_1}{\partial x} + \frac{\partial f_2}{\partial y}$$

d.h. bei bekannten Geschwindigkeiten erhält man eine Poissongleichung für den Druck.

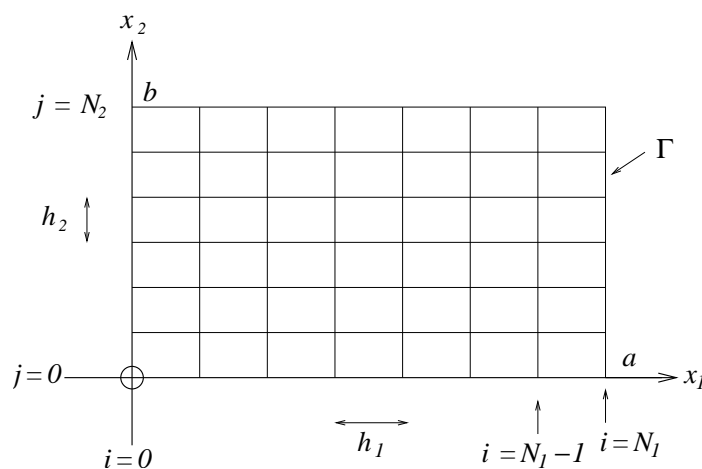
§ 12 Die erste RWA für die Poissongleichung im Rechteck

Für die Aufgabe

$$(12.1) \quad \begin{array}{ll} \Delta u + f = 0 & \text{in } \Omega = (0, a) \times (0, b) \quad (\text{Poissongleichung}) \\ u = g & \text{auf } \Gamma, \quad \Gamma = \delta\Omega \quad (\text{Dirichletwerte}) \end{array}$$

haben wir die einfachste Diskretisierung schon im § 10 beschrieben.

Wir wiederholen die Bezeichnungen



$$\begin{aligned} h_1 &= \frac{a}{N_1}, & h_2 &= \frac{b}{N_2}, & h &= (h_1, h_2) \\ \omega_h &= \text{Menge der inneren Gitterpunkte} \\ \gamma_h &= \text{Menge der Randpunkte} \\ \overline{\omega_h} &= \omega_h \cup \gamma_h = \text{Menge aller Gitterpunkte} \end{aligned}$$

Wir benutzen folgende **Abkürzungen**:

$\|\mathbf{f}\|_{C(\omega_h)} = \max_{x \in \omega_h} |\mathbf{f}(x)|$ Maximumnorm der von f erzeugten Gitterfunktion \mathbf{f} auf den inneren Gitterpunkten des Gitters ω_h mit den Maschenweiten $h = (h_1, h_2)$.

$\|\mathbf{f}\|_{C(\overline{\omega_h})}$, $\|\mathbf{g}\|_{C(\gamma_h)}$ entsprechend.

y_{ij} , $(ij) \in \omega_h$ Doppelindizierung der Punkte zu ihrer Festlegung im Gitter und „ $(ij) \in \omega_h$ “ (z.B.) als Zugehörigkeit zu den inneren Gitterpunkten.

w_{ij} als Komponente einer durch eine Funktion \mathbf{w} erzeugten Gitterfunktion.

Zur Aufstellung der Diskretisierungsmatrizen werden die Gitterpunkte zeilenweise durchnummeriert, beginnend mit $j = 0, i = 0, \dots, N_1, j = 1, i = 0, \dots, N_1$ usw. Für die Näherungswerte der exakten Lösung u ist als Bezeichnung üblich

$$u(x_i, x_j) \approx y_{ij}, \quad \text{wobei } x_i = i h_1, \quad x_j = j h_2.$$

Die Diskretisierung von (12.1) in inneren Gitterpunkten lautet

$$(12.2) \quad \begin{aligned} y_{x_1 x_1, ij} &+ y_{\bar{x}_2 \bar{x}_2, ij} &+ f_{ij} &= 0 \\ \frac{1}{h_1^2} (y_{i-1, j} - 2y_{ij} + y_{i+1, j}) &+ \frac{1}{h_2^2} (y_{i, j-1} - 2y_{ij} + y_{i, j+1}) &+ f_{ij} &= 0 \end{aligned}$$

wobei für die Ableitungen gemäß (2.5) gilt (bzgl. jeder Variablen in jedem Punkt)

$$(12.3) \quad u_{\bar{x}x} = \begin{cases} u'' + \frac{h^2}{24} (u_+^{(4')} + u_-^{(4')}) & \text{falls } u \in C^4 \\ u'' + \frac{h^2}{12} u^{(4')} + \frac{h^4}{720} (u_+^{(6')} + u_-^{(6')}) & \text{falls } u \in C^6 \end{cases}$$

Die Restglieder kann man wie folgt abschätzen: Ist $u \in C^4$, also stetig, so folgt aus dem Zwischensatz: $\exists \xi \in [\tilde{x}_1, \tilde{x}_2]: u^{(4')}(\tilde{x}_1) + u^{(4')}(\tilde{x}_2) = 2u^{(4')}(\xi)$.

Damit erhält man

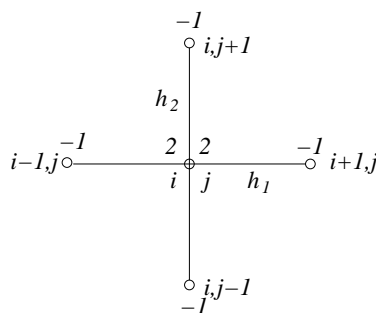
$$(12.4) \quad \left| \frac{h^2}{24} (u_+^{(4')} + u_-^{(4')}) \right| \leq \frac{h^2}{12} M_4, \quad M_4 = \max_x |u^{(4)}(x)|.$$

Analog kann man im Fall $u \in C^6$ verfahren.

Wir wollen die Differenzgleichungen für alle Gitterpunkte in Matrixform aufschreiben unter Benutzung des 5-Punkte-Sterns (vgl.(12.2)):

$$(12.5) \quad \mathbf{A}_h \mathbf{y} = \boldsymbol{\varphi}$$

Beachte die Vorzeichen: In der Differentialgleichung stehen Δu und f auf derselben Seite. In der Matrixdarstellung stehen sie auf verschiedenen Seiten.



Konvergenz und Apriori Schranke für die Lösung von $A_h \mathbf{y} = \varphi$.

Satz 12.1
 Für die Lösung u von (12.1)

$$\Delta u + f = 0 \text{ in } \Omega = (0, a) \times (a, b) \subset \mathbb{R}^2, \quad u|_{\partial\Omega} = g,$$

betrachten wir das Verfahren $A_h \mathbf{y} = \varphi$, (12.5), mit

$$y_{\bar{x}_1 x_1, ij} + y_{\bar{x}_1 x_2, ij} + f_{ij} = 0, \quad (i, j) \in w_h$$

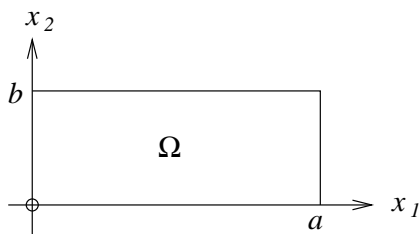
Dann gilt:

- a) Die Diskretisierungsmatrix A_h aus (12.5) ist eine M -Matrix.
- b) Das Verfahren konvergiert quadratisch ($O(h_1^2 + h_2^2)$)
- c) Es gilt die Apriori Abschätzung

$$(12.6) \quad \|\mathbf{y}\|_{C(\bar{w}_h)} \leq \left(1 + \frac{a^2 + b^2}{16}\right) \max(\|\mathbf{g}\|_{C(\gamma_h)}, \|\mathbf{f}\|_{C(\bar{w}_h)})$$

Beweis a)

Die Vorzeichenverteilung einer M -Matrix wurde für A_h schon gezeigt. Zur Konstruktion eines Vektors $\mathbf{w} > \mathbf{0}$ mit $A_h \mathbf{w} > \mathbf{0}$ (elementweise vgl. Definition 6.2) machen wir den Ansatz



$$w(x_1, x_2) = 4 + x_1(a - x_1) + x_2(b - x_2) \text{ auf } \bar{\Omega}.$$

Offensichtlich gilt: $w \geq 4$ in $[0, a] \times [0, b]$.

Nun ist

$$(12.7) \quad (A_h \mathbf{w})_{ij} = \begin{cases} w_{ij}, & (i, j) \in \gamma_h & \text{(Randwerte)} \\ -(w_{\bar{x}_1 x_1} + w_{\bar{x}_2 x_2})_{ij}, & (i, j) \in w_h & \text{(innere Punkte)} \end{cases}$$

Gemäß (12.3): $u_{\bar{x}_1 x_1} = \frac{\partial^2}{\partial x_1^2} u + \frac{h^2}{12} \frac{\partial^4}{\partial x_1^4} u(x_1 + \theta h_1, x_2)$, entsprechend für x_2 , gilt für die quadratische Funktion w :

$$w_{\bar{x}_1 x_1} = \frac{\partial^2}{\partial x_1^2} w, \quad w_{\bar{x}_2 x_2} = \frac{\partial^2}{\partial x_2^2} w, \quad \text{da } \frac{\partial^4}{\partial x_1^4} w = \frac{\partial^4}{\partial x_2^4} w = 0.$$

Weiter ist $w_{\bar{x}_1 x_1} + w_{\bar{x}_2 x_2} = -4$, somit folgt aus (12.7)

$$(12.8) \quad (\mathbf{A}_h \mathbf{w})_{ij} \begin{cases} \geq 4 & (i, j) \in \gamma_h, \\ = 4 & (i, j) \in \omega_h \end{cases}$$

Bezeichnet \mathbf{w} die durch w erzeugte Gitterfunktion, so gilt $\mathbf{w} \geq 4\mathbf{e}$, $\mathbf{e} = (1, 1, \dots, 1)^T$, und nach (12.8) $(\mathbf{A}_h \mathbf{w})_{ij} \geq 4 > 0$, bzw. $\mathbf{A}_h \mathbf{w} > \mathbf{0}$ (elementweise).

$\implies \mathbf{A}_h$ ist eine M -Matrix

Beweis b)

$$\begin{aligned} \text{Es ist } \|\mathbf{w}\|_{C(\bar{\omega})} &= 4 + \frac{a^2 + b^2}{4} \quad (\text{Ableiten, Maximum berechnen}) \\ \min_i (\mathbf{A} \mathbf{w})_i &= 4 \quad \text{laut (12.8)} \end{aligned}$$

Deshalb liefert Satz 6.3

$$\|\mathbf{A}_h^{-1}\|_{\infty} \leq \frac{\|\mathbf{w}\|_{\infty}}{\min_i (\mathbf{A} \mathbf{w})_i} \leq \frac{4 + \frac{a^2 + b^2}{4}}{4} = 1 + \frac{a^2 + b^2}{16}.$$

Der Verfahrensfehler $\mathbf{z} = \mathbf{u} - \mathbf{y}$ hat Nullrandwerte. Man kann das System (12.5) also beschränken auf die Matrix \mathbf{A}_h^0 , die man aus \mathbf{A}_h durch Streichen der Randterme erhält. \mathbf{A}_h^0 wirkt dann nur auf den Gittervektor $\overset{\circ}{\mathbf{y}}$ der inneren Punkte. Der Verfahrensfehler genügt dann dem Differenzenschema (mit $\mathbf{f}^0 = \mathbf{f}|_{\omega_h}$)

$$\mathbf{A}_h^0 \mathbf{z} = \mathbf{A}_h^0 \mathbf{u} - \underbrace{\mathbf{A}_h^0 \mathbf{y}}_{=\mathbf{f}^0} = \mathbf{u}_{\bar{x}_1 x_1} + \mathbf{u}_{\bar{x}_2 x_2} - \mathbf{f}^0 = \underbrace{(\Delta u - f)}_{=0} + \mathcal{O}(h_1^2 + h_2^2) =: \boldsymbol{\psi}$$

Aus $\mathbf{A}_h \mathbf{z} = \boldsymbol{\psi}$ ($\boldsymbol{\psi}, \mathbf{z}$ sind hier um die Nullkomponenten der Randterme angereichert) folgt somit (siehe oben)

$$\|\mathbf{z}\|_{\infty} \leq \left(1 + \frac{a^2 + b^2}{16}\right) \mathcal{O}(h_1^2 + h_2^2).$$

Beweis c)

Aus $\mathbf{A}_h \mathbf{y} = \boldsymbol{\varphi}$, $\boldsymbol{\varphi}$ gemäß (12.5), erhält man ebenso

$$\begin{aligned} \|\mathbf{y}\|_{C(\bar{\omega}_h)} &\leq \|\mathbf{A}_h^{-1}\|_{\infty} \max(\|\mathbf{g}\|_{C(\gamma_h)}, \|\mathbf{f}\|_{C(\bar{\omega}_h)}) \\ &\leq \left(1 + \frac{a^2 + b^2}{16}\right) \max(\|\mathbf{g}\|_{C(\gamma_h)}, \|\mathbf{f}\|_{C(\bar{\omega}_h)}). \end{aligned} \quad \blacksquare$$

Bemerkung

Diese Abschätzung bedeutet auch die Stabilität für das Verfahren (12.5), denn die Konstante auf der rechten Seite, also die Abschätzung von $\|\mathbf{A}_h^{-1}\|_{\infty}$, ist unabhängig von der Diskretisierung.

Schärfere Abschätzungen kann man mit Hilfe des diskreten Maximumprinzips erhalten, das wir nun beweisen.

Das diskrete Maximumprinzip

Definition 12.2 Maximumprinzip für Matrizen

Für eine Matrix \mathbf{A} und einen Vektor $\mathbf{y} \in \mathbb{R}^n$ definieren wir

$$N^0(\mathbf{A}\mathbf{y}) = \{i; (\mathbf{A}\mathbf{y})_i = 0, i \in \{1, \dots, n\}\}$$

$$N^\neq(\mathbf{A}\mathbf{y}) = \{i; (\mathbf{A}\mathbf{y})_i \neq 0, i \in \{1, \dots, n\}\}$$

\mathbf{A} genügt dem (strengen) Maximumprinzip, falls gilt

a) $N^\neq(\mathbf{A}\mathbf{y}) = \emptyset \implies \mathbf{y} = \mathbf{0} \forall \mathbf{y} \in \mathbb{R}^n$ (d.h. \mathbf{A} ist regulär)

b) $\max_{i \in N^0(\mathbf{A}\mathbf{y})} |y_i| < (\text{bzw. } \leq) \max_{j \in N^\neq(\mathbf{A}\mathbf{y})} |y_j| \forall \mathbf{y} \in \mathbb{R}^n, \mathbf{y} \neq \mathbf{0}$

d.h. in Worten: Dort wo die rechte Seite $\neq 0$ ist, wird das Maximum angenommen. Das Maximum der restlichen Komponenten ist kleiner oder (im nichtstrengen Fall) höchstens gleich groß.

Satz 12.3

Das strenge Maximumprinzip gilt genau dann, wenn die Hauptdiagonale in \mathbf{A} streng überwiegt (streng diagonaldominant), d.h.

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{i,j}| \quad \forall i$$

Beweis „ \Leftarrow “

a) Wir zeigen zuerst $\exists \mathbf{A}^{-1}$.

Aus der Diagonaldomonaz folgt $a_{ii} \neq 0; \forall i$

$$\implies \mathbf{D} = \text{diag}(a_{11}, \dots, a_{nn}) \text{ ist invertierbar.}$$

$$\implies \mathbf{D}^{-1}\mathbf{A} \text{ hat in der Hauptdiagonalen nur Einsen.}$$

$$\implies \mathbf{B} := (b_{ij}) = \mathbf{I} - \mathbf{D}^{-1}\mathbf{A} \text{ hat in der Hauptdiagonalen nur Nullen und}$$

$$b_{ij} = -\frac{a_{ij}}{a_{ii}} \quad \forall j \neq i.$$

Für die Zeilensummennorm $\|\mathbf{B}\|_\infty$ folgt aus der Diagonaldominanz

$$\|\mathbf{B}\|_\infty < 1.$$

Nach Satz 3.1 (Neumansche Reihe) $\exists (\mathbf{I} - \mathbf{B})^{-1} = (\mathbf{D}^{-1}\mathbf{A})^{-1}$

$$\implies \mathbf{A} \text{ ist invertierbar, also a).}$$

Beweisvariante:

In keinem Gershgorinkreis von \mathbf{A} ist die Null enthalten auf Grund der starken Diagonaldominanz, also sind alle Eigenwerte $\lambda(\mathbf{A}) \neq 0 \implies \exists \mathbf{A}^{-1}$.

- b) Sei \mathbf{y} der Lösungsvektor von $\mathbf{Ax} = \mathbf{b} \neq \mathbf{0}$, $\implies \exists i : y_i \neq 0$ und $|y_i| = \|\mathbf{y}\|_\infty$.
(der Lösungsvektor ist eindeutig bestimmt, da \mathbf{A} regulär ist nach a).

Wir zeigen indirekt: $i \in N^\neq(\mathbf{Ay})$.

Annahme: $i \in N^0(\mathbf{Ay})$, d.h. $(\mathbf{Ay})_i = 0$ d.h.

$$\begin{aligned} a_{ii}y_i + \sum_{j=1, j \neq i}^n a_{ij}y_j &= 0 \implies \\ 0 &= |a_{ii}y_i + \sum_{j=1, j \neq i}^n a_{ij}y_j| \\ &\geq |a_{ii}y_i| - \sum_{j=1, j \neq i}^n |a_{ij}| |y_j| \quad \text{und wegen } |y_i| = \|\mathbf{y}\|_\infty \\ &\geq \underbrace{|y_i|}_{>0} \underbrace{\left(|a_{ii}| - \sum_{j=1, j \neq i}^n |a_{ij}| \right)}_{>0} \implies \text{W!} \end{aligned}$$

$|y_i|$ nimmt also sein Maximum nicht für $i \in N^\neq(\mathbf{Ay})$ an, d.h. b).

„ \implies “

Wir zeigen zuerst $a_{ii} \neq 0 \forall i$.

Wende das Maximumprinzip an auf $\mathbf{y} = \mathbf{e}^i$ (i -ter Einheitsvektor).

$\mathbf{e}^i \neq \mathbf{0}$ und da nach a) $\exists \mathbf{A}^{-1} \implies N^\neq(\mathbf{Ae}^i) \neq \emptyset$, d.h. $\mathbf{Ae}^i = \mathbf{a}^i \neq \mathbf{0}$ ($\mathbf{a}^i = i$ -te Spalte von \mathbf{A}).

Wegen $e_i^i = 1$, $e_j^i = 0$ für $i \neq j$, also $\max_j |e_j^i| = e_i^i = 1$, muß nach b) gelten:

$$i \in N^\neq(\mathbf{Ae}^i), \text{ d.h. } (\mathbf{Ae}^i)_i = a_{ii} \neq 0.$$

Wir zeigen die strenge Diagonaldominanz.

Für beliebiges, aber festes i definiere

$$\begin{aligned} y_i &:= - \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{a_{ii}}, & y_j &:= \overline{\text{sgn}} a_{ij} \text{ für } i \neq j, \text{ wo } \overline{\text{sgn}}(t) = \begin{cases} 1 & \text{für } t \geq 0 \\ -1 & \text{für } t < 0 \end{cases} \\ &\implies \mathbf{y} \neq \mathbf{0}, \end{aligned}$$

denn selbst wenn alle $a_{ij} = 0$ wären, folgt das aus der Definition von $\overline{\text{sgn}}$.

Nun ist

$$(\mathbf{A}\mathbf{y})_i = \sum_{j=1}^n a_{ij}y_j = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| + a_{ii} \left(- \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{a_{ii}} \right) = 0,$$

also $i \in N^0(\mathbf{A}\mathbf{y})$, d.h. $|y_j|$ nimmt sein Maximum nicht für $j = i$ an. (strenges Maximumprinzip!).

Damit folgt

$$|y_i| = \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|a_{ii}|} < \max_j |y_j| = 1 \quad \text{laut Definition von } y.$$

Das ist die starke Diagonaldominanz. ■

Satz 12.4

\mathbf{A} sei invertierbar und schwach diagonal dominant ($|a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \quad \forall i$)

Dann gilt das schwache Maximumprinzip.

Beweis:

Eigenschaft a) von Definition 12.2 ist erfüllt, da \mathbf{A} invertierbar.

Da \mathbf{A} schwach diagonal dominant und invertierbar ist, sind alle $a_{ii} \neq 0$, denn \mathbf{A} kann keine Nullzeile haben.

Sei $\mathbf{A} = (a_{ij})$, so gilt für ein beliebiges $\varepsilon > 0$ mit $\tilde{\varepsilon} = \varepsilon(\text{sgn}(a_{11}), \dots, \text{sgn}(a_{nn}))^T$:

$\mathbf{A}_\varepsilon := \mathbf{A} + \tilde{\varepsilon}\mathbf{I}$ ist $\forall \varepsilon > 0$ streng diagonal dominant, erfüllt also das strenge Maximumprinzip und ist invertierbar.

Für ein beliebiges $\mathbf{y} \neq \mathbf{0}$ sei $\boldsymbol{\varphi} := \mathbf{A}\mathbf{y}$.

Dadurch sind $N^0(\mathbf{A}\mathbf{y})$ und $N^\neq(\mathbf{A}\mathbf{y})$ festgelegt. Sei $\mathbf{y}^{(\varepsilon)}$ Lösung von

$$(*) \quad \mathbf{A}_\varepsilon \mathbf{y}^{(\varepsilon)} = \boldsymbol{\varphi} \\ \implies N^0(\mathbf{A}_\varepsilon \mathbf{y}^{(\varepsilon)}) = N^0(\mathbf{A}\mathbf{y}), \quad N^\neq(\mathbf{A}_\varepsilon \mathbf{y}^{(\varepsilon)}) = N^\neq(\mathbf{A}\mathbf{y}).$$

dann gilt nach dem starken Maximumprinzip für \mathbf{A}_ε .

$$(**) \quad \max_{i \in N^0(\mathbf{A}\mathbf{y})} |y_i^{(\varepsilon)}| < \max_{j \in N^\neq(\mathbf{A}\mathbf{y})} |y_j^{(\varepsilon)}|$$

In (*) kann der Grenzübergang $\varepsilon \rightarrow 0$ ausgeführt werden, da $\mathbf{A}_\varepsilon^{-1}$ existiert für $\varepsilon \geq 0$. Dann folgt aus (**) die Behauptung b). ■

Bemerkung: Für schwach diagonal dominante M-Matrizen gilt das schwache Maximumprinzip, denn M-Matrizen sind invertierbar.

Satz 12.5 Diskretes Maximumprinzip für die 1. RWA der Poissongl.

Für die Aufgabe

$$\Delta u = f \text{ in } \Omega = (0, a) \times (0, b), \quad u|_{\partial\Omega} = g$$

sei durch (12.5) die diskrete Form

$$\mathbf{A}_h \mathbf{y} = \boldsymbol{\varphi}$$

gegeben. Dann gilt für beliebiges $h = (h_1, h_2)$

$$(12.9) \quad \text{falls } f = 0 : \quad \|\mathbf{y}\|_\infty \leq \|\mathbf{g}\|_{C(\gamma_h)}$$

$$(12.10) \quad \text{falls } g = 0 : \quad \|\mathbf{y}\|_\infty \leq \frac{a^2 + b^2}{16} \|\mathbf{f}\|_{C(\bar{\omega}_h)}$$

$$(12.11) \quad \text{allgemein : } \quad \|\mathbf{y}\|_\infty \leq \|\mathbf{g}\|_{C(\gamma_h)} + \frac{a^2 + b^2}{16} \|\mathbf{f}\|_{C(\bar{\omega}_h)}$$

Bemerkung: Die Aussage des Satzes kann direkt für parabolische Aufgaben übernommen werden. Wesentlich ist nur, daß $\mathbf{A}_h \mathbf{y} = \boldsymbol{\varphi}$ eine Diskretisierung der parabolischen Aufgabe ist, welche die Randwerte mit einschließt und über alle Zeitschichten geht..

Beweis des Satzes: Wir zerlegen die allgemeine Aufgabe $\mathbf{A}_h \mathbf{y} = \boldsymbol{\varphi}$ in zwei Teilaufgaben

$$1. \quad \mathbf{A}_h \mathbf{y}^{(1)} = \boldsymbol{\varphi}_{\gamma_h}, \quad \boldsymbol{\varphi}_{\gamma_h} = \begin{cases} 0 & \text{für } (i, j) \in \omega_h \\ g_{i,j} & \text{für } (i, j) \in \gamma_h \end{cases} \quad (\text{Fall (12.9)})$$

$$2. \quad \mathbf{A}_h \mathbf{y}^{(2)} = \boldsymbol{\varphi}_{\omega_h}, \quad \boldsymbol{\varphi}_{\omega_h} = \begin{cases} f_{i,j} & \text{für } (i, j) \in \omega_h \\ 0 & \text{für } (i, j) \in \gamma_h \end{cases} \quad (\text{Fall (12.11)})$$

Die Linearität der Aufgabe bewirkt (Superposition)

$$\mathbf{A} \mathbf{y} = \mathbf{A} \mathbf{y}^{(1)} + \mathbf{A} \mathbf{y}^{(2)} = \boldsymbol{\varphi}_{\gamma_h} + \boldsymbol{\varphi}_{\omega_h} = \boldsymbol{\varphi}.$$

Dies liefert (12.11), wenn (12.9) und (12.10) bewiesen sind.

Fall (12.9): $f = 0$

\mathbf{A}_h ist eine schwach diagonal dominante Matrix gemäß (12.5) und nach Satz 12.1 eine M -Matrix, also gilt das schwache Maximumprinzip (Satz 12.4).

Die Nullenverteilung von $\boldsymbol{\varphi}_{\gamma_h}$ liefert (interpretiere γ_h und ω_h als Indexmengen:)

$$N^\neq(\mathbf{A}_h \mathbf{y}^{(1)}) \subset \gamma_h, \quad N^0(\mathbf{A}_h \mathbf{y}^{(1)}) \supset \omega_h.$$

Damit folgt

$$\max_{i \in \omega_h} |y_i^{(1)}| \leq \max_{i \in N^0(\mathbf{A}_h \mathbf{y}^{(1)})} |y_i^{(1)}| \leq \max_{j \in N^\neq(\mathbf{A}_h \mathbf{y}^{(1)})} |y_j^{(1)}| \leq \max_{j \in \gamma_h} |y_j^{(1)}|,$$

oder kurz

$$\|\mathbf{y}^{(1)}\|_{C(\omega_h)} \leq \|\mathbf{y}^{(1)}\|_{C(\gamma_h)} = \|\mathbf{g}\|_{C(\gamma_h)}, \quad \text{also (12.9).}$$

Fall(12.10): $g = 0$

Nach Satz 6.3 c) gilt für $\mathbf{A}_h \mathbf{y}^{(2)} = \varphi_{\omega_h}$

$$\|\mathbf{y}^{(2)}\|_{C(\bar{\omega}_h)} \leq \frac{\|\mathbf{p}\|_{C(\bar{\omega}_h)}}{\min_{i: (\varphi_{\omega_h})_i \neq 0} (\mathbf{A}_h \mathbf{p})_i} \|\varphi_{\omega_h}\|_{C(\bar{\omega}_h)} \quad \forall \mathbf{p} > 0 \text{ und } \mathbf{A}\mathbf{p} > 0 \text{ (komponentenweise)}$$

In diesem Fall genügt für p der Ansatz

$$w = \varepsilon + x_1(a - x_1) + x_2(b - x_2), \quad (x_1, x_2) \in \Omega, \quad \mathbf{p} = \mathbf{w}.$$

Dann ist

$$(\mathbf{A}_h \mathbf{p})_i \begin{cases} \geq \varepsilon > 0 & \text{auf } \gamma_h \\ = 4 & \text{auf } \omega_h \end{cases}$$

$$\|\mathbf{p}\|_{C(\bar{\omega}_h)} = \varepsilon + \frac{a^2 + b^2}{4},$$

$$\min_{i: (\varphi_{\omega_h})_i \neq 0} (\mathbf{A}_h \mathbf{p})_i = 4,$$

also

$$\|\mathbf{y}^{(2)}\|_{C(\bar{\omega}_h)} \leq \left(\frac{\varepsilon}{4} + \frac{a^2 + b^2}{16} \right) \|\mathbf{f}\|_{C(\bar{\omega}_h)} \quad \forall \varepsilon > 0,$$

und da diese Ungleichung von ε unabhängig ist, gilt sie auch für $\varepsilon = 0$. Daraus folgt (12.10).

Fall (12.11) durch Superposition. ■

Bemerkungen

1. Die Abschätzungen (12.9) - (12.11) sind besser als die jeweiligen Abschätzungen aus Satz 12.1 c), die nur die M -Matrixeigenschaft verwendeten.
2. Für die Aufgabe 2) am Beweisanfang (Nullrandwerte) liefert das schwache Maximumprinzip (Satz 12.4)

$$\|\mathbf{y}^{(2)}\|_{C(\gamma_h)} \leq \|\mathbf{y}^{(2)}\|_{C(\omega_h)},$$

das Maximum wird also im Innern angenommen, was nicht verwunderlich ist, da die Randwerte $= 0$ sind.

3. In der praktischen Anwendung werden in $\mathbf{A}_h \mathbf{y} = \varphi$ die Randwerte immer auf die rechte Seite gebracht. Man erhält dann eine Matrix \mathbf{A}_h^0 (vgl. (12.5)), die symmetrisch ist (die Übergangselemente $\frac{1}{h_1^2}$ entfallen) und daher bessere numerische Eigenschaften hat. Man löst dann ein System $\mathbf{A}_h^0 \mathbf{y}^0 = \varphi^0$, in dem \mathbf{y}^0 nur innere Punkte enthält. Man beachte jedoch, daß nun in den Komponenten von φ^0 additiv zu den Komponenten von f_{ij} auch Komponenten von \mathbf{g} auftreten. Für den Punkt $i = 1, j = 1$ (zum Beispiel) lautet die rechte Seite $f_{11} + \frac{1}{h_1^2} g_{10}$ (vgl. dazu auch (2.12)).

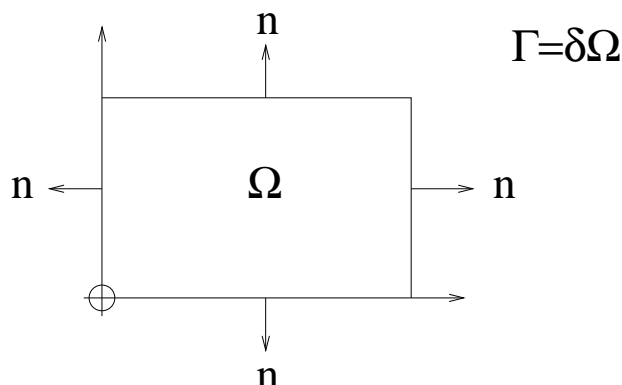
§ 13 Die 3. RWA für die Poissongleichung

Physikalische Herkunft der 3. RWA:

$$\Delta u + f = 0, \quad x \in \Omega$$

$$\frac{\partial u}{\partial n} + \sigma(u - u_0) = 0 \text{ auf } \delta\Omega$$

$$\sigma \geq \sigma_0 > 0, \quad u_0 = \text{Außentemperatur.}$$

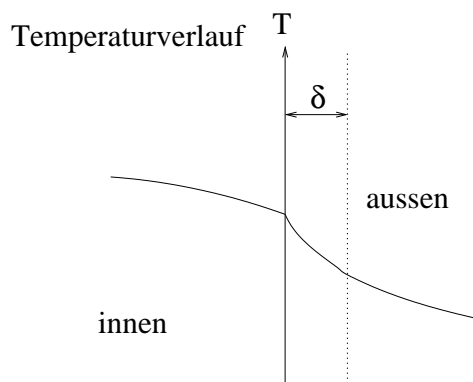


Diese Randbedingung ist bei einem stationären Wärmeleitungsprozess (z.B. in Ω wird geheizt, außerhalb ist es kälter) die richtige physikalische Formulierung dafür, daß der Wärmestrom $k \frac{\partial u}{\partial n}$ (wobei u = Körpertemperatur, u_0 = Außentemperatur, n = äußere Normale) stetig durch die Trennfläche geht. Sie wird bestimmt durch das

Fick'sche Gesetz: $\exists \alpha$:

$$\left(k \frac{\partial u}{\partial n} \right) \Big|_{\Gamma-\varepsilon} = \left(\alpha \frac{\partial u}{\partial n} \right) \Big|_{\Gamma+\varepsilon}$$

\uparrow \uparrow
 innere äußere
 Randgrenzwerte



Probleme in der Praxis: Der Diffusionskoeffizient k kann sich beim Durchgang durch Γ (unstetig) ändern. Dies bewirkt einen Knick in der Temperaturableitung. In der Grenzschicht (Dicke = δ , $\delta = ?$) herrschen komplizierte Verhältnisse. Ist $u(\Gamma)$ die Körpertemperatur in Ω (genauer: der Grenzwerte von innen zum Rand Γ) und u_0 eine Schätzung für die Außentemperatur (, die nicht notwendig konstant sein muß,) so wird eine grobe Schätzung der rechten Seite des Fick'schen Gesetzes gegeben durch

$$\left(\alpha \frac{\partial u}{\partial n} \right) \Big|_{\Gamma+\varepsilon} = \alpha \frac{u_0 - u(\Gamma)}{\delta} \quad \left(= \left(k \frac{\partial u}{\partial n} \right) \Big|_{\Gamma-\varepsilon} \right)$$

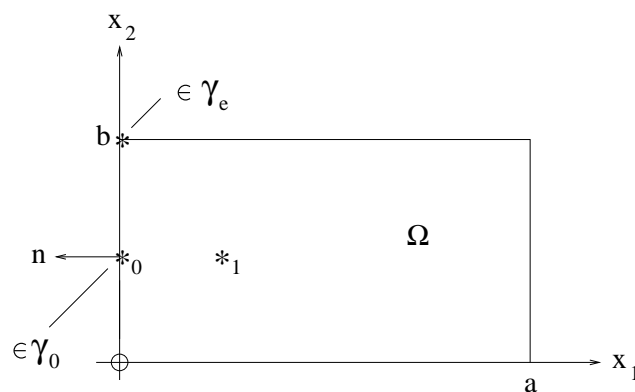
woraus in der Randbedingung, die durch das Fick'sche Gesetz beschrieben wird, sich der Faktor σ ergibt zu $\sigma = \frac{\alpha}{\delta k}$; α , δ und k sind in der Anwendung oft schwierig zu erhalten.

Grenzfall: Ist der Außenraum ein Isolator, so fließt keine Wärme von innen nach außen und man erhält als Randbedingung: $\frac{\partial u}{\partial n} = 0$.

Es sind nun folgende Probleme zu untersuchen:

1. Herstellung einer geeigneten Differenzenapproximation und Abschätzung des Diskretisierungsfehlers (Taylorabgleich).
2. Nachweis, daß die Diskretisierungsmatrix eine M -Matrix ist ($\Rightarrow \exists \mathbf{A}^{-1}$).
3. Konvergenzabschätzung.

1) Differenzenapproximation für die 3. Randbedingung: $\frac{\partial u}{\partial n} + \sigma(u - u_0) = 0$
 Beachte: Die Diskretisierung muß für jeden Gitterpunkt, also auch jeden Randpunkt, eine Gleichung liefern weil keine Dirichletrandwerte gegeben sind.



In Ω wird wie bekannt approximiert

$$(\mathbf{A}_h y)_i = -(y_{\bar{x}_1 x_1} + y_{\bar{x}_2 x_2})_i, \quad i \in \omega_h.$$

Diese Approximation ist von 2. Ordnung. Wir suchen deshalb für die Randpunkte ebenfalls eine Approximation 2. Ordnung. Dazu muß der Rand γ in zwei Teile zerlegt werden.

$$\gamma = \gamma_0 + \gamma_e, \quad \gamma_0 = \text{innere Randpunkte}, \quad \gamma_e = \text{Eckpunkte.}$$

Wir konstruieren zunächst die

Approximation am linken Rand

In obiger Abbildung bezeichnen wir mit Index 0 einen inneren Randpunkt $\in \gamma_0$, mit Index 1 seinen rechten Gitternachbarn. Die Taylorentwicklung im Randpunkt liefert

$$u_1 = u_0 + h_1 \frac{\partial u}{\partial x_1} \Big|_{\Gamma} + \frac{h_1^2}{2} \frac{\partial^2 u}{\partial x_1^2} \Big|_{\Gamma} + \frac{h_1^3}{6} \frac{\partial^3 u}{\partial x_1^3} \Big|_{\Gamma}$$

$$(13.1) \quad u_{x_1,0} = \frac{u_1 - u_0}{h_1} = \frac{\partial u}{\partial x_1} \Big|_{\Gamma} + \frac{h_1}{2} \frac{\partial^2 u}{\partial x_1^2} \Big|_{\Gamma} + \frac{h_1^2}{6} \frac{\partial^3 u}{\partial x_1^3} \Big|_{\Gamma} \implies$$

$$(13.2) \quad \frac{\partial u}{\partial n} \Big|_{\Gamma} = -\frac{\partial u}{\partial x_1} \Big|_{\Gamma} = -u_{x_1,0} + \frac{h_1}{2} \frac{\partial^2 u}{\partial x_1^2} \Big|_{\Gamma} + O(h_1^2)$$

Nun kann $\frac{\partial^2 u}{\partial x_1^2}$ im Randpunkt nicht (oder nur schlecht) approximiert werden. Zwar wäre eine Approximation von $\frac{\partial^2 u}{\partial x_1^2}$ von entsprechender Ordnung durch Taylorabgleich unter Hinzunahme entsprechend vieler x_1 -Werte auf dem Level x_2 möglich. Dadurch wird jedoch die Struktur der Matrix beeinflusst und ein Nachweis, daß sie eine M -Matrix ist, erschwert, wenn nicht unmöglich gemacht. Deshalb wird $\frac{\partial^2 u}{\partial x_1^2}$ mit Hilfe der Differentialgleichung ersetzt. Laut Differentialgleichung gilt

$$(13.3) \quad \frac{\partial^2 u}{\partial x_1^2} = -f - \frac{\partial^2 u}{\partial x_2^2}$$

und analog zu (12.3), falls $u \in C^3$ (Taylorreihe, vgl. z.B. (12.2)-(12.3))

$$(13.4) \quad \frac{\partial^2 u}{\partial x_2^2} \Big|_{\Gamma} = u_{\bar{x}_2 x_2, 0} - \underbrace{\left(\frac{h_2}{6} \frac{\partial^3 u_+}{\partial x_2^3} + \frac{h_2}{6} \frac{\partial^3 u_-}{\partial x_2^3} \right)}_{|\leq \frac{h_2}{3} M_3}$$

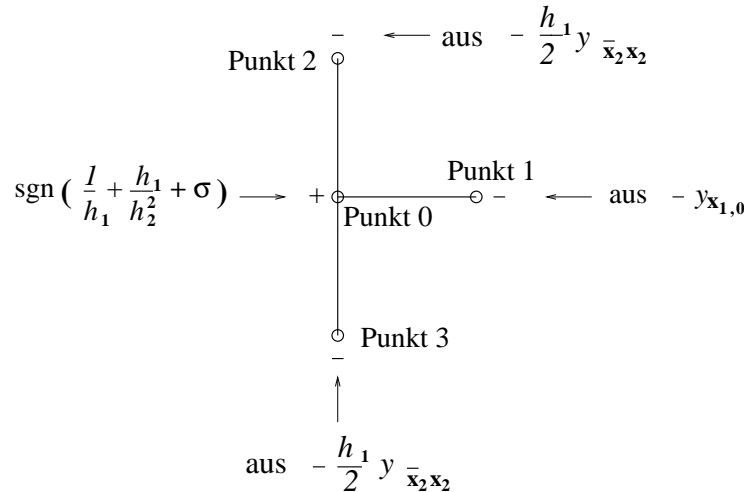
Bemerkung: In (13.2) ist $\frac{\partial^2 u}{\partial x_1^2}$ schon mit einem Faktor h_1 versehen, deshalb genügt hier eine in h_2 lineare Abschätzung um gemäß $h_1 h_2 = \frac{1}{2}(h_1^2 + h_2^2)$ eine quadratische Abschätzung zu erhalten, insgesamt also: $\frac{h_1}{2} \frac{h_2}{3} \leq M_3 \frac{h_1^2 + h_2^2}{6}$.
Wir ersetzen in (13.2) die 2.te Ableitung gemäß (13.3), (13.4) und erhalten damit insgesamt die Approximation

$$(13.5) \quad \begin{aligned} \frac{\partial u}{\partial n} \Big|_{\Gamma} + \sigma(u - u_0) &= -y_{x_1, 0} + \frac{h_1}{2}(-f - y_{\bar{x}_2 x_2, 0}) + \sigma(y_0 - u_0) + R = 0 \\ \text{mit } |R| &\leq \left(\underbrace{\frac{h_2^1}{6}}_{\text{aus (13.1)}} + \underbrace{\frac{h_1^2 + h_2^2}{6}}_{\text{vgl. oben}} \right) M_3 = M_3 \left(\frac{2h_1^2 + h_2^2}{6} \right) \leq \frac{M_3}{3}(h_1^2 + h_2^2). \end{aligned}$$

Die Diskretisierung am linken Rand lautet also (bekannte Daten auf die rechte Seite und Verwendung der Punktnummern (vgl. Abb.) als Indizes)

$$(L) \quad \begin{aligned} (\mathbf{A}_h \mathbf{y})_0 &:= -y_{x_1, 0} + \sigma y_0 - \frac{h_1}{2} y_{\bar{x}_2 x_2, 0} = \sigma u_0 + \frac{h_1}{2} f, \quad \text{bzw.} \\ -\frac{y_1 - y_0}{h_1} + \sigma y_0 - \frac{h_1}{2} \frac{y_2 - 2y_0 + y_3}{h_2^2} &= \sigma u_0 + \frac{h_1}{2} f \end{aligned}$$

mit folgender Vorzeichenverteilung in \mathbf{A}_h



Beachte: Die Vorzeichenverteilung ist die richtige für eine M -Matrix, denn die y -Komponente des Punkt 0 steht in der Hauptdiagonalen, die Komponenten der y -Werte der Punkte 1 und 2 stehen rechts neben der Hauptdiagonalen, der Koeffizient zur y -Komponente von Punkt 3 steht links der Hauptdiagonalen.

Entsprechend findet man am rechten Rand (beachte: $\frac{\partial}{\partial n} = +\frac{\partial}{\partial x_1}$ und $y_{\bar{x}_1, N_1}$)

$$(R) \quad (\mathbf{A}_h y)_0 := +y_{\bar{x}_1, N_1} + \sigma y_{N_1} - \frac{h_1}{2} y_{\bar{x}_2 x_2, 0} = \sigma u_0 + \frac{h_1}{2} f.$$

$\left. \begin{array}{c} - \\ \hline 0 \\ \hline - \end{array} \right\} +$

\uparrow
 beachte \bar{x}_1 : rückwärts genommener Differenzenquotient

Die Randbedingungen am unteren bzw. oberen Rand ergeben sich aus (L) bzw. (R) durch Vertauschung von $x_1 \leftrightarrow x_2$ und $h_1 \leftrightarrow h_2$ mit den Vorzeichenverteilungen



Behandlung der 4 Eckpunkte

Wir betrachten zunächst den Fall

(die Ziffern indizieren die Punkte)

Ziel:
 Man konstruiert als Approximation für $\frac{\partial}{\partial n}$ eine Konvexkombination von $\frac{\partial}{\partial x_1}$ und

einen Vektor $\mathbf{p} > \mathbf{0}$ mit $\mathbf{A}\mathbf{p} > \mathbf{0}$ mit $\mathbf{p} = \mathbf{w}$ (Gittervektor).

Für die inneren Punkte ($\in \omega_h$) gilt (vgl. (12.7)- (12.8)): $(\mathbf{A}\mathbf{p})_l = 4, l \in \omega_h$.

Für die linken Randpunkte ($\in \gamma_0$) folgt aus (L) durch Taylorabgleich:
(beachte: Im linken Randpunkt hat der rechte Nachbar den Wert $x_1 = h_1$)

$$w_1 = w_0 + h_1 \frac{\partial w|_0}{\partial x_1} + \frac{h_1^2}{2} \frac{\partial^2 w|_0}{\partial x_1^2} + \underbrace{\frac{h_1^3}{6} \frac{\partial^3 w|_0}{\partial x_1^3}}_{=0}$$

$$\frac{w_1 - w_0}{h_1} = a - 2 \cdot 0 - h_1 = a - h_1 \quad \text{und damit nach (L)}$$

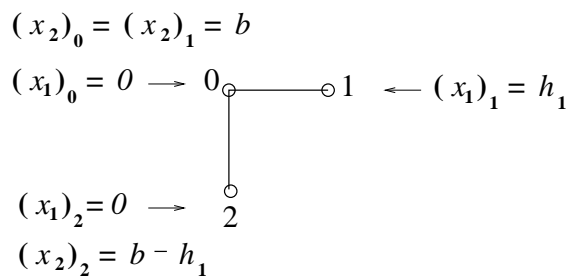
$$\begin{aligned} (\mathbf{A}_h \mathbf{w})_0 &:= -w_{x_1,0} + \sigma w_0 - \frac{h_1}{2} w_{\bar{x}_2 x_2,0} \\ &\geq -(a - h_1) + \sigma c - \frac{h_1}{2}(-2), \text{ da } x_2(b - x_2) \geq 0 \text{ und } \frac{\partial^2}{\partial x_2^2}[x_2(b - x_2)] = -2 \\ &- a + \sigma c + 2h_1 \stackrel{!}{\geq} 4 \iff c \stackrel{!}{\geq} \frac{4 + a - 2h_1}{\sigma}. \end{aligned}$$

Die letzte Abschätzung hätte man gerne im Hinblick auf die Abschätzung von $\|\mathbf{A}_h^{-1}\| : \min_l (\mathbf{A}\mathbf{p})_l$ soll nicht kleiner werden als 4, was durch die Abschätzung für die inneren Punkte schon vorgegeben ist. Mit $\sigma \geq \sigma_0 > 0$ ist sie erfüllt, falls verschärft gilt

$$c \geq \frac{4 + a}{\sigma_0}.$$

Für die rechten Randpunkte folgt analog: $c \geq \frac{4 + b}{\sigma_0}$, und entsprechend für die oberen bzw. unteren Randpunkte $c \geq \frac{4 + b}{\sigma_0}$ bzw. $c \geq \frac{4 + a}{\sigma_0}$.

Eckpunkte (zunächst links oben)



Laut (E) gilt

$$(\mathbf{A}_h \mathbf{y})_0 := -\frac{H}{h_1^2}(y_1 - y_0) + \frac{H}{h_2^2}(y_0 - y_2) + \sigma y_0.$$

Aus $w = c + x_1(a - x_1) - x_2(b - x_2)$ folgt

$$\left. \begin{aligned} w_0 &= c + 0 \\ w_1 &= c + h_1(a - h_1) \\ w_2 &= c + (b - h_2)(b - (b - h_2)) = c + h_2(b - h_2) \end{aligned} \right\} \left. \begin{aligned} w_1 - w_0 &= h_1(a - h_1) \\ w_0 - w_2 &= -h_2(b - h_2) \end{aligned} \right\}$$

Damit erhält man

$$\begin{aligned}
(\mathbf{A}_h \mathbf{w})_0 &:= -\frac{H}{h_1^2}(h_1(a-h_1)) + \frac{H}{h_2^2}(-h_2(b-h_2)) + \sigma c \\
&= -\frac{H}{h_1}(a-h_1) - \frac{H}{h_2}(b-h_2) + \sigma c \\
&= -\frac{h_2}{h_1+h_2}(a-h_1) - \frac{h_1}{h_1+h_2}(b-h_2) + \sigma c \\
&= -\underbrace{\frac{h_2}{h_1+h_2}}_{\leq 1} a - \underbrace{\frac{h_1}{h_1+h_2}}_{\leq 1} b + \underbrace{\frac{2h_1h_2}{h_1+h_2}}_{\geq 0} + \sigma c \\
&\geq -a - b + \sigma c \stackrel{!}{\geq} 4 \quad (\text{möchte man wieder!})
\end{aligned}$$

Dies führt zu der hinreichenden Bedingung

$$(13.7) \quad c \geq \frac{4+a+b}{\sigma_0}, \quad (\text{analog für die anderen Eckpunkte})$$

Diese Forderung impliziert die vorigen, ist also insgesamt ausreichend. Also ist \mathbf{A}_h eine M -Matrix.

3) Stabilität und Konvergenz

Nach obiger Wahl von $w = c + x_1(a-x_1) + x_2(b-x_2)$, $\mathbf{p} = \mathbf{w}$ (Gitterfunktion) erhält man (vgl. Satz 6.3 und Beweis von Satz 12.1)

$$\begin{aligned}
\|\mathbf{A}_h^{-1}\|_{C(\omega_h)} &\leq \frac{\|\mathbf{p}\|_\infty}{\min_l(\mathbf{A}\mathbf{p})_l} \leq \frac{c + \frac{a^2+b^2}{4}}{4}, \quad \text{also mit (13.7)} \\
(13.8) \quad \|\mathbf{A}_h^{-1}\|_{C(\omega_h)} &\leq \frac{4+a+b}{4\sigma_0} + \frac{a^2+b^2}{16}.
\end{aligned}$$

Diese Abschätzung ist unabhängig von der Diskretisierung, d.h. auch die Stabilität ist gesichert.

Der Verfahrensfehler $\mathbf{z} = \mathbf{y} - \mathbf{u}$, (\mathbf{u} = Gitterfunktion der exakten Lösung,) genügt $\mathbf{A}_h \mathbf{z} = \boldsymbol{\psi}$, $\boldsymbol{\psi}$ = Diskretisierungsfehler.

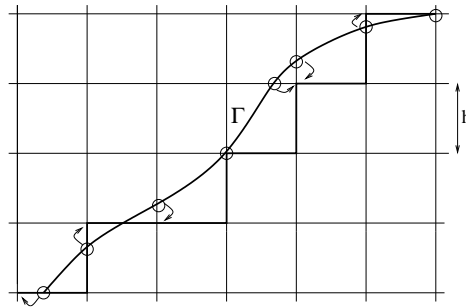
$\boldsymbol{\psi}$ = enthält den Fehler der Differentialgleichung: $|\cdot| \leq \frac{h_1^2+h_2^2}{12} M_4$ (vgl. (12.3),(12.4)) und den Fehler des Randes $|\cdot| \leq \frac{h_1^2+h_2^2}{3} M_3$ (vgl. (13.5)). Insgesamt also

$$\begin{aligned}
\|\mathbf{z}\|_\infty &\leq \underbrace{\left(\frac{4+a+b}{4\sigma_0} + \frac{a^2+b^2}{16} \right)}_{=: C_1} \max \left(\frac{h_1^2+h_2^2}{12} M_4, \frac{h_1^2+h_2^2}{3} M_3 \right) \\
\|\mathbf{z}\|_\infty &\leq C_1 \max(4M_3, M_4) \frac{h_1^2+h_2^2}{12} \quad \text{falls } u \in C^4,
\end{aligned}$$

also quadratische Konvergenz.

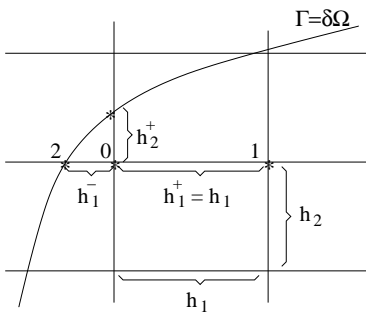
§ 14 Die 1. RWA der Poissongl. in allgemeineren Gebieten

Was macht man z.B. mit dem Rand der Nordsee? Er ist nicht eindeutig.
 Ein einfaches, stabiles Verfahren erhält man, indem man die echten Randpunkte auf die benachbarten Gitterpunkte verschiebt.



Dies liefert einen Fehler 1. Ordnung. Die Ergebnisse sind jedoch, in einem gewissen Abstand vom Rand brauchbar. Wir beschreiben im Folgenden eine Diskretisierung in den randnahen inneren Punkten, die zunächst von 1. Ordnung ist, jedoch zu einem quadratisch konvergenten Verfahren für den Gesamtbereich führt.

Shortley-Wellers: Approximation und Verfahren



$h_{1,2}^{\pm}$ bezeichnen die Abstände des Bezugspunktes (hier der mit 0 bezeichnete Punkt) zu den benachbarten Gitterpunkten. Dies können innere Gitterpunkte sein (wie hier der Punkt 1), als auch Randpunkte, die durch den Schnitt des Randes von Ω mit den Gitterlinien entsteht (hier z.B. der Punkt 2).

Randnahe Gitterpunkte sind solche, für welche mindestens einer der Werte $h_{1,2}^{\pm}$ von der Maschenweite verschieden ist. Wir untersuchen, zunächst nur für die x_1 -Richtung, die Approximation von $\frac{\partial^2}{\partial x_1^2}$ im Punkt 0.

Für die Differenzenquotienten gilt (Taylorabgleich)

$$\frac{y_1 - y_0}{h_1^+} = y_0' + \frac{h_1^+}{2} y_0'' + \frac{(h_1^+)^2}{6} y_0^{(3')} + \frac{(h_1^+)^3}{24} y_0^{(4')} + \dots$$

$$\frac{y_0 - y_2}{h_1^-} = y_0' - \frac{h_1^-}{2} y_0'' + \frac{(h_1^-)^2}{6} y_0^{(3')} - \frac{(h_1^-)^3}{24} y_0^{(4')} + \dots$$

Beide Gleichungen weisen einen Fehler 1. Ordnung auf. Durch Subtraktion der Gleichungen und Multiplikation mit $\frac{2}{h_1^+ + h_1^-}$ folgt

$$\frac{2}{h_1^+ + h_1^-} \left(\frac{y_1 - y_0}{h_1^+} - \frac{y_0 - y_2}{h_1^-} \right) = y_0'' + \underbrace{\frac{2}{h_1^+ + h_1^-} \left(\frac{(h_1^+)^2}{6} - \frac{(h_1^-)^2}{6} \right)}_{\frac{h_1^+ - h_1^-}{3}} y_0^{(3')} + O((h_1^+)^2 + (h_1^-)^2)$$

$$(14.1) \quad y_0'' = \frac{2}{h_1^+ + h_1^-} \left(\frac{y_1 - y_0}{h_1^+} - \frac{y_0 - y_2}{h_1^-} \right) - \frac{h_1^+ - h_1^-}{3} y_0^{(3')} + O((h_1^+)^2 + (h_1^-)^2)$$

Das Fehlerglied folgt mit Hilfe Abschätzung

$$\frac{(h_1^+)^3}{h_1^+ + h_1^-} + \frac{(h_1^-)^3}{h_1^+ + h_1^-} \leq \frac{(h_1^+)^3}{h_1^+} + \frac{(h_1^-)^3}{h_1^-} \leq (h_1^+)^2 + (h_1^-)^2.$$

Beachtet man weiter

$$|h_1^+ - h_1^-| \leq h_1 \quad \text{und} \quad (h_1^+)^2 + (h_1^-)^2 \leq 2h_1^2$$

so läßt sich (14.1) auch schreiben als

$$(14.2) \quad y_0'' = \frac{2}{h_1^+ + h_1^-} \left(\frac{y_1 - y_0}{h_1^+} - \frac{y_0 - y_2}{h_1^-} \right) + O(h_1) + O(h_1^2)$$

mit $|O(h_1)| \leq \frac{h_1}{3} M_3, \quad |O(h_1^2)| \leq \frac{h_1^2}{12} M_4, \quad M_i = \max_{\Omega} |y^{(i)}|.$

Dies ist eine Approximation 1. Ordnung, die auch bei nichtäquidistanten Gittern verwendet werden kann.

Mögliche Abhilfe: Ist h_{max} die maximale Maschenweite, so kann man durch die

$$\text{Forderung:} \quad |h_1^+ - h_1^-| \leq M h_{max}^2, \quad \text{falls} \quad M \geq \frac{1}{h_{max}}$$

erreichen, daß die Approximation (14.1) von 2. Ordnung ist.

Man kann Bedingungen angeben, wie ein gegebenes, nicht äquidistantes Gitter verfeinert werden kann, damit diese Forderung erhalten bleibt.

Obwohl (14.2) nur eine Abschätzung 1. Ordnung ist, werden wir zeigen, daß sie zu einem Verfahren 2. Ordnung führt.

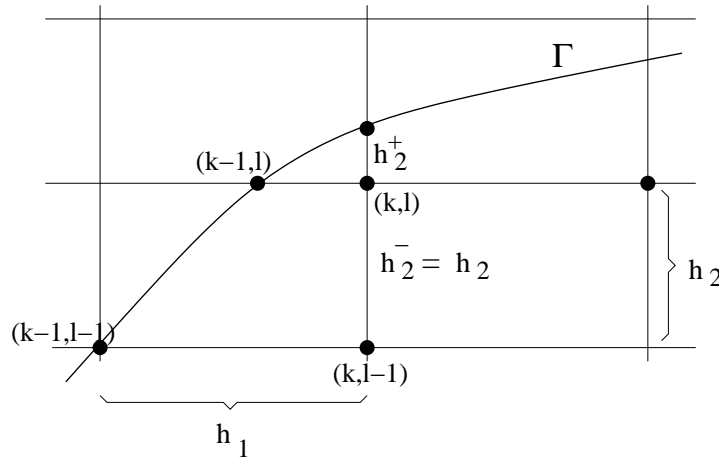
Wir legen nun über Ω ein nicht notwendig äquidistantes Gitter mit den Maschenweiten $h = (h_1, h_2)$.

Die Menge der Gitterpunkte, die wir bei der Diskretisierung benutzen, setzt sich wie folgt zusammen aus

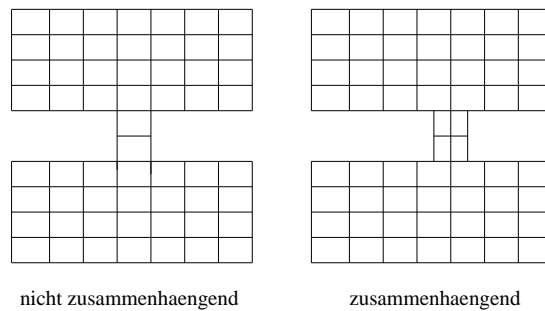
$$(14.3) \quad \bar{\omega}_h = \gamma_h \cup \omega_h^* \cup \omega_h^0.$$

Dabei bedeuten (vgl. Abbildung)

- γ_h = Randpunkte: Schnittpunkte von $\Gamma = \delta\Omega$ mit den Gitterlinien
(z.B. $(k-1, l), (k-1, l-1)$)
- ω_h^* = irreguläre Punkte: Innere Gitterpunkte, die mindestens einen Randpunkt als Nachbar haben (z.B. $(k, l), (k, l-1)$).
- ω_h^0 = reguläre Punkte: Innere Gitterpunkte, deren Nachbarn alle innere Punkte sind.



Das Gitter soll so fein sein (zusammenhängend), daß sich zwei beliebige innere Gitterpunkte durch einen Polygonzug aus inneren Gitterlinien verbinden lassen. Sonst zerfällt das Gleichungssystem der diskretisierten Gleichungen in separate Teilsysteme und die Kopplung zwischen den Teilsystemen geht verloren.



Bei nichtzusammenhängenden Gebieten kann man durch Verschiebung des Gitters oder Gitterverdichtung den Zusammenhang herstellen.

Wir bezeichnen nun mit $\tilde{\mathbf{y}}$ der Vektor, der alle Punkte aus $\bar{\omega}_h$ enthält. $\tilde{\mathbf{A}}_h$ sei die zugehörige Diskretisierungsmatrix.

Diskretisierung:
in ω_h^0 durch den bekannten 5-Punktstern von 2. Ordnung,

in ω_h^* gemäß (14.1) durch die Shortley-Wellers-Approximation

$$(14.4) \quad (\tilde{A}_h \tilde{y})_{(k,l)} = \begin{aligned} & -\frac{2}{h_1^+ + h_1^-} \left(\frac{y_{k+1,l} - y_{k,l}}{h_1^+} - \frac{y_{k,l} - y_{k-1,l}}{h_1^-} \right) & \text{mit } h_1 \geq h_1^\pm, \\ & -\frac{2}{h_2^+ + h_2^-} \left(\frac{y_{k,l+1} - y_{k,l}}{h_2^+} - \frac{y_{k,l} - y_{k,l-1}}{h_2^-} \right) & h_2 \geq h_2^\pm \end{aligned}$$

Beachte: Für $h_1^+ = h_1^- = h_1$, $h_2^+ = h_2^- = h_2$, wie das in inneren Punkten der Fall ist, ist das der übliche 5-Punktstern, d.h. in den entsprechenden Matrixkomponenten ergibt das die gewohnte Vorzeichenverteilung.

Dann gilt gemäß (14.2)

$$(14.5) \quad (\tilde{A}_h \tilde{u})_{(k,l)} = -(\Delta u)_{(k,l)} + \underbrace{O(h_1 + h_2)}_{|\leq \frac{h_1+h_2}{3} M_3 \text{ vgl. (14.3)}} + \underbrace{O(h_1^2 + h_2^2)}_{|\leq \frac{h_1^2+h_2^2}{12} M_4, \text{ vgl. (12.3)}}$$

mit $M_i = \max_{\Omega} \left(\left| \frac{\partial^i u}{\partial x_1^i} \right|, \left| \frac{\partial^i u}{\partial x_2^i} \right| \right)$

Dies ist offenbar eine Approximation 1. Ordnung. Trotzdem werden wir zeigen, daß sie ein konvergentes Verfahren 2. Ordnung liefert. Wir zeigen zuerst

Satz 14.1

\tilde{A}_h gemäß (14.4) ist eine M -Matrix und es gilt

$$\|\tilde{A}_h^{-1}\|_{C(\bar{\omega}_h)} \leq 1 + \frac{1}{4}(\text{diam}(\Omega))^2, \quad (\text{diam}(\Omega) := \max_{\mathbf{x}, \mathbf{y} \in \Omega} \|\mathbf{x} - \mathbf{y}\|_2).$$

Beweis: Im Gleichungssystem

$$\tilde{A}_h \tilde{\mathbf{y}} = \tilde{\boldsymbol{\varphi}} = \begin{pmatrix} \mathbf{g} \\ \mathbf{f} \end{pmatrix}, \quad \begin{aligned} \mathbf{g} &= \text{Werte auf } \gamma_h \\ \mathbf{f} &= \text{rechte Seite} \end{aligned}$$

↑

bedeutet keine Reihenfolge der Anordnung

ist die Vorzeichenverteilung in den irregulären Punkten die gleiche wie in den regulären Punkten. Die echten Randpunkte liefern nur eine Eins in der Diagonalen. Die Vorzeichenverteilung in \tilde{A}_h genügt also der einer M -Matrix.

Zu konstruieren ist also ein Vektor $\tilde{\mathbf{p}} > \mathbf{0}$ mit $\tilde{A}_h \tilde{\mathbf{p}} > \mathbf{0}$.

Dazu machen wir einen quadratischen Ansatz

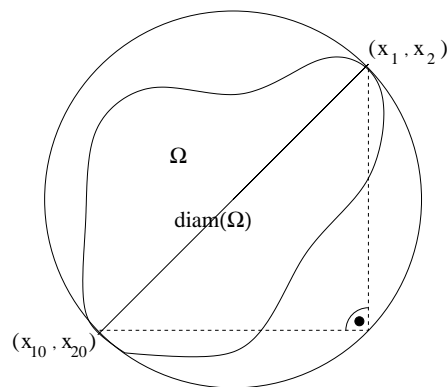
$$(14.6) \quad \begin{aligned} w(x_1, x_2) &= 4 + c - ((x_1 - x_{10})^2 + (x_2 - x_{20})^2), \quad (x_1, x_2) \in \Omega \\ &(x_{10}, x_{20}) \text{ innerhalb des Umkreises um } \Omega, \quad c = (\text{diam}(\Omega))^2. \end{aligned}$$

Begründung und Bemerkungen

1. Wir setzen $\tilde{\mathbf{p}} = \mathbf{w}$. Die Vorzeichenverteilung der quadratischen Glieder in w ist die gleiche wie in der Ansatzfunktion (12.7). Das benötigt man für eine M -Matrix.
2. $w \geq 4$ wird gebraucht im Hinblick auf die Abschätzung aus Satz 6.3, damit der Nenner nicht kleiner wird. Dies erfordert einen Korrekturterm c , der so beschaffen sein muß, daß " c -quadratische Glieder" ≥ 0 ausfällt. c geht in $\|\tilde{\mathbf{p}}\|_\infty$ ein und sollte deshalb möglichst klein sein. Die quadratischen Glieder werden deshalb so angesetzt, daß sie der Vorzeichenverteilung genügen, aber nicht zu groß ausfallen.
3. Liegt der Punkt (x_{10}, x_{20}) im Umkreis um Ω , so ist $(x_1 - x_{10})^2 + (x_2 - x_{20})^2 \leq (\text{diam}(\Omega))^2$ (Phytagoras.) Das " \leq " kann zum " $=$ " werden, wenn (x_{10}, x_{20}) auf dem Umkreis liegt und $(x_{10}, x_{20}) \in \delta\Omega$ (vgl. Zeichnung).
4. Die Lage von (x_{10}, x_{20}) und die Wahl von c bewirken deshalb, daß

$$c - (x_1 - x_{10})^2 - (x_2 - x_{20})^2 \geq 0 \quad \text{und} \quad w(x_1, x_2) \geq 4 \text{ in } \bar{\Omega}.$$

Extreme Lage von (x_0, x_1) und (x_{10}, x_{20})



Mit $\tilde{\mathbf{p}} = \mathbf{w}$ gilt nach (14.4)

$$(14.7) \quad (\tilde{\mathbf{A}}_h \tilde{\mathbf{p}})_{k,l} = \begin{cases} \geq 4 & \text{auf } \gamma_h \text{ Randpunkte, vgl. Bemerkung d)} \\ 4 & \text{auf } \omega_h^0 \text{ reguläre innere Punkte, vgl. (12.9),} \\ 4 & \text{auf } \omega_h^* \text{ weil } w \text{ quadratisch ist und } \frac{\partial^3 w}{\partial x_i^3} = 0. \end{cases}$$

Begründung zum Fall $(k, l) \in \omega_h^*$:

$(\tilde{\mathbf{A}}_h \tilde{\mathbf{p}})_{k,l}$ wird gemäß (14.4) dargestellt und (14.1) zeigt, daß der Fehler der Darstellung erst in der 3. Ableitung nicht verschwindet. Das tut nicht weh, da w quadratisch ist.

Mit $\tilde{\mathbf{p}} = \mathbf{w}$ gilt nach (14.4)

$$(14.8) \quad (\tilde{\mathbf{A}}_h \tilde{\mathbf{p}})_{k,l} = \begin{cases} \geq 4 & \text{auf } \gamma_h \text{ Randpunkte, vgl. Bemerkung 4)} \\ 4 & \text{auf } \omega_h^0 \text{ reguläre innere Punkte, vgl. (12.9),} \\ 4 & \text{auf } \omega_h^* \text{ weil } w \text{ quadratisch ist und } \frac{\partial^3 w}{\partial x_i^3} = 0. \end{cases}$$

Also ist $\tilde{\mathbf{A}}_h$ eine M -Matrix und wir erhalten nach Satz 6.3 b) die Abschätzung

$$\|\tilde{\mathbf{A}}_h^{-1}\|_{C(\bar{\omega}_h)} \leq \frac{\|\tilde{\mathbf{p}}\|_\infty}{\min_i (\tilde{\mathbf{A}}_h \tilde{\mathbf{p}})_i} \leq \frac{4+c}{4} = 1 + \frac{c}{4} = 1 + \frac{1}{4}(\text{diam}(\Omega))^2 \quad \text{für } c = (\text{diam}(\Omega))^2.$$

■

Bemerkung:

Für den Verfahrensfehler $\tilde{\mathbf{z}} = \tilde{\mathbf{y}} - \tilde{\mathbf{u}}$, ($\tilde{\mathbf{u}}$ = Gitterfunktion der exakten Lösung u in allen Gitterpunkten inclusive Randpunkten), liefert diese Abschätzung wegen $\tilde{\mathbf{A}}_h \tilde{\mathbf{z}} = \tilde{\boldsymbol{\psi}} = O(h_1 + h_2)$ (vgl. (14.5))

$$(14.9) \quad \|\tilde{\mathbf{z}}\|_\infty \leq \|\tilde{\mathbf{A}}_h^{-1}\|_{C(\bar{\omega}_h)} \|\tilde{\boldsymbol{\psi}}\|_\infty = O(h_1 + h_2).$$

Dieses Ergebnis können wir jedoch verbessern mit Hilfe der schärferen Abschätzung Satz 6.3 c)

Satz 14.2 Shortly-Wellers für die 1. RWA der Poissongleichung

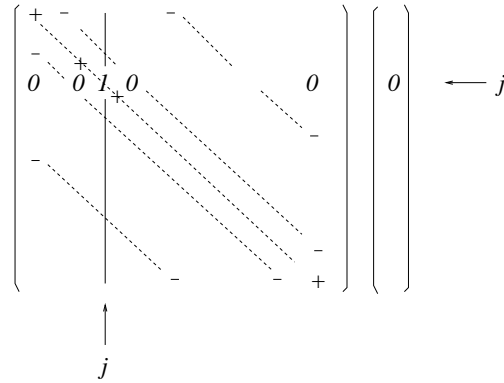
Falls $u \in C^4(\bar{\Omega})$ konvergiert das Verfahren gemäß (14.4) quadratisch.

Beweis: Idee:

Es ist $\bar{\omega}_h = \gamma_h \cup \omega_h^* \cup \omega_h^0$. Die Konvergenzabschätzung erhält man aus dem Gleichungssystem $\tilde{\mathbf{A}}_h \tilde{\mathbf{z}} = \tilde{\boldsymbol{\psi}}$. Nun sind bei der 1. RWA die Komponenten von $\tilde{\mathbf{z}}$ und $\tilde{\boldsymbol{\psi}}$, die zu den Komponenten $\in \gamma_h$ gehören, =0. Daher kann man die Randwerte aus dem Gleichungssystem eliminieren, sodaß nach der Elimination keine Randwerte mehr auftauchen.

Elimination: (vgl. Zeichnung)

Zu jedem Randpunkt gehört in $\tilde{\mathbf{A}}_h$ eine Zeile (z.B. j), in der die 1 in der Diagonale das einzige Element $\neq 0$ ist. Die Zeile j im Gleichungssystem $\tilde{\mathbf{A}}_h \tilde{\mathbf{z}} = \tilde{\boldsymbol{\psi}}$ und die Spalte j in $\tilde{\mathbf{A}}_h$ werden gestrichen, letztere weil sie mit $\tilde{z}_j = 0$ multipliziert wird. Auf diese Weise entsteht aus $\tilde{\mathbf{A}}_h$ die Matrix \mathbf{A}_h . Durch das Streichen der j -ten Spalte werden nur Nebendiagonalelemente von $\tilde{\mathbf{A}}_h$ und \mathbf{A}_h entfernt, was in den entsprechenden Zeilen die Hauptdiagonale stärkt. Das bedeutet: Geht \mathbf{p} aus $\tilde{\mathbf{p}}$ hervor durch Streichen der Komponenten, die zu den Randwerten gehören, ist $(\mathbf{A}_h)_i$ eine Zeile, in der ein Nebendiagonalelement gestrichen wurde, $\tilde{\mathbf{A}}_{h\tilde{i}}$ die entsprechende Zeile von $\tilde{\mathbf{A}}_h$, so gilt $(\mathbf{A}_h)_i \mathbf{p} > (\tilde{\mathbf{A}}_h)_{\tilde{i}} \tilde{\mathbf{p}}$, denn die Nebendiagonalelemente sind ≤ 0 . Da $\tilde{\mathbf{A}}_h$ eine M -Matrix war, gilt dies auch für \mathbf{A}_h , denn die Vorzeichenverteilung bleibt dieselbe.



Für \mathbf{A}_h untersuchen wir nun für die Fehlerfunktion \mathbf{z} das System

$$\mathbf{A}_h \mathbf{z} = \boldsymbol{\psi},$$

das wir in zwei Teile zerlegen

$$(14.10) \quad \mathbf{A}_h \mathbf{z}^0 = \boldsymbol{\psi}^0, \quad \boldsymbol{\psi}^0 = \begin{cases} \psi_{kl}, & (k, l) \in \omega_h^0 \\ 0, & (k, l) \in \omega_h^* \end{cases} \implies \|\boldsymbol{\psi}^0\|_{C(\omega_h^0)} = O(h_1^2 + h_2^2)$$

$$(14.11) \quad \mathbf{A}_h \mathbf{z}^* = \boldsymbol{\psi}^*, \quad \boldsymbol{\psi}^* = \begin{cases} \psi_{kl}, & (k, l) \in \omega_h^* \\ 0, & (k, l) \in \omega_h^0 \end{cases} \xrightarrow{\text{zeige}} \|\boldsymbol{\psi}^*\|_{C(\omega_h^*)} = O(h_1^2 + h_2^2)$$

Dann gilt

$$(14.12) \quad \mathbf{A}_h \mathbf{z} = \mathbf{A}_h (\mathbf{z}^0 + \mathbf{z}^*) = \boldsymbol{\psi}^0 + \boldsymbol{\psi}^*.$$

Wir betrachten beide Aufgaben getrennt.

Zu Aufgabe (14.10) benutzen wir die verschärfte Abschätzung Satz 6.3 c)

$$(14.13) \quad \|\mathbf{z}\|_\infty \leq \frac{\|\mathbf{p}\|_\infty}{\min_{(i,k) \in \omega_h^0} (\mathbf{A}\mathbf{p})_{ik}} \|\boldsymbol{\psi}^0\|_\infty \quad \forall \mathbf{p} > \mathbf{0} \text{ mit } (\mathbf{A}\mathbf{p})_{ik} > \mathbf{0} \quad \forall (i, k) \in \omega_h^0.$$

Es genügt hier $w(x_1, x_2) = c - (x_1 - x_{10})^2 - (x_2 - x_{20})^2$, $\mathbf{p} = \mathbf{w}$ zu wählen mit $c \geq (\text{diam}(\Omega))^2$, denn $\min w$ wird nur auf γ_h angenommen (vgl. Bemerkung c) nach (14.6)), also ist $\mathbf{w} > \mathbf{0}$ in $\omega_h = \omega_h^* \cup \omega_h^0$ und $\mathbf{A}_h \mathbf{w} = 4$ in ω_h^0 (vgl. (14.7)). Man erhält also

$$(14.14) \quad \|\mathbf{z}^0\|_{C(\omega_h)} \leq \frac{c}{4} \underbrace{\|\boldsymbol{\psi}^0\|_\infty}_{O(h_1^2 + h_2^2)} \leq \frac{c}{4} M_4 (h_1^2 + h_2^2) \quad \text{nach (12.3)}.$$

Das ist die gewöhnliche Abschätzung für Δu in den inneren Punkten.

Zu Aufgabe (14.11) benutzen wir das schwache Maximumprinzip (vgl. dazu die Sätze 12.3 - 12.4) \mathbf{A}_h ist eine invertierbare, zumindest schwach diagonaldominante M -Matrix

(vgl. (14.16)). Wir wenden es auf $\mathbf{A}_h \mathbf{z}^* = \boldsymbol{\psi}^*$ und erhalten wegen $\omega_h^0 \subset N^0(\mathbf{A}_h \mathbf{z}^*)$, $N^\neq(\mathbf{A}_h \mathbf{z}^*) \subset \omega_h^*$ (identifiziere die Gitterpunkte mit ihren Indizes)

$$\max_{i \in \omega_h^0} |z_i^*| \leq \max_{i \in N^0(\mathbf{A}_h \mathbf{z}^*)} |z_i^*| \leq \max_{j \in N^\neq(\mathbf{A}_h \mathbf{z}^*)} |z_j^*| \leq \max_{j \in \omega_h^*} |z_j^*|.$$

Die maximale z -Komponente gehört also zu einem Gitterpunkt $\in \omega_h^*$.

Sei $\mathbf{A}_h = (a_{ij}) \in \mathbb{R}^{n \times n}$, dann gilt für $i \in \omega_h^*$ mit $|z_i^*| = \max_j |z_j^*|$ wegen $\mathbf{A}_h \mathbf{z}^* = \boldsymbol{\psi}^*$

$$\begin{aligned} a_{ii}|z_i^*| &= |a_{ii}z_i^*| = |\psi_i^* - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij}z_j^*| \leq |\psi_i^*| + \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| |z_j^*| \leq |\psi_i^*| + \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| |z_i^*| \\ \implies |z_i^*| |a_{ii} - \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|| &\leq |\psi_i^*| \end{aligned}$$

Wir zeigen, daß die Zeilensumme $\neq 0$ ist. Dann folgt

$$(14.15) \quad |z_i^*| \leq \frac{1}{|a_{ii} - \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}||} |\psi_i^*|.$$

Eine Abschätzung der Zeilensumme nach unten wird für $|z_i^*|$ die gewünschte Abschätzung liefern.

Dazu betrachten wir zunächst für die volle Matrix $\tilde{\mathbf{A}}_h$ und den dazugehörigen längeren Vektor $\tilde{\mathbf{e}} = (1, \dots, 1)^T$ die zu einem Punkt $(k, l) \in \omega_h^*$ gehörige Komponente $(\tilde{\mathbf{A}}_h \tilde{\mathbf{e}})_{k,l}$ in der Diskretisierung (14.4) (für zeilenweise Nummerierung)

$$(14.16) \quad \begin{aligned} (\tilde{\mathbf{A}}_h \tilde{\mathbf{e}})_{k,l} &= - \left(\frac{2}{h_2^+ + h_2^-} \frac{1}{h_2^-} \right) \tilde{e}_{k,l-1} - \left(\frac{2}{h_1^+ + h_1^-} \frac{1}{h_1^-} \right) \tilde{e}_{k-1,l} \\ &+ \left(\frac{2}{h_1^+ + h_1^-} \frac{1}{h_1^-} + \frac{2}{h_1^+ + h_1^-} \frac{1}{h_1^+} + \frac{2}{h_2^+ + h_2^-} \frac{1}{h_2^+} + \frac{2}{h_2^+ + h_2^-} \frac{1}{h_2^-} \right) \tilde{e}_{k,l} \\ &- \left(\frac{2}{h_1^+ + h_1^-} \frac{1}{h_1^+} \right) \tilde{e}_{k+1,l} - \left(\frac{2}{h_2^+ + h_2^-} \frac{1}{h_2^+} \right) \tilde{e}_{l,k+1} \end{aligned}$$

Hierdurch wird die Zeilensumme der Zeile (k, l) der Matrix dargestellt.

In der 1. Zeile stehen die Elemente links der Hauptdiagonalen, in der 2. Zeile die Hauptdiagonalelemente und in der 3. Zeile die Elemente rechts der Hauptdiagonalen.

Diese Darstellung gilt in ω_h^* und in ω_h^0 (wobei im letzteren Fall alle $h_{1,2}^\pm = h_{1,2}$).

Beachte: $\tilde{\mathbf{A}}_h$ ist schwach diagonaldominant. Die Zeilensummen, die zu Elementen von $\omega_h^0 \cup \omega_h^*$ gehören, sind =1.

Beim Übergang zur Matrix \mathbf{A}_h werden die Randpunkte $\in \gamma_h$ gestrichen, d.h. in \mathbf{A}_h fehlt in den Punkten $\in \omega_h^*$ mindestens ein Summand in (14.16), *mindestens*, weil ein irregulärer Punkt als Nachbarn mehr als einen Randpunkt haben kann.

In \mathbf{A}_h fehlen mindestens

1. bei einem unteren Randpunkt $\tilde{e}_{k,l-1}$ eine frühere Komponente von \mathbf{A}_h
2. bei einem linken Randpunkt $\mathbf{e}_{k-1,l}$ eine frühere Komponente von \mathbf{A}_h
3. bei einem rechten Randpunkt $\tilde{e}_{k+1,l}$ eine spätere Komponente von \mathbf{A}_h
4. bei einem oberen Randpunkt $\tilde{e}_{l,k+1}$ eine spätere Komponente von \mathbf{A}_h

Um den jeweils gestrichenen Term überwiegt der entsprechende Term in der Hauptdiagonalen. Also ist $(\mathbf{A}_h \mathbf{e})_{k,l}$, $(k, l) \in \omega_h^*$ größer um den jeweils fehlenden Term, den wir abschätzen.

$$\frac{2}{h_2^+ + h_2^-} \frac{1}{h_2^-} \geq \frac{2}{2h_2} \frac{1}{h_2} = \frac{1}{h_2^2}, \quad \frac{2}{h_1^+ + h_1^-} \frac{1}{h_1^-} \geq \frac{2}{2h_1} \frac{1}{h_1} = \frac{1}{h_1^2}$$

entsprechend die anderen beiden Terme \implies

$$(\tilde{\mathbf{A}}_h \tilde{\mathbf{e}})_{kl} \geq \min \left(\frac{1}{h_1^2}, \frac{1}{h_2^2} \right) \implies \frac{1}{(\tilde{\mathbf{A}}_h \tilde{\mathbf{e}})_{kl}} \leq \frac{1}{\min \left(\frac{1}{h_1^2}, \frac{1}{h_2^2} \right)} \leq h_1^2 + h_2^2 \quad \forall (k, l) \in \omega_h^*$$

Wir tragen diese Abschätzung in (14.14) ein und erhalten

$$(14.17) \quad |z_i^*| = \|\mathbf{z}^*\|_\infty \leq (h_1^2 + h_2^2) |\psi_i^*| = O(h_1^2 + h_2^2) \cdot O(h_1 + h_2)$$

■.

Bemerkungen zum Verfahren

1. Dieses Verfahren arbeitet gut bei Randwertaufgaben.
2. Bei Eigenwertaufgaben ist es problematisch. Die kontinuierliche Aufgabe ist oft selbstadjungiert, zumindest immer formal selbstadjungiert. Die Matrix $\tilde{\mathbf{A}}_h$ ist aber nicht symmetrisch, hat also ggf. komplexe Eigenwerte. Die Asymmetrie von $\tilde{\mathbf{A}}_h$ stört.

Aufgabe:

- (a) Zeige, daß \mathbf{A}_h (Randwerte auf die rechte Seite gebracht) symmetrisch ist bei beliebigen, einfach zusammenhängenden Gebieten.
- (b) Zeige, daß $\tilde{\mathbf{A}}_h$ nicht symmetrisch ist.
3. Die Abschätzung (14.17) läßt für die Diskretisierung in den irregulären Randpunkten Spielraum ohne die quadratische Konvergenz zu gefährden. Man kann also “unfaire” Approximationen für Δu benutzen, oder “gar keine”, d.h. einen Fehler von $O(1)$, man verliert dann den Faktor $O(h_1 + h_2)$, bzw. muß ihn durch eine Konstante ersetzen. Das tut der quadratischen Konvergenz nichts. Solche “unfairen” Approximationen werden benutzt.

Insbesondere: Ersetzt man in (14.4) alle $h_{1,2}^\pm$ (, die auch noch vom irregulären

Punkt abhängen) durch $h_{1,2}$, so geht die Ordnung in (14.1) verloren und somit auch der Faktor $O(h_1 + h_2)$ in (14.15). Die Abschätzungen von (14.15) zu (14.17) bleiben jedoch richtig. Der Faktor $O(h_1^2 + h_2^2)$ in (14.17) bleibt erhalten, und somit auch die quadratische Konvergenz.

Daß der Einfluß der Diskretisierung in den irregulären Randpunkten auf die Konvergenzgeschwindigkeit gering ist, kann mit Hilfe der Greenschen Funktion erklärt werden, die in den randnahen Punkten klein ausfällt.

4. Statt Δu in irregulären Punkten zu diskretisieren, kann man für die Funktionswerte in inneren Punkten auch einen Wert durch lineare Interpolation der Funktionswerte in den benachbarten Punkten erhalten.

§ 15 Jacobi und Gauss-Seidel

Direkte Verfahren zur Lösung großer linearer Gleichungssysteme sind im allgemeinen zu aufwendig und zu ungenau. Deshalb werden iterative Verfahren vorgezogen. Jacobi-Verfahren (Gesamtschrittverfahren) zur Lösung von

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b}, \mathbf{A} \in \mathbb{R}^{n \times n}, \mathbf{A} = (\mathbf{L} + \mathbf{D} + \mathbf{R}) \\ \mathbf{L} &= \text{untere Dreiecksmatrix von } \mathbf{A} \\ \mathbf{D} &= \text{Diagonalmatrix von } \mathbf{A} \\ \mathbf{L} &= \text{obere Dreiecksmatrix von } \mathbf{A} \end{aligned}$$

$$\begin{aligned} \mathbf{x}^{m+1} &= -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{R})\mathbf{x}^m + \mathbf{D}^{-1}\mathbf{b} = -\mathbf{D}^{-1}(\mathbf{A} - \mathbf{D})\mathbf{x}^m + \mathbf{D}^{-1}\mathbf{b} \\ &= (\mathbf{I} - \mathbf{D}^{-1}\mathbf{A})\mathbf{x}^m + \mathbf{D}^{-1}\mathbf{b} = \mathbf{x}^m + \mathbf{D}^{-1}(\mathbf{b} - \mathbf{Ax}^m) \end{aligned}$$

Entsprechend lautet das ‘‘gedämpfte’’ Jacobi-Verfahren (der Defekt wird gedämpft)

$$(15.1) \quad \mathbf{x}^{m+1} = (\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A})\mathbf{x}^m + \omega\mathbf{D}^{-1}\mathbf{b} = \mathbf{x}^m + \omega\mathbf{D}^{-1}(\mathbf{b} - \mathbf{Ax}^m), \quad \omega > 0.$$

Aus (15.1) liest man ab, daß alle Fixpunkte Lösungen von $\mathbf{Ax} = \mathbf{b}$ sind. Subtrahiert man die Fixpunktgleichung

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

so erhält man die Fehlerdarstellung

$$(15.2) \quad \begin{aligned} \mathbf{e}^{m+1} &:= \mathbf{x}^{m+1} - \mathbf{x} = (\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A})\mathbf{e}^m = \dots = (\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A})^{m+1}\mathbf{e}^0 \\ \|\mathbf{e}^{m+1}\| &\leq \|\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A}\|^{m+1} \|\mathbf{e}^0\|. \end{aligned}$$

Die Iterationsmatrix $\|\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A}\|$ heißt auch Konvergenzrate von (15.1) (normabhängig). Konvergenz liegt vor, wenn die Konvergenzrate < 1 ausfällt (hinreichend), bzw. wenn der Spektralradius der Iterationsmatrix $\rho(\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A}) < 1$ ist (notwendig und hinreichend). Letzteres wird bewiesen durch den

Satz

$\forall \varepsilon > 0 \wedge \mathbf{A} \in \mathbb{R}^{n \times n} \exists$ Vektornorm $\|\cdot\|_V$ und eine zugeordnete Matrixnorm $\|\cdot\|_M$ mit $\|\mathbf{A}\|_M \leq \rho(\mathbf{A}) + \varepsilon$.

Bemerkungen:

1. Ist $\mathbf{A} = \mathbf{A}^T > 0$ und $\mathbf{D}^{-1}(\mathbf{A}) = c\mathbf{I}$ (dies ist gegeben, wenn \mathbf{A} die Diskretisierungsmatrix von Δu ist), so ist $\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A}$ symmetrisch und es gilt

$$(15.3) \quad \rho(\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A}) = \max |\lambda(\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A})| = \|\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A}\|_S = \|\mathbf{I} - \omega c\mathbf{A}\|_S.$$

2. Übung: Man zeige, daß $\lambda \in \sigma(\mathbf{A}) \iff (1 - \lambda) \in \sigma(\mathbf{I} - \mathbf{A})$.
 $\sigma(\mathbf{A})$ bezeichnet das Spektrum von \mathbf{A} .

Satz 15.1

Für das gedämpfte Jacobi-Verfahren zur Lösung von $\mathbf{Ax} = \mathbf{b}$ mit $\mathbf{A} = \mathbf{A}^T > 0$

$$\mathbf{x}^{m+1} = (\mathbf{I} - \omega \mathbf{D}^{-1} \mathbf{A}) \mathbf{x}^m + \omega \mathbf{D}^{-1} \mathbf{b}, \quad \mathbf{D} = \text{diag}(\mathbf{A}), \quad \omega > 0, \quad \mathbf{x}^0 \text{ gegeben}$$

gilt

Alle Eigenwerte von $\mathbf{D}^{-1} \mathbf{A}$ sind positiv.

Sind λ_{\min} bzw. λ_{\max} der minimale bzw. maximale Eigenwert von $\mathbf{D}^{-1} \mathbf{A}$, so wird für den Dämpfungsfaktor

$$(15.4) \quad \omega_{\text{opt}} = \frac{2}{\lambda_{\min} + \lambda_{\max}}$$

der Spektralradius der Iterationsmatrix $\rho(\omega) := \rho(\mathbf{I} - \omega \mathbf{D}^{-1} \mathbf{A})$ minimal und

$$(15.5) \quad \rho(\omega_{\text{opt}}) = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\min} + \lambda_{\max}} < 1.$$

Das mit ω_{opt} gedämpfte Verfahren konvergiert für jeden Anfangswert gegen eine Lösung von $\mathbf{Ax} = \mathbf{b}$.

Bemerkungen:

1. Für das ungedämpfte Verfahren ($\omega = 1$) ist der Satz falsch (vgl. Übungen).
2. Ist $\mathbf{D} = c\mathbf{I}$, so können in (15.5) für λ_{\min} und λ_{\max} die entsprechenden Eigenwerte von \mathbf{A} eingetragen werden ($\sigma(c\mathbf{A}) = c\sigma(\mathbf{A})$).
3. Zur praktischen Anwendung benötigt man (zumindest Schätzungen für) λ_{\min} und λ_{\max} .

Beweis des Satzes:

$$\begin{aligned} \lambda \in \sigma(\mathbf{D}^{-1} \mathbf{A}) &\iff \exists \mathbf{x} \neq 0 : \mathbf{D}^{-1} \mathbf{Ax} = \lambda \mathbf{x} \\ &\iff \mathbf{Ax} = \lambda \mathbf{Dx} \\ &\implies \mathbf{x}^* \mathbf{Ax} = \lambda \mathbf{x}^* \mathbf{Dx} \\ &\iff \lambda = \frac{\mathbf{x}^* \mathbf{Ax}}{\mathbf{x}^* \mathbf{Dx}} \end{aligned}$$

Da \mathbf{A} und \mathbf{D} positiv definit sind, folgt $\lambda > 0 \quad \forall \lambda \in \sigma(\mathbf{D}^{-1} \mathbf{A})$.

$$0 < \lambda_{\min} \leq \lambda(\mathbf{D}^{-1}\mathbf{A}) \leq \lambda_{\max} \xrightarrow{\omega > 0} 1 - \omega\lambda_{\max} \leq 1 - \omega\lambda_{\min}$$

$\rho(\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A})$ wird in Abhängigkeit von ω minimal, falls

$$|1 - \omega\lambda_{\max}| = |1 - \omega\lambda_{\min}| \quad \text{bzw.} \quad \omega\lambda_{\max} - 1 = 1 - \omega\lambda_{\min} \quad \text{woraus folgt}$$

$$\omega_{\text{opt}} = \frac{2}{\lambda_{\min} + \lambda_{\max}} \quad \text{und} \quad 1 - \omega\lambda_{\min} = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} < 1. \quad \blacksquare$$

Für nichtsymmetrische Matrizen gilt

Satz 15.2

Ist $\mathbf{A} \in \mathbb{C}^{n \times n}$ irreduzibel und diagonaldominant, d.h.

$$\sum_{j=1, j \neq i}^n |a_{ij}| \leq |a_{ii}| \quad \forall i, \quad < \text{für mindestens ein } i,$$

so konvergiert das Jacobi-Verfahren für jeden Startwert gegen die Lösung von $\mathbf{Ax} = \mathbf{b}$.

Definition 15.3

$\mathbf{A} \in \mathbb{C}^{n \times n}$ heißt *reduzibel*, wenn gilt

1. $\exists N_1, N_2 \subset N = \{1, 2, \dots, n\}, N_1, N_2 \neq \emptyset$
2. $N_1 \cup N_2 = N, N_1 \cap N_2 = \emptyset$
3. $a_{ij} = 0 \quad \forall (i, j) \in N_1 \times N_2.$

\mathbf{A} heißt *irreduzibel* $\iff \mathbf{A}$ ist nicht reduzibel.

Beweis von Satz (15.2): Übung.

Bemerkung zur Wahl der Startwerte:

Eine gute Wahl der Startwerte ist im Allgemeinen schwierig.

Oft wird empfohlen: $\mathbf{x}^0 = \mathbf{D}^{-1}\mathbf{b}$. Dies ist nichts anderes als der Iterationswert \mathbf{x}^1 für $\mathbf{x}^0 = \mathbf{0}$ und $\omega = 1$ mit dem Jacobi-Verfahren.

Beim (noch zu besprechenden) Mehrgitter-Verfahren wird eine Näherung mitgeliefert.

Abbruchkriterien

In der Praxis wird oft empfohlen: $\|\mathbf{x}^{m+1} - \mathbf{x}^m\| \leq \text{eps}$ (Maschinenengenauigkeit). Dies ist problematisch, insbesondere bei schlechter Kondition, weil dieses Kriterium nichts

über die Genauigkeit aussagen muß.

Empfehlung: Defektabschätzung über die Kondition.

Es ist $\mathbf{A}(\underbrace{\mathbf{x}^m - \mathbf{x}}_{\mathbf{e}^m}) = \mathbf{A}\mathbf{x}^m - \mathbf{b} =: \mathbf{d}^m$ Defektvektor, also

$$\|\mathbf{e}^m\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{d}^m\|.$$

Aus $\|\mathbf{b}\| = \|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$ folgt $\|\mathbf{x}\| \geq \frac{\|\mathbf{b}\|}{\|\mathbf{A}\|}$ also gilt für den relativen Fehler

$$\frac{\|\mathbf{e}^m\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\mathbf{d}^m\|}{\|\mathbf{b}\|} = \text{cond}(\mathbf{A}) \frac{\|\mathbf{d}^m\|}{\|\mathbf{b}\|}.$$

Ein vernünftiges Abbruchkriterium

$$(15.6) \quad \frac{\|\mathbf{e}^m\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{d}^m\|}{\|\mathbf{b}\|} \stackrel{!}{\leq} \varepsilon$$

ist nur möglich, wenn man Werte oder Schätzungen für $\text{cond}(\mathbf{A})$ hat. Gute Formelpakete liefern das.

Schätzung einer Iterationszahl zur Erreichung einer vorgegebenen Genauigkeit

Satz 15.4

Für das gedämpfte Jacobi-Verfahren zur Lösung von $\mathbf{A}\mathbf{x} = \mathbf{b}$ gelte $\text{diag}(\mathbf{A}) = c\mathbf{I}$ und die Voraussetzungen von Satz (15.1) seien erfüllt (insbesondere gilt dann $\rho := \rho(\omega_{opt}) < 1$).

Dann gilt für den absoluten Fehler nach m Iterationen

$$\|\mathbf{e}^m\|_2 \leq \varepsilon,$$

falls

$$m = m(\varepsilon) \geq \frac{\log \|\mathbf{e}^0\|_2 + \log(1/\varepsilon)}{\log(1/\rho)} \quad \text{und}$$

$$\frac{\log \|\mathbf{e}^0\|_2 + \log(1/\varepsilon)}{\log(1/\rho)} \leq (\log \|\mathbf{e}^0\|_2 + \log(1/\varepsilon)) \frac{1}{2} \text{cond}(\mathbf{A}),$$

wobei \mathbf{e}^0 = Fehler der Ausgangsnäherung.

Bemerkungen:

1. $\text{diag}(\mathbf{A}) = c\mathbf{I}$ ist nur technisch. Es ist immer $\text{diag}(\mathbf{D}^{-1}\mathbf{A}) = \mathbf{I}$.
2. Die notwendige Iterationszahl m wächst (in etwa) linear mit $\text{cond}(\mathbf{A})$.

3. Für $\|\mathbf{e}^0\|_2 \approx 1$ ist $\log\|\mathbf{e}^0\|_2 \approx 0$, deshalb lautet eine Schätzung (für relativ gute Anfangsnäherungen)

$$(15.7) \quad m(\varepsilon) \approx \frac{1}{2}(\log(1/\varepsilon)) \operatorname{cond}(\mathbf{A})$$

Beweis:

Wir haben gezeigt (vgl. (15.2))

$$\|\mathbf{e}^{m+1}\|_2 \leq \|\mathbf{I} - \omega \mathbf{D}^{-1} \mathbf{A}\|_2^m \|\mathbf{e}^0\|_2 = \rho(\omega_{opt})^m \|\mathbf{e}^0\|_2.$$

Aus der Forderung $\rho^m \|\mathbf{e}^0\|_2 \leq \varepsilon$ folgt (beachte: $\log \rho < 0$)

$$(15.8) \quad m \geq \frac{\log(\varepsilon) - \log(\|\mathbf{e}^0\|_2)}{\log(\rho)} = \frac{\log(\|\mathbf{e}^0\|_2 + \log(1/\varepsilon))}{\log(1/\rho)}$$

Wir zeigen $\frac{1}{\log(1/\rho)} \leq \frac{1}{2} \operatorname{cond}_2(\mathbf{A})$.

Bekannt ist für die Eigenwerte von $\mathbf{D}^{-1} \mathbf{A}$ (vgl. (15.5))

$$\rho(\omega_{opt}) = \frac{\lambda_{max} - \lambda_{min}}{\lambda_{max} + \lambda_{min}} = \frac{\frac{\lambda_{max}}{\lambda_{min}} - 1}{\frac{\lambda_{max}}{\lambda_{min}} + 1} = \frac{\kappa - 1}{\kappa + 1},$$

dabei ist wegen $\mathbf{D} = \operatorname{diag} \mathbf{A} = c \mathbf{I}$

$$\kappa = \operatorname{cond}_2(\mathbf{D}^{-1} \mathbf{A}) = \operatorname{cond}_2(c^{-1} \mathbf{A}) = \|c^{-1} \mathbf{A}\|_2 \|c \mathbf{A}^{-1}\|_2 = \operatorname{cond}_2(\mathbf{A}) = \frac{\lambda_{max}(A)}{\lambda_{min}(A)}.$$

In (15.8) schätzen wir $\log(1/\rho)$ nach unten ab. Nun ist

$$\log(1/\rho) = \log\left(\frac{\kappa + 1}{\kappa - 1}\right) = \log\left(1 + \frac{2}{\kappa - 1}\right)$$

Mit der Abschätzung $\log(1 + x) \geq \frac{2x}{x+2}$ (Beweis: richtig für $x = 0$, und für $x \geq 0$ wächst die linke Seite stärker als die rechte) folgt

$$\log(1/\rho) = \log\left(1 + \underbrace{\frac{2}{\kappa - 1}}_{=x}\right) \geq \frac{2x}{x+2} = \frac{2 \frac{2}{\kappa - 1}}{\frac{2}{\kappa - 1} + 2} = \frac{4}{2(\kappa - 1) + 2} = \frac{2}{\kappa}.$$

Trägt man dies in (15.8) ein, so folgt die Behauptung. ■

Anwendung auf die Diskretisierungsmatrix $\mathbf{A}_h^0 = \frac{1}{h^2} \operatorname{tridiag}(-1, 2, -1)$.

Für diese Matrix (Diskretisierung von u_{xx} in $[0, 1]$) wurde gezeigt (vgl. (2.11), (3.16), (3.17))

$$8 \leq \lambda_{min} \leq \lambda_{max} \leq \frac{4}{h^2}.$$

Mit den Schätzungen

$$\|\mathbf{A}^{-1}\|_S = \frac{1}{\lambda_{\min}(\mathbf{A})} \approx \frac{1}{8}, \quad \|\mathbf{A}\|_S = \lambda_{\max}(\mathbf{A}) \approx \frac{4}{h^2} \quad \text{ist} \quad \text{cond}_2(\mathbf{A}) \approx \frac{1}{2h^2}.$$

Die Schätzung (15.7) für die Iterationszahl $m(\varepsilon)$ liefert also

$$m \approx \frac{1}{2} \left(\log\left(\frac{1}{\varepsilon}\right) \right) \text{cond}\mathbf{A} = \frac{1}{4h^2} \log\left(\frac{1}{\varepsilon}\right) \stackrel{h=1/N}{=} N^2 \frac{\log\left(\frac{1}{\varepsilon}\right)}{4}.$$

Da eine Jacobiiteration mindestens N^2 Rechenoperationen braucht (N =Zahl der Unbekannten), folgt

$$m \approx N^4 \frac{\log\left(\frac{1}{\varepsilon}\right)}{4} \text{Iterationsschritte.}$$

Zu einem vorgegebenen absoluten Fehler ε , der sich natürlich an einem sinnvollen relativen Fehler ausrichten wird, wächst also auch die Zahl der notwendigen Rechenoperationen in etwa proportional zu N^4 .

Fazit: Die Lösung einer solchen Gleichung durch reine Iteration ist zu aufwendig.

Das Gauß-Seidel-Verfahren

(Einzelschrittverfahren (ESV))

Mit der Matrixzerlegung

$$(15.9) \quad \mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{R}$$

erhält man zur Lösung von $\mathbf{A}\mathbf{x} = \mathbf{b}$ wegen $(\mathbf{L} + \mathbf{D})\mathbf{x} = -\mathbf{R}\mathbf{x} + \mathbf{b}$ falls \mathbf{D} regulär ist, das

$$(15.10) \quad \text{ESV:} \quad \mathbf{x}^{m+1} = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{R}\mathbf{x}^m + (\mathbf{L} + \mathbf{D})^{-1}\mathbf{b}$$

mit der Iterationsmatrix

$$(15.11) \quad \mathbf{S} = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{R} = -(\mathbf{L} + \mathbf{D})^{-1}(\mathbf{A} - (\mathbf{L} + \mathbf{D})) = \mathbf{I} - (\mathbf{L} + \mathbf{D})^{-1}\mathbf{A}.$$

Beachte: Das Verfahren ist abhängig von der Anordnung der Gleichungen.

Satz 15.5

Ist $\mathbf{A} = \mathbf{A}^T > 0$ (symmetrisch, positiv definit), so konvergiert das ESV für jede Ausgangsnäherung gegen die Lösung von $\mathbf{A}\mathbf{x} = \mathbf{b}$.

Beweis: Wir zeigen für eine geeignete Matrixnorm $\|\mathbf{S}\| < 1$, genauer:

Durch $\|\mathbf{y}\|_A^2 := (\mathbf{A}\mathbf{y}, \mathbf{y}) = (\mathbf{y}, \mathbf{A}\mathbf{y})$ wird eine Vektornorm definiert, die zugeordnete Matrixnorm (Energienorm) ist

$$(15.12) \quad \|\mathbf{S}\|_A = \sup_{\mathbf{y} \neq 0} \frac{\|\mathbf{S}\mathbf{y}\|_A}{\|\mathbf{y}\|_A} = \max_{\|\mathbf{x}\|_A=1} \|\mathbf{S}\mathbf{x}\|_A.$$

Das Maximum wird angenommen, da $\|\mathbf{x}\|_A = 1$ eine kompakte Menge ist. Wir zeigen

$$(15.13) \quad \|\mathbf{S}\mathbf{y}\|_A < \|\mathbf{y}\|_A \quad \forall \mathbf{y} \neq 0.$$

Gemäß der Definition (15.12) gibt es dann ein $\tilde{\mathbf{y}}$ mit $\|\mathbf{S}\tilde{\mathbf{y}}\|_A = \|\mathbf{S}\|_A \|\tilde{\mathbf{y}}\|_A$. Deshalb folgt aus (15.13): $\|\mathbf{S}\|_A < 1$.

Beweis (15.13): (der Einfachheit halber schreiben wir $\|\cdot\| = \|\cdot\|_A$).

Mit $\mathbf{z} = \mathbf{S}\mathbf{y} = (\mathbf{I} - (\mathbf{L} + \mathbf{D})^{-1}\mathbf{A})\mathbf{y}$, also $\|\mathbf{z}\|^2 = \mathbf{z}^T \mathbf{A} \mathbf{z}$ folgt

$$\begin{aligned} \|\mathbf{y}\|^2 - \|\mathbf{S}\mathbf{y}\|^2 &= \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{z}^T \mathbf{A} \mathbf{z} \\ &= \mathbf{y}^T \mathbf{A} \mathbf{y} - \mathbf{y}^T (\mathbf{I} - \mathbf{A}^T (\mathbf{L} + \mathbf{D})^{-1})^{-T} \mathbf{A} (\mathbf{I} - (\mathbf{L} + \mathbf{D})^{-1} \mathbf{A}) \mathbf{y} \\ &\quad \text{mit der Abkürzung } \mathbf{B} = (\mathbf{L} + \mathbf{D}) \\ &= \mathbf{y}^T \{ \mathbf{A} - (\mathbf{I} - \mathbf{A}^T \mathbf{B}^{-T}) \mathbf{A} (\mathbf{I} - \mathbf{B}^{-1} \mathbf{A}) \} \mathbf{y} \\ &= \mathbf{y}^T \{ \mathbf{A} - (\mathbf{A} - \mathbf{A}^T \mathbf{B}^{-T} \mathbf{A}) (\mathbf{I} - \mathbf{B}^{-1} \mathbf{A}) \} \mathbf{y} \quad \stackrel{\mathbf{A}=\mathbf{A}^T}{\Longrightarrow} \\ &= \mathbf{y}^T \mathbf{A}^T \{ \mathbf{I} - (\mathbf{I} - \mathbf{B}^{-T} \mathbf{A}) (\mathbf{I} - \mathbf{B}^{-1} \mathbf{A}) \} \mathbf{y} \\ &= \mathbf{y}^T \mathbf{A}^T \{ \mathbf{I} - (\mathbf{I} - \mathbf{B}^{-T} \mathbf{A} - \mathbf{B}^{-1} \mathbf{A} + \mathbf{B}^{-T} \mathbf{A} \mathbf{B}^{-1} \mathbf{A}) \} \mathbf{y} \\ &= \mathbf{y}^T \mathbf{A}^T \{ \mathbf{B}^{-T} + \mathbf{B}^{-1} - \mathbf{B}^{-T} \mathbf{A} \mathbf{B}^{-1} \} \mathbf{A} \mathbf{y}. \end{aligned}$$

Wegen $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{L}^T = \mathbf{B} + (\mathbf{L} + \mathbf{D})^T - \mathbf{D} = \mathbf{B} + \mathbf{B}^T - \mathbf{D}$ folgt

$$\begin{aligned} \mathbf{B}^{-T} \mathbf{A} \mathbf{B}^{-1} &= \mathbf{B}^{-T} (\mathbf{B} + \mathbf{B}^T - \mathbf{D}) \mathbf{B}^{-1} \\ &= \mathbf{B}^{-T} (\mathbf{I} + \mathbf{B}^T \mathbf{B}^{-1}) - \mathbf{D} \mathbf{B}^{-1} \\ &= \mathbf{B}^{-T} + \mathbf{B}^{-1} - \mathbf{B}^{-T} \mathbf{D} \mathbf{B}^{-1}. \end{aligned}$$

Setzt man dies ein, so folgt

$$\|\mathbf{y}\|^2 - \|\mathbf{S}\mathbf{y}\|^2 = \mathbf{y}^T \mathbf{A}^T \mathbf{B}^{-T} \mathbf{D} \mathbf{B}^{-1} \mathbf{A} \mathbf{y} = (\mathbf{B}^{-1} \mathbf{A} \mathbf{y})^T \mathbf{D} (\mathbf{B}^{-1} \mathbf{A} \mathbf{y}) > 0,$$

denn \mathbf{D} hat wegen $\mathbf{A} = \mathbf{A}^T > 0$ und $(\mathbf{e}^i, \mathbf{A} \mathbf{e}^i) = a_{ii} > 0$ nur positive Diagonalelemente, ist also positiv definit.

Bemerkungen

1. Die Diskretisierungsmatrizen für Δu und \mathbf{A}_h erfüllen die Voraussetzungen von Satz (15.5).
2. Der Iterationsmatrix (15.11) der Gauß-Seidel-Iteration fehlen Symmetrieeigenschaften. Dies wird sich im Zusammenhang mit dem Mehrgitterverfahren als Nachteil erweisen (vgl. § 17 f). Wir behandeln deshalb das symmetrische Gauß-Seidel-Verfahren.

Das symmetrische Gauß-Seidel-Verfahren

Man kann Einzelschrittverfahren durch verschiedene Zerlegungen der Matrix \mathbf{A} erhalten.

$$\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{R}, \quad \mathbf{B}_1 = \mathbf{L} + \mathbf{D}, \quad \mathbf{B}_2 = \mathbf{D} + \mathbf{R}.$$

Das gewöhnliche Gauß-Seidel-Verfahren benutzt \mathbf{B}_1 und führt zur Iteration (vgl. (15.9),(15.10))

$$(15.14) \quad \mathbf{x}^{m+1} = (\mathbf{I} - \mathbf{B}_1^{-1}\mathbf{A})\mathbf{x}^m + \mathbf{B}_1^{-1}\mathbf{b}.$$

das rückwärtsgenommene ESV benutzt \mathbf{B}_2 und man erhält analog

$$(15.15) \quad \mathbf{x}^{m+1} = (\mathbf{I} - \mathbf{B}_2^{-1}\mathbf{A})\mathbf{x}^m + \mathbf{B}_2^{-1}\mathbf{b}.$$

Natürlich gilt Satz 15.5 auch für (15.15).

Das symmetrische ESV faßt je einen Schritt von (15.13) und (15.4) zu einem Iterationsschritt zusammen.

$$(15.16) \quad \begin{aligned} \mathbf{x}^{m+\frac{1}{2}} &= (\mathbf{I} - \mathbf{B}_1^{-1}\mathbf{A})\mathbf{x}^m + \mathbf{B}_1^{-1}\mathbf{b} \\ \mathbf{x}^{m+1} &= (\mathbf{I} - \mathbf{B}_2^{-1}\mathbf{A})\mathbf{x}^{m+\frac{1}{2}} + \mathbf{B}_2^{-1}\mathbf{b}. \end{aligned}$$

Für das symmetrische ESV erhält man deshalb

$$(15.17) \quad \mathbf{x}^{m+1} = \underbrace{(\mathbf{I} - \mathbf{B}_2^{-1}\mathbf{A})(\mathbf{I} - \mathbf{B}_1^{-1}\mathbf{A})}_{\mathbf{S}} \mathbf{x}^m + \underbrace{(\mathbf{B}_1^{-1} - \mathbf{B}_2^{-1}\mathbf{A}\mathbf{B}_1^{-1} + \mathbf{B}_2^{-1})}_{\mathbf{Q}} \mathbf{b}$$

Die Iterationsmatrix hat die Darstellung

$$(15.18) \quad \mathbf{S} = \mathbf{I} - \underbrace{(\mathbf{B}_2^{-1} + \mathbf{B}_1^{-1} - \mathbf{B}_2^{-1}\mathbf{A}\mathbf{B}_1^{-1})}_{=: \mathbf{W}^{-1}} \mathbf{A}$$

Satz 15.6 Eigenschaften von \mathbf{W}

Sei $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{R}$, $\mathbf{B}_1 = \mathbf{L} + \mathbf{D}$, $\mathbf{B}_2 = \mathbf{D} + \mathbf{R}$, \mathbf{D} regulär.

a) Für die Iterationsmatrix des symmetrischen Gauß-Seidel-Verfahrens gilt

$$(15.19) \quad \mathbf{S} = \mathbf{I} - (\mathbf{B}_2^{-1} + \mathbf{B}_1^{-1} - \mathbf{B}_2^{-1}\mathbf{A}\mathbf{B}_1^{-1})\mathbf{A} =: \mathbf{I} - \mathbf{W}^{-1}\mathbf{A} \text{ mit } \mathbf{W} = \mathbf{A} + \mathbf{L}\mathbf{D}^{-1}\mathbf{R}.$$

b) Aus $\mathbf{A} = \mathbf{A}^T > 0$ folgt $0 < \mathbf{A} \leq \mathbf{W} = \mathbf{W}^T$ und $\|\mathbf{S}\|_S < 1$ (Spektralnorm)

Bemerkungen:

1. \mathbf{W} ist symmetrisch, falls \mathbf{A} symmetrisch ist.
Beachte auch $\mathbf{B}_1^T = \mathbf{B}_2$, $\mathbf{B}_2^T = \mathbf{B}_1$.

2. Die Bezeichnung

$$\mathbf{W}^{-1} := (\mathbf{B}_2^{-1} + \mathbf{B}_1^{-1} - \mathbf{B}_2^{-1}\mathbf{A}\mathbf{B}_1^{-1})$$

rührt daher, daß (15.17) durch Multiplikation mit \mathbf{W} auch dargestellt werden kann als

$$(15.20) \quad \mathbf{W}(\mathbf{x}^{m+1} - \mathbf{x}^m) = \mathbf{b} - \mathbf{A}\mathbf{x}^m \quad (= \mathbf{d}^m)$$

\mathbf{W} ist dann eine Präkonditionsmatrix.

Hinweis: Eine Möglichkeit ein Iterationsverfahren zu beschleunigen wird untersucht durch Einführen sog. Präkonditionsmatrizen. Ist z.B. die Matrixdiagonale von \mathbf{A} schon auf \mathbf{I} normiert, so lautet das Gesamtschritt-Verfahren (vgl. (15.1) mit $\omega = 1$)

$$\mathbf{x}^{m+1} - \mathbf{x}^m = \mathbf{b} - \mathbf{A}\mathbf{x}^m.$$

Es wird versucht durch Multiplikation der linken Seite mit einer geeigneten (Präkonditions-)Matrix $\tilde{\mathbf{W}}$ ein schnelleres Verfahren der Art

$$\tilde{\mathbf{W}}(\mathbf{x}^{m+1} - \mathbf{x}^m) = \mathbf{b} - \mathbf{A}\mathbf{x}^m$$

zu erhalten. Für die Wahl $\tilde{\mathbf{W}} = \omega^{-1}\mathbf{I}$ mit einer geeigneten Konstanten ω erhält man z.B. das gedämpfte Gesamtschritt-Verfahren. Wählt man $\tilde{\mathbf{W}} = \mathbf{W}$, so erhält man gerade das symmetrische Gauß-Seidel-Verfahren.

3. Man kann das symmetrische Gauß-Seidel-Verfahren so programmieren, daß der Aufwand derselbe ist wie für das Gesamtschritt-Verfahren. (vgl. Hanke-Bourgeois: Grundlagen der Numerischen Mathematik und des Wissenschaftlichen Rechnens. Teubner 2002, §II,8, S. 83).

Beweis a): Wir berechnen (vgl. (15.18))

$$\begin{aligned} \mathbf{B}_2^{-1}\mathbf{A}\mathbf{B}_1^{-1} &= \mathbf{B}_2^{-1}(\mathbf{B}_1 + \mathbf{R})\mathbf{B}_1^{-1} = \mathbf{B}_2^{-1}(\mathbf{B}_1 + \mathbf{B}_2 - \mathbf{D})\mathbf{B}_1^{-1} \\ &= (\mathbf{B}_2^{-1}\mathbf{B}_1 + \mathbf{I} - \mathbf{B}_2^{-1}\mathbf{D})\mathbf{B}_1^{-1} = \mathbf{B}_2^{-1} + \mathbf{B}_1^{-1} - \mathbf{B}_2^{-1}\mathbf{D}\mathbf{B}_1^{-1}. \end{aligned}$$

Eingesetzt in \mathbf{W}^{-1} folgt (vgl. (15.19))

$$\begin{aligned} \mathbf{W}^{-1} &= \mathbf{B}_2^{-1}\mathbf{D}\mathbf{B}_1^{-1} \\ \mathbf{W} &= (\mathbf{L} + \mathbf{D})\mathbf{D}^{-1}(\mathbf{D} + \mathbf{R}) = (\mathbf{L} + \mathbf{D})(\mathbf{I} + \mathbf{D}^{-1}\mathbf{R}) \\ &= \mathbf{L} + \mathbf{D} + \mathbf{L}\mathbf{D}^{-1}\mathbf{R} + \mathbf{R} = \mathbf{A} + \mathbf{L}\mathbf{D}^{-1}\mathbf{R}. \end{aligned}$$

Beweis b): Aus $\mathbf{A} = \mathbf{A}^T$, $\mathbf{L}^T = \mathbf{R}$, $\mathbf{R}^T = \mathbf{L}$ folgt

$$\mathbf{W}^T = \mathbf{A}^T + \mathbf{R}^T\mathbf{D}^{-1}\mathbf{L}^T = \mathbf{A} + \mathbf{L}\mathbf{D}^{-1}\mathbf{R} = \mathbf{W}.$$

Aus $(\mathbf{W}\mathbf{x}, \mathbf{x}) = (\mathbf{A}\mathbf{x}, \mathbf{x}) + (\mathbf{L}\mathbf{D}^{-1}\mathbf{R}\mathbf{x}, \mathbf{x}) = (\mathbf{A}\mathbf{x}, \mathbf{x}) + (\mathbf{D}^{-1}\mathbf{R}\mathbf{x}, \mathbf{R}\mathbf{x})$

folgt wegen $\mathbf{D}^{-1} > 0$ sofort

$$\mathbf{W} \geq \mathbf{A}.$$

Beweis $\|\mathbf{S}\|_2 < 1$ als Übung (Hinweis: Satz (15.5) verwenden).

Kapitel III

Das Mehrgitterverfahren (MGV)

§ 16 Motivation und Grobstruktur

Ausgangssituation: Löse

$$(16.1) \quad \mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$$

ein lineares Gleichungssystem, das durch Diskretisierung einer Aufgabe (h = Diskretisierungsparameter, Maschenweite) entstanden ist, die ihre Struktur jeder Diskretisierung überträgt. Die Gleichungssysteme können aus elliptischen, hyperbolischen, parabolischen oder auch Integralgleichungen stammen. Auf alle diese Typen läßt sich das MGV anwenden.

Um die Dimension der zu lösenden Probleme zu veranschaulichen, hier ein bescheidenes Beispiel: Poissongleichung im Einheitsquadrat mit 40×40 inneren Punkten, also 1600 Unbekannten, $1600^2 = 2.56 \cdot 10^6$ Matrixelementen. Pro Zeile sind nur 5 Elemente $\neq 0$, also 8000 Matrixelemente $\neq 0$.

Man kann sich überlegen, daß zur Lösung (z.B. für Bandmatrizen) ein Rechenaufwand von Q arithmetischen Operationen mit $Q \sim n^2$ (\sim proportional) nötig ist, mit $n = (N - 1)^2$ (hier $n = 1600$.)

Probleme und Fakten:

1. der Aufwand für exaktes Lösen (sowohl an Operationen als auch an Speicherplatz) steigt ungeheuerlich.
2. Eine exakte Lösung ist gar nicht sinnvoll. Es genügt eine Approximationsgenauigkeit der Lösung in der Größenordnung des Diskretisierungsfehlers.
3. Lösung der Gleichungen durch reine Iteration erfordert ebenfalls viel Aufwand, aber
4. Die "Effektivität" der Iterationsverfahren (z.B. Jacobi und Gauß-Seidel) ist in den ersten Iterationsschritten viel größer, als die lineare Konvergenzgeschwindigkeit erwarten läßt (vgl. Übungen.)

Dies führte nacheinander zu folgenden Entwicklungen.

1: Grundidee: Nachiteration (schon länger bekannt.)

Löse das Gleichungssystem $\mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$ "billig" (als erste Möglichkeit etwa LR-Zerlegung und einfache Genauigkeit, später werden wir noch billigere Möglichkeiten kennen lernen.) Man erhält dann eine Näherung \mathbf{y}_h^0 .

Berechne den Defekt \mathbf{d}_h exakt (z.B. mit doppelter Genauigkeit)

$$(16.2) \quad \mathbf{d}_h = \mathbf{b}_h - \mathbf{A}_h \mathbf{y}_h^0.$$

Löse die Korrekturgleichung exakt (was immer das numerisch bedeutet)

$$(16.3) \quad \mathbf{A}_h \mathbf{v}_h = \mathbf{d}_h.$$

Gewinne eine neue Näherung durch

$$(16.4) \quad \mathbf{y}_h^1 = \mathbf{y}_h^0 + \mathbf{v}_h.$$

Dies bewirkt:

$$\mathbf{A}_h \mathbf{y}_h^1 = \mathbf{A}_h \mathbf{y}_h^0 + \mathbf{A}_h \mathbf{v}_h = \mathbf{A}_h \mathbf{y}_h^0 + \mathbf{d}_h \stackrel{(16.2)}{=} \mathbf{A}_h \mathbf{y}_h^0 + \mathbf{b}_h - \mathbf{A}_h \mathbf{y}_h^0 = \mathbf{b}_h,$$

d.h. die neue Näherung ist die exakte Lösung.

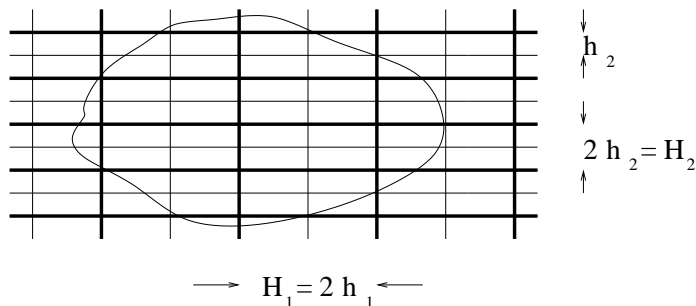
Ist die Lösung von (16.3) nicht exakt, so ist mindestens zu erwarten, daß die Näherung \mathbf{y}_h^1 besser ist, als \mathbf{y}_h^0 . Naheliegender ist daher der Gedanke (16.3) "billiger" zu berechnen und diese Korrektur iterativ zu wiederholen. Daraus ergibt sich die

2. Grundidee: Löse $\mathbf{A}_h \mathbf{v}_h = \mathbf{d}_h$ "billig" (und dafür öfter) durch Ausnutzen der Aufgabenstruktur,

d.h. wir wissen, daß $\mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$ aus einer Diskretisierung mit der Schrittweite h stammt.

Idee: Löse die Korrekturgleichung (16.3) nur auf einem gröberen Gitter mit der Schrittweite H . (Standart: $H = 2h$). Löse also

$$(16.5) \quad \mathbf{A}_H \mathbf{v}_H = \mathbf{d}_H$$



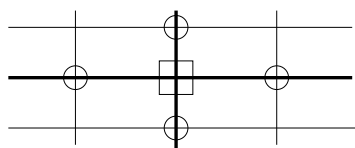
Ausgangsvoraussetzung ist also: Man hat eine Ausgangsnäherung \mathbf{y}_h^0 für (16.1). Also ist der Defekt $\mathbf{d}_h = \mathbf{b}_h - \mathbf{A}_h \mathbf{y}_h^0$ bekannt.

Problem 1: Wie macht man aus dem “großen” Vektor \mathbf{d}_h (mit den vielen Komponenten) den “kleinen” Vektor \mathbf{d}_H ?

1. Möglichkeit: Man könnte die Defektkomponenten den einzelnen Gitterpunkten zuordnen und für die gemeinsamen Punkte beider Gitter (also für die Punkte des groben Gitters) $\mathbf{d}_h = \mathbf{d}_H$ setzen.

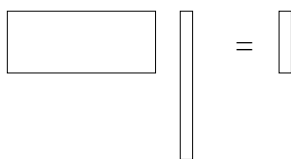
Diese Variante steigt gelegentlich aus (zur Erklärung vgl. Problem 3).

2. Möglichkeit: Man gewinnt \mathbf{d}_H durch konvexe Mittelbildung der “umgebenden \mathbf{d}_h ”.



Dies kann man in Matrixschreibweise angeben mit dem Restriktionsoperator \mathbf{R}_h^H

$$(16.6) \quad \mathbf{R}_h^H \mathbf{d}_h = \mathbf{d}_H.$$



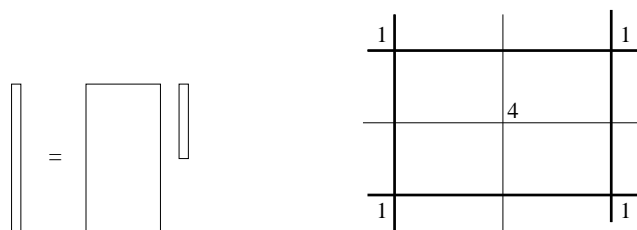
Die Matrix \mathbf{R}_h^H wird nicht gespeichert, sondern im Programm berechnet. Dann ist

$$(16.7) \quad \mathbf{A}_H \mathbf{v}_H = \mathbf{R}_h^H \mathbf{d}_h \text{ lösbar} \implies \mathbf{v}_h \text{ berechenbar}$$

Problem 2: Die Lösung \mathbf{v}_H ist zu “klein” (zu wenige Komponenten), also in $\mathbf{y}_h^1 = \mathbf{y}_h^0 + \mathbf{v}_h$ nicht einsetzbar. Man braucht eine

Prolongation \mathbf{P}_H^h :

$$(16.8) \quad \mathbf{v}_h = \mathbf{P}_H^h \mathbf{v}_H$$



Beim Übergang vom groben auf's feine Gitter wird linear interpoliert. Die 4 in der Mitte bedeutet nur daß der Funktionswert in der Mitte aus den 4 umgebenden, mit “1” bezeichneten Punkten gewonnen wird. (vgl. § 17)

Bemerkungen:

1. Die Indizierungen bei \mathbf{R} und \mathbf{P} zeigen die Richtung an, in welcher die Umformung vorgenommen wird (jeweils von unten nach oben), also vom feineren zum gröberem Gitter oder umgekehrt. Wenn man das weiß, kann man die Indizes auch weglassen.
2. Wenn das Verfahren so liefere, wie beschrieben beim Vorgehen zur Lösung von Problem 1 oder 2, wäre das vom Rechenaufwand her (für die Lösung der Korrekturgleichung) schon lohnend, denn liegt der Rechenaufwand (arithmetische Operationen) beim feinen Gitter bei $Q \approx n^2$ (größenordnungsmäßig bei n Unbekannten), so liegt er beim gröberem Gitter nur bei $Q_H \approx (\frac{1}{4}n)^2 = \frac{1}{16}Q_h$ (im 2D-Fall bei $H = 2h$, asymptotisch), d.h. man sparte mehr als eine Größenordnung.
3. Natürlich wäre mit diesem Verfahren ein Genauigkeitsverlust verbunden. Wie ernst das wäre in Anbetracht des Diskretisierungsfehlers, wäre zu überlegen.

Problem 3

Leider läuft das Verfahren so, wie beschrieben nicht, denn \mathbf{R} hat – notwendigerweise – einen nichttrivialen Nullraum $N(\mathbf{R})$.

$$\text{Gilt } \mathbf{d}_h \in N(\mathbf{R}) \xrightarrow{(16.7)} \mathbf{v}_h = \mathbf{0} \xrightarrow{(16.4)} \mathbf{y}_h^1 = \mathbf{y}_h^0.$$

Die Iteration steht.

Abhilfeüberlegungen

1. Man muß an Hand einer Fehleranalyse überlegen, welche Art von Fehlern mit einer ‘‘Grobgitterkorrektur’’ überhaupt korrigiert werden können und nicht ein $\mathbf{d} \in N(\mathbf{R})$ liefern.
2. Dann wäre es wünschenswert, eine Ausgangsnäherung \mathbf{y}_h^0 so zu wählen, oder so abzuändern, daß sie nur – oder zumindest überwiegend – Fehler enthält, die man mit der Grobgitterkorrektur abbauen kann. Dazu untersuchen wir das Fehlerverhalten bei der Lösung eines Gleichungssystems $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit Hilfe der Entwicklung des Fehlers nach Eigenvektoren. Dies hängt natürlich wesentlich davon ab, daß man zur Lösung ein Verfahren benutzt, das eine solche Analyse ermöglicht. Zur Erklärung des Sachverhalts benutzen wir hier das gedämpfte Jacobi-Verfahren. In der Praxis nimmt man lieber Gauß-Seidel, aber das Prinzip läßt sich hier einfacher erklären.

Fehleranalyse:

Zur Lösung von $\mathbf{A}\mathbf{x} = \mathbf{b}$ setzen wir voraus $\mathbf{A} = \mathbf{A}^T > \mathbf{0}$. Dann lautet das gedämpfte Jacobi-Verfahren (vgl. (15.1))

$$\mathbf{x}^{m+1} = (\mathbf{I} - \omega \mathbf{D}^{-1} \mathbf{A}) \mathbf{x}^m + \omega \mathbf{D}^{-1} \mathbf{b}, \quad \mathbf{D} = \text{diag}(\mathbf{A}),$$

2. daß sie langsamer $\rightarrow 0$ gehen, als die Anteile mit mittleren k -Werten, für die $|q_k(\omega)|$ kleiner ausfällt.

Problem 4

Das gröbere Gitter produziert Matrizen \mathbf{A}_H mit sehr viel weniger Eigenwerten (und damit weniger Eigenvektoren) als das feinere Gitter. Nur die kleinen Eigenwerte des groben Gitters (und damit auch die "kleinen" Eigenvektoren) kann man als Näherungen für die Eigenwerte und Eigenvektoren des feineren Gitters betrachten.

Die großen Eigenwerte und Eigenvektoren des feineren Gitters kann man aus dem groben Gitter gar nicht approximieren

Folgerungen

1. Die Fehleranteile der "großen" Eigenvektoren (das sind die Eigenvektoren, die zu großen λ_k gehören,) kann man durch Korrekturen, die aus dem groben Gitter kommen (d.h. iterative Verbesserung auf dem groben Gitter), gar nicht korrigieren.
2. Man müßte \mathbf{y}_k^0 dahingehend korrigieren, daß die Fehleranteile in \mathbf{y}_k^0 , die von "großen" Eigenvektoren herrühren, möglichst klein sind und dann versuchen, durch Korrekturen, die aus dem groben Gitter errechnet werden, die Fehleranteile der kleinen Eigenvektoren zu reduzieren.
3. Man wird also nach einem (Iterations-) Verfahren zur Verbesserung und Konstruktion von \mathbf{y}_k^0 suchen, das die Fehleranteile der "großen" Eigenvektoren klein macht.

Glättung

Im 1D-Fall werden die Eigenwerte und Eigenvektoren von $\mathbf{D}^{-1}\mathbf{A}_h^0$, $\mathbf{D} = \text{diag}(\mathbf{A}_h^0)$ gegeben durch (vgl. Lemma 3.7 und beachte: $\mathbf{D} = \text{diag}(\frac{2}{h^2})$)

$$\lambda_k^h = 2 \sin^2\left(\frac{k\pi h}{2}\right), \quad \mathbf{y}_i^k = \sin(k\pi x_i), \quad i = 0, \dots, N, \quad \mathbf{y}_0^k = \mathbf{y}_N^k = 0.$$

Für große k sind die \mathbf{y}^k hochfrequente Schwingungen. Werden die Koeffizienten dieser Eigenvektoren in der Fehlerdarstellung (16.9) betragsmäßig verkleinert, so entspricht dies optisch einer Glättung des Fehlers, daher der Name für diesen Prozess.

Im 2D-Fall, $\Omega = (1, 0)^2$, $h_1 = h_2 = h$ hat man für $(\text{diag}\mathbf{A}_h^0)^{-1}\mathbf{A}_h^0$ die Eigenwerte

$$\lambda_k^h = \lambda_{k_1 k_2}^h = \lambda_{k_1}^{(1)} + \lambda_{k_2}^{(2)} = \sin^2\left(\frac{k_1\pi h}{2}\right) + \sin^2\left(\frac{k_2\pi h}{2}\right)$$

mit den Eigenvektoren

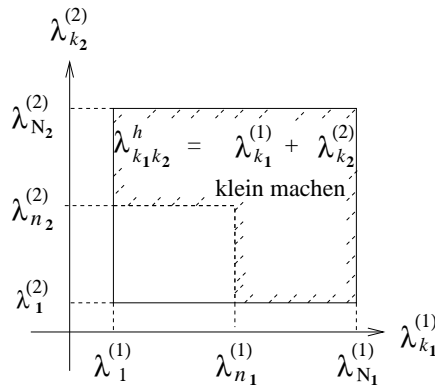
$$\mathbf{v}_{k_1 k_2} = 2 \sin(k_1\pi x_1) \sin(k_2\pi x_2), \quad (\text{komponentenweise definiert für } (x_1, x_2) \in \omega_h^0).$$

Will man wieder das gedämpfte Jacobi-Verfahren benutzen, für den Glättungsprozess, (in der Praxis wird oft Gauß-Seidel vorgezogen, darauf kommen wir später zurück), so entspricht die vorige Wahl von ω_{opt} dieser Absicht nicht.

Ziel: $|1 - \omega\lambda_k|$ soll klein sein für die großen Eigenwerte, d.h. im

$$\text{1D-Fall: } \frac{N}{2} \leq k \leq N,$$

$$\text{2D-Fall: } \frac{N}{4} \leq k \leq N, \text{ (genauer } k_1 \geq \frac{N_1}{2} =: n_1, k_2 \geq \frac{N_2}{2} =: n_2 \text{.)}$$



Man kann zeigen (später): Es ist möglich, die Werte $|1 - \omega\lambda_k|$, die zu $\frac{N}{4} \leq k \leq N$ gehören (2D-Fall), unter eine von der Schrittweite unabhängige Konstante zu drücken, entsprechend im 1D-Fall.

Wir demonstrieren die Glättung an einem einfachen (1D-) Beispiel mit der Diskretisierungsmatrix, die schon im Zusammenhang mit den parabolischen Differentialgleichungen untersucht wurde. (Wegen 1D wird $\frac{N}{4}$ durch $\frac{N}{2}$ ersetzt, vgl. dazu die Eigenvektoren des 2D-Falles.)

Einfaches (1D-) Demonstrationsbeispiel zur Glättung
(unter Verwendung des gedämpften Jacobi-Verfahrens)

$$y_{\bar{x}x} + f = 0, \quad y(0) = y(1) = 1, \quad h = \frac{1}{N}.$$

Die Eigenwerte von $\mathbf{D}^{-1}\mathbf{A}$, (\mathbf{A} = Diskretisierungsmatrix) und die zugehörigen Eigenvektoren sind (beachte: $\mathbf{D} = \text{diag}(\frac{2}{h^2})$)

$$\lambda_k = 2 \sin^2\left(\frac{k\pi h}{2}\right), \quad \mathbf{v}^{(k)} = \begin{pmatrix} \sin(k\pi x_1) \\ \vdots \\ \sin(k\pi x_{N-1}) \end{pmatrix}, \quad (\text{vgl. Lemma 3.7})$$

wobei $N = 2n$, n = Matrixdimension fürs grobe Gitter,
 N = Matrixdimension fürs feine Gitter,
 $h = \frac{1}{2n}$ Gitterweite des feinen Gitters.

Ziel: Für $n \leq k \leq N-1$ sollen die Eigenwerte der Iterationsmatrix, also die Werte von $|1 - \omega\lambda_k|$, durch geeignete Wahl von $\omega > 0$ möglichst klein ausfallen (vgl. § 15). Da alle $\lambda_k > 0$ sind, wird dies erreicht durch die Forderung

$$\begin{array}{c} \lambda_1 \quad \lambda_n \quad \lambda_{N-1} \\ \oplus \quad | \quad | \quad | \\ \hline \end{array}$$

$$|1 - \omega\lambda_n| = |1 - \omega\lambda_{N-1}| \quad \begin{array}{c} \hline | \quad | \quad | \\ \oplus \quad | \quad | \quad | \\ \hline 1-\omega\lambda_{N-1} \quad 1-\omega\lambda_n \quad 1-\omega\lambda_1 \end{array}$$

Hieraus folgt

$$\omega = \frac{2}{\lambda_n + \lambda_{N-1}} \quad \text{und} \quad q = 1 - \omega\lambda_n = \frac{\lambda_{N-1} - \lambda_n}{\lambda_{N-1} + \lambda_n}.$$

Nun gilt für die Eigenwerte des feinen Gitters

$$\begin{aligned} \lambda_n &= 2 \sin^2\left(\frac{n\pi h}{2}\right) \stackrel{nh=1/2}{=} 2 \sin^2\left(\frac{\pi}{4}\right) = 1, \quad \left(\sin\left(\frac{\pi}{4}\right) = \frac{\sqrt{2}}{2}\right) \\ \lambda_{N-1} &= 2 \sin^2\left(\frac{(N-1)\pi h}{2}\right) = 2 \sin^2\left(\frac{\pi h}{2} - \frac{Nh\pi}{2}\right) \stackrel{Nh=1}{=} 2 \underbrace{\cos^2\left(\frac{\pi h}{2}\right)}_{\approx 1 \text{ für kleine } h} \leq 2, \\ \omega &\approx \frac{2}{3}, \end{aligned}$$

also folgt für das Betragsmaximum der Eigenwerte von $(I - \omega \mathbf{D}^{-1} \mathbf{A})$ für die Eigenwerte λ_k , $\frac{N}{2} \leq k \leq N$ von $\mathbf{D}^{-1} \mathbf{A}$

$$q = 1 - \omega\lambda_n = \frac{\lambda_{N-1} - \lambda_n}{\lambda_{N-1} + \lambda_n} = \frac{2 \cos^2\left(\frac{\pi h}{2}\right) - 1}{2 \cos^2\left(\frac{\pi h}{2}\right) + 1} \leq \frac{1}{3}.$$

Diese, bzw. eine entsprechende Wahl im höherdimensionalen Fall, wird dem Mehrgitterverfahren zur Konvergenz verhelfen.

Folge: Man kann also eine Ausgangsnäherung \mathbf{y}_h^0 des Gleichungssystems $\mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$ durch einige Schritte des beschriebenen Glättungsverfahrens so abändern, daß man Näherungen erhält, deren Fehler weniger hochfrequente Anteile erhalten.

Bemerkung: Eine Glättungiteration muß im allgemeinen Fall nicht mehr konvergieren. Es sollen ja nur die Fehleranteile der "großen" Vektoren vermindert werden. In diesem Fall konvergiert sie "gerade noch".

$$\rho(\mathbf{I} - \omega \mathbf{D}^{-1} \mathbf{A}) = 1 - \omega\lambda_1 = 1 - \frac{2}{3} 2 \sin^2\left(\frac{\pi h}{2}\right) \approx 1 - \frac{4}{3} \frac{\pi^2 h^2}{4} \approx 1 - \pi h^2.$$

Beachte jedoch: Die Konvergenz wird schlechter je kleiner h wird.

Zwei-Gitter-Verfahren

Mit der Bezeichnung

$$\mathcal{S} = \mathcal{S}(\mathbf{A}_h, \mathbf{b}_h) \quad \text{für den Glättungsoperator} \quad (\mathcal{S} \hat{=} \text{smoothing})$$

kann man also ein 2-Gitter-Verfahren wie folgt beschreiben:

Verlauf eines 2-Gitter-Verfahrens zur Lösung von $\mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$

1. Billige Beschaffung einer Ausgangsnäherung \mathbf{y}_h^0
2. Glättung der Ausgangsnäherung (durch i Schritte eines einfachen Verfahrens: Gauß-Seidel ist beliebter als Jacobi)

$$\mathbf{y}^{(i)} = \mathcal{S}^i(\mathbf{A}_h, \mathbf{b}_h) \mathbf{y}_h^0 \quad \text{es verbleiben hauptsächlich niederfrequente Fehleranteile}$$
3. $\mathbf{d}_h = \mathbf{b}_h - \mathbf{A}_h \mathbf{y}_h^{(i)}$ Defektberechnung
4. $\mathbf{R}_h^H \mathbf{d}_h = \mathbf{d}_H$ Restriktion auf "kleineren" Vektor
5. $\mathbf{A}_H \mathbf{v}_H = \mathbf{d}_H$ Korrekturgleichung lösen
6. $\mathbf{y}_h^{(1)} = \mathbf{y}_h^{(i)} + \mathbf{P}_H^h \mathbf{v}_H$ Prolongation und Verbesserung
7. \longrightarrow 2)

Dieses Verfahren konvergiert schon.

Beschaffung einer Ausgangsnäherung: Löse $\mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$ auf einem ganz groben Gitter, "fahre es hoch" durch Prolongation bis aufs Gitter h , und benutze diese Approximation als Ausgangsnäherung.

Bemerkungen:

1. Dieser Zyklus wird drei bis vier mal durchlaufen. (Es gibt Varianten mit nur 2 Zyklen.)
Beachte: Es muß ja nur eine Lösungsgenauigkeit erreicht werden, die mit der Genauigkeit des Diskretisierungsoperators vergleichbar ist.
2. Der Schritt 2) heißt *Vorglättung*. Man könnte als Schritt 6a) nochmals eine Glättung einfügen (*Nachglättung*). Es gibt Varianten ohne Vorglättung aber mit Nachglättung.
3. Die Glättungsstrategie bewirkt, daß bei der Restriktion der Defekt $\mathbf{d}_H \neq 0$ ausfällt, in 5) also tatsächlich eine Verbesserung erreicht wird (genauerer muß der Konvergenzbeweis zeigen).
4. Das Verfahren klappt auch bei Finite-Element-Rechnungen, ist aber schneller beim Differenzenverfahren.

5. Selbst für nichtlineare Probleme wird es mit Erfolg angewendet (vgl. § 20). Man benötigt dafür sehr gute Lösungsverfahren auf dem untersten Gitter.
6. Überschlägt man den Aufwand an Rechenoperationen für einen Zyklus, so erkennt man leicht, daß er in allen Schritten, außer 5), linear mit der Anzahl der Unbekannten wächst.

Zur weiteren Beschreibung ändern wir die Indizierung. Da wir mehr als nur zwei Gitter einsetzen wollen, reichen zu ihrer Bezeichnung die Indizes h und H nicht mehr aus. Wir bezeichnen das Gitter, auf dem wir die Gleichung $\mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$ lösen wollen, mit der Gitternummer l , das gröbere Gitter (üblich ist Schrittweitenverdopplung bei der Gitterverfeinerung) mit $l-1$; $l=0$ bezeichnet das größte Gitter. Entsprechend werden die Gleichungsgrößen indiziert, also z.B. $\mathbf{A}_l \mathbf{y}_l = \mathbf{b}_l$.

Algorithmisch können wir das Zweigitter-Verfahren (ZGV, englisch TGM $\hat{=}$ two-grid-method) in einer programmierähnlichen Schreibweise wie folgt beschreiben: (vgl. Hackbusch: Multigrid Methods and Applications)

Mit den Abkürzungen

$l \hat{=}$ Gitternummer

$\mathbf{u} \hat{=}$ als Eingabevektor die Anfangsnäherung \mathbf{y}_h^0 ,

als Ausgabevektor für die durch Glättung verbesserte Approximation

$\mathbf{b} \hat{=}$ rechte Seite von $\mathbf{A}\mathbf{x} = \mathbf{b}$

$\mathbf{d} \hat{=}$ Defekt

$\mathbf{v} \hat{=}$ Lösung der Korrekturgleichung

lautet die Programmieranweisung

```
procedure ZGV( $l, \mathbf{u}, \mathbf{b}$ ); integer  $l$ ; array  $\mathbf{u}, \mathbf{b}$ ;
```

```
if  $l = 0$  then  $\mathbf{u} := \mathbf{A}_0^{-1} * \mathbf{b}_0$  Lösung auf dem größten Gitter
```

```
else
```

```
begin array  $\mathbf{v}, \mathbf{d}$ ; Korrekturvektor und Defekt
```

```
 $\mathbf{u} := \mathcal{S}_l^\nu(\mathbf{u}, \mathbf{b});$   $\nu$  mal glätten auf Gitter  $l$  (Vorglättung)
```

```
 $\mathbf{d} := \mathbf{R} * (\mathbf{b} - \mathbf{A}_l * \mathbf{u});$  Restriktion und Defektberechnung
```

```
 $\mathbf{v} := \mathbf{A}_{l-1}^{-1} \mathbf{d};$  Lösung der Korrekturgleichung auf dem gröberen Gitter
```

```
 $\mathbf{u} := \mathbf{u} + \mathbf{P} * \mathbf{v};$  Prolongation und Korrektur
```

```
 $\mathbf{u} := \mathcal{S}^{\nu_2}(\mathbf{u}, \mathbf{b});$  ggf. Nachglättung ( $\nu_2$  Iterationen)
```

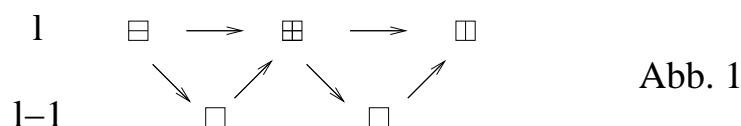
```
end ZGV
```

Der Aufruf $\text{ZGV}(l, \mathbf{u}, \mathbf{b})$; mit $l > 0$ beschreibt einen Schritt des ZGV.

Graphisch kann man das mit Symbolen wie folgt veranschaulichen:

- \boxplus Glatte
- \searrow Defektberechnung und Restriktion des Defekt auf das groebere Gitter
- \square Exakte Loesung der Gleichung
- \nearrow Prolongation auf das feinere Gitter
- \boxminus Korrektur auf dem feineren Gitter
- \longrightarrow zeigt an, welcher Wert korrigiert wird

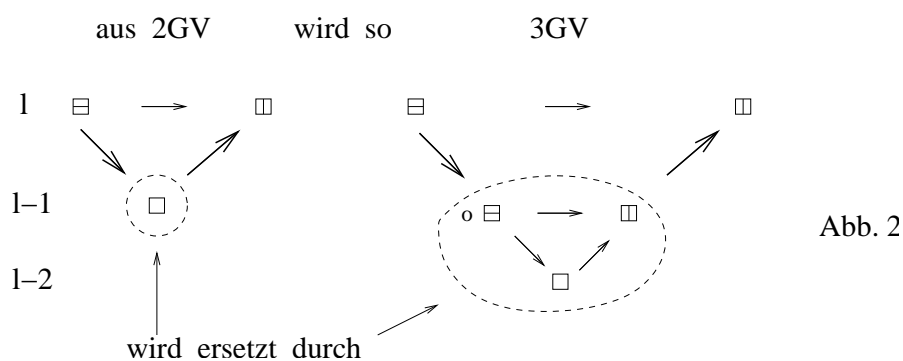
So lassen sich z.B. zwei Zyklen des ZGV veranschaulichen als



Die waagrechten Pfeile werden oft weggelassen. Korrigiert wird dann immer die Lösung, die links daneben auf derselben Höhe steht.

Mehrgitter-Verfahren

In Anbetracht dessen, daß exakte Lösung der Defektgleichung auf dem gröberen Gitter + Prolongation auf das feinere Gitter + Korrektur auf dem feineren Gitter nur eine neue Näherung \mathbf{y}_l^0 liefert, kann man sich überlegen, ob die exakte Lösung auf dem gröberen Gitter nicht durch eine billigere, approximative Lösung ersetzt werden kann. Beachtet man, daß \mathbf{A}_{l-1} ja dieselbe Struktur wie \mathbf{A}_l besitzt, so liegt es nahe die Korrekturgleichung auf dem Gitter $l-1$ zu lösen durch Anwendung eines neuerlichen ZGV unter Benutzung eines nochmals vergrößerten Gitters $l-2$. Als Ausgangsnäherung wird $\mathbf{v} = \mathbf{0}$ benutzt (graphisch durch o angedeutet).



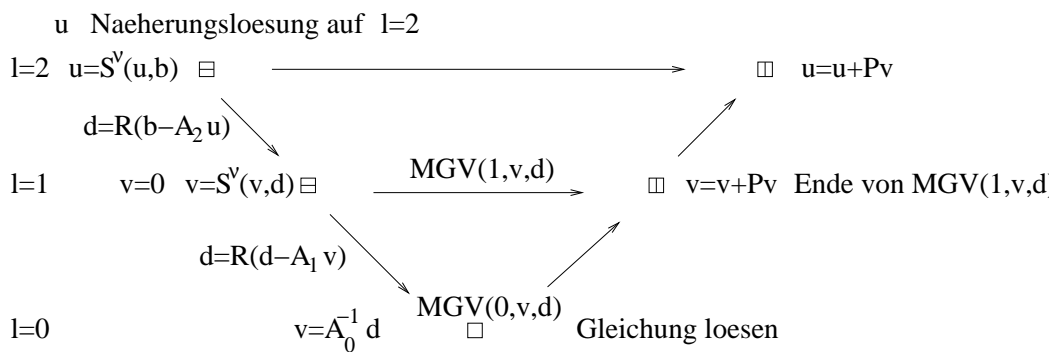
Algorithmisch bedeutet das, daß innerhalb des ZGV nochmals ein ZGV aufgerufen wird (rekursiver Aufruf des ZGV) zur approximativen Lösung der Korrekturgleichung.

Algorithmisch kann man das wie folgt als Mehrgitterverfahren fassen durch mehrmaligen rekursiven Aufruf des ZGV.

```

procedure MGV(l, u, b); integer l; array u, b;
if l = 0 then u := A0-1 * b0
else
begin array v, d;
    u := Slv(u, b);
    d := R * (b - Al * u);      ↓ Änderung gegenüber ZGV
    v := 0; MGV(l - 1, v, d);   Approximative Lösung der Grobgitterkorrektur
    u = u + P * v;
    u := Sv2(u, b);    als Variante gelegentlich noch eine Nachglättung
end;
    
```

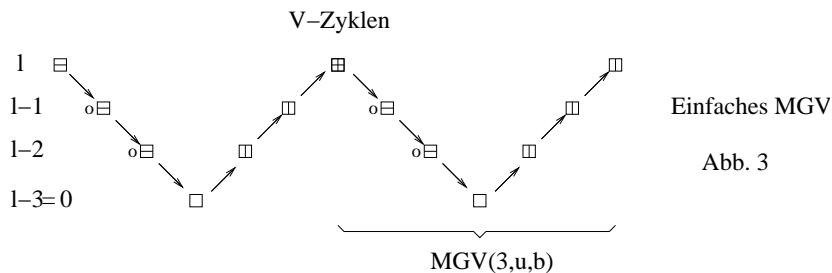
Wir verfolgen den Programmverlauf an einem Beispiel für *l*=2.



3GV sind unüblich, weil auch die Lösung auf Gitter *l* - 3 zu aufwendig ist. Man löst also die Korrekturgleichung auf *l* - 3 approximativ durch ein weiteres ZGV, also MGV(3, **u**, **b**) vgl. nächstes Beispiel.

Praktisch verwendet werden *l* ≥ 4 Gitter.

Graphisch sehen zwei Zyklen eines solchen 4GV wie folgt aus



In diesem Beispiel wird zur approximativen Lösung der Korrekturgleichung auf Gitter $l - 1$ ein 3GV benutzt.

Der erste und zweite V-Zyklus aus Abb 3 wird beschrieben durch ein MGV $(3, \mathbf{u}, \mathbf{b})$ (beachte: das ist ein 4GV), d.h. die 1. Grobgitterkorrekturgleichung wird durch ein 3GV gelöst. Erst auf dem größten Gitter ($l = 0$) wird die dortige Gleichung exakt gelöst. Für gutartige Probleme läuft dieses Verfahren schon ganz gut.

Verfahrensvariante: In den aufsteigenden Linien wird nicht nur eine Korrektur durchgeführt, sondern auch geglättet, d.h. ersetze $\boxed{|}$ durch $\boxed{+}$

Ausbau:

Bei empfindlichen Problemen, (z.B. schlechte Kondition, zu viele Gitterverfeinerungen) kann es vorkommen, daß die Verbesserung, die man durch die approximative Lösung der Korrekturgleichung erhält, auf Grund der vielen Restriktionen und Prolongationen, nicht gut genug ist. Dem kann man begegnen, indem man den im MGV eingerahmten Prozess zur approximativen Lösung der Grobgittergleichung mehrfach ausführt, (d.h. die Lösung der Grobgitterkorrektur wird genauer), z.B. γ mal ($\gamma = 2$ wird oft benutzt), bevor man nach Prolongation die Korrektur auf dem obersten Gitter ausführt. Algorithmisch geschieht das, indem man für den eingerahmten Teil des MGV eine Wiederholungsschleife einführt. Dies bewirkt u.a., daß auf dem untersten Gitter auch das Gleichungslösen mehrfach (γ -fach im folgenden Programm) ausgeführt wird.

Wir bezeichnen diese Prozedur mit $\text{MGV } \gamma$, um anzudeuten, daß γ mal eine Nachiteration auf dem größeren Gitter stattfindet. Mit dieser Bezeichnung ist $\text{MGV} = \text{MGV}1$.

```

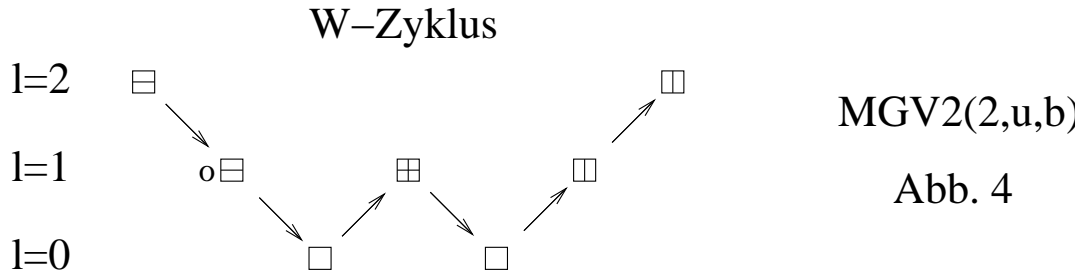
procedure MGV  $\gamma(l, \mathbf{u}, \mathbf{b})$ ; integer  $l$ ; array  $\mathbf{u}, \mathbf{b}$ ;
if  $l = 0$  then  $\mathbf{u} := \mathbf{A}_0^{-1} * \mathbf{b}_0$ 
else
begin  $\boxed{\text{integer } j;}$  array  $\mathbf{v}, \mathbf{d}$ ;
 $\mathbf{u} := \mathcal{S}_l^\nu(\mathbf{u}, \mathbf{b})$ ;
 $\mathbf{d} := \mathbf{R} * (\mathbf{b} - \mathbf{A}_l * \mathbf{u})$ ;
 $\mathbf{v} := 0$ 
 $\boxed{\text{for } j = 1 \text{ step } 1 \text{ until } \gamma \text{ do}}$  MGV( $l - 1, \mathbf{v}, \mathbf{d}$ );
 $\mathbf{u} := \mathbf{u} + \mathbf{P} * \mathbf{v}$ ;
end;
```

Die eingerahmten Teile sind neu im Vergleich zum einfachen Verfahren MGV, das für $\gamma = 1$ in MGV1 enthalten ist.

Wir erläutern dieses Programm durch eine Reihe von Beispielen graphisch, jeweils für

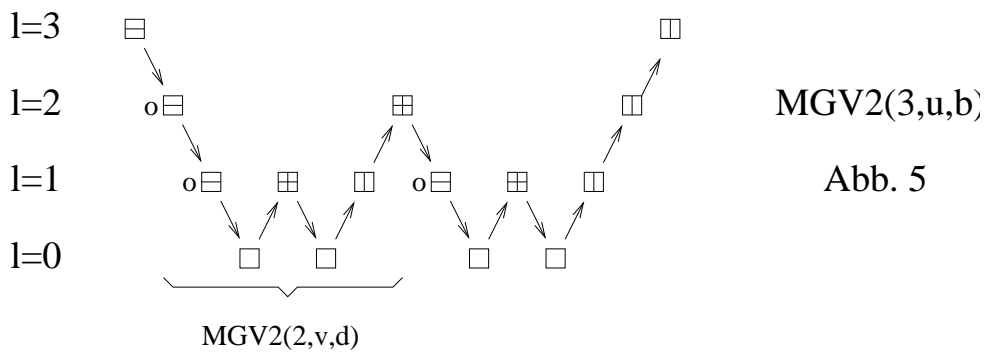
$\gamma=2$ und verschiedene l .

Beispiel: Für $l=2$ wird ein 2-Gitter-Verfahren γ mal zur Konstruktion der Grobgitterkorrektur benutzt.

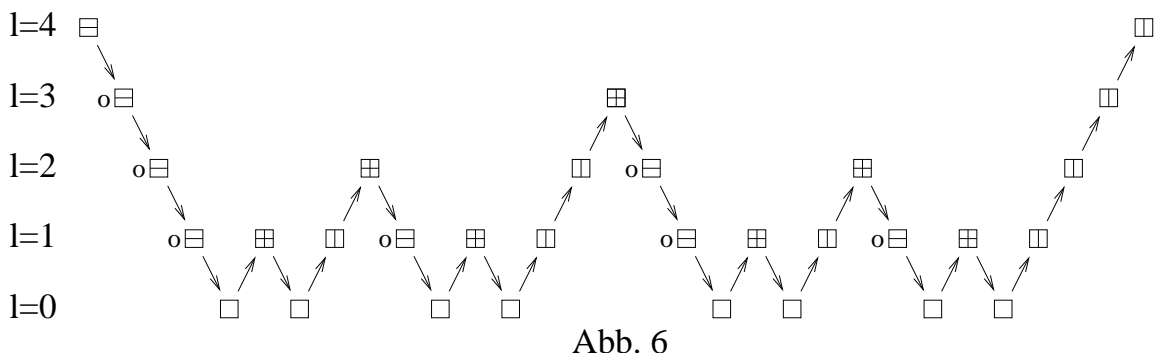


Es wurden $\gamma = 2$ Schritte zur Konstruktion der Approximation für die Grobgitterlösung benutzt.

Beispiel: $l=3$: Die Lösung der Korrekturgleichung auf $l = 2$ wird durch $\gamma = 2$ Schritte des obigen 3-Gitter-Verfahrens approximiert.



Beispiel: $l=4$: Die Korrekturgleichung auf $l = 3$ wird approximativ gelöst durch $\gamma = 2$ Schritte des vorigen 4-Gitter-Verfahrens.



Bekannt seien $(\mathbf{A}_l)_{l=0}^{l_{max}}$, $(\mathbf{b}_l)_{l=0}^{l_{max}}$, gesucht sei die Lösung von $\mathbf{A}_{l_{max}} \mathbf{u} = \mathbf{b}_{l_{max}}$.

$$\mathbf{u}_0 := \mathbf{A}_0^{-1} \mathbf{b}_0$$

for $l = 1$ step l until l_{max} do

begin $\mathbf{u}_l := \tilde{\mathbf{P}} \mathbf{u}_{l-1}$;

$\mu :=$ if $l = l_{max}$ then μ_2 else μ_1 ;

for $j = 1$ step 1 until μ do MGVB $\gamma(l, \tilde{\mathbf{u}}_l, \mathbf{b}_l)$;

end;

Man handelt sich also mit jeweils μ_1 Schritten MGVB hoch bis l_{max} und rechnet dann μ_2 Zyklen von MGVB. Ist l die aktuelle Stufe, so wird γ mal die Korrekturgleichung (also auf Stufe $l - 1$) durch ein MGVB1 gelöst.

Bemerkungen zu den Iterationszahlen

Üblich sind $\nu = \nu_1 + \nu_2 \leq 3$, wenn man die Anzahl ν der Glättungen in Vor- und Nachglättung aufteilt. Vernünftig sind 1-2 Vorglättungen und eine Nachglättung zum Glätten höherer Frequenzen, die von der Prolongation kommen könnten.

$\mu_1 \leq 2 \leq \mu_2 \leq 4$, die beste Wahl ist problemabhängig, $\gamma = 2$ ist normal.

Die Nachglättung wird in MGVB1 eingebaut als letzter Schritt.

Geschichtlicher Überblick

(genauer bei Hackbusch 2.6.5)

- 1961 Federenko: 1) Glättung, 2) Grobgitterkorrektur
- 1964 Federenko: Konvergenzbeweis für Laplace im Einheitsquadrat mit Hilfe der Fourier-Analyse (Entwicklung nach Eigenvektoren)
- 1966 Bachvalov: kompletter Konvergenzbeweis für das volle MGVB für allgemeine elliptische Operatoren inklusive Nachweis der Optimalität bzgl. des Rechenaufwandes
- 1974-1975 Achi Brand: Erfuhr auf Umweg über Amerika von Schülern der Moskauer Schule die Idee des Verfahrens, hat es weiterentwickelt und auf nichttriviale Probleme der amerikanischen Navy angewandt, ohne auf Konvergenzbetrachtungen einzugehen.
- 1976 Hackbusch: Unabhängige Neuerfindung des MGVB. Von ihm stammt das erste Buch über MGVB.
- 1989 erstes russisches Buch über MGVB.

Zum Paragraphenabschluß zitieren wir aus dem Buch von A. Iserless: A first course in the numerical Analysis for Differential Equations (Cambridge University Press 1996): The number of different strategies and implementations of multigrid is a source of major preoccupation to professionals, although it might be at times baffling to other numerical analysts and to users of computational algorithms.

§ 17 Glättung, Restriktion, Prolongation

Wir befassen uns erst mit dem Problem, den Glättungseffekt zu "messen". In der Literatur werden dazu zwei Möglichkeiten betrachtet.

1) Die Glättungsrate (eingeführt von A.Brand 1977)

Unsere Glättungsüberlegungen beruhen auf der in § 16 geschilderten Grundidee.

Sei n_l die Anzahl der Gitterpunkte (=Anzahl der Eigenwerte) des l -ten Gitters und n_{l-1} die des nächst größeren Gitters.

Bei der Verdoppelung der Maschenweite gilt (vgl. S. 125)

im 1D-Fall: $n_{l-1} \approx \frac{n_l}{2}$,

im 2D-Fall: $n_{l-1} \approx \frac{n_l}{4}$.

Erinnerung: Die Eigenvektorentwicklung der Fehlerdarstellung (vgl. (15.2),(16.9)) war eine Entwicklung nach den Eigenvektoren der Iterationsmatrix \mathbf{S} . Die Iterationsmatrizen haben die Gestalt $\mathbf{S} = \mathbf{I} - \omega \mathbf{M} \mathbf{A}$, wo \mathbf{A} die Diskretisierungsmatrix ist und z.B. $M = \mathbf{D}^{-1}$. Sind λ_i die Eigenwerte von $\mathbf{M} \mathbf{A}$ zu den Eigenvektoren $\mathbf{v}^{(i)}$, so hat die Iterationsmatrix die Eigenwerte $(1 - \omega \lambda_i)$ zu den Eigenvektoren $\mathbf{v}^{(i)}$.

Die Eigenvektoren $\mathbf{v}_l^{(\mu)}$ der Iterationsmatrixmatrix wurden eingeteilt in

niedere Frequenzen: $1 \leq \mu \leq n_{l-1}$ und

hohe Frequenzen $n_{l-1} + 1 \leq \mu \leq n_l$.

Sei σ_μ der μ -te Eigenwert der Iterationsmatrix \mathbf{S}_l des Dämpfungsverfahrens, also $\mathbf{S}_l \mathbf{v}_l^{(\mu)} = \sigma_\mu \mathbf{v}_l^{(\mu)}$, so definiert man für ν -fache Glättung (wenn $l=1$ das größte Gitter bezeichnet) die

$$(17.1) \quad \text{Glättungsrate} \quad \rho_B = \sup_{l \geq 1} \max\{|\sigma_\mu|^\nu; n_{l-1} + 1 \leq \mu \leq n_l\}$$

Beachte: Es ist (sollte sein): $|\sigma_\mu| < 1$

Im praktischen Fall ist ihre Bestimmung naturgemäß mit einigem Aufwand verbunden. Einen etwas anderen Zugang liefert die

2) Die Glättungsnummer (Hackbusch)

Die Glättung einer Näherungslösung \mathbf{u}_l^j für $\mathbf{A}_l \mathbf{u}_l = \mathbf{b}_l$ wird durch eine Glättungsiteration beschrieben (z.B. gedämpfter Jacobi oder Gauß-Seidel)

$$\bar{\mathbf{u}}_l = \mathbf{S} \mathbf{u}_l^j + \mathbf{T} \mathbf{b}_l, \quad \mathbf{S} = \text{Iterationsmatrix}$$

$\mathbf{T} \mathbf{b}_l$ bezeichnet in der Iteration den Restsummanden, der nicht von der Iterationsmatrix beeinflusst wird.

(Beispiel: gedämpfter Jacobi, falls Diagonale zu 1 normiert: $\mathbf{S} = (\mathbf{I} - \omega \mathbf{A})$.)

Die exakte Lösung ist immer ein Fixpunkt der Iteration

$$\mathbf{u}_l = \mathbf{S} \mathbf{u}_l^j + \mathbf{T} \mathbf{b}_l.$$

Nur Verfahren mit dieser Eigenschaft werden zur Glättung benutzt.

Geglättet wird nicht die Näherungslösung \mathbf{u}^j , sondern ihr Fehler \mathbf{e}^j (vgl. (16.9))

$$\bar{\mathbf{e}} := \bar{\mathbf{u}}_l - \mathbf{u}_l = \mathbf{S}(\underbrace{\mathbf{u}_l^j - \mathbf{u}_l}_{\mathbf{e}^j}) =: \mathbf{S}\mathbf{e}^j$$

bzw. bei mehrfacher Glättung

$$\bar{\mathbf{e}}^j = \mathbf{S}^\nu \mathbf{e}^j.$$

Ein Fehlervektor $\mathbf{v} = (v_1, \dots, v_{n_l})^T$ kann im 1D-Fall anschaulich als glatt bezeichnet werden (“glatt” im Sinne “geringes Oszillieren”), wenn eine Norm der Differenz “benachbarter Steigungen”: $\frac{v_{j+1} - v_j}{h} - \frac{v_j - v_{j-1}}{h}$ klein ausfällt.

Ein spezieller Differenzenoperator 2. Ordnung, der diese Differenzen beschreibt, wird durch die Diskretisierungsmatrix von Δu gegeben. Wir bezeichnen sie hier mit \mathbf{L}_l .

$$(\mathbf{L}_l \mathbf{v})_i = \frac{1}{h} \left(\frac{v_{j+1} - v_j}{h} - \frac{v_j - v_{j-1}}{h} \right) \quad \text{im 1D-Fall.}$$

Zum 2D-Fall vgl. (12.2).

Die Norm der Steigungsdifferenzen des geglätteten Fehlervektors

$$\|\mathbf{L}_l \bar{\mathbf{e}}^j\| = \|\mathbf{L}_l \mathbf{S}^\nu \mathbf{e}^j\| \leq \|\mathbf{L}_l \mathbf{S}^\nu\| \|\mathbf{e}^j\|$$

sollte klein ausfallen im Vergleich zur Norm der Steigungsdifferenzen des ungeglätteten Fehlers

$$\|\mathbf{L}_l \mathbf{e}^j\| \leq \|\mathbf{L}_l\| \|\mathbf{e}^j\|$$

Wir motivieren so (wieder sei $l=1$ das grösste Gitter) die

$$(17.2) \quad \text{Glättungsnummer} \quad \rho_L(\nu) = \sup_{l \geq 1} \frac{\|\mathbf{L}_l \mathbf{S}^\nu\|}{\|\mathbf{L}_l\|},$$

die möglichst klein ausfallen sollte. (Beispiele werde wir in den Übungen kennen lernen.) Die Glättungsnummer ist normabhängig. Das wird im Konvergenzbeweis eine Rolle spielen.

Glättungsiterationen

Als Glättungsverfahren werden benutzt

1. Gauß-Seidel (Standard),
2. gedämpftes Jacobi-Verfahren: Iterationsmatrix: $\mathbf{I} - \omega \mathbf{D}^{-1} \mathbf{A}$,
3. Tschebyscheff-Verfahren, Iterationsmatrix: $\mathbf{I} - \omega_i \mathbf{A}$ (in den einzelnen Iterationsschritten werden verschiedene ω_i benutzt),

4. SOR klappt nicht, da nur auf Konvergenz ausgerichtet, Dämpfung nicht möglich
5. Konjugierte Gradienten-Verfahren: Da der Aufwand relativ hoch ist, ist es nur sinnvoll als äußere Iteration mit MGW als Präkonditionierung. Ansonsten ist es gut (insbesondere bei dünn besetzten Matrizen, wenn man keine Glättungsiteration finden kann (z.B. bei Navier-Stokes-Gleichungen)).

Wir kennen bereits das **gedämpfte Jacobi-Verfahren**.

Beachte:

1. Dieses Verfahren ist unabhängig von der Nummerierung der Unbekannten.
2. Auf Grund seiner Struktur läßt es sich leicht parallelisieren (z.B. pro Zeile ein Prozessor).

Dieses Verfahren läßt sich schreiben als

$$(17.3) \quad \mathbf{u}_l^{j+1} = \mathbf{u}_l^j - \omega \mathbf{D}^{-1}(\mathbf{A}_l \mathbf{u}_l^j - \mathbf{b}_l)$$

Die Bestimmung eines optimalen Dämpfungsparameters ist im allgemeinen Fall schwierig. In vielen Anwendungen sind die Diagonaleinträge von \mathbf{A} von der Art $\frac{k}{h^{2m}}$, $k = \text{const}$, wo $2m$ die Ordnung des diskreten Differentialoperators ist. Man untersucht leichter (im Hinblick auf die Glättungseigenschaften) statt (17.3) ein Verfahren der Art

$$\mathbf{u}_l^{j+1} = \mathbf{u}_l^j - \omega \frac{h_l^{2m}}{k} (\mathbf{A}_l \mathbf{u}_l^j - \mathbf{b}_l)$$

Ist $\text{diag } \mathbf{A}_l = \text{diag}(\frac{k}{h^{2m}})$, so ist ein beliebiges, weil bequem zu berechnendes ω , das gute Dämpfungseigenschaften liefert,

$$(17.4) \quad \omega = \omega_l = \frac{k}{\|h_l^{2m} \mathbf{A}_l\|_S} \quad (\text{Spektralnorm})$$

Wir zeigen dies im nächsten Paragraphen (vgl. Satz 18.2 und (18.29))

Das am meisten benutzte Verfahren ist

die Gauß-Seidel-Iteration (Einzelschritt-Verfahren)

Man kann zeigen, daß sie für eine große Klasse von Matrizen (nicht für alle) schneller konvergiert als die Jacobi-Iteration (vgl. Satz 18.4, Satz von Stein-Rosenberg (zu finden in Varga: Matrix Iterative Analysis, Prentice Hall 1962, vgl. dazu auch die Konvergenzbeispiele in den Übungen)

Zudem sind die Glättungseigenschaften besser als die des gedämpften Jacobi (vgl. Übungen).

Ein Iterationsschritt des Verfahrens zur Lösung der Gleichung $\mathbf{Ax} = \mathbf{b}$, $\mathbf{A} = (a_{ij})$, kann algorithmisch beschrieben werden durch

$$(17.5) \quad \text{for } i = 1 \text{ step } 1 \text{ until } n_l \text{ do } x_i := \frac{-1}{a_{ii}} \left(\sum_{j=1, j \neq i}^{n_l} a_{ij} x_j - b_j \right)$$

Man beachte, daß in jedem Schritt die “alten” Komponenten des Vektors durch die neuen überschrieben werden (Speicherersparnis gegenüber Jacobi).

Beachte:

1. Dieses Verfahren ist abhängig von der Reihenfolge, in der die Unbekannten iterativ verbessert werden (Anordnung der Gitterpunkte).
2. Ein Iterationsschritt durchläuft “sequentiell” die Zeilen des Systems und kann deshalb nicht in dieser Form parallelisiert werden, **jedoch**, eine geschickte Anordnung der Zeilen des Systems, d.h. eine vorgeschriebene Ordnung, in der die Gitterpunkte abgearbeitet werden, kann die Glättungseigenschaften gewährleisten, sogar verbessern, und abhängig von der Diskretisierungsmatrix, eine gewisse Parallelisierung ermöglichen.

Parallelisierung

Wir besprechen einige der gebräuchlichsten und erfolgreichsten Anordnungen der Gitterpunkte, die in Verbindung mit dem Gauß-Seidel-Verfahren für die Diskretisierung von Δu gute Dämpfungseigenschaften liefern und eine Parallelisierung ermöglichen. Wir beschränken uns auf den 2D-Fall.

1. Schachbrettanordnung: schwarz-weiß oder weiß-schwarz (Chequer bord ordering, red-black ordering),
2. 4-Farben-Ordnung (Four-colour-ordering),
3. zeilen- oder spaltenweise Anordnung (lexicographical ordering oder rotated lexicographical ordering),

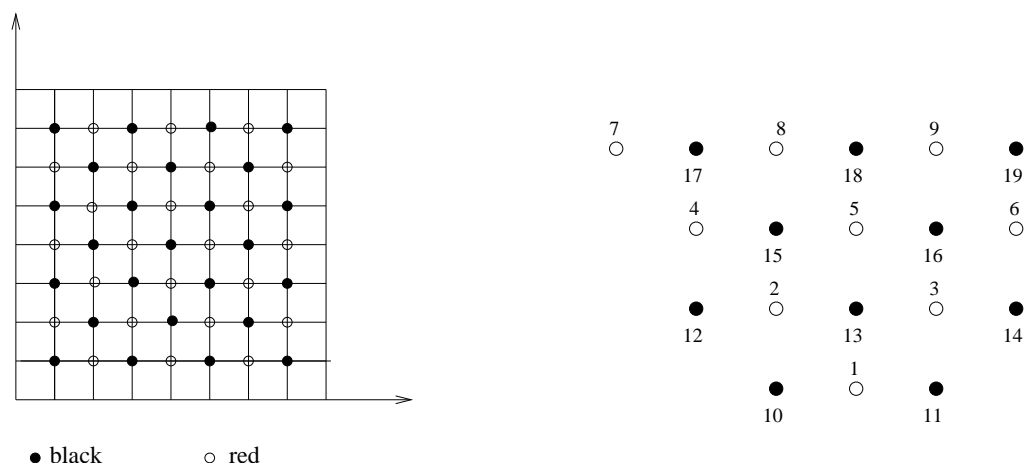
Für weitere Möglichkeiten vgl. etwa Hackbusch: Multigrid Methods, § 3.3 und 6.2; oder

Großmann/Roos: Numerik partieller Differentialgleichungen, § 5.2, 5.6

Diese Anordnungen betreffen jeweils eine Einteilung der Gitterpunkte in Klassen, die im Iterationsverfahren einzeln abgearbeitet werden. Die Klasseneinteilung ist abhängig von der Diskretisierungsmatrix, wie wir sehen werden.

Schachbrettanordnung (red-black): Gut geeignet für den 5-Punktstern für Δu . Wir illustrieren die Klasseneinteilung für das Einheitsquadrat und die Nummerierung

der Punkte am Beispiel eines beliebigen Gebiets.



Zuerst werden die roten, dann die schwarzen Punkte zeilenmäßig von unten nach oben durchgezählt. Anfangspunkt ist der erste rote Punkt (red-black).

Parallelisierung

Wir führen das GS-Verfahren (17.5) zunächst für die roten Punkte durch. Die iterative Verbesserung durch das GS für einen "roten" Punkt (rote Unbekannte) benötigt nur schwarze Nachbarpunkte (vgl. obige Abbildungen und den 5-Punktstern (12.5)), d.h. innerhalb der roten Punkteklasse ist die Reihenfolge der Punkte beim GS-Verfahren beliebig.

Dies kann zur Parallelisierung benutzt werden: z.B. einen extra Prozessor für jede Zeile zur Verbesserung der Punkte der roten Klasse, die in dieser Zeile enthalten sind.

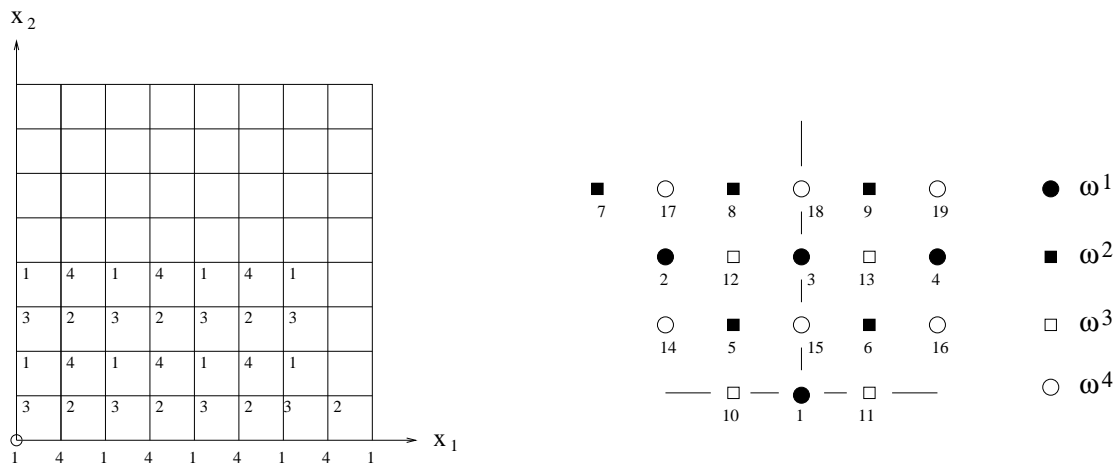
Nachdem alle roten Komponenten verbessert wurden, arbeitet man die schwarzen Komponenten ab, deren Verbesserung nur von den (schon verbesserten) roten Punkten abhängt. Für die Parallelisierung gilt das oben Gesagte.

Welche der Anordnungen bessere Glättungseigenschaften zeigt (red-black oder black-red), ist problemabhängig.

Beachte: Wir haben hier die Parallelisierung für Δu besprochen. Nicht alle Diskretisierungsmatrizen müssen eine solche Einteilung erlauben. Muß z.B. eine Ableitung höherer Ordnung diskretisiert werden, so benötigt dies üblicherweise mehr als nur die beiden unmittelbaren Gitternachbarn. Dann ist die beschriebene Parallelisierung nicht mehr möglich.

Sehr beliebt (auch für den 5-Punkte-Stern) ist auch die

Vierfarben-Ordnung: Gut geeignet für den 9-Punkte-Stern für Δu .



Die Punkte werden in 4 Klassen $\omega^1, \dots, \omega^4$, eingeteilt. Im Einheitsquadrat bezeichnen wir die Klasseneinteilung der Punkte durch Zahlen 1 bis 4.

Beim allgemeinen Gebiet führen wir zur Beschreibung ein zweites Koordinatenkreuz ein, dessen Ursprung bei zeilenweiser Nummerierung von unten nach oben im 2.ten Gitterpunkt liegt. Die einzelnen Klassen sind durch verschiedene Symbole gekennzeichnet. Die Nummerierung gibt die Reihenfolge an, in der die Punkte verbessert werden.

Wir beschreiben die Klassen wie folgt:

- $\omega^1 = \{(\nu h, \mu h); \nu, \mu \text{ gerade}\}$
- $\omega^2 = \{(\nu h, \mu h); \nu, \mu \text{ ungerade}\}$
- $\omega^3 = \{(\nu h, \mu h); \nu \text{ gerade, } \mu \text{ ungerade}\}$
- $\omega^4 = \{(\nu h, \mu h); \nu \text{ ungerade, } \mu \text{ gerade}\}$

Parallelisierung:

Innerhalb ω^1 werden nur die zu ω^1 gehörigen Komponenten verbessert. Die Verbesserung der Punkte $\in \omega^1$ mit dem 9-Punkte Stern benötigt nur Punkte aus $\omega^2, \omega^3, \omega^4$. Die Reihenfolge der Nummerierung innerhalb ω^1 ist dann bedeutungslos. Dasselbe gilt für die anderen Klassen.

Beispiel: Zuerst werden alle ω^1 -Punkte verbessert (z.B. ein Prozessor für jede Zeile, in der Punkte $\in \omega$ vorkommen), dann (wieder parallel) alle ω^2 -Punkte mit den verbesserten Werten aus ω^1 . Die Verbesserung der Punkte $\in \omega^2$ benötigt nur Punkte aus $\omega^1, \omega^3, \omega^4$, usw.

Die Vierfarben-Ordnung wird auch in Verbindung mit einem 5-Punkte-Stern angewandt (vgl. Hackbusch: Abschnitt 3.3.3 und Exercise 3.9.2)

Wir erwähnen weiter die

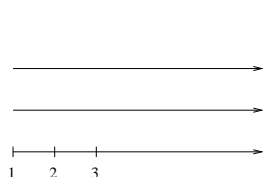
Lexikographische Ordnung: In jeder Zeile von links nach rechts und zeilenweise von unten nach oben oder

rotierte Lexikographische Ordnung: In jeder Spalte von unten nach oben und spaltenweise von links nach rechts.

Beispiel: stationäre Strömungsaufgaben. Bei der Gleichung

$$\nu \Delta u - v \operatorname{grad} u + f = 0, \quad v > 0.$$

läuft die Strömung von links nach rechts.



Zeilenweise Numerierung: gegen die Strömung zu rechnen oder senkrecht zur Strömung ist hier nicht sinnvoll. In diesem Beispiel verteilt sich die Information von links nach rechts.

Die Anordnung der Punkte in der Zeichnung ist sinnvoll, da der Zustand in “2” nur vom Zustand in “1” beeinflusst wird.

Ist der Diffusionskoeffizient $\nu = 0$, so liefert das Verfahren die exakte Lösung, da es auf die Transportgleichung reduziert ist.

Man erkennt die Abhängigkeit der Wahl von der Anwendung.

Ein Spezialfall der lexikographischen Anordnung im Hinblick auf Parallelisierung beim obigen Beispiel wäre etwa die

zebra-line-ordering: 1.te Klasse: die Zeilen 1,3,5,...;

2.te Klasse: die Zeilen 2,4,6,...;

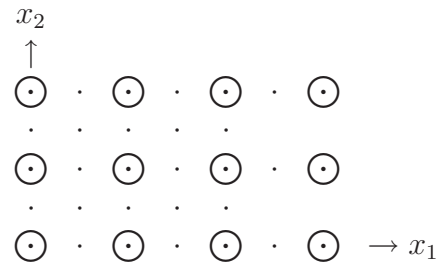
Auch red-black mit dem 5-Punktstern und zentralen Differenzen für die $\frac{\partial u}{\partial x_i}$ (und dann wieder zeilenweise) wäre möglich.

Für weitere Anordnungen siehe Hackbusch: Multigrid Methods.

Prolongation und Restriktion

Vernünftigerweise gibt es genau ein Vorgehen für die Prolongation, nämlich die Interpolation. Dabei hat sich gezeigt, daß für Gleichungen vom Grad $2m$ eine Interpolation der Ordnung m genügt. Wir beschränken uns hier auf den 2D-Fall und Δu , d.h. also, lineare Interpolation. Man beachte, daß auf Grund des GSV immer Defekte restringiert und prolongiert werden. Defekte haben (bei vorgegebenen Dirichlet-Randwerten, auf die wir uns hier beschränken,) immer Nullrandwerte. Dies muß bei der Prolongation berücksichtigt werden. Wir beschränken uns hier, der einfacheren Schreibweise wegen, auf den Fall gleicher Maschenweiten $h_1 = h_2 = h$. Das Vorgehen für $h_1 \neq h_2$ ist völlig analog.

Im Einheitsquadrat bezeichnen wir
 Grobgitterpunkte durch \odot ,
 Feingitterpunkte nur durch einen Punkt,
 $2h \hat{=}$ Grobgittermaschenweite,
 $h \hat{=}$ Feingittermaschenweite.
 Es seien
 $(0, 0), (0, 2h), (2h, 2h), (2h, 0) \in \omega_l \cap \omega_{l-1}$



Dann kann man die Prolongation einer auf ω_{l-1} (=Gitterpunkte auf Gitter $l-1$) definierten Funktion v (Lösung der Korrekturgleichung) nach ω_l wie folgt beschreiben:

9-Punkte-Prolongation

Auf dem Grobgitter bleiben die Werte unverändert.

$$(17.6) \quad \begin{aligned} v_l(0, 0) &= v_{l-1}(0, 0), & v_l(0, 2h) &= v_{l-1}(0, 2h); \\ v_l(2h, 0) &= v_{l-1}(2h, 0), & v_l(2h, 2h) &= v_{l-1}(2h, 2h); \end{aligned}$$

Für die Feingitterpunkte, die zwischen 2 Grobgitterpunkten liegen, liefert die lineare Interpolation

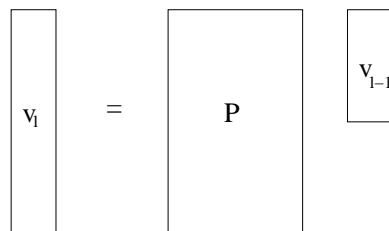
$$(17.7) \quad \begin{aligned} v_l(0, h) &= \frac{1}{2}v_{l-1}(0, 0) + \frac{1}{2}v_{l-1}(0, 2h) \\ v_l(h, 0) &= \frac{1}{2}v_{l-1}(0, 0) + \frac{1}{2}v_{l-1}(2h, 0) \\ v_l(2h, h) &= \frac{1}{2}v_{l-1}(2h, 0) + \frac{1}{2}v_{l-1}(2h, 2h) \\ v_l(h, 2h) &= \frac{1}{2}v_{l-1}(0, 2h) + \frac{1}{2}v_{l-1}(2h, 2h) \end{aligned}$$

und schließlich

$$(17.8) \quad \begin{aligned} v_l(h, h) &= \frac{1}{2}v_{l-1}(0, h) + \frac{1}{2}v_{l-1}(2h, h) \\ &= \frac{1}{4}v_{l-1}(0, 0) + \frac{1}{4}v_{l-1}(0, 2h) + \frac{1}{4}v_{l-1}(2h, 0) + \frac{1}{4}v_{l-1}(2h, 2h) \end{aligned}$$

Die Formeln (17.6) – (17.8) heißen 9-Punkte-Prolongation. Natürlich kann man diese Formeln in Matrixnotation aufschreiben.

$$v_l = P v_{l-1}$$



Billiger und Speicherplatz sparender ist jedoch, die prolongierten Werte direkt zu berechnen.

Wir demonstrieren dies für den Fall des Einheitsquadrats $\bar{\Omega}$. Die Randpunkte müssen für die Interpolation mitgeführt werden. Dann haben wir das

$$\text{Grobgridter:} \quad \omega_{l-1} = \{(\nu h_{l-1}, \mu h_{l-1}) \in \bar{\Omega}; 0 \leq \nu, \mu \leq \frac{1}{h_{l-1}}\},$$

$$\text{Feingitter:} \quad \omega_l = \{(\nu h_l, \mu h_l) \in \bar{\Omega}; 0 \leq \nu, \mu \leq \frac{1}{h_l}\}, \quad h_{l-1} = 2h_l.$$

Beachte: Bei Dirichlet-Randwerten sind die Defekte in den Randwerten = 0, deshalb werden die Randpunkte mit $v = 0$ vorbesetzt.

Dann wird, unter der Voraussetzung, daß v_{l-1} bekannt ist, die Prolongation beschrieben durch

for $x := 0$ step $2h_l$ until 1 do for $y = h_l$ step $2h_l$ until $1 - h_l$ do

$$v(x, y) := [v(x, y - h_l) + v(x, y + h_l)]/2;$$

for $y := 0$ step h_l until 1 do for $x = h_l$ step $2h_l$ until $1 - h_l$ do

$$v(x, y) := [v(x - h_l, y) + v(x + h_l, y)]/2;$$

d.h. zuerst werden die Feingitterpunkte unter- und oberhalb der Grobgitterpunkte (also in vertikaler Richtung) besetzt, danach die restlichen.

Natürlich sollen Interpolation und Restriktion sich "irgendwie entsprechen". Der Konvergenzbeweis zeigt, daß dies der Fall ist, wenn sie zueinander adjungiert sind. Natürlich kann auch funktionalanalytisch begründet werden, warum das sinnvoll ist.

Definition 17.1 Adjungierte Operatoren

Seien H_H und H_h Hilberträume mit den inneren Produkten $(\cdot, \cdot)_H$ und $(\cdot, \cdot)_h$, und

$$\mathbf{R}: H_h \xrightarrow{\text{auf}} H_H$$

eine lineare Abbildung.

Die Abbildung

$$\mathbf{P}: H_H \xrightarrow{\text{linear}} H_h$$

heißt adjungiert zu \mathbf{R} , falls gilt

$$(17.9) \quad (\mathbf{P}\mathbf{u}^H, \mathbf{v}^h)_h = (\mathbf{u}^H, \mathbf{R}\mathbf{v}^h)_H \quad \forall \mathbf{u}^H \in H, \quad \mathbf{v}^h \in H_h.$$

In unserem Fall sind $H_H = \mathbb{R}^{n_{l-1}}$, $H_h = \mathbb{R}^{n_l}$, dabei sind n_{l-1} bzw. n_l die Anzahl der Grobgitter- bzw. Feingitterpunkte im Einheitsquadrat inclusive der Randpunkte.

$$(\mathbf{u}^h, \mathbf{v}^h)_h = \sum_{j=1}^{n_l} \mathbf{u}_j^h \mathbf{v}_j^h h^{(2)}, \quad h^{(2)} = h_1 h_2, \quad h_i = \text{Maschenweite in } x_i\text{-Richtung,}$$

$$(\mathbf{u}^H, \mathbf{v}^H)_H = \sum_{i=1}^{n_{l-1}} \mathbf{u}_i^H \mathbf{v}_i^H H^{(2)} \quad H^{(2)} = H_1 H_2, \quad H_i = \text{Maschenweite in } x_i\text{-Richtung.}$$

Die 2 im Exponenten ist ein Hinweis auf den 2-dimensionalen Fall.

\mathbf{P} und \mathbf{R} sind Prolongation und Restriktion, wir setzen $H_i = 2h_i$. (Das ist Standard, aber kein "Muß".) Dann lautet (17.9) in Matrixschreibweise

$$(\mathbf{P}\mathbf{u}^H)^T \mathbf{v}^h h^{(2)} = (\mathbf{u}^H)^T \mathbf{P}^T \mathbf{v}^h h^{(2)} \stackrel{(17.9)}{=} (\mathbf{u}^H)^T \mathbf{R}\mathbf{v}^h \cdot H^{(2)} = (\mathbf{u}^H)^T \mathbf{R}\mathbf{v}^h \cdot 4h^{(2)}$$

oder

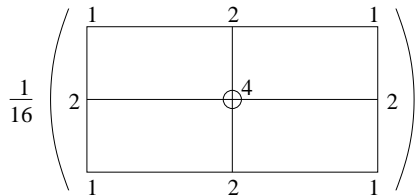
$$(\mathbf{u}^H)^T \mathbf{P}^T \mathbf{v}^h = (\mathbf{u}^H)^T \mathbf{R}\mathbf{v}^h \cdot 4.$$

Dies liefert als Konstruktionsvorschrift für eine Restriktionsmatrix

$$(17.10) \quad \mathbf{R} = \frac{1}{4} \mathbf{P}^T.$$

Benutzt man die 9-Punkte-Prolongation, so erhält man als Restriktion im 2D-Fall die

9-Punkte-Restriktion (Standardvariante) für den Defekt, (folgt aus (17.10) mit der Prolongation (17.6)-(17.8)) (Übung)



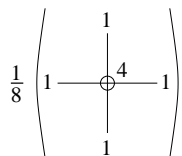
○ ist der Grobgitterpunkt, die Zahlen zeigen die Gewichte der Defekte d_h in den einzelnen Punkten. (Konvexkombination der beteiligten Punkte)

$$(17.11) \quad d_{ij}^H = \frac{1}{16} \{ d_{i+1,j+1}^h + d_{i-1,j+1}^h + d_{i-1,j-1}^h + d_{i+1,j-1}^h + 2(d_{i,j+1}^h + d_{i,j-1}^h + d_{i+1,j}^h + d_{i-1,j}^h) + 4d_{i,j}^h \}$$

Bemerkung: Beim FEM-Verfahren (Finite Element Method) sind Restriktion und Prolongation vorgeschrieben. Beim Differenzenverfahren hat man Wahlmöglichkeiten.

Aufgabe:

1. Man beschreibe die 9-Punkte-Prolongationsmatrix und leite daraus die zugehörige 9-Punkte-Restriktion ab. (inclusive Aufstellung der Restriktionsmatrix)
2. Wie sieht die Restriktionsmatrix für die 5-Punkte-Restriktion aus?



Wie sieht die zugehörige Prolongationsmatrix aus? Beschreibt sie eine stückweis lineare Interpolation?

3. Statt (17.8) könnte man auch wählen

$$v_l(h, h) = \frac{1}{2} v_{l-1}(0, 0) + \frac{1}{2} v_{l-1}(2h, 2h)$$

Wie sehen die zugehörige Prolongations- und Restriktionsmatrix aus?

Bemerkung: Es liegt nahe, als Restriktion die triviale Injektion zu wählen, d.h.

$$(\mathbf{R}d_l)(\mathbf{x}) = d_l(\mathbf{x}) \quad \mathbf{x} \in \omega_{l-1} \subset \omega_l.$$

Man kann zeigen (vgl. Hackbusch 3.5), daß für diese Restriktion in Verbindung mit der red-black-Ordnung, dem GS-Verfahren und dem 5-Punkte-Stern für Δu die Iteration zum Stehen kommen kann.

§ 18 Konvergenz des ZGV

Idee: Zur approximativen Lösung von $\mathbf{A}_k \mathbf{y}_k = \mathbf{b}_k$ mit dem ZGV leiten wir (wie beim Jacobi-Verfahren) für den Fehler \mathbf{e}^m der m -ten Iteration mit der Ausgangsnäherung \mathbf{y}_h^0 eine Fehlergleichung her mit einer noch aufzustellenden Iterationsmatrix \mathbf{K}_h

$$\mathbf{e}_h^{m+1} = \mathbf{K}_h \mathbf{e}_h^m, \quad \mathbf{e}_h^m = \mathbf{y}_h^m - \mathbf{y}_h^0, \quad \mathbf{K}_h = \text{Iterationsmatrix des ZGV.}$$

Konvergenz liegt vor, wenn der Spektralradius $\rho(\mathbf{K}_h)$ kleiner als 1 ausfällt.

Zur Herleitung der Fehlerdarstellung betrachten wir folgendes ZGV (mit Vor- und Nachglättung)

Gegeben: Ausgangsnäherung \mathbf{y}_h^0 für $\mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$. \mathbf{S} sei die Iterationsmatrix des Glättungsverfahrens, welches durch $\bar{\mathbf{y}}_h = \mathbf{S} \mathbf{y}_h^0 + \mathbf{T} \mathbf{b}_h$ beschrieben wird. Dann ist für die ν -malige Glättung die Restmatrix \mathbf{T}_ν aus der Iterationsvorschrift berechenbar.

$$(18.1) \quad \text{Vorglättung:} \quad \bar{\mathbf{y}}_h = \mathcal{S}(\mathbf{A}_h, \mathbf{b}_h) \mathbf{y}_h^0 = \mathbf{S}^{\nu_1} \mathbf{y}_h^0 + \mathbf{T}_{\nu_1} \mathbf{b}_h$$

$$(18.2) \quad \text{Defektberechnung:} \quad \mathbf{d}_h = \mathbf{b}_h - \mathbf{A}_h \bar{\mathbf{y}}_h$$

$$(18.3) \quad \text{Restriktion und Grobgitterlösung:} \quad \mathbf{v}_H = \mathbf{A}_H^{-1} (\mathbf{R} \mathbf{d}_h)$$

$$(18.4) \quad \text{Prolongation und Korrektur:} \quad \bar{\bar{\mathbf{y}}}_h = \bar{\mathbf{y}}_h + \mathbf{P} \mathbf{v}_H$$

$$(18.5) \quad \text{Nachglättung:} \quad \mathbf{y}_h^1 = \mathbf{S}^{\nu_2} \bar{\bar{\mathbf{y}}}_h + \mathbf{T}_{\nu_2} \mathbf{b}_h$$

Wir konstruieren zunächst die

Iterationsmatrix des ZGV

Als Glättungsiteration werden nur Verfahren verwendet, welche die exakte Lösung \mathbf{y}_h von $\mathbf{A}_h \mathbf{y}_h = \mathbf{b}_h$ als Fixpunkt haben (vgl. S 135), d.h. $\mathbf{y}_h = \mathbf{S}^{\nu_1} \mathbf{y}_h + \mathbf{T}_{\nu_1} \mathbf{b}_h$.

Abziehen dieser Gleichung von (18.1) liefert

$$(18.6) \quad \bar{\mathbf{e}}_h := \bar{\mathbf{y}}_h - \mathbf{y}_h = \mathbf{S}^{\nu_1} \mathbf{e}_h^0, \quad \text{mit} \quad \mathbf{e}_h^0 = \mathbf{y}_h^0 - \mathbf{y}_h.$$

Mit $\mathbf{d}_h \stackrel{(18.2)}{=} \mathbf{b}_h - \mathbf{A}_h \bar{\mathbf{y}}_h = \mathbf{A}_h (\mathbf{y}_h - \bar{\mathbf{y}}_h) = -\mathbf{A}_h \bar{\mathbf{e}}_h$ folgt aus (18.3)

$$(18.7) \quad \mathbf{v}_H \stackrel{(18.3)}{=} \mathbf{A}_H^{-1} \mathbf{R} \mathbf{d}_h \stackrel{(18.2)}{=} \mathbf{A}_H^{-1} \mathbf{R} (\mathbf{b}_h - \mathbf{A}_h \bar{\mathbf{y}}_h) = -\mathbf{A}_H^{-1} \mathbf{R} (\mathbf{A}_h \bar{\mathbf{y}}_h - \mathbf{A}_h \mathbf{y}_h) = -\mathbf{A}_H^{-1} \mathbf{R} \mathbf{A}_h \bar{\mathbf{e}}_h$$

Für den Fehler nach Prolongation und Korrektur erhält man

$$(18.8) \quad \begin{aligned} \bar{\bar{\mathbf{e}}}_h &:= \bar{\bar{\mathbf{y}}}_h - \mathbf{y}_h \stackrel{(18.4)}{=} \bar{\mathbf{y}}_h + \mathbf{P} \mathbf{v}_H - \mathbf{y}_h \\ &\stackrel{(18.6)}{=} \mathbf{S}^{\nu_1} \mathbf{e}_h^0 + \mathbf{P} \mathbf{v}_H \stackrel{(18.7)}{=} \mathbf{S}^{\nu_1} \mathbf{e}_h^0 - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R} \mathbf{A}_h \bar{\mathbf{e}}_h \\ &\stackrel{(18.6)}{=} \mathbf{S}^{\nu_1} \mathbf{e}_h^0 - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R} \mathbf{A}_h \mathbf{S}^{\nu_1} \mathbf{e}_h^0 = (\mathbf{I} - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R} \mathbf{A}_h) \mathbf{S}^{\nu_1} \mathbf{e}_h^0 \\ &= (\mathbf{A}_h^{-1} - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R}) \mathbf{A}_h \mathbf{S}^{\nu_1} \mathbf{e}_h^0 \end{aligned}$$

Die Nachglättung führt unter Verwendung der Fixpunktgleichung zu

$$\begin{aligned} e_h^1 &:= \mathbf{y}_h^1 - \mathbf{y}_h \stackrel{(18.5)}{=} \mathbf{S}^{\nu_2} \bar{\mathbf{y}}_h + \mathbf{T}_{\nu_2} \mathbf{b}_h - (\mathbf{S}^{\nu_2} \mathbf{y}_h + \mathbf{T}_{\nu_2} \mathbf{b}_h) \\ &= \mathbf{S}^{\nu_2} (\bar{\mathbf{y}}_h - \mathbf{y}_h) = \mathbf{S}^{\nu_2} \bar{\mathbf{e}}_h \quad (\text{Definition von } \bar{\mathbf{e}}_h) \\ &\stackrel{(18.8)}{=} \mathbf{S}^{\nu_2} \underbrace{(\mathbf{A}_h^{-1} - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R})}_{=: \mathbf{C}_h} \mathbf{A}_h \mathbf{S}^{\nu_1} \mathbf{e}_h^0 \end{aligned}$$

Damit erhalten wir die Iterationsmatrix des ZGV

$$(18.9) \quad \mathbf{K}_h(\nu_1, \nu_2) = \mathbf{S}^{\nu_2} (\mathbf{A}_h^{-1} - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R}) \mathbf{A}_h \mathbf{S}^{\nu_1} = \mathbf{S}^{\nu_2} \mathbf{C}_h \mathbf{S}^{\nu_1}.$$

Für die Iterationsmatrix $\mathbf{K}_h(\nu_1, \nu_2)$ des ZGV muß gezeigt werden

$$\|\mathbf{K}_h(\nu_1, \nu_2)\| < 1 \quad \text{bzw.} \quad \rho(\mathbf{K}_h(\nu_1, \nu_2)) < 1.$$

Führt man μ Schritte des ZGV aus, so erhält man mit der Abkürzung

$$(18.10) \quad \begin{aligned} \mathbf{C}_h &:= (\mathbf{A}_h^{-1} - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R}) \\ \mathbf{K}_h^\mu &= \mathbf{S}^{\nu_2} \mathbf{C}_h \mathbf{A}_h \mathbf{S}^{\nu_1} \mathbf{S}^{\nu_2} \mathbf{C}_h \mathbf{A}_h \mathbf{S}^{\nu_1} \dots \mathbf{S}^{\nu_2} \mathbf{C}_h \mathbf{A}_h \mathbf{S}^{\nu_1} \end{aligned}$$

Wir wollen Vor- und Nachglättung zusammenfassen, $\nu = \nu_1 + \nu_2$, und dazu in (18.10) geeignet klammern

$$(18.11) \quad \mathbf{K}_h^\mu = \underbrace{\mathbf{S}^{\nu_2}}_{\mathbf{F}} \underbrace{(\mathbf{C}_h \mathbf{A}_h \mathbf{S}^\nu)^{\mu-1} \mathbf{C}_h \mathbf{A}_h \mathbf{S}^{\nu_1}}_{\mathbf{U}}$$

Nun gilt für quadratische Matrizen \mathbf{F}, \mathbf{U} gleichen Formats: $\rho(\mathbf{F}\mathbf{U}) = \rho(\mathbf{U}\mathbf{F})$ (Übung, vgl. unten), sodaß

$$(18.12) \quad \begin{aligned} \rho(\mathbf{K}_h^\mu) &= \rho[(\mathbf{C}_h \mathbf{A}_h \mathbf{S}^\nu)^\mu] = [\rho(\mathbf{C}_h \mathbf{A}_h \mathbf{S}^\nu)]^\mu \quad \text{und} \\ \|\mathbf{K}_h^\mu\|_S &= \rho(\mathbf{K}_h)^\mu \quad \forall \mu \in \mathbb{N} \end{aligned}$$

Bemerkung: Für ein ZGV ohne Nachglättung gilt (18.12) mit $\nu = \nu_1 + 0$, d.h.

$$\mathbf{K}_h^\mu(\nu_1, 0) =: \mathbf{K}_h^\mu(\nu) = [\mathbf{K}_h(\nu)]^\mu = [(\mathbf{A}_h^{-1} - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R}) \mathbf{A}_h \mathbf{S}^\nu]^\mu = [\mathbf{C} \mathbf{A}_h \mathbf{S}^\nu]^\mu$$

Konvergenz liegt vor, wenn

$$(18.13) \quad \|\mathbf{K}_h(\nu)\| < 1 \quad \text{bzw.} \quad \|\mathbf{K}_h(\nu_1, \nu_2)\| < 1.$$

Der Vergleich zeigt, daß ein Nachweis von $\rho((\mathbf{A}_h^{-1} - \mathbf{P} \mathbf{A}_H^{-1} \mathbf{R}) \mathbf{A}_h \mathbf{S}^\nu) < 1$ gleichermaßen die Konvergenz von Verfahren mit und ohne Nachglättung sichert.

Nötig dazu ist, wie oben erwähnt, folgende

Aufgabe: Zeige für quadratische Matrizen \mathbf{A}, \mathbf{B} : $\rho(\mathbf{A}\mathbf{B}) = \rho(\mathbf{B}\mathbf{A})$.

Hinweise:

(i) Zeige es zuerst unter der Voraussetzung, daß eine der Matrizen, z.B. \mathbf{B} regulär ist.

(ii) Verallgemeinere (i) durch ein Stetigkeitsargument, angewandt auf eine "leicht gestörte" Matrix \mathbf{B} (Jordan-Normalform).

Konvergenzbeweise für das ZGV werden üblicherweise geführt, indem man die Iterationsmatrix

$$\mathbf{K}_h(\nu) = (\mathbf{A}_h^{-1} - \mathbf{P}\mathbf{A}_H^{-1}\mathbf{R})\mathbf{A}_h\mathbf{S}^\nu$$

aufteilt in

$$(18.14) \quad \text{Glättungsteil} \quad \mathbf{A}_h\mathbf{S}^\nu \quad \text{und}$$

$$(18.15) \quad \text{Approximationsteil} \quad \mathbf{C}_h = (\mathbf{A}_h^{-1} - \mathbf{P}\mathbf{A}_H^{-1}\mathbf{R}).$$

Definition 18.1 (Hackbusch)

Das ZGV besitzt die Glättungseigenschaft, falls

$$(18.16) \quad \|\mathbf{A}_h\mathbf{S}^\nu\| \leq c_S(\nu) \quad \text{mit} \quad c_S(\nu) \rightarrow 0 \text{ für } \nu \rightarrow \infty$$

mit einer von ν abhängigen Konstanten $c_S = c_S(\nu)$.

Das ZGV besitzt die Approximationseigenschaft, falls

$$(18.17) \quad \|\mathbf{A}_h^{-1} - \mathbf{P}\mathbf{A}_H^{-1}\mathbf{R}\| \leq c_a(h) \quad (= \text{const}).$$

Bemerkungen:

1. Ziel eines Konvergenzbeweises ist die Ungleichung

$$\begin{aligned} \|\mathbf{K}_h(\nu)\| &= \|(\mathbf{A}_h^{-1} - \mathbf{P}\mathbf{A}_H^{-1}\mathbf{R})\mathbf{A}_h\mathbf{S}^\nu\| \\ &\stackrel{(*)}{\leq} \|(\mathbf{A}_h^{-1} - \mathbf{P}\mathbf{A}_H^{-1}\mathbf{R})\| \|\mathbf{A}_h\mathbf{S}^\nu\| \leq c_a c_S \stackrel{!}{<} 1 \quad \text{für kleines } \nu. \end{aligned}$$

Die Abschätzungen (18.16),(18.17) müssen also möglichst scharf ausfallen. Sie werden üblicherweise gesondert bewiesen, müssen beide aber bzgl. Normen ausgeführt werden, für welche die Abschätzung (*) gilt

2. c_a in (18.17) ist eine feste Zahl, die < 1 sein kann, aber nicht muß. Im letzteren Fall hängt die Konvergenz von \mathbf{S} ab. Eine optimale Glättung durch \mathbf{S} muß also nicht notwendigerweise mit einer optimalen Konvergenz äquivalent sein. (vgl. Bemerkung auf S. 126)
3. Bei der Durchführung des MGV hat man immer zwei benachbarte Aufgaben zu untersuchen (auf dem feineren und auf dem gröberen Gitter)

$$\mathbf{A}_h\mathbf{y}_h = \mathbf{d}_h, \quad \mathbf{A}_H\mathbf{y}_H = \mathbf{d}_H \quad (\text{Erinnerung: } \mathbf{d}_H = \mathbf{R}\mathbf{d}_h)$$

Die Lösungen dieser Defektgleichungen werden verglichen (Nachiteration). Es gilt

$$(18.18) \quad \|\mathbf{y}_h - \mathbf{P}\mathbf{y}_H\| = \|(\mathbf{A}_h^{-1} - \mathbf{P}\mathbf{A}_H^{-1}\mathbf{R})\mathbf{d}_h\| = \|\mathbf{C}_h\mathbf{d}_h\| \leq \|\mathbf{C}_h\| \|\mathbf{d}_h\|.$$

Hieraus erklärt sich der Name Approximationseigenschaft. Bei vorgegebenem Gitter kann $\|\mathbf{C}_h\|$ nicht gegen Null gehen.

4. Es ist einleuchtend, dass man Abschätzungen der Art

$$(18.19) \quad \|\mathbf{y}_h - \mathbf{P}\mathbf{y}_H\| \leq c_\alpha h^\alpha \quad (\text{Standart: } \alpha = 2 \text{ oder } 0)$$

zeigen kann. Solche Abschätzungen wurden u.a. von Dryja (polnischer Mathematiker), Hackbusch und Bachvalov bewiesen. Sie sind sehr aufwendig, weshalb wir hier darauf verzichten. (Wir verweisen z.B. auf das Hackbusch-Buch, Abschnitt 6.3.2)

Dies führt dazu, dass man in den Konvergenzbeweisen die Eigenschaften (18.16),(18.17) gelegentlich abändert zu

$$(18.16') \quad \|\mathbf{A}_h \mathbf{S}^\nu\| \leq \frac{1}{h^\alpha} c_S(\nu), \quad c_S(\nu) \xrightarrow{\nu \rightarrow \infty} 0,$$

$$(18.17') \quad \|\mathbf{C}_h\| \leq c_\alpha h^\alpha.$$

In unseren Beispielen zum Beweis der Glättungseigenschaft werden wir sehen, daß sich Abschätzungen der Art (18.16') "natürlich" ergeben. $\alpha = 0$ entspricht (18.16),(18.17), $\alpha = 0$ ist üblich für Δu .

Wir zeigen im nächsten Abschnitt die Glättungseigenschaften einiger Verfahren. Beliebte sind dabei Normen, die sich (unter gewissen Symmetrievoraussetzungen) auf die Spektralnorm stützen. Dann kann nämlich mit dem Spektralradius (und damit mit den Eigenwerten) statt mit Normen gerechnet werden.

Glättungseigenschaften des gedämpften Jacobi-Verfahrens

Die Iterationsmatrix des gedämpften Jacobi-Verfahrens lautet

$$\mathbf{S} = \mathbf{I} - \omega \mathbf{D}_h^{-1} \mathbf{A}_h, \quad \mathbf{D}_h = \text{Diagonalmatrix von } \mathbf{A}_h.$$

Die Dämpfungseigenschaft ist also neben dem Verfahren auch von \mathbf{A}_h abhängig.

Wir beschränken unsere Untersuchungen auf Matrizen \mathbf{A}_h mit

$$(18.20) \quad \mathbf{A}_h = \mathbf{A}_h^T > 0, \quad \mathbf{D}_h = \frac{k}{h^{2m}} \mathbf{I}, \quad k = \text{const.}$$

\mathbf{A} hat dann nur positive Eigenwerte. Die Iterationsmatrix \mathbf{S} des Jacobi-Verfahrens ist symmetrisch läßt sich nun schreiben als (vgl. (18.4))

$$(18.21) \quad \mathbf{S} =: \mathbf{I} - \tilde{\omega} \mathbf{A}_h \quad \text{mit } \tilde{\omega} = \frac{h^{2m}}{k} \omega.$$

Zum Beweis der Glättungseigenschaft ist mit einer geeigneten Norm der Ausdruck $\|\mathbf{A}\mathbf{S}^\nu\|$ abzuschätzen. Wegen (18.20) gilt mit der Spektralnorm $\|\cdot\|_S$ (Beachte: \mathbf{S}^ν ist ein Polynom in \mathbf{S} , also ist $\mathbf{A}_h \mathbf{S}^\nu$ symmetrisch)

$$(18.22) \quad \|\mathbf{A}_h \mathbf{S}^\nu\|_S = \max_{\lambda \in \sigma(\mathbf{A}_h)} |\lambda(1 - \tilde{\omega}\lambda)^\nu|,$$

Um mit Eigenwerten rechnen zu können, legen wir deshalb für die weiteren Untersuchungen die Spektralnorm zu Grunde.

Ziel: Gesucht ist eine Abschätzung dieses Ausdrucks durch eine in ν fallende Schranke, die auch Auskunft über zulässige $\tilde{\omega}$ -Werte gibt.

(Bemerkung: In Hackbusch (6.2.3) wird eine $\tilde{\omega}$ -Schranke nicht hergeleitet, vgl. jedoch Exercise 6.6.3)

Satz 18.2

Seien $\mathbf{A}_h = \mathbf{A}_h^T > 0$, $\mathbf{D}_h := \text{diag}(\mathbf{A}_h) = \frac{k}{h^{2m}} \mathbf{I}$, $k = \text{const.}$, und

$\mathbf{S} = \mathbf{I} - \tilde{\omega} \mathbf{A}_h$, $\tilde{\omega} = \frac{h^{2m}}{k} \omega$ die Iterationsmatrix des gedämpften Jacobi-Verfahrens

Dann gilt für

$$\omega \leq \frac{1.2}{\|\mathbf{A}_h\|_S} \frac{k}{h^{2m}} \quad (\text{Spektralnorm})$$

die Glättungseigenschaft

$$\|\mathbf{A}_h \mathbf{S}^\nu\|_S \leq \frac{\|\mathbf{A}_h\|_S}{1.2} \frac{1}{2(\nu + 1)}.$$

Beweis

Wir setzen $t = \tilde{\omega} \lambda$ und erhalten für ein noch zu bestimmendes t_{max}

$$(18.23) \quad \|\mathbf{A}_h \mathbf{S}^\nu\|_S = \frac{1}{\tilde{\omega}} \max_t |t(1-t)^\nu| \quad \text{für } 0 \leq t \leq t_{max}.$$

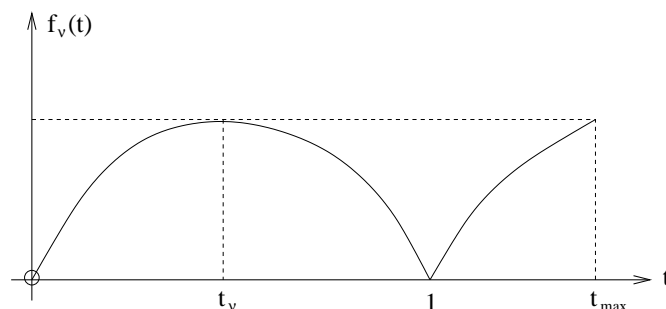
Die Untersuchung dieses Ausdrucks ist ein Kernstück vieler Glättungseigenschaftsnachweise (vgl. z.B. symmetrisches GSV und GSV mit Schachbrettanordnung.)

Vorgehen: Durch Abschätzung des Ausdrucks

$$f_\nu(t) = |t(1-t)^\nu|$$

nach oben leiten wir eine Ungleichung her, die uns in Verbindung mit $\lambda_{max}(\mathbf{A}_h) = \|\mathbf{A}_h\|_S$ eine Schranke für $\tilde{\omega}$ liefern wird (beachte: alle Eigenwerte von \mathbf{A} sind positiv).

$f_\nu(t) = |t(1-t)^\nu|$ hat für alle $\nu \geq 1$ qualitativ das folgende Aussehen (vgl. Abbildung), wobei für $\nu = 1$ in $t = 1$ eine Spitze auftritt.



In $0 \leq t \leq 1$ gilt $f_\nu(t) = t(1-t)^\nu$. Gesucht: $\max f_\nu(t)$ in $[0, 1]$

$$\begin{aligned} f'_\nu(t) &= (1-t)^\nu - t\nu(1-t)^{\nu-1} = (1-t)^{\nu-1}[1-t-t\nu] \stackrel{!}{=} 0 \\ 1-t(1+\nu) &= 0 \implies t_\nu = \frac{1}{(\nu+1)} \implies \\ f_\nu(t_\nu) &= \frac{1}{(\nu+1)} \left(1 - \frac{1}{\nu+1}\right) = \frac{1}{(\nu+1)} \left(\frac{1}{(1+\frac{1}{\nu})^\nu}\right) \end{aligned}$$

wegen $(1 + \frac{1}{\nu})^\nu \nearrow e$ für $\nu \rightarrow \infty$ gilt für $\nu = 1$

$$f_\nu(t) = \frac{1}{(\nu+1)} \left(\frac{1}{(1+\frac{1}{\nu})^\nu}\right) \leq \frac{1}{(\nu+1)} \frac{1}{2} \quad \forall t \in [0, 1].$$

Wir bestimmen, zunächst für $\nu = 1$, das maximale Intervall $(0, t_{max}]$ in dem diese Ungleichung gilt. Da $f_\nu(t) = t(t-1)$ für $t \geq 1$, ergibt sich t_{max} aus

$$\begin{aligned} t(t-1) &\leq \frac{1}{2} \frac{1}{(\nu+1)} \stackrel{\nu=1}{=} \frac{1}{4} \quad \text{zu} \\ (18.24) \quad t_{max} &= \frac{1+\sqrt{2}}{2}. \end{aligned}$$

Damit gilt in $[0, t_{max}]$ für $\nu = 1$: (beachte $f_\nu(t)$ ist monoton wachsend für $t \geq 1$)

$$(18.25) \quad |t(1-t)^\nu| \leq t_{max}(t_{max}-1) = \frac{1+\sqrt{2}}{2} \left(\frac{1+\sqrt{2}}{2} - 1\right)^\nu \leq \frac{1}{2(\nu+1)}$$

Mit $\frac{1+\sqrt{2}}{2} \approx 1.2$ ist die zweite Ungleichung äquivalent zu

$$(18.26) \quad 2.4(\nu+1) \cdot 0.2^\nu \leq 1, \quad \text{für } \nu = 1.$$

Die linke Seite dieser Ungleichung ist für $\nu \geq 1$ monoton fallend in ν (Ableitung ausrechnen). Deshalb gilt die Ungleichung (18.26) für alle $\nu \geq 1$, also gelten auch die Ungleichungen (18.25) für alle $\nu \geq 1$ und $0 \leq t \leq 1.2$.

Von Bedeutung für t sind die Werte $t_i = \tilde{\omega} \lambda_i(\mathbf{A}_h)$. Nun gilt $\lambda_i(\mathbf{A}_h) \leq \lambda_{max} = \|\mathbf{A}_h\|_S$. Man beachte, daß (18.25) bis auf den Faktor $\frac{1}{\tilde{\omega}}$ eine Abschätzung für den Glättungsterm liefert (vgl. (18.23)). Um die Glättungseigenschaft zu gewährleisten (d.h. um die Ungleichung (18.25) zu erhalten), fordern wir deshalb die Ungleichung

$$\tilde{\omega} \lambda_i \leq \tilde{\omega} \lambda_{max} = \tilde{\omega} \|\mathbf{A}_h\|_S \stackrel{!}{\leq} t_{max} = 1.2,$$

woraus sich für $\tilde{\omega}$ folgende Schranke ergibt

$$(18.27) \quad \tilde{\omega} \leq \frac{1.2}{\|\mathbf{A}_h\|_S} \quad \text{bzw. } \omega = \tilde{\omega} \frac{k}{h^{2m}}. \quad (\text{vgl. Satz (18.2)})$$

Die Ungleichung (18.25) gilt für t_{max} . Aus (18.22) erhält man deshalb die Glättungseigenschaft

$$(18.28) \quad \|\mathbf{A}_h \mathbf{S}^\nu\|_S = \frac{1}{\tilde{\omega}} \max_{t \in (0, t_{max}]} |t(1-t)| \leq \frac{1}{\tilde{\omega}} \frac{1}{2(\nu+1)} = \frac{\|\mathbf{A}_h\|_S}{1.2} \cdot \frac{1}{2(\nu+1)}.$$



Bemerkungen:

1. Satz 18.2 zeigt, daß der in (17.4) angegebene Dämpfungsparameter zulässig ist.

$$(18.29) \quad \omega_l = \frac{k}{h^{2m} \|A_h\|_S} \leq \tilde{\omega} \frac{k}{h^{2m}} = \frac{1.2 k}{h^{2m} \|A_h\|_S}.$$

2. Ist A_h die Diskretisierungsmatrix von $-\Delta u$ so gilt im Fall

$$1D: \quad \|A_h\|_S \leq \frac{4}{h^2} \quad (\text{vgl. Lemma 3.7})$$

$$2D: \quad \|A_h\|_S \leq \frac{8}{h^2}$$

Wird dies in (18.28) eingesetzt, erhält man eine Abschätzung vom Typ (18.16') (z.B. für 1D).

$$\|A_h S^\nu\|_S \leq \frac{\|A_h\|_S}{1.2} \cdot \frac{1}{2(\nu+1)} \leq \frac{4}{1.2h^2} \cdot \frac{1}{2(\nu+1)}$$

3. Im 1D-Fall ist $h^{2m} = h^2$, $k = 2$ (vgl. (18.21)). Dann folgt aus (18.28):

$$\omega \leq \frac{1.2k}{h^{2m} \|A_h\|_S} = \frac{2 \cdot 1.2}{h^2 \frac{4}{h^2}} = 0.6$$

Als optimalen Dämpfungsparameter hatten wir in § 15 im Demonstrationsbeispiel $\omega_{opt} = \frac{2}{3} = 0.6\dots$ ermittelt. ω_{opt} fällt also nicht in die zulässige Schranke. Nun kann man in (18.24) natürlich eine größere Schranke für $\max |t(t-1)|$ zulassen und dadurch (vgl. die Zeichnung) das Intervall $(0, t_{max}]$ so erweitern, daß ω_{opt} auch noch in den zulässigen Bereich fällt. Rechnet man dies nach, so stellt man fest, daß sich die Abschätzung für $\|A_h S^\nu\|_S$ insgesamt verschlechtert. Um die Konvergenzbedingung zu erfüllen, muß ein größeres ν gewählt werden.

Man erkennt also: Optimal Glättung und optimale Konvergenz sind nicht äquivalent.

4. Beachte: Der Beweis von Satz 18.2 beruht auf der Gleichung (18.22). Diese Darstellung hat als wesentliche Voraussetzung die Symmetrie der Iterationsmatrix S des Jacobi-Verfahrens und damit des Glättungsteils $A_h S^\nu$. Sonst ist eine Normdarstellung (18.22) nicht möglich und damit auch kein Rechnen mit dem Spektralradius, bzw. mit der Spektralnorm.)

Die Iterationsmatrix des GSV ist nicht symmetrisch, weshalb für die Glättungseigenschaft dieses Verfahrens bei zeilenweiser Nummerierung der Punkte (bisher) kein Beweis existiert.

Man kann die Glättungseigenschaft für das GSV retten unter Zuhilfenahme zusätzlicher Struktureigenschaften der Iterationsmatrix in Zusammenhang mit einer geeigneten Nummerierung der Giterpunkte (ohne Beweis).

Wir werden sehen, daß für das symmetrische GSV der Glättungsteil wieder symmetrisch ist.

Satz 18.3

Sei $\mathbf{A}_h = \mathbf{L} + \mathbf{D} + \mathbf{U}$, $\mathbf{A}_h = \mathbf{A}_h^T > 0$ und $\mathbf{S} = \mathbf{I} - \mathbf{W}^{-1}\mathbf{A}_h$ die Iterationsmatrix des symmetrischen Gauß-Seidel-Verfahrens mit $\mathbf{W} = \mathbf{A}_h + \mathbf{LD}^{-1}\mathbf{U}$ (vgl. Satz 15.6). Dann gilt (mit der Spektralnorm) folgende Glättungseigenschaft

$$\|\mathbf{W}_h^{-\frac{1}{2}}\mathbf{A}_h\mathbf{S}^\nu\mathbf{W}_h^{-\frac{1}{2}}\|_S \leq \frac{1}{2(\nu+1)}$$

Beweis

Idee: Wir formen den Glättungsteil \mathbf{G} so um, daß wir (wie beim Jacobi-Verfahren) eine zu (18.23) analoge Darstellung erhalten.

Für den Glättungsteil gilt

$$\begin{aligned} \mathbf{G} &= \underbrace{\mathbf{W}_h^{-\frac{1}{2}}\mathbf{A}_h\mathbf{W}_h^{-\frac{1}{2}}}_{=: \mathbf{X}_h} \mathbf{W}_h^{\frac{1}{2}}\mathbf{S}^\nu\mathbf{W}_h^{-\frac{1}{2}}, & \mathbf{X}_h &:= \mathbf{W}_h^{-\frac{1}{2}}\mathbf{A}_h\mathbf{W}_h^{-\frac{1}{2}} \\ &= \mathbf{X}_h \mathbf{W}_h^{\frac{1}{2}}\mathbf{S}\mathbf{W}_h^{-\frac{1}{2}}\mathbf{W}_h^{\frac{1}{2}}\mathbf{S}\mathbf{W}_h^{-\frac{1}{2}} \dots \mathbf{W}_h^{\frac{1}{2}}\mathbf{S}\mathbf{W}_h^{-\frac{1}{2}} \\ &= \mathbf{X}_h(\mathbf{W}_h^{\frac{1}{2}}\mathbf{S}\mathbf{W}_h^{-\frac{1}{2}})^\nu. \end{aligned}$$

Nun ist gemäß (15.19) (Darstellung der Iterationsmatrix \mathbf{S} des symmetrischen Gauß-Seidel)

$$\mathbf{W}_h^{\frac{1}{2}}\mathbf{S}\mathbf{W}_h^{-\frac{1}{2}} = \mathbf{W}_h^{\frac{1}{2}}(\mathbf{I} - \mathbf{W}_h^{-1}\mathbf{A}_h)\mathbf{W}_h^{-\frac{1}{2}} = \mathbf{I} - \mathbf{W}_h^{-\frac{1}{2}}\mathbf{A}_h\mathbf{W}_h^{-\frac{1}{2}} = \mathbf{I} - \mathbf{X}_h$$

also (vgl. die 3. Formelzeile des Beweises)

$$(18.32) \quad \mathbf{G} = \mathbf{X}_h(\mathbf{I} - \mathbf{X}_h)^\nu, \quad \text{Darstellung des Glättungsteils.}$$

Da \mathbf{W} und \mathbf{A}_h symmetrisch sind, gilt dies auch für $\mathbf{W}_h^{-\frac{1}{2}}$ und \mathbf{X}_h .

Die Matrix des Glättungsteils ist also symmetrisch. (beachte: $(\mathbf{I} - \mathbf{X}_h)^\nu$ ist ein Polynom in \mathbf{X}_h .)

Wir zeigen

$$(18.33) \quad 0 < \mathbf{X}_h \leq \mathbf{I}, \quad \text{positive Definitheit.}$$

Beweis von (18.33):

$$(\mathbf{X}_h\mathbf{x}, \mathbf{x}) = (\mathbf{W}_h^{-\frac{1}{2}}\mathbf{A}_h\mathbf{W}_h^{-\frac{1}{2}}\mathbf{x}, \mathbf{x}) = (\mathbf{A}_h\mathbf{W}_h^{-\frac{1}{2}}\mathbf{x}, \mathbf{W}_h^{-\frac{1}{2}}\mathbf{x}) > 0 \implies \mathbf{X}_h > 0, \text{ da } \mathbf{A}_h > 0.$$

$$\text{Wegen } (\mathbf{A}_h\mathbf{y}, \mathbf{y}) \leq (\mathbf{W}_h\mathbf{y}, \mathbf{y}) \quad \forall \mathbf{y} \quad (\text{vgl. Satz 15.6 b)) \text{ ist}$$

$$(\mathbf{A}_h\mathbf{W}_h^{-\frac{1}{2}}\mathbf{x}, \mathbf{W}_h^{-\frac{1}{2}}\mathbf{x}) \leq (\mathbf{W}_h\mathbf{W}_h^{-\frac{1}{2}}\mathbf{x}, \mathbf{W}_h^{-\frac{1}{2}}\mathbf{x}) = (\mathbf{W}_h^{\frac{1}{2}}\mathbf{x}, \mathbf{W}_h^{-\frac{1}{2}}\mathbf{x}) = (\mathbf{I}\mathbf{x}, \mathbf{x}) \quad \text{daraus folgt}$$

$$\mathbf{X}_h \leq \mathbf{I}.$$

Wegen (18.33) gilt für die Eigenwerte $\lambda \in \sigma(\mathbf{X}_h)$: $0 < \lambda \leq 1$, denn

$$0 \leq \frac{\mathbf{x}^T \mathbf{X}_h \mathbf{x}}{\|\mathbf{x}\|_2^2} \leq \frac{(\mathbf{x}, \mathbf{x})}{\|\mathbf{x}\|_2^2} \leq 1 \quad (\text{Rayleighquotient})$$

Deshalb folgt für die Abschätzung des Glättungsteils aus (18.32)- (18.33) (Spektralnorm)

$$(18.34) \quad \|\mathbf{G}\|_S \leq \max_{\lambda \in [0,1]} |\lambda(1-\lambda)^\nu| \leq \frac{1}{2(1+\nu)} \quad (\text{vgl. (18.23) - (18.25) mit } \tilde{\omega} = 1)$$

Dies ist die gewünschte Glättungseigenschaft. Eine Abschätzung des Approximationsanteils (vgl. (18.31)) durch eine Konstante c_α (das ist schwächer als durch $c_\alpha h^\alpha$) beendet den Konvergenzbeweis für hinreichend großes ν . ■

Nachbemerkungen

1. Angenehm ist, daß die Abschätzung
 - a) unabhängig ist von h (vgl. dazu die Nachbemerkung 3)) und
 - b) daß man sich nicht um Dämpfungsparameter kümmern muß.
2. Wesentliche Grundvoraussetzung für beide Glättungsbeweise war jeweils die Symmetrie von $\mathbf{D}^{-1}\mathbf{A}$ und \mathbf{W}^{-1} . Sie führte beim gedämpften Jacobi-Verfahren zur Gleichung (18.22), beim GSV zu (18.32) und (18.34). Da diese Symmetrie beim gewöhnlichen GSV nicht vorliegt, konnte eine Glättungseigenschaft für dieses Verfahren bisher nicht gezeigt werden. Trotzdem wird der Gauß-Seidel mit Erfolg benutzt.
3. Beim Jacobi-Verfahren war die Glättungseigenschaft in natürlicher Weise von h abhängig über die Diskretisierungsmatrix (vgl. Bemerkung 2 nach dem Beweis von Satz 18.2). Symmetrisches GSV und Jacobi-Verfahren lassen sich in einer gemeinsamen Form schreiben:

$$\mathbf{x}^{m+1} = (\mathbf{I} - \mathbf{W}^{-1}\mathbf{A}_h)\mathbf{x}^m + \mathbf{W}^{-1}\mathbf{b}, \quad (\text{symmetrisches GSV, vgl. (15.17),(15.18)})$$

$$\mathbf{x}^{m+1} = (\mathbf{I} - \omega\mathbf{D}^{-1}\mathbf{A}_h)\mathbf{x}^m + \omega\mathbf{D}^{-1}\mathbf{b}, \quad (\text{Jacobi-Verfahren, vgl. (15.1)})$$

Dadurch erklären sich auch die Ähnlichkeiten in den Glättungsbeweisen. Der Glättungsbeweis des symmetrischen GSV läßt sich damit auch auf das Jacobi-Verfahren übertragen. Daß die Glättungsabschätzung für das symmetrische GSV unabhängig von h ausfällt, liegt an der verwendeten Norm, welche die h -Abhängigkeit auf den Approximationsteil überträgt.

4. Glättungseigenschaftern für nichtsymmetrische Verfahren stecken noch in den Kinderschuhen.
5. Die Konvergenzeigenschaften bei der Verwendung des symmetrischen GSV werden durch den Beweis sicherlich unterschätzt. Sie fallen für Jacobi-Verfahren und symmetrisches GSV etwa gleich aus, was nicht ganz einsichtig ist, wenn man sich überlegt, daß ein Schritt symmetrisches GSV zwei Schritten GSV entspricht und daß das GSV üblicherweise schneller konvergiert als der Jacobi.

6. Zum Vergleich der Konvergenzgeschwindigkeit der Verfahren zitieren wir, ohne Beweis (vgl. Varga: Matrix Iterative Analysis, Kap. 3.3)

Satz 18.4 Stein-Rosenberg

Sei $\mathbf{A} \in C^{n \times n}$, $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{R}$,

$\mathbf{S} = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{R})$ die Iterationsmatrix des Jacobi-Verfahrens mit
 $s_{ij} \leq 0 \forall i \neq j, s_{ii} = 0 \forall i$,

$\mathbf{G} = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{R}$ die Iterationsmatrix des Gauß-Seidel-Verfahrens.

Dann gilt genau eine der folgenden Aussagen

- (a) $\rho(\mathbf{S}) = \rho(\mathbf{G}) = 0$
- (b) $0 < \rho(\mathbf{G}) < \rho(\mathbf{S}) < 1$
- (c) $1 = \rho(\mathbf{G}) = \rho(\mathbf{S})$
- (d) $1 < \rho(\mathbf{S}) < \rho(\mathbf{G})$

Setzt man $\mathbf{L} = \mathbf{D}\tilde{\mathbf{L}}$ und $\mathbf{R} = \mathbf{D}\mathbf{U}$, so erhält man die Darstellung aus Varga.

Beachte: Für unsere Diskretisierungen sind alle diese Voraussetzungen erfüllt, denn alle Nebendiagonalelemente von \mathbf{A} sind ≤ 0 , alle Diagonalelemente von \mathbf{A} sind positiv. Es liegt Fall b) vor, denn die Konvergenzeigenschaft wird durch Satz 15.5 gesichert.

Insbesondere: Jacobi ist langsamer als Gauß-Seidel und Gauß-Seidel braucht weniger Speicherplatz als Jacobi.

Letzteres folgt aus der Darstellung (17.5). Die alten Komponenten werden beim Gauß-Seidel sofort durch die neuen überschrieben. Beim Jacobi müssen sie aufbewahrt werden, bis alle Komponenten berechnet sind.

§ 19 Konvergenz des Mehrgitterverfahrens

Der Konvergenzbeweis beruht natürlich auf einer Untersuchung der Iterationsmatrix des MGV. Zu ihrer Beschreibung benötigen wir einige

Vorbetrachtungen über Iterationsverfahren zur Lösung der Korrekturgleichungen

Jedes lineare Iterationsverfahren zur Lösung von

$$(19.1) \quad \mathbf{A}\mathbf{v} = \mathbf{d}, \quad \mathbf{A} \in \mathbb{R}^{n \times n}$$

kann in der Form

$$(19.2) \quad \mathbf{v}^{j+1} = \varphi(\mathbf{v}^j, \mathbf{d}) := \mathbf{M}\mathbf{v}^j + \mathbf{N}\mathbf{d}, \quad \mathbf{M}, \mathbf{N} \in \mathbb{R}^{n \times n}$$

dargestellt werden. Als Minimalvoraussetzung wird vom Verfahren verlangt, daß die Lösung von $\mathbf{A}\mathbf{v} = \mathbf{d}$ ein Fixpunkt des Verfahrens ist. Diese Forderung erlaubt es \mathbf{N} durch \mathbf{M} auszudrücken, d.h.

$$\mathbf{v} = \varphi(\mathbf{v}, \mathbf{d}) := \mathbf{M}\mathbf{v} + \mathbf{N}\mathbf{d} = \mathbf{M}\mathbf{v} + \mathbf{N}\mathbf{A}\mathbf{v} \quad \text{und damit} \\ \mathbf{I} - \mathbf{M} = \mathbf{N}\mathbf{A}.$$

Ist \mathbf{A} invertierbar, was wir immer voraussetzen, so kann man \mathbf{N} ausdrücken durch

$$(19.3) \quad \mathbf{N} = (\mathbf{I} - \mathbf{M})\mathbf{A}^{-1}.$$

Das Iterationsverfahren ist dann durch seine Iterationsmatrix schon eindeutig bestimmt und hat dann die Gestalt

$$(19.4) \quad \mathbf{v}^{j+1} = \mathbf{M}\mathbf{v}^j + (\mathbf{I} - \mathbf{M})\mathbf{A}^{-1}\mathbf{d}.$$

Ist $\mathbf{v}^0 = \mathbf{0}$ Ausgangsnäherung des Verfahrens, was typisch ist als Anfangsnäherung für die Lösung der Korrekturgleichung, so folgt

$$\begin{aligned} \mathbf{v}^1 &= \mathbf{M}\mathbf{v}^0 + \mathbf{N}\mathbf{d} = \mathbf{N}\mathbf{d} \\ \mathbf{v}^2 &= \mathbf{M}\mathbf{v}^1 + \mathbf{N}\mathbf{d} = \mathbf{M}\mathbf{N}\mathbf{d} + \mathbf{N}\mathbf{d} \\ \mathbf{v}^3 &= \mathbf{M}\mathbf{v}^2 + \mathbf{N}\mathbf{d} = (\mathbf{M}^2\mathbf{N} + \mathbf{M}\mathbf{N} + \mathbf{N})\mathbf{d} \end{aligned}$$

und durch vollständige Induktion

$$(19.5) \quad \mathbf{v}^\gamma = \sum_{j=0}^{\gamma-1} \mathbf{M}^j \mathbf{N} \mathbf{d} \stackrel{(19.3)}{=} \sum_{j=0}^{\gamma-1} \mathbf{M}^j (\mathbf{I} - \mathbf{M}) \mathbf{A}^{-1} \mathbf{d} = \sum_{j=0}^{\gamma-1} (\mathbf{M}^j - \mathbf{M}^{j+1}) \mathbf{A}^{-1} \mathbf{d} \\ \mathbf{v}^\gamma = (\mathbf{I} - \mathbf{M}^\gamma) \mathbf{A}^{-1} \mathbf{d},$$

und falls $\|\mathbf{M}\| < 1 \implies \mathbf{v}^\gamma \xrightarrow{\gamma \rightarrow \infty} \mathbf{A}^{-1} \mathbf{d},$

d.h. wenn das Verfahren konvergiert, dann gegen die Lösung von (19.1), denn

$$\|\mathbf{M}\| < 1 \implies \|\mathbf{M}^\gamma\| \leq \|\mathbf{M}\|^\gamma < 1 \xrightarrow{\gamma \rightarrow \infty} 0 \implies \mathbf{M} \rightarrow \mathbf{0}.$$

Beim MGV werden nur Glättungsiterationen benutzt $\tilde{\mathbf{u}} = \mathcal{S}(\mathbf{u}, \mathbf{d})$, welche die Lösung von $\mathbf{A}\mathbf{v} = \mathbf{d}$, bzw. $\mathbf{A}\mathbf{y} = \mathbf{b}$ als Fixpunkt haben (z.B. Jakobi, GSV, symmetrisches GSV). Hieraus folgt

(19.6) Auf jedem Gitterlevel ist die Lösung \mathbf{v}_l von $\mathbf{A}_l\mathbf{v}_l = \mathbf{d}_l$ bzw. $\mathbf{A}_l\mathbf{y} = \mathbf{b}_l$ ein Fixpunkt der Multigriditeration,

denn:

\mathbf{v}_l ist Fixpunkt der Vorglättung	\implies
Für den Defekt gilt:	$\mathbf{d}_l - \mathbf{A}_l\mathbf{v}_l = \mathbf{0}$
Restriktion des Defekts	$= \mathbf{0}$
Grobgridterlösung der Defektgleichung	$= \mathbf{0}$
Prolongation der GrobgridterLösung	$= \mathbf{0}$
Korrektur von \mathbf{v}_l :	$\mathbf{v}_l = \mathbf{v}_l + \mathbf{0}$
\mathbf{v}_l ist Fixpunkt der Nachglättung	

Man kann also für die Konstruktion der Iterationsmatrix des MGV, das ein lineares Verfahren ist, die Eigenschaft (19.3) verwenden und weiß dann, daß das MGV gegen die Lösung der Aufgabe konvergiert, falls $\|\mathbf{M}\| < 1$ ist.

Konstruktion der Iterationsmatrix \mathbf{M}_l des MGV auf Level l

Auf Grund der Eigenschaften des MGV ist die Konstruktion notwendigerweise rekursiv. Wir konstruieren \mathbf{M}_l in Abhängigkeit von der Iterationsmatrix \mathbf{M}_{l-1} des MGV auf Gitter $l-1$, indem wir von Gitter l auf Gitter $l-1$ heruntersteigen. Dort ist die Korrekturgleichung $\mathbf{A}_{l-1}\mathbf{v}_{l-1} = \mathbf{d}_{l-1}$ mit der Ausgangsnäherung $\mathbf{v}_{l-1}^0 = \mathbf{0}$ zu lösen. Dies geschieht wieder durch das MGV indem man pro Iterationsschritt bis zum größten Gitter $l=0$ heruntersteigt und dann wieder bis zum Level $l-1$ hochkommt. Dabei benutzen wir für \mathbf{M}_{l-1} die Darstellung (19.5). Schließlich steigen wir wieder zum Gitter l auf.

Wir listen die einzelnen Schritte auf. Hierbei ist das Aussehen des Terms der jeweiligen Iterationen, der nur von der rechten Seite der zu lösenden Gleichung abhängt, zwar konstruierbar, aber für die Konstruktion der Iterationsmatrix ohne Interesse. Wir konstruieren die Iterationsmatrix für das MGV.

$$\begin{aligned} \mathbf{y}_l^0 &= \text{Ausgangsnäherung auf Gitter } l \text{ zur Lösung von } \mathbf{A}_l\mathbf{y}_l = \mathbf{b}_l \\ \bar{\mathbf{y}}_l &= \mathbf{S}^{\nu_1}\mathbf{y}_l^0 + \mathbf{Q}_{\nu_1}(\mathbf{b}_l) \quad \text{Vorglättung} \\ \bar{\bar{\mathbf{y}}}_l &= \bar{\mathbf{y}}_l + \mathbf{P}\mathbf{v}_{l-1}^\gamma \end{aligned}$$

Dabei ist \mathbf{v}_{l-1}^γ die Näherung für die Lösung der Korrekturgleichung, die man, ausgehend von $\mathbf{v}_{l-1}^0 = \mathbf{0}$ nach γ Schritten des MGV auf dem Level $l-1$ erhält.

$$\begin{aligned} \mathbf{y}_l^1 &= \mathbf{S}^{\nu_2} \bar{\mathbf{y}}_l + \mathbf{Q}_{\nu_2}(\mathbf{b}_l) \quad \text{Nachglättung} \\ &= \mathbf{S}^{\nu_2} (\bar{\mathbf{y}}_l + \mathbf{P} \mathbf{v}_{l-1}^\gamma) + \mathbf{Q}_{\nu_2}(\mathbf{b}_l). \end{aligned}$$

Mit $\bar{\mathbf{v}}_{l-1}^\gamma \stackrel{(19.5)}{=} (\mathbf{I} - \mathbf{M}_{l-1}^\gamma) \mathbf{A}_{l-1}^{-1} \mathbf{d}_{l-1}$ und $\mathbf{d}_{l-1} = \mathbf{R}(\mathbf{b}_l - \mathbf{A}_l \bar{\mathbf{y}}_l)$ folgt

$$\begin{aligned} \mathbf{y}_l^1 &= \mathbf{S}^{\nu_2} \{ \bar{\mathbf{y}}_l + \mathbf{P} (\mathbf{I} - \mathbf{M}_{l-1}^\gamma) \mathbf{A}_{l-1}^{-1} \mathbf{R} (\mathbf{b}_l - \mathbf{A}_l \bar{\mathbf{y}}_l) \} + \mathbf{Q}_{\nu_2}(\mathbf{b}_l) \\ &= \mathbf{S}^{\nu_2} \{ \bar{\mathbf{y}}_l - \mathbf{P} (\mathbf{I} - \mathbf{M}_{l-1}^\gamma) \mathbf{A}_{l-1}^{-1} \mathbf{R} \mathbf{A}_l \bar{\mathbf{y}}_l \} + \mathbf{Q}_3(\mathbf{b}_l) \end{aligned}$$

$\mathbf{Q}_3(\mathbf{b}_l), \mathbf{Q}_4(\mathbf{b}_l)$ bezeichnen Restausdrücke, die nur von den jeweiligen rechten Seiten abhängig sind

Ausmultiplizieren, Ausklammern und Einsetzen von $\bar{\mathbf{y}}_l$ ergibt

$$\begin{aligned} \mathbf{y}_l^1 &= \mathbf{S}^{\nu_2} \{ \mathbf{I} - \mathbf{P} \mathbf{A}_{l-1}^{-1} \mathbf{R} \mathbf{A}_l + \mathbf{P} \mathbf{M}_{l-1}^\gamma \mathbf{A}_{l-1}^{-1} \mathbf{R} \mathbf{A}_l \} \underbrace{(\mathbf{S}^{\nu_1} \mathbf{y}_l^0 + \mathbf{Q}_{\nu_1}(\mathbf{b}_l))}_{\bar{\mathbf{y}}_l} + \mathbf{Q}_3(\mathbf{b}_l) \\ &= \underbrace{\{ \mathbf{S}^{\nu_2} [\mathbf{I} - \mathbf{P} \mathbf{A}_{l-1}^{-1} \mathbf{R} \mathbf{A}_l] \mathbf{S}^{\nu_1} + \mathbf{S}^{\nu_2} \mathbf{P} \mathbf{M}_{l-1}^\gamma \mathbf{A}_{l-1}^{-1} \mathbf{R} \mathbf{A}_l \mathbf{S}^{\nu_1} \}}_{\mathbf{K}_l(\nu_1, \nu_2)} \underbrace{\mathbf{y}_l^0}_{=: \mathbf{E}_{l-1}} + \mathbf{Q}_4(\mathbf{b}_l) \\ &\quad \underbrace{\hspace{10em}}_{\text{Iterationsmatrix des MGV}} \end{aligned}$$

$\mathbf{K}_l(\nu_1, \nu_2)$ ist die Iterationsmatrix des ZGV (vgl. (18.9)) auf Level l , $l \geq 1$, wenn $l=0$ das größte Gitter ist.

Wir haben damit

Satz 19.1 Iterationsmatrix des MGV für $l \geq 1$

Bezeichnet $l=0$ das größte Gitter, so gilt für die Iterationsmatrix \mathbf{M}_l des MGV auf Gitter l

$$\mathbf{M}_l = \mathbf{K}_l(\nu_1, \nu_2) + \mathbf{S}^{\nu_2} \mathbf{P} \mathbf{M}_{l-1}^\gamma \mathbf{E}_{l-1}, \quad l \geq 1, \quad \mathbf{M}_0 = \mathbf{0}$$

mit der Iterationsmatrix des ZGV

$$\mathbf{K}_l(\nu_1, \nu_2) = \mathbf{S}^{\nu_2} [\mathbf{I} - \mathbf{P} \mathbf{A}_{l-1}^{-1} \mathbf{R} \mathbf{A}_l] \mathbf{S}^{\nu_1}$$

und

$$\mathbf{E}_{l-1} = \mathbf{A}_{l-1}^{-1} \mathbf{R} \mathbf{A}_l \mathbf{S}^{\nu_1}$$

Beachte zum Anfangswert: Für $l=1$ liegt ein ZGV vor, also gilt

$$\mathbf{S}^{\nu_2} \mathbf{P} \mathbf{M}_{l-1}^\gamma \mathbf{E}_{l-1} = \mathbf{0}.$$

Dies wird erfüllt durch $\mathbf{M}_0 = \mathbf{0}$.

Vorbemerkungen zum Konvergenzsatz

1. Man kann die Iterationsmatrix des MGV betrachten als die mit einer Störung versehene Iterationsmatrix des ZGV. Den Konvergenzbeweis kann man führen, indem man durch Voraussetzungen dafür sorgt, daß die Summe beider Terme normmäßig < 1 ausfällt.
2. Ausschlaggebend für die Konvergenz ist der Spektralradius ρ der Iterationsmatrix. Er hängt, soweit es die Glättungen betrifft, wegen $\rho(\mathbf{UF}) = \rho(\mathbf{FU})$ (vgl. die Bemerkung nach (18.12)) nur von der absoluten Zahl $\nu = \nu_1 + \nu_2$ der Glättungen ab und nicht davon, ob sie auf Vor- oder Nachglättung verteilt sind. Wir führen den Konvergenzbeweis für die Variante $\nu = \nu_1$ ($\nu_2 = 0$, ohne Nachglättung).
3. Sind H_1, H_2 endlichdimensionale Hilberträume mit den inneren Produkten $(\cdot, \cdot)_{H_1}$ und $(\cdot, \cdot)_{H_2}$ und den von ihnen erzeugten Normen und ist

$$\mathbf{F} : H_1 \xrightarrow{\text{linear}} H_2$$

eine lineare Abbildung, so wird die Operatornorm definiert durch

$$(19.7) \quad \|\mathbf{F}\|_{H_1 \rightarrow H_2} := \sup_{\mathbf{x} \in H_1, \mathbf{x} \neq 0} \frac{\|\mathbf{F}\mathbf{x}\|_{H_2}}{\|\mathbf{x}\|_{H_1}}.$$

Man beachte, daß dadurch auch für nicht quadratische Matrizen eine Norm erklärt wird. (Übung: Normeigenschaften nachweisen)

Wir verwenden diese Normdefinition für nichtquadratische Matrizen.

Ist $H_1 = H_2 = \mathbb{R}^\mu$, $\mu \in \mathbb{N}$, so wird durch (19.7) die Spektralnrm definiert.

Zur Definition der Norm der Prolongationsmatrix gemäß (19.7) verwenden wir für benachbarte Gitter folgende folgende innere Produkte und die von ihnen erzeugten Normen (vgl. dazu Definition (17.1) und die anschließende Anwendung):

$$(19.8) \quad \begin{cases} (\cdot, \cdot)_{l,2} &= \text{euklidisches Produkt in } \omega_l \\ (\cdot, \cdot)_l &= (\cdot, \cdot)_{l,2} h_1 h_2 \\ (\cdot, \cdot)_{l-1} &= (\cdot, \cdot)_{l-1,2} 4h_1 h_2 \end{cases} \quad \text{im 2D-Fall}$$

Ansonsten verwenden wir durchgehend, ohne Kennzeichnung, auf allen Gittern die Euklidische Vektornorm und die zugeordnete Operatornorm (Spektralnrm).

Satz 19.2 Konvergenzsatz für das MGV (ohne Nachglättung)

Es seien folgende Voraussetzungen erfüllt

$$(19.9) \quad \|\mathbf{P}\| \leq 1 \quad \mathbf{P} = \text{Prolongationsmatrix (vgl. (19.7), (19.8))}$$

$$(19.10) \quad \|\mathbf{P}\mathbf{x}\| \geq c_p \|\mathbf{x}\| \quad \forall \mathbf{x}, \text{ für ein } c_p > 0 \quad (c_p \leq 1 \text{ wegen (19.9)})$$

$$(19.11) \quad \|\mathbf{S}\| \leq 1, \quad \mathbf{S} = \text{Iterationsmatrix der Glättung}$$

$$(19.12) \quad \|\mathbf{K}_l(\nu)\| \leq \xi(\nu) < \frac{\gamma - 1}{\gamma} \left(\frac{c_p}{2\gamma} \right)^{\frac{1}{\gamma-1}} \quad \text{für}$$

$$(19.13) \quad \gamma \geq 2, \quad \gamma = \text{Wiederholungszahl des MGV auf den Gittern } l \leq l_{max} - 1.$$

Dabei ist $\xi(\nu)$ eine Normschränke für die Iterationsmatrix \mathbf{K}_l des ZGV auf Level l $1 \leq l \leq l_{max}$, ($\nu = \nu_1, \nu_2 = 0$).

Dann gilt für die Iterationsmatrix \mathbf{M}_l des MGV für $l \geq 1$

$$(19.14) \quad \zeta_l := \|\mathbf{M}_l\| \leq \frac{\gamma}{\gamma - 1} \xi(\nu) < \left(\frac{c_p}{2\gamma} \right)^{\frac{1}{\gamma-1}} \quad (< 1 \text{ gemäß (19.10)})$$

Bemerkungen zu den Voraussetzungen:

zu (19.9)

$\|\mathbf{P}\| \leq 1$ ist für die Interpolation als Prolongation erfüllt. Für die lineare Interpolation (9-Punkte-Stern) rechnet man das im Einheitsquadrat (etwas langwierig, aber ohne Schwierigkeiten) nach.

Hinweise: Gemäß (19.7), (19.8) ist

$$\|\mathbf{P}\| = \sup_{\mathbf{x} \neq 0} \sqrt{\frac{(\mathbf{P}\mathbf{x}, \mathbf{P}\mathbf{x})_l}{(\mathbf{x}, \mathbf{x})_{l-1}}}$$

und damit

$$\|\mathbf{P}\| = \sup_{\mathbf{x} \neq 0} \sqrt{\frac{(\mathbf{P}\mathbf{x}, \mathbf{P}\mathbf{x})_{l,2} h_1 h_2}{(\mathbf{x}, \mathbf{x})_{l-1,2} 4 h_1 h_2}} = \sup_{\mathbf{x} \neq 0} \sqrt{\frac{1}{4} \frac{(\mathbf{P}\mathbf{x}, \mathbf{P}\mathbf{x})_{l,2}}{(\mathbf{x}, \mathbf{x})_{l-1,2}}}$$

Unter Benutzung der Interpolationsformeln (17.6)-(17.8) schätzt man ab

$$(\mathbf{P}\mathbf{x}, \mathbf{P}\mathbf{x})_{l,2} \leq 4 (\mathbf{x}, \mathbf{x})_{l-1,2} \quad \implies \quad \|\mathbf{P}\| \leq 1.$$

Hinweis: Verwende bei der Abschätzung die Ungleichung $2x_i x_{i+1} \leq (x_i^2 + x_{i+1}^2)$ für aufeinanderfolgende Variablenwerte im groben Gitter.

zu (19.10)

Dies ist eine Invertierbarkeitsvoraussetzung für \mathbf{P} auf dem Bildraum von \mathbf{P} . Man sieht, (zumindest bei allen uns bekannten Prolongationen), daß unsere Prolongationsmatrizen Höchststrang haben, woraus folgt

$$\text{Ker } \mathbf{P} = \{0\} \quad \implies \quad \min_{\|\mathbf{x}\|_{l-1,2}=1} \|\mathbf{P}\mathbf{x}\|_{l,2} =: c_p > 0.$$

So ein c_p existiert also. Werte für c_p erhält man aus der Konstruktion von \mathbf{P} im Zusammenhang mit dem Beweis der Approximationseigenschaft. (vgl. etwa Hackbusch: Multigrid Methods Lemma 6.3.13)

zu (19.11)

1. Das gedämpfte Jacobi-Verfahren:

In der Anwendung haben wir uns auf Matrizen $\mathbf{A} = \mathbf{A}^T > 0$ beschränkt, für die die Diagonalelemente alle konstant waren. Die Iterationsmatrix \mathbf{S} des gedämpften Jacobi-Verfahrens war

$$\mathbf{S} = \mathbf{I} - \tilde{\omega}\mathbf{A} \quad (\text{vgl. (18.21)})$$

und es galt für die Spektralnorm

$$\|\mathbf{S}\|_S = \|\mathbf{I} - \tilde{\omega}\mathbf{A}\|_S = \max_{i=1,\dots,n} |1 - \tilde{\omega}\lambda_i|, \quad \lambda_i \in \sigma(\mathbf{A})$$

Nun ist $\|\mathbf{S}\| \leq 1$ (19.11) erfüllt, falls $|1 - \tilde{\omega}\lambda_i| \leq 1 \quad \forall i$.

Wegen $\mathbf{A} = \mathbf{A}^T > 0$ ist $\lambda_i > 0 \quad \forall i$, also ist die Forderung wegen $\tilde{\omega} > 0$ äquivalent zu

$$\tilde{\omega}\lambda_i \leq 2 \quad \forall i \quad \text{bzw.} \quad \tilde{\omega} \leq \frac{2}{\lambda_i} \quad \forall i,$$

was erfüllt ist, falls

$$\tilde{\omega} \leq \frac{2}{\|\mathbf{A}\|_S}.$$

Da die Glättungseigenschaft $\tilde{\omega} \leq \frac{1.2}{\|\mathbf{A}\|_S}$ verlangte, gilt (19.11) für den gedämpften Jacobi.

2. Das Gauß-Seidel-Verfahren

erfüllt $\|\mathbf{S}\| \leq 1$ (vgl. Satz (15.6)). Damit ist die Voraussetzung bei der Schachbrettanordnung erfüllt.

3. Das symmetrische Gauß-Seidel-Verfahren

Sind \mathbf{S}_1 , bzw. \mathbf{S}_2 die Iterationsmatrizen des vorwärts- bzw. rückwärtsgenommenen Gauß-Seidel-Verfahrens, so ist die Iterationsmatrix des symmetrischen Verfahrens $\mathbf{S} = \mathbf{S}_1\mathbf{S}_2$ (vgl. (15.17)).

$\|\mathbf{S}_1\| < 1$ wurde in Satz (15.5) gezeigt, analog zeigt man $\|\mathbf{S}_2\| < 1$ und damit folgt $\|\mathbf{S}\| = \|\mathbf{S}_1\| \|\mathbf{S}_2\| < 1$.

zu (19.12)

Dies ist eine verschärfte Forderung an die Iterationsmatrix des ZGV auf Level l , die durch hinreichend großes $\nu = \nu_1 + \nu_2$ immer erfüllt werden kann. Glättungs- und Approximationseigenschaften sind in (19.12) enthalten.

Da $c_p \leq 1$, erhält man folgende Schranken für $\xi(\nu)$

$$\frac{\gamma-1}{\gamma} \left(\frac{c_p}{2\gamma} \right)^{\frac{1}{\gamma-1}} \leq \begin{cases} \frac{1}{8} = 0.125c_p \leq 0.125 & \text{für } \gamma = 2 \\ \frac{2}{3} \sqrt{\frac{1}{6}} \sqrt{c_p} \approx 0.272 \sqrt{c_p} \leq 0.272 & \text{für } \gamma = 3 \\ \frac{4}{5} \sqrt[3]{\frac{1}{8}} \sqrt[3]{c_p} = 0.4 \sqrt[3]{c_p} \leq 0.4 & \text{für } \gamma = 4 \end{cases}$$

zu (19.13)

$\gamma \geq 2$ ist eine beweistechnische Voraussetzung. Für $\gamma = 1$ muß ein anderer Beweis gemacht werden, der wesentlich aufwendiger ist.

Der **Beweis des Konvergenzsatzes** gliedert sich in zwei Teile

1. Gemäß der rekursiven Definition der Iterationsmatrix wird eine Rekursionsformel für $\zeta_l = \|\mathbf{M}_l\|$ hergeleitet.
2. Beweis des Konvergenzsatzes induktiv mit Hilfe von 1: Zeige $\|\mathbf{M}_l\| < \frac{\gamma}{\gamma-1} \xi(\nu)$.

Beweis 1:

Aus der Gestalt von \mathbf{M}_l (Satz (19.1)) folgt für $l \geq 1$ und $\nu = \nu_1, \nu_2 = 0$

$$(19.15) \quad \begin{aligned} \zeta_l := \|\mathbf{M}_l\| &\leq \xi(\nu) + \|\mathbf{P}\mathbf{M}_{l-1}^\gamma \mathbf{A}_{l-1}^{-1} \mathbf{R}\mathbf{A}_l \mathbf{S}^\nu\| \\ &\stackrel{\|\mathbf{P}\| \leq 1}{\leq} \xi(\nu) + \|\mathbf{M}_{l-1}^\gamma\| \underbrace{\|\mathbf{A}_{l-1}^{-1} \mathbf{R}\mathbf{A}_l \mathbf{S}^\nu\|}_{\mathbf{E}_{l-1}} \\ &\leq \xi(\nu) + \zeta_{l-1}^\gamma \|\mathbf{E}_{l-1}\|. \end{aligned}$$

Wegen $\mathbf{M} = \mathbf{0}$ ist $\zeta_0 = 0$.

\mathbf{E}_{l-1} ist im Wesentlichen aus der Darstellung der Iterationsmatrix $\mathbf{K}_l(\nu)$ des ZGV bekannt (vgl. (18.9)). Dies wird zur Abschätzung benutzt.

$$\mathbf{K}_l(\nu) = \mathbf{S}^\nu - \mathbf{P} \underbrace{\mathbf{A}_{l-1}^{-1} \mathbf{R}\mathbf{A}_l \mathbf{S}^\nu}_{\mathbf{E}_{l-1}}.$$

Hieraus folgt mit der Dreiecksungleichung (rückwärts)

$$\|\mathbf{P}\mathbf{E}_{l-1}\| \leq \underbrace{\|\mathbf{K}_l(\nu)\|}_{<1} + \underbrace{\|\mathbf{S}^\nu\|}_{\leq 1} < 2. \quad (\text{Brutalabschätzung})$$

Wegen $\|\mathbf{P}\mathbf{x}\| \geq c_p \|\mathbf{x}\|$ (vgl. (19.10)), folgt (indirekter Beweis unter Beachtung von (19.7))

$$\begin{aligned} c_p \|\mathbf{E}_{l-1}\| &\leq \|\mathbf{P}\mathbf{E}_{l-1}\| < 2, \\ \|\mathbf{E}_{l-1}\| &< \frac{2}{c_p}, \end{aligned}$$

Aus (19.15) folgt damit für $l \geq 1$

$$(19.16) \quad \zeta_l \leq \xi(\nu) + \frac{2}{c_p} \zeta_{l-1}^\gamma.$$

Beweis 2:

Wir beweisen durch vollständige Induktion: (Abkürzung: $\xi := \xi(\nu)$)

$$(19.17) \quad \zeta_l \leq \xi(\nu) + \frac{2}{c_p} \zeta_{l-1}^\gamma \leq \frac{\gamma}{\gamma-1} \xi(\nu).$$

$l = 1$: richtig nach (19.15), da $\zeta_0 = 0$.

Induktionsvoraussetzung: $\zeta_{l-1} \leq \frac{\gamma}{\gamma-1} \xi$.

Aus (19.16) folgt durch Eintragen der Induktionsvoraussetzung

$$\zeta_l \leq \xi + \frac{2}{c_p} \zeta_{l-1}^\gamma \leq \xi + \frac{2}{c_p} \left(\frac{\gamma}{\gamma-1} \xi \right)^\gamma.$$

Wir zeigen, daß in Verschärfung von (19.17) sogar gilt

$$\xi + \frac{2}{c_p} \left(\frac{\gamma}{\gamma-1} \xi \right)^\gamma \leq \frac{\gamma}{\gamma-1} \xi,$$

wie sich aus der folgenden, äquivalenten Umformung ergibt

$$\begin{aligned} \xi + \frac{2}{c_p} \left(\frac{\gamma}{\gamma-1} \xi \right)^\gamma &\leq \frac{\gamma}{\gamma-1} \xi \\ 1 + \frac{2}{c_p} \left(\frac{\gamma}{\gamma-1} \right)^\gamma \xi^{\gamma-1} &\leq \frac{\gamma}{\gamma-1} \\ \left(\frac{\gamma}{\gamma-1} \right)^\gamma \xi^{\gamma-1} &\leq \frac{c_p}{2(\gamma-1)} \\ \xi^{\gamma-1} &\leq \frac{c_p (\gamma-1)^{\gamma-1}}{2 \gamma^\gamma} = \frac{c_p}{2\gamma} \left(\frac{\gamma-1}{\gamma} \right)^{\gamma-1} \\ \xi &\leq \frac{\gamma-1}{\gamma} \left(\frac{c_p}{2\gamma} \right)^{\frac{1}{\gamma-1}} \end{aligned}$$

Die letzte Ungleichung ist aber die Voraussetzung (19.12). ■

Überlegungen zur Iterationszahl des einfachen MGV

(d.h. ohne Nachglättung)

zur Erreichen der vorgegebenen Genauigkeit

Sei

\mathbf{M}_l = die Iterationsmatrix des MGV auf dem l -ten Gitter

\mathbf{y}_l = exakte Lösung von $\mathbf{A}_l \mathbf{y}_l = \mathbf{b}_l$

$\mathbf{y}_l^{(0)}$ = Ausgangsnäherung für \mathbf{y}_l

$\mathbf{y}_l^{(m)}$ = m -te Näherung für \mathbf{y}_l nach m Schritten des MGV, ausgehend von $\mathbf{y}_l^{(0)}$.

Da \mathbf{y}_l ein Fixpunkt des MGV ist, da auf allen Gittern nur Iterationen verwendet werden, welche die Lösung der jeweiligen Gleichung als Fixpunkt haben (vgl. u.a. (19.6)), gilt für den Fehler

$$\begin{aligned} \mathbf{y}_l^{(m)} - \mathbf{y}_l &= \mathbf{M}_l(\mathbf{y}_l^{(m)} - \mathbf{y}_l), \quad \text{also mit (19.14) und } \xi := \xi(\nu) \\ \|\mathbf{y}_l^{(m)} - \mathbf{y}_l\| &\leq \left(\frac{\gamma}{\gamma - 1} \xi \right) \|\mathbf{y}_l^{(m-1)} - \mathbf{y}_l\| \\ &\leq \left(\frac{\gamma}{\gamma - 1} \xi \right)^m \|\mathbf{y}_l^{(0)} - \mathbf{y}_l\|. \end{aligned}$$

Eine naheliegende Forderung für eine ε -Fehlergenauigkeit könnte sein

$$(19.18) \quad \left(\frac{\gamma - 1}{\gamma} \xi \right)^m \leq \varepsilon \approx \text{Approximierungsfehler}, \quad \|\mathbf{y}_l^{(0)} - \mathbf{y}_l\| \approx 1.$$

Für die notwendige Iterationszahl folgt hieraus

$$m = \left\lceil \frac{\ln \varepsilon}{\ln \left(\frac{\gamma - 1}{\gamma} \xi \right)} \right\rceil,$$

wobei $\lceil \cdot \rceil$ die nächst größere, ganze Zahl bedeutet. Beachte: $\ln(\varepsilon) < 0$, es ist an $\varepsilon < 1$ gedacht, ebenso im Nenner $\ln \left(\frac{\gamma - 1}{\gamma} \xi \right) < 0$.

Die Abschätzung ist unbefriedigend, weil der Diskretisierungsfehler bestenfalls asymptotisch bekannt ist und man somit keine Aussage über die notwendige Iterationszahl erhält, und weil auf dem Gitter l die Differenz $\|\mathbf{y}_l^{(0)} - \mathbf{y}_l\|$ auch nicht so ohne weiteres abschätzbar ist, \mathbf{y}_l ist ja unbekannt.

Abhilfe kann hier die Untersuchung des vollen MGV (also mit Anlaufiteration, vgl. § 16, Volles MGV) schaffen, die ja die Anfangsnäherungen mitliefert.

Gewünscht wird eine Abschätzung $\|\mathbf{y}_l^{(0)} - \mathbf{y}_l\| = O(h_l^\kappa)$, wobei $O(h_l^\kappa)$ der Diskretisierungsfehler auf Level l ist, unter Beachtung des Ausgangsfehlers einer Mehrgitteriteration. Wir beweisen eine entsprechende Abschätzung unter den

Voraussetzungen (19.19)-(19.21):

$$(19.19) \quad h_l = 2 h_{l-1}.$$

Dies beschreibt den Standardfall der Schrittweitenhalbierung von einem Gitter zum nächst feineren.

$$(19.20) \quad \|\tilde{\mathbf{P}}\| \leq 1, \quad \|\tilde{\mathbf{P}}\mathbf{y}_l - \mathbf{y}_{l+1}\| \leq c_0 h_{l+1}^\kappa. \quad (\text{Approx.-Eigenschaft, vgl. (18.18),(18.19)})$$

Hierbei ist $\tilde{\mathbf{P}}$ die ggf. bessere Prolongation, die bei der Anlaufiteration des vollen MGW benutzt wird.

h^κ beschreibt die Diskretisierungsordnung, $\mathbf{y}_l, \mathbf{y}_{l+1}$ sind die exakten Lösungen der Gleichungen auf Level l , bzw. $l+1$. (vgl. dazu die Abschätzungen zum Approximationsteil)

$c_0 = c_0(u)$ sei eine Konstante, die unabhängig vom Gitter ist, aber abhängig sein darf von der exakten Lösung u der Randwertaufgabe.

Von der Qualität her besagt diese Voraussetzung, daß $\tilde{\mathbf{P}}$ so gut sein soll, dass die Approximation $\tilde{\mathbf{P}}\mathbf{y}_k$ die Approximationsgenauigkeit der Diskretisierung nicht verdirbt.

$$(19.21) \quad \|\mathbf{M}_l\| \leq \zeta_l, \quad \zeta := \max_l \zeta_l,$$

d.h. ζ , unabhängig von l_{max} , ist eine gleichmäßige obere Schranke für die Kontraktionszahlen der Iterationsmatrizen des MGW auf allen Gittern. (vgl. (19.14)). Wenn man die Abschätzungszahlen für das ZGV kennt, kann man aus Satz 19.2 Werte für ζ erhalten, natürlich in Abhängigkeit von γ .

Wir beweisen den

Satz 19.3

Die Voraussetzungen (19.19)-(19.21) seien erfüllt. Sei \mathbf{y}_0 die exakte Lösung von $\mathbf{A}_0\mathbf{y}_0 = \mathbf{d}_0$ auf dem größten Gitter $l=0$.

Das volle MGW, das auf jedem Gitterniveau μ ($\mu = \mu_1 = \mu_2$) einfache MGW-Schritte verwendet, liefert, ausgehend von \mathbf{y}_0 , die Näherungslösung $\tilde{\mathbf{y}}_l$ für $\mathbf{A}_l\mathbf{y}_l = \mathbf{b}_l$.

Sei $\zeta = \max_l \zeta_l = \max_l \|\mathbf{M}_l\|$ und μ genüge der Forderung

$$(19.22) \quad \zeta^\mu \leq \frac{1}{1 + 2^\kappa}, \quad (\kappa = \text{Diskretisierungsordnung})$$

so gilt

$$(19.23) \quad \|\tilde{\mathbf{y}}_l - \mathbf{y}_l\| \leq c_0 h^\kappa \quad (c_0 \text{ laut (19.20)})$$

Beachte: (19.22) ist eine Bedingung an die Normen $\|\mathbf{M}_l\|$, die durch entsprechend großes ν (Glättungsanzahl) erfüllt werden kann.

Wir verdeutlichen die Bedeutung der Konstanten μ (siehe obiger Satz) und γ (vgl. Satz 19.2 und Beweis des Satzes 19.3) durch Beispiele: **Eine** Mehrgitteriteration mit der Iterationsmatrix \mathbf{M}_l und $\gamma = 2$ ist für $l = 2, 3, 4$ in den Abb. 4-6 auf S. 132

dargestellt. μ gibt an, wie oft eine Iteration mit \mathbf{M}_l auf jedem Gitterlevel ausgeführt werden soll.

Beweis: Wir zeigen induktiv

$$(19.24) \quad \|\tilde{\mathbf{y}}_l - \mathbf{y}_l\| \leq c_0 \frac{\zeta^\mu}{1 - 2^\kappa \zeta^\mu} h_l^\kappa.$$

Beachte hierzu: Der Bruch in (19.24) ist in ζ monoton wachsend. Wird in (19.24) die Abschätzung (19.22) eingetragen: $\zeta^\mu \leq \frac{1}{1+2^\kappa}$, also $2^\kappa \zeta^\mu \leq \frac{2^\kappa}{1+2^\kappa} < 1$, so ist der Nenner in (19.24) erklärt und es folgt

$$\frac{\zeta^\mu}{1 - 2^\kappa \zeta^\mu} \leq \frac{\frac{1}{1+2^\kappa}}{1 - \frac{2^\kappa}{1+2^\kappa}} = 1, \quad \text{also die Behauptung (19.23).}$$

Für $l = 0$ ist (19.24) erfüllt, da auf $l = 0$ exakt gelöst wird.

Sei $\tilde{\mathbf{y}}_l$ die Näherungslösung für $\mathbf{A}\mathbf{y}_l = \mathbf{b}_l$ nach μ Iterationsschritten, ausgehend von der Näherungslösung $\tilde{\mathbf{P}}\tilde{\mathbf{y}}_{l-1}$, so folgt (beachte: Die exakte Lösung \mathbf{y}_l von $\mathbf{A}_l\mathbf{y}_l = \mathbf{d}$, bzw. \mathbf{b} ist Fixpunkt der Iterationsverfahren auf dem jeweiligen Gitter

$$\begin{aligned} \tilde{\mathbf{y}}_l - \mathbf{y}_l &= (\mathbf{M}_l)^\mu (\tilde{\mathbf{P}}\tilde{\mathbf{y}}_{l-1} - \mathbf{y}_l) \\ \|\tilde{\mathbf{y}}_l - \mathbf{y}_l\| &\leq \zeta_l^\mu \|\tilde{\mathbf{P}}\tilde{\mathbf{y}}_{l-1} - \mathbf{y}_l\| \implies \\ &\leq \zeta^\mu (\|\tilde{\mathbf{P}}\tilde{\mathbf{y}}_{l-1} - \tilde{\mathbf{P}}\mathbf{y}_{l-1}\| + \|\tilde{\mathbf{P}}\mathbf{y}_{l-1} - \mathbf{y}_l\|) \xrightarrow{(19.20), (19.21)} \\ &\leq \zeta^\mu (\|\tilde{\mathbf{y}}_l - \mathbf{y}_l\| + c_0 h_l^\kappa) \xrightarrow{\text{Induktionsvoraussetzung}} \\ &\leq \zeta^\mu (c_0 h_{l-1}^\kappa \frac{\zeta^\mu}{1 - 2^\kappa \zeta^\mu} + c_0 h_l^\kappa) \xrightarrow{(19.19)} \\ &\leq \zeta^\mu (c_0 h_l^\kappa \frac{2^\kappa \zeta^\mu}{1 - 2^\kappa \zeta^\mu} + c_0 h_l^\kappa) \\ &\leq c_0 h_l^\kappa \zeta^\mu (\frac{2^\kappa \zeta^\mu}{1 - 2^\kappa \zeta^\mu} + 1) \\ &= c_0 h_l^\kappa \zeta^\mu \frac{1}{1 - 2^\kappa \zeta^\mu}. \quad \blacksquare \end{aligned}$$

Beachte zur **Anwendung von Satz** (19.3):

Es hat sich als numerisch sinnvoll erwiesen zu wählen:

$$\begin{aligned} \nu &\approx 3 \\ \gamma &\geq 3 \quad \text{im 2-dimensionalen Fall} \\ \gamma &\geq 7 \quad \text{im 3-dimensionalen Fall} \end{aligned}$$

Wir machen eine Überschlagrechnung zur Anwendung von Satz (19.3) für $\gamma = 3$ und

$\kappa = 2$ (quadratische Diskretisierungsordnung). Wir erinnern dazu an das Herkommen der verschiedenen Konstanten.

$$\|\mathbf{P}x\| \geq c_p \|x\|, \quad 0 < c_p \leq 1 \quad \text{vgl. (19.10)}$$

$$\xi(\nu) < \frac{\gamma-1}{\gamma} \left(\frac{c_p}{2\gamma}\right)^{\frac{1}{\gamma-1}} \stackrel{c_p \leq 1}{\leq} \frac{\gamma-1}{\gamma} \left(\frac{1}{2\gamma}\right)^{\frac{1}{\gamma-1}} \quad \text{vgl. Satz 19.2 und}$$

$$\frac{\gamma-1}{\gamma} \left(\frac{1}{2\gamma}\right)^{\frac{1}{\gamma-1}} \leq 0.272 \sqrt[3]{c_p} \quad \text{für } \gamma = 3 \quad (\text{vgl. Satz 19.3})$$

$$\zeta_l \leq \zeta \leq \frac{\gamma}{\gamma-1} \xi(\nu) \quad (\text{vgl. (19.20)}) \quad \implies$$

$$\zeta \leq \frac{3}{2} 0.272 \sqrt[3]{c_p} = 0.408 \sqrt[3]{c_p} \quad \text{und mit (19.22)}$$

$$(0.408 \sqrt[3]{c_p})^\mu \leq \frac{1}{1+2^2} = \frac{1}{5}$$

$$\mu \geq \frac{|\ln 5|}{|\ln(0.408) + \ln(\sqrt[3]{c_p})|} \stackrel{c_p \approx 1}{\approx} \frac{|\ln 5|}{|\ln(0.408)|} \approx 1.8,$$

d.h. mit $\mu = 2$ ist man gut beraten.

Bemerkung: Eine analoge Berechnung für $\gamma = 2$ zeigt, daß auch in diesem Fall $\mu = 2$ ausreichend ist. Man beachte jedoch, daß dann die Konvergenzbedingung von

Satz 19.3: $\xi(\nu) < \frac{\gamma-1}{\gamma} \left(\frac{c_p}{2\gamma}\right)^{\frac{1}{\gamma-1}}$ für $\nu = 3$ nicht erfüllbar ist.

§ 20 MGW für nichtlineare Probleme

Gelöst werden soll

$$(20.1) \quad \mathbf{A}_l(\mathbf{y}_l) = \mathbf{f}_l,$$

ein nichtlineares Gleichungssystem, das z.B. aus der Diskretisierung der folgenden Aufgabe kommen könnte

$$(20.2) \quad \sum_{i=1}^3 \frac{\partial}{\partial x_i} \left(k(u) \frac{\partial u}{\partial x_i} \right) + f(u) = 0 \quad \text{in } \Omega \subset \mathbb{R}^3, \quad u|_{\delta\Omega} = g.$$

k und f können nichtlinear sein.

Das \mathbf{f}_l aus (20.1) muß nicht aus dem f aus (20.2) stammen. Man könnte (auf dem feinsten Gitter) das \mathbf{f}_l auch auf die linke Seite bringen und ein System der Art

$$(20.3) \quad \tilde{\mathbf{A}}_{l_{max}}(\mathbf{y}_{l_{max}}) := \mathbf{A}_{l_{max}}(\mathbf{y}_{l_{max}}) - \mathbf{f}_{l_{max}} = \mathbf{0}$$

lösen.

Naheliegende Idee zur Lösung von (20.1) ist die Anwendung des Newton-Verfahrens. Statt (20.1) löst man eine Folge von linearisierten Problemen (Newton-Iterationsvorschrift) mit Hilfe des linearen MGW. Mehrere Anwendungsvarianten dieser Art wurden schon untersucht (vgl. Hackbusch, Kap. 9)

Man kann jedoch die **Verwendung von Ableitungen** vermeiden durch Anwenden eines nichtlinearen MGW, von dem wir nun eine Variante, die von A. Brand vorgeschlagen wurde, beschreiben.

Für einen Konvergenzbeweis verweisen wir auf die Literatur (vgl. Hackbusch: Multigrid Methods and Applications, und die dort zitierte Literatur).

Zu lösen ist also (20.1), wobei $\mathbf{f}_l = \mathbf{0}$ erlaubt ist.

Glättungsiteration

Üblich ist ein nichtlinearer Gauß-Seidel
Grobbeschreibung:

$$\text{Löse } \mathbf{f}(\mathbf{x}) = \mathbf{0}, \quad \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{f} = (f_1, \dots, f_n)^T.$$

Dann lautet der i -te Teilschritt des nichtlinearen GS-Verfahrens innerhalb des k -ten Gesamtschritts

$$0 = f_i(x_1^{(i)}, \dots, x_i^{(i)}, x_{i+1}^{(i-1)}, \dots, x_n^{(i-1)}), \quad i = 1, \dots, n,$$

d.h. im i -ten Schritt wird aus der i -ten Komponente f_i von \mathbf{f} das $x_i^{(i)}$ berechnet (eine Gleichung mit einer Unbekannten).

Empfohlen wird: Statt der Lösung der Gleichung führe man einen Iterationsschritt nach dem Newtonverfahren aus und zwar nach einem diskretisierten Newton-Verfahren (ohne Verwendung von Ableitungen.)

Beachte: Der nichtlineare GS braucht ohne Voraussetzungen (z.B. Zeilensummenkriterium und Irreduzibilität) nicht zu konvergieren.

Restriktionen und Prolongation: wie im linearen Fall

Grobgitterkorrektur

Auf Level l sei eine Ausgangsnäherung \mathbf{y}_l^0 gegeben. Gelöst werden soll auf Gitter l (vgl. die Nachiteration in (§ 7))

$$(20.4) \quad \mathbf{A}(\mathbf{y}_l^0 + \mathbf{v}_l) = \mathbf{f}_l$$

Gesucht wird also die Korrektur \mathbf{v}_l , die \mathbf{y}_l^0 zur exakten Lösung ergänzt.

Beachte: Auch wenn $\mathbf{f}_l = \mathbf{0}$ ist für $l = l_{\max}$, muß das, wie wir sehen werden, für $l < l_{\max}$ nicht gelten. So ist etwa nicht \mathbf{f}_{l-1} die Restriktion von \mathbf{f}_l . Man erinnere sich, dass beim MGVS nicht die rechten Seiten restringiert werden, sondern die Defekte. Wenn schon (20.4) nicht exakt gelöst wird, so soll doch eine Verbesserung $\mathbf{y}_l^1 = \mathbf{y}_l^0 - \mathbf{v}_l$ konstruiert werden, wobei \mathbf{y}_l^1 eine bessere Näherung für die Lösung von (20.4) sein soll als \mathbf{y}_l^0 . Dabei will man die Verbesserung \mathbf{v}_l als $\mathbf{P}\mathbf{v}_{l-1}$ aus einer Rechnung auf Level $l-1$ haben.

Auf $l = l_{\max}$ seien gegeben seien gegeben $\mathbf{A}_l, \mathbf{y}_l^0$ und \mathbf{f}_l (z.B. $\mathbf{f}_{l_{\max}} = \mathbf{0}$). Die Aufgabe besteht darin, auf Level $l-1$ eine geeignete Gleichung zu finden, aus der \mathbf{v}_{l-1} berechnet werden kann. Als Ansatz soll die zu konstruierende Gleichung folgende Form haben

$$(20.5) \quad \mathbf{A}_{l-1}(\mathbf{R}\mathbf{y}_l^0 + \mathbf{v}_{l-1}) = \mathbf{f}_{l-1}.$$

$\mathbf{A}_{l-1}, \mathbf{y}_l^0, \mathbf{R}$ sind bekannt, \mathbf{f}_{l-1} muß geeignet bestimmt werden. Dann ist \mathbf{v}_{l-1} berechenbar. Aus (20.5) erhält man durch Linearisierung (Taylor)

$$(20.6) \quad \mathbf{A}_{l-1}(\mathbf{R}\mathbf{y}_l^0 + \mathbf{v}_{l-1}) = \mathbf{A}_{l-1}(\mathbf{R}\mathbf{y}_l^0) + \underbrace{\mathbf{A}'_{l-1}(\mathbf{R}\mathbf{y}_l^0)}_{\text{Jakobimatrix}} \mathbf{v}_{l-1} + \dots$$

Die rechte Seite dieser Gleichung könnte ein geeignetes \mathbf{f}_{l-1} sein (unter Vernachlässigung von "+..."), wäre da nicht die Ableitung.

Idee: Ersetze (20.4) durch die linearisierte Fassung

$$(20.7) \quad \mathbf{f}_l = \mathbf{A}_l(\mathbf{y}_l^0 + \mathbf{v}_l) \approx \mathbf{A}_l(\mathbf{y}_l^0) + \underbrace{\mathbf{A}'_l(\mathbf{y}_l^0)}_{\text{Jacobi-Matrix}} \mathbf{v}_l,$$

berechne hieraus $\mathbf{A}'_l(\mathbf{y}_l^0)\mathbf{v}_l \approx \mathbf{f}_l - \mathbf{A}_l(\mathbf{y}_l^0)$ und verwende dessen Restriktion

$$\mathbf{R}(\mathbf{A}'_l(\mathbf{y}_l^0)\mathbf{v}_l) \approx \mathbf{R}(\mathbf{f}_l - \mathbf{A}_l(\mathbf{y}_l^0))$$

als Approximation für $\mathbf{A}'_{l-1}(\mathbf{R}\mathbf{y}_l^0)\mathbf{v}_{l-1}$. Damit folgt aus (20.6) die approximierte Fassung

$$(20.8) \quad \mathbf{A}_{l-1}(\underbrace{\mathbf{R}\mathbf{y}_l^0 + \mathbf{v}_{l-1}}_{\mathbf{w}_{l-1}}) = \mathbf{A}_{l-1}(\mathbf{R}\mathbf{y}_l^0) + \mathbf{R}(\mathbf{f}_l - \mathbf{A}_l(\mathbf{y}_l^0)) := \mathbf{f}_{l-1}.$$

Die rechte Seite betrachten wir als Definitionsgleichung für \mathbf{f}_{l-1} .

Beachte: Auch wenn $\mathbf{f}_l = \mathbf{0}$ wäre, müßte nicht $\mathbf{f}_{l-1} = \mathbf{0}$ gelten.

Wir lösen nun (näherungsweise) das System (20.8), (ein nichtlineares System mit viel weniger Variablen als in (20.4)), erhalten $\mathbf{w}_{l-1} = \mathbf{R}\mathbf{y}_l^0 + \mathbf{v}_{l-1}$, setzen

$$(20.9) \quad \mathbf{v}_{l-1} = \mathbf{w}_{l-1} - \mathbf{R}\mathbf{y}_l^0$$

und definieren

$$(20.10) \quad \mathbf{y}_l^1 = \mathbf{y}_l^0 + \mathbf{P}\mathbf{v}_{l-1}.$$

Bemerkungen: (20.8), (20.9), (20.10) sind Schritte zu einer verbesserten Näherung \mathbf{y}_l^0 . Ein Konvergenzbeweis stammt von Hackbusch.

Statt (20.4) wird (20.8) (näherungsweise) gelöst, ein System mit vielweniger Variablen. (Bis hierher wäre das ein ZGV). Auf die beschriebene Weise steigt man ab bis zum größten Gitter, auf dem dann wirklich die Gleichung $\mathbf{A}_0(\mathbf{R}\mathbf{y}_1^0 + \mathbf{v}_0) = \mathbf{f}_0$ (mit entsprechend konstruiertem \mathbf{f}_0) gelöst werden muß.

Dies ist a) die entscheidende Klippe, liefert

b) aber auch Chancen.

zu a): Die Lösung auf dem größten Gitter muß gut sein, sonst erhält man möglicherweise Divergenz.

zu b): Das System (20.4) kann mehrere Lösungen haben. Es kann leichter sein, auf dem größten Gitter die richtige Lösung rauszusuchen und dann hochzutransportieren (Beispiel bei Bifurkation), als die Auswahl erst auf dem feinen Gitter zu treffen.

Die Eindeutigkeit der Lösung kann bei konkavem \mathbf{f} (oft) gezeigt werden. Bei nichtlineare Aufgaben kann man oft mittels Monotonie Ober- und Unterfunktionen konstruieren und dazwischen die Eindeutigkeit beweisen.

Wesentliche Punkte sind also:

1. Auf $l = 1$ wird ein gutes Verfahren für wenige Punkte benötigt. Es muß nicht die exakte Lösung sein, sollte den Defekt aber deutlich mindern.
2. Im Unterschied zum linearen Verfahren wird auch die Näherungslösung restringiert. (vgl. (20.6),(20.8))
3. Die γ Iterationen werden, wie bekannt, durchgeführt.
4. Dieses so beschriebene Verfahren ist eine direkte Verallgemeinerung des linearen Verfahrens. Ist (20.4) eine lineare Gleichung, so fällt es mit dem bekannten linearen MGW zusammen. Dies erkennt man wie folgt:
Betrachte die Korrekturgleichung (20.8).
Im linearen Fall ist $\mathbf{A}_{l-1}(\mathbf{R}\mathbf{y}_l^0 + \mathbf{v}_{l-1})$ ein Produkt aus der Matrix \mathbf{A}_{l-1} und dem Vektor $\mathbf{R}\mathbf{y}_l^0 + \mathbf{v}_{l-1}$. Entsprechendes gilt für $\mathbf{R}(\mathbf{f}_l - \mathbf{A}_l(\mathbf{y}_l^0))$. Man kann (20.8)

mit \mathbf{A}_{l-1}^{-1} multiplizieren und erhält

$$\mathbf{R}\mathbf{y}_l^0 + \mathbf{v}_{l-1} = \mathbf{R}\mathbf{y}_l^0 + \mathbf{A}_{l-1}^{-1} \underbrace{\mathbf{R}(\underbrace{\mathbf{f}_l - \mathbf{A}_l(\mathbf{y}_l^0)}_{\mathbf{d}_l})}_{\mathbf{d}_{l-1}}$$

also $\mathbf{v}_{l-1} = \mathbf{A}_{l-1}^{-1} \mathbf{d}_{l-1}$.

Die ist die Korrekturgleichung, die aus dem linearen Fall bekannt ist. $\mathbf{R}\mathbf{y}_l^0$ fällt weg.

Beachte jedoch: Im nichtlinearen Fall kann man (in dieser Version) nicht aufsteigen. *Man kommt von oben.* Eine Näherungslösung wird restringiert.

Nachbemerkungen: Neuere Entwicklungen zum MGVS gibt es bei GMD: Chemnitz; Inria, Kiel bzw Bremen (Hackbusch); USA, Conrad Zuse Institut (Deuffhard).

Von Deuffhard stammt eine Anwendung auf Differentialgleichungen, wo nicht immer zum größeren Gitter heruntergestiegen wird, sondern vom größten aufs feinste Gitter aufgestiegen wird. Wesentliche Punkte sind dazu

1. Konjugierte Gradienten-Verfahren
2. Genau dosierte Iterationszahlen auf dem jeweiligen Gitter für das Konjugierte Gradienten-Verfahren
3. $l \rightarrow l + 1$ durch geeignete Interpolation

Shaidurov hat, (Anfang der 90er Jahre einen Beweis dazu gemacht.)

Literaturempfehlung: Shaidurov/Marshuk: Difference Methods and their Extrapolation, Springer 1993.

MGVS zur Lösung von Integralgleichungen

Ein Beispiel: Löse

$$y(x) = \int_{\Omega} K(x, \xi) y(\xi) d\xi = f.$$

Diskretisiere das Integral nach ξ . Dies liefert

$$y(x) = \sum_{i=1}^n a_i(x, \xi_i) y(\xi_i) h, \quad h = \xi_{i+1} - \xi_i,$$

mit von x abhängigen Gewichten a_i .

Löse diese Gleichung an den Punkten $x = \xi_j$, $j = 1, \dots, n$.

Die ergibt ein Gleichungssystem

$$\begin{aligned}(\mathbf{I} - \mathbf{A})\mathbf{y} &= \mathbf{f} \\ \text{mit } \mathbf{y} &= (y_1(\xi), \dots, y_n(\xi))^T \\ \mathbf{A} &= (a_{i,j}), \quad a_{i,j} = a(\xi_i, \xi_j).\end{aligned}$$

Dies kann mit einem MGV behandelt werden.

Für Integralgleichungen wurde schon vor Aufkommen des MGV gezeigt, daß es solche Verfahren gibt.

Kapitel IV

Hyperbolische Differentialgleichungen

§ 21 Die Wellengleichung

Wir beginnen unsere Untersuchungen mit einem (bekanntem) Satz aus der Theorie der Partiellen Differentialgleichungen.

Satz 21.1 (d'Alembert, eigentlich Euler)

Die Anfangswertaufgabe

$$\begin{aligned}u_{tt} &= a^2 u_{xx} + f, & a \in \mathbb{R}, \quad a \neq 0, \quad u \in C^2(\mathbb{R}^2), f \in C^1(\mathbb{R}^2) \\u(x, 0) &= u_0(x), & u_0 \in C^2(\mathbb{R}), \\u_t(x, 0) &= u_1(x), & u_1 \in C^1(\mathbb{R}).\end{aligned}$$

ist sachgemäß.

Ihre Lösung wird gegeben durch die

D'Alembert'sche Lösungsformel

$$u(x, t) = \frac{1}{2}(u_0(x + at) + u_0(x - at)) + \frac{1}{2a} \int_{x-at}^{x+at} u_1(\xi) d\xi + \frac{1}{2a} \int_0^t \int_{x-a(t-\tau)}^{x+a(t-\tau)} f(\xi, \tau) d\xi d\tau$$

Abhängigkeitsverhältnisse nach der Lösungsformel (für $t \geq 0$)

Zur Bestimmung von $u(x_0, t_0)$ werden abgefragt

Anfangspositionen in $x_0 \pm at_0$,

Anfangsgeschwindigkeit im Intervall $[x_0 - at_0, x_0 + at_0]$

f im gesamten Dreieck.

Für die Randwertaufgabe gilt

Satz 21.2

Die 1. Randwertaufgabe

$$\begin{aligned} u_{tt} &= a^2 u_{xx}, & a \in \mathbb{R}, \quad a \neq 0, & \quad u \in C^2([0, \ell] \times [0, \infty)), \\ u(x, 0) &= u_0(x), & u_0 \in C^2([0, \ell]), \\ u_t(x, 0) &= u_1(x), & u_1 \in C^1([0, \ell]), \\ u(0, t) &= g_1(t), \\ u(\ell, t) &= g_2(t), & g_i \in C^2([0, \infty)), \end{aligned}$$

besitzt genau dann eine eindeutige Lösung $\in C^2([0, \ell] \times [0, \infty))$, wenn in jedem Randpunkt α_i , ($\alpha_1 = 0$, $\alpha_2 = \ell$), folgende Verträglichkeitsbedingungen erfüllt sind

$$(21.1) \quad \begin{cases} g_i(0) = u_0(\alpha_i) & \text{(Stetigkeit der Funktion } u) \\ g_i'(0) = u_1(\alpha_i) & \text{(Stetigkeit von } u_t) \\ g_i''(0) = a^2 u_0''(\alpha_i) & \text{(Erfülltheit der Differentialgleichung)}. \end{cases}$$

Numerische Differenzenapproximation

Wir benutzen äquidistante Gitter in Raum und Zeit. Mit den inneren Gitterpunkten

$$\omega_h := \{x_i = +ih; i = 1, \dots, N-1, h = \frac{1}{n}\}$$

$$\omega_\tau := \{t_j = j\tau; j = 1, \dots, m, m = \frac{T}{\tau}\}$$

approximieren wir (vgl. (2.4),(2.7))

mit $u(x_i, t_j) \approx y_i^j, \quad \mathbf{y}^{j+1} \hat{=} \hat{\mathbf{y}}, \quad \mathbf{y}^j \hat{=} \mathbf{y}, \quad \mathbf{y}^{j-1} \hat{=} \bar{\mathbf{y}}$

$$\frac{\partial^2 u}{\partial t^2} \approx \frac{\hat{\mathbf{y}} - 2\mathbf{y} + \bar{\mathbf{y}}}{h^2} =: \mathbf{y}_{\hat{t}t} \quad \text{für beliebige Zeitschichten, d.h.}$$

(21.2)

$$\frac{\partial^2 u(x_i, t_j)}{\partial t^2} \approx \frac{y_i^{j+1} - 2y_i^j + y_i^{j-1}}{\tau^2} =: y_{\hat{t}t,i}^j = \frac{\partial^2 u}{\partial t^2} + \frac{\tau^2}{12} \frac{\partial^4 u}{\partial t^4}(t_j + \theta_j \tau) \quad \text{falls } u \in \mathbb{C}^{0,4}$$

(21.3)

$$\frac{\partial^2 u(x_i, t_j)}{\partial x^2} \approx \frac{y_{i+1}^j - 2y_i^j + y_{i-1}^j}{h^2} =: y_{\hat{x}x,i}^j = \frac{\partial^2 u}{\partial x^2} + \frac{h^2}{12} \frac{\partial^4 u}{\partial x^4}(x_i + \eta_i h) \quad \text{falls } u \in \mathbb{C}^{4,0}$$

Da in der Differentialgleichung eine 2. Zeitableitung vorliegt, brauchen wir 3 Zeitschichten zur Diskretisierung:

$$\mathbf{y}^\sigma := \sigma \hat{\mathbf{y}} + (1 - 2\sigma)\mathbf{y} + \sigma \bar{\mathbf{y}} = (\mathbf{y} + \sigma \tau^2 \mathbf{y}_{\hat{t}t}) = \mathbf{y} + \mathbf{y}_{\hat{t}t} + O(\tau^2).$$

Die Diskretisierung der Differentialgleichung liefert damit

$$\mathbf{y}_{\hat{t}t} = a^2 (\mathbf{y}^{(\sigma)})_{\hat{x}x} + \boldsymbol{\varphi} \quad \text{in } \omega_h \times \omega_\tau \quad \text{und z.B. } \boldsymbol{\varphi} = \mathbf{f}.$$

Es kann jedoch opportun sein als $\boldsymbol{\varphi}$ das \mathbf{f} auf einer Zwischenschicht zu wählen. Diskretisierung der Anfangswerte

1. Funktionswerte: $y_i^0 = u_0(x_i)$ ist trivial.
2. Ableitungswerte: Die Taylorentwicklung von $u(x, t_1)$ liefert

$$u(x, t_1) = u(x, 0) + \tau \dot{u}(x, 0) + \frac{\tau^2}{2} \ddot{u}(x, 0) + O(\tau^3).$$

Unter Benutzung der Differentialgleichung und der Anfangswerte

$$\begin{aligned} u(x, t_1) &= u_0(0) + \tau u_1(0) + \frac{\tau^2}{2} (a^2 u_{xx}(x, 0) + f(x, 0)) + O(\tau^3) \\ &= u_0(0) + \tau u_1(0) + \frac{\tau^2}{2} (a^2 y_{\bar{x}\bar{x}}^0 + f(x, 0) + O(h^2)) + O(\tau^3) \\ &= u_0(0) + \tau u_1(0) + \frac{\tau^2}{2} (a^2 y_{\bar{x}\bar{x}}^0 + f(x, 0)) + \frac{\tau^2}{2} O(h^2) + O(\tau^3) \end{aligned}$$

erhält man mit $\mathbf{y}^1 \approx u(x, t_1)$ als 2. Anfangsbedingung

$$\begin{aligned} \frac{\mathbf{y}^1 - \mathbf{y}^0}{\tau} &= \mathbf{u}_1 + \frac{\tau}{2} (a^2 (\mathbf{u}_0)_{\bar{x}\bar{x}} + \mathbf{f}^0) + \frac{\tau}{2} O(h^2), \quad \left(\frac{\tau}{2} O(h^2) = O(\tau^2 + h^4)\right) \quad \text{bzw.} \\ \mathbf{y}^1 &= \mathbf{y}^0 + \tau \mathbf{u}_1 + \frac{\tau^2}{2} (a^2 (\mathbf{u}_0)_{\bar{x}\bar{x}} + \mathbf{f}^0). \end{aligned}$$

Diskretisierung der Randwerte: $y_0^{j+1} = g_0(t_{j+1})$, $y_N^{j+1} = g_1(t_{j+1})$.

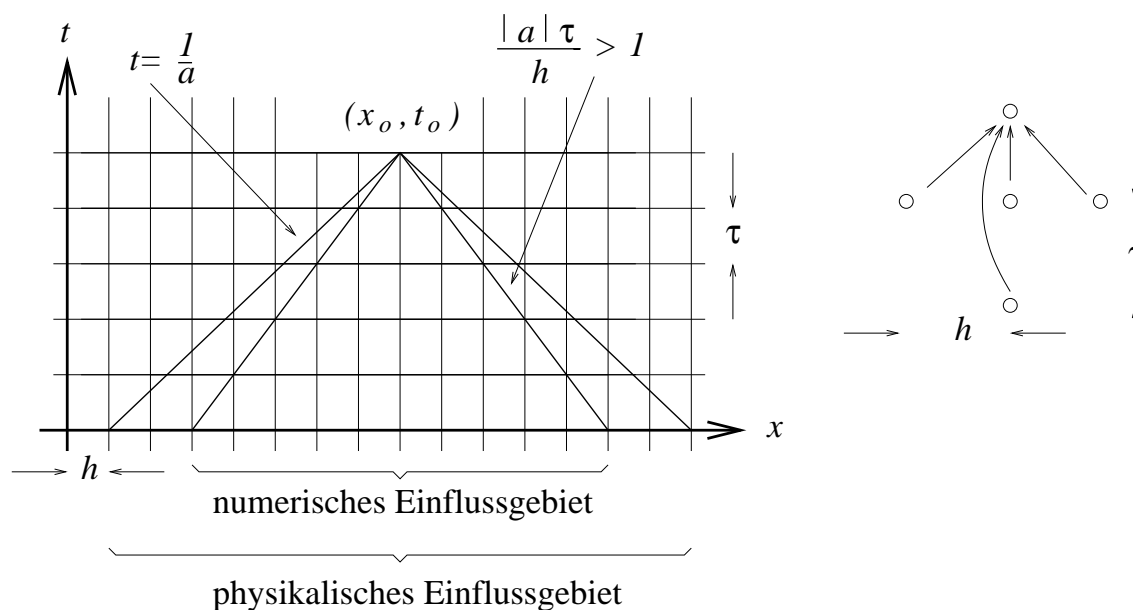
Das komplette Differenzschema lautet damit (nach Zeitschichten geordnet)

$$\begin{aligned} \frac{1}{\tau^2} y^{j+1} - a^2 \sigma y_{\bar{x}\bar{x}}^{j+1} &= a^2 (1 - 2\sigma) y_{\bar{x}\bar{x}}^j + a^2 \sigma y_{\bar{x}\bar{x}}^{j-1} + \frac{2}{\tau^2} y^j - \frac{1}{\tau^2} y^{j-1} + f_i^j \\ & \quad i = 1, \dots, N-1, \quad j = 1, \dots, m-1, \\ y_0^{j+1} &= g_0(t_{j+1}), \quad y_N^{j+1} = g_1(t_{j+1}) \quad j = 0, \dots, m-1, \\ y_i^0 &= u_0(x_i), \quad i = 0, \dots, N, \\ y_i^1 &= y_i^0 + \tau u_1(x_i) + \frac{\tau^2}{2} (a^2 y_{\bar{x}\bar{x}}^0 + f_i^0), \quad i = 1, \dots, N-1. \end{aligned}$$

Folge: $\sigma = 0$ explizites Verfahren

$\sigma > 0$ implizites Verfahren

Wir untersuchen zunächst den expliziten Fall und vergleichen im Punkt (x_0, t_0) den numerischen und den physikalischen Einflußbereich.



Verfolge den Differenzenstern auf dem Gitter.

Gemäß dem zur Diskretisierung verwendeten Differenzenstern beschreiben Geraden durch die Gitterpunkte mit der Gitter-Steigung $\frac{\tau}{h}$ die Grenzen des numerischen Einflußbereichs. Ideal ist, wenn die Gittersteigung mit der Steigung der Charakteristiken (der physikalischen Steigung) $\frac{1}{a}$ übereinstimmt. (Beachte dazu: $t = \pm \frac{1}{a}x + k$), also

$$\frac{\tau}{h} = \frac{1}{|a|} \iff \gamma := \frac{|a|\tau}{h} = 1.$$

Ist die Gittersteigung größer als die physikalische Steigung ($\frac{\tau}{h} > \frac{1}{a} \iff \gamma := \frac{|a|\tau}{h} > 1$), so ist der numerische Einflußbereich kleiner als der physikalische Einflußbereich.

Wir bezeichnen mit $\frac{h}{\tau}$ die Gittergeschwindigkeit
 $|a|$ die Wellengeschwindigkeit
 $\gamma := \frac{|a|\tau}{h}$ die **Courant-Levy-Friedrich-Zahl**, CLF-Zahl.

Für die CLF-Zahl haben wir also folgende 3 Fälle:

$$\gamma = \frac{|a|\tau}{h} \begin{cases} = 1 & \text{idealer Fall: numerisches = physikalisches Einflußgebiet} \\ > 1 & \text{physikalisches Einflußgebiet} > \text{numerisches Einflußgebiet} \\ < 1 & \text{physikalisches Einflußgebiet} < \text{numerisches Einflußgebiet} \end{cases}$$

Man erkennt sofort (vgl. Zeichnung, heuristische Überlegung):

Ist $\gamma > 1$, so ändert sich die physikalische (exakte) Lösung, wenn die Anfangswerte außerhalb des numerischen, aber innerhalb des physikalischen Einflußgebiets geändert werden. Die numerische Lösung bleibt jedoch unverändert.

Dieser Fall muß ausgeschlossen werden.

$\frac{|a|\tau}{h} \leq 1$ ist ein notwendiges, kein hinreichendes Stabilitätskriterium (vgl. Beispiel 4)
(CLF-Kriterium, Courant-Levy-Friedrich).

Diese Bezeichnung muß natürlich mathematisch gerechtfertigt werden.

Wir zeigen am einfachen Beispiel der homogenen Wellengleichung, daß die Verletzung dieses Kriteriums zu Instabilitäten führt.

Die Differentialgleichung

$$u_{tt} = a^2 u_{xx}$$

hat die Diskretisierung

$$\frac{y_k^{j+1} - 2y_k^j + y_k^{j-1}}{h^2} = a^2 \frac{y_{k+1}^j - 2y_k^j + y_{k-1}^j}{h^2}.$$

Durch den (für Differenzgleichungen üblichen) Ansatz

$$y_k^j = \mu^k \lambda^j$$

bestimmen wir eine Lösung der Differenzgleichung an der man ablesen können wird, daß keine Stabilität vorliegt.

Bemerkung: Eigentlich müßten die Zeitindizes j bei den y_k^j -Werten in Klammern eingeschlossen werden, besonders bei skalaren Werten, um Verwechslungen mit Potenzen zu vermeiden. Man hat sich jedoch daran gewöhnt (vermutlich aus Bequemlichkeitsgründen) diese Klammern wegzulassen. Bei y_k^j -Werten bedeutet deshalb das j im Exponenten grundsätzlich nur die Zeitschicht, bzw. bei Folgen (in Zeitrichtung, vgl. Satz 22.1) die Laufvariable der Folge.

Setzt man den obigen Ansatz in die Differenzgleichung ein, so erhält man

$$\frac{\mu^k \lambda^j (\lambda - 1)^2}{\tau^2 \lambda} = a^2 \frac{\mu^k \lambda^j (\mu - 1)^2}{h^2 \mu}$$

woraus mit $\gamma = \frac{|a|\tau}{h}$ folgt

$$\frac{(\lambda - 1)^2}{\lambda} = \gamma^2 \frac{(\mu - 1)^2}{\mu}.$$

Setzen wir in dieser Gleichung $\mu = -1$ und ersetzen λ durch $-\lambda$, so erhält man

$$\begin{aligned} (\lambda + 1)^2 &= 4\gamma^2 \lambda \quad \text{bzw.} \\ \lambda^2 - (4\gamma^2 - 2)\lambda + 1 &= 0 \end{aligned}$$

mit den Wurzeln

$$\lambda_{1,2} = 2\gamma^2 - 1 \pm 2\sqrt{\gamma^2(\gamma^2 - 1)}.$$

Wir untersuchen die bekannten 3 Fälle für γ :

$\gamma < 1$: $\lambda_{1,2}$ komplex, $|\lambda_{1,2}| = 1$, Betrag ausrechnen

$\gamma = 1$: $\lambda_1 = \lambda_2 = 1$,

$\gamma > 1$: $\lambda_1 = 2\gamma^2 - 1 + 2\sqrt{\gamma^2(\gamma^2 - 1)} > 1$.

Da wir eine Differenzgleichung 2. Ordnung (in der Zeit) haben, muß man für eine eindeutige Lösung Anfangswerte auf 2 Zeitschichten vorgeben. Wir geben beschränkte Anfangswerte vor

$$|y_k^0| = 1, \quad |y_k^1| = \lambda_1.$$

Sie werden durch die Lösung $y_k^j = \mu^k \lambda_1^j$ angenommen. (Wir erinnern uns an $\lambda_1 > 1$ und $\mu = -1$).

In einem Punkt (x_k, \tilde{t}) hat die exakte Lösung den Wert $u(x_k, \tilde{t})$. Für ein gegebenes Gitter sei $\tilde{t} = t_j$. Eine Lösung der Differenzgleichung wird durch $(-1)^k \lambda_1^j$ gegeben. Verfeinert man das Gitter, indem man γ einfriert, d.h. das Verhältnis $\frac{h}{\tau}$ bleibt konstant, so ändert sich das numerische Abhängigkeitsgebiet nicht.

Wird die Gitterweite z.B. halbiert, so wird $\tilde{t} = t^{2j}$ und die numerische Lösung für (x_k, \tilde{t}) gegeben durch $(-1)^k \lambda_1^{2j}$, d.h. bei weitergehender Verfeinerung gilt wegen $\lambda_1 > 1$ für ein beliebiges, festes, z.B. geradzahliges k

1. für k : $y_k^j \xrightarrow{j \rightarrow \infty} +\infty$
2. für $k + 1$: $y_{k+1}^j \xrightarrow{j \rightarrow \infty} -\infty$,

d.h. die Lösung geht mit k oszillierend gegen $\pm\infty$, ist also sicher nicht stabil im Sinne der alten Stabilitätsdefinition (vgl. parabolische Gleichungen).

Wir erhalten also :

$\gamma > 1$ ist eine hinreichende Instabilitätsbedingung.

Man beachte, daß diese Stabilitätsaussage zumindest nicht unmittelbar mit dem vorliegenden Anfangswertproblem verknüpft ist.

§ 22 Die Neumann'sche Stabilitätsanalyse

Diese Analyse untersucht die Stabilität bzgl. der Anfangswerte. Wir beschreiben hier die Stabilitätsanalyse für lineare Differentialgleichungen mit konstanten Koeffizienten. Für weitere Untersuchungen, insbesondere für mehrdimensionale Differentialgleichungen und Systeme von Gleichungen, verweisen wir etwa auf das Buch von Godunov/Ryabenki: Difference Schemes, North Holland.

Diese Stabilitätsanalyse ist nicht auf hyperbolische Gleichungen beschränkt. Eine erste Anwendung haben wir schon zu Beginn der Vorlesung (§ 4) bei der Untersuchung der Wärmeleitungsgleichung kennengelernt. Ausgangspunkt unserer Überlegungen ist das vorige Beispiel ebenso wie das Beispiel aus § 4.

Wird eine Differentialgleichung diskretisiert, so erhält man bei festgehaltenem Ort (k bzw. x_k fest) eine Differenzgleichungen (lineare) in y^j der Ordnung N (=Anzahl der Zeitschichten - 1)

Über die Lösung solcher, homogener Differenzgleichungen gilt folgender Satz (vgl. z.B. Henrici: Elemente der Numerischen Analysis I, BI 1964, Satz 6.8)

Satz 22.1

Das charakteristische Polynom einer linearen Differenzgleichung N -ter Ordnung besitze die verschiedenen Nullstellen $\lambda_1, \dots, \lambda_k$, $k \leq N$, mit der Vielfachheit $m_i + 1$, ($\sum m_i = N - k$) von λ_k . Dann bilden die Folgen

$$(22.1) \quad y^{(n)} = \left(\prod_{\mu=0}^{m_i-1} (n - \mu) \right) \lambda_i^{n-m}, \quad m = 0, \dots, m_i, \quad i = 1, \dots, k$$

ein System von N linear unabhängigen Lösungen.

$$(22.2) \quad y^{(n)} = n^m \lambda_i^n, \quad m = 0, \dots, m_i, \quad i = 1, \dots, k$$

sind ebenfalls ein System von N linear unabhängigen Lösungen.

Man kann also in jeder Ortsschicht alle partiellen Lösungen der Differenzgleichung (alle!) in der genannten Form darstellen mit einem konstanten Faktor als Koeffizienten (- im vorigen Beispiel μ^k -) zur Beschreibung der Ortsabhängigkeit), der bei einer homogenen Differenzgleichung keine Rolle spielt. Wir bekommen auf jeder Ortsschicht dann die Zeitabhängigkeit in Potenzform von λ (bei einfachen Nullstellen, bei mehrfachen Nullstellen mit einer natürlichen Zahl als Koeffizienten).

Bei Stabilitätsuntersuchungen ist man daran interessiert, Aussagen zu gewinnen, wann die Lösungen der Differenzgleichungen mit beschränkten Anfangswerten in Zeitrichtung bei Gitterverfeinerung beschränkt bleiben.

Falls ein $|\lambda_i| > 1$ ausfällt, ist dies offensichtlich nicht der Fall. Dann kann das vorgelegte Differenzschema ausgeschlossen werden. Bei mehrfachen Nullstellen, z.B. Lösungen

der Art $n\lambda^n$, muß in Abhängigkeit von der Größe von λ untersucht werden, ob der Ausdruck beschränkt bleibt. Man erhält auf diese Weise Bedingungen für Instabilität. Bedingungen für Stabilität eines Verfahrens kann man nur erhalten, wenn man **alle** Lösungen des Differenzschemas (nicht nur bei festgehaltenem Index k) in dieser “potenzförmigen” Zeitabhängigkeit beschreiben kann, und wenn alle Lösungen beschränkt bleiben.

Dies ist insbesondere der Fall, wenn alle Lösungen λ_i des charakteristischen Polynoms $|\lambda_i| < 1$ erfüllen und einfach sind.

Der Fall mehrfacher Nullstellen muß gesondert untersucht werden. Wir können also festlegen:

Definition 22.2

Ein Mehrschichtschema (j – Index für die Zeitschichten, k – Index für die Punkte in x -Richtung) heißt instabil, falls unter den Lösungen

$$y_k^j := \mu^k \lambda^j$$

solche sind, für die $|\lambda| > 1$ und $|\mu| = 1$ (oder umgekehrt $|\lambda| = 1$ und $|\mu| > 1$) gilt.

Es muß also untersucht werden, wann alle Lösungen des Differenzschemas in der oben genannten Form dargestellt werden können.

Dazu normieren wir das Intervall, das wir untersuchen auf $[0, 1]$. Man kann sich bei hyperbolischen Differentialgleichungen darauf beschränken auf Grund der Abhängigkeitsgebiete.

Wir benutzen ein Ortsgitter

$$x_k = k \cdot h, \quad k = 0, 1, \dots, n+1, \quad \text{also } (n+1)h = 1.$$

Die Feinheit des Gitters wird über $h = \frac{1}{n+1}$ bestimmt.

Bei Trennungsansätzen der Art $y_k^j := \mu^k \lambda^j$ muß $|\mu| = 1$ sein, weil sonst bei Gitterverfeinerung ($n \rightarrow \infty$) die Anfangswerte auf der Zeitschicht $j = 0$ gegen ∞ wachsen oder gegen Null fallen und somit eine Untersuchung des Verhaltens der Lösungen illusorisch wird. Sinnvoll sind für μ also Ausdrücke der Art $e^{i2\pi l x}$. (Dabei kann l als “äußerer” Index dienen), die wir als Gitterfunktion (in x) auffassen können, also $x = x_k$ einsetzen. Wir betrachten die Vektoren

$$\mathbf{v}^{(\ell)} := (1, e^{2\pi i \ell h}, e^{2\pi i \ell 2h}, \dots, e^{2\pi i \ell n h})^T \quad \text{für } l = 0, \dots, n, \quad h = \frac{1}{n+1}.$$

Man beachte, daß die Funktionen $e^{2\pi i \ell x}$ in x periodisch sind ($x = kh$):

$$e^{2\pi i \ell \cdot 0} = 1 = e^{2\pi i \ell \cdot (n+1)h} = e^{2\pi i \ell \cdot 1} = \cos(2\ell\pi) + i \sin(2\ell\pi) = 1.$$

Dies überträgt sich natürlich auf die $\mathbf{v}^{(\ell)}$, d.h. für $\mathbf{v}^{(\ell)} = \mathbf{v}^{(\ell)}(h)$ gilt somit $\mathbf{v}^{(\ell)}(0) = \mathbf{v}^{(\ell)}((n+1)h)$.

Wir zeigen zunächst, daß die $\mathbf{v}^{(\ell)}$, $\ell = 0, \dots, n$ ein Orthonormalssystem bzgl. des diskreten $L_2(\omega_h)$ -Skalarprodukts

$$(22.3) \quad (\mathbf{u}, \mathbf{w})_{(0,h)} := \sum_{k=0}^N u_k \bar{w}_k h$$

bilden (\bar{w}_k ist der konjugiert komplexe Wert), denn

$$(\mathbf{v}^{(\ell)}, \mathbf{v}^{(m)})_{(0,h)} = \sum_{k=0}^n e^{2\pi i(\ell-m)kh} h =: \sum_{k=0}^n w^k h \quad \text{mit } w = e^{2\pi i(\ell-m)h} h.$$

Hieraus folgt

$$(\mathbf{v}^{(\ell)}, \mathbf{v}^{(m)})_{(0,h)} = \begin{cases} 1, & \ell = m, \\ \frac{w^{n+1}-1}{w-1} = 0, & \ell \neq m, \end{cases} \quad \begin{array}{l} \text{weil } \sum_{k=0}^n h = (n+1)h = 1, \\ \text{endliche geometrische Reihe und Periodizität.} \end{array}$$

Daraus folgt: Ist y_k^0 periodisch in k mit der Periode $k = n + 1$, so kann man jede Anfangswert(gitter)funktion \mathbf{y}^0 darstellen als

$$(22.4) \quad \mathbf{y}^0 = \sum_{l=0}^n c_l^{(0)} \mathbf{v}^{(l)}, \quad \text{bzw. komponentenweise} \quad y_k^0 = \sum_{l=0}^n c_l^{(0)} v_k^{(l)},$$

und jede Gitterfunktion auf Level j durch

$$(22.5) \quad \mathbf{y}^j = \sum_{l=0}^n c_l^{(j)} \mathbf{v}^{(l)}, \quad \text{bzw. komponentenweise} \quad y_k^j = \sum_{l=0}^n c_l^{(j)} v_k^{(l)}.$$

Wir bezeichnen mit $\mathbf{y}_{(i)}^j$ eine partielle Lösung der Differenzgleichung. Dabei ist j der Laufindex der Folge und i die Anzeige, zu welcher Nullstelle der charakteristischen Gleichung die Folge gehört, i hat nichts mit der Ortsschicht zu tun.

Wir setzen nun voraus, daß alle Lösungen der charakteristischen Gleichung einfach sind.

Dann hat jede solche partielle Lösung $\mathbf{y}_{(i)}^j$ mit der Anfangswertfunktion \mathbf{y}^0 nach Satz 22.1 die Darstellung (wobei wir den Index i bei $\mathbf{y}_{(i)}^j$ zunächst unterdrücken)

$$(22.6) \quad \mathbf{y}^j = \sum_{l=0}^n c_l^{(0)} \mathbf{v}^{(l)} \lambda_i^j.$$

Der Vergleich von (22.5) und (22.6) liefert

$$(22.7) \quad c_l^{(j)} = c_l^{(0)} \lambda_i^j.$$

Insbesondere erhält man aus (22.7)

$$(22.8) \quad |c_l^{(j)}| \leq |c_l^{(j-1)}| \iff |\lambda_i| \leq 1.$$

Da die $\mathbf{v}^{(\ell)}$ ein Orthonormalsystem bilden, gilt auf jeder Zeitschicht

$$\begin{aligned} (\mathbf{y}^j, \mathbf{y}^j)_{(0,h)} &= \left(\sum_{l=0}^n c_l^{(j)} \mathbf{v}^{(l)}, \sum_{m=0}^n c_m^{(j)} \mathbf{v}^{(m)} \right) = \sum_{m,l=0}^n c_l^{(j)} \overline{c_m^{(j)}} (\mathbf{v}^{(l)}, \mathbf{v}^{(m)}) \\ &= \sum_{l=0}^n |c_l^{(j)}|^2 = \|\mathbf{y}^j\|_{(0,h)}^2 \end{aligned}$$

Mit $|\lambda_i| \leq 1$ folgt aus (22.8) also

$$(22.9) \quad \|\mathbf{y}^{j+1}\|_{(0,h)}^2 = \sum_{l=0}^n |c_l^{j+1}|^2 \leq \sum_{l=0}^n |c_l^j|^2 = \|\mathbf{y}^j\|_{(0,h)}^2 \leq \dots \leq \|\mathbf{y}^0\|_{(0,h)}^2.$$

Mit diesem Stabilitätsbegriff arbeitet Neumann.

Die allgemeine Lösung der Differenzengleichung in jeder Ortsschicht x_k wird durch eine Linearkombination der Lösungen $\sum_{i=1}^n a_{i,(k)} \lambda_i^j$ gegeben, die an Stelle von λ_i in (22.6)-(22.7) eingetragen werden muß. Auf Grund der Homogenität der Differenzengleichung spielt ein konstanter Faktor für die Lösung keine Rolle, weswegen jeweils einer der Koeffizienten $a_{i,(k)}$ zu 1 normiert werden kann. Bei einer Differentialgleichung n -ter Ordnung in der Zeitableitung, sowie der zugehörigen Differenzengleichung, wird die eindeutige Lösung durch $n - 1$ Anfangswerte bestimmt, (u, u_t, \dots) bei der Differentialgleichung, bzw. den Vorgaben auf $n - 1$ Zeitschichten bei der Differenzengleichung. Die letzteren werden in jedem Punkt x_k durch die ersteren bestimmt und legen damit in jedem Punkt x_k die Koeffizienten $a_{i,(k)}$ fest.

Wir halten fest: Sind die Anfangswerte der Differentialgleichung periodisch und alle λ einfach, so kann man auf die beschriebene Weise die allgemeine Lösung darstellen.

Der für die Stabilität schlimmste Fall tritt ein, wenn in (22.7) das betragsgrößte λ_i eingetragen wird. Sind alle $|\lambda_i| \leq 1$, so bleibt (22.9) unberührt. Andernfalls explodiert die Lösung. Wir haben also folgendes Ergebnis:

Satz 22.3

1. Sind die Anfangswerte einer Aufgabe periodisch und sind alle Lösungen λ_i der charakteristischen Gleichung einfach, so gilt:

$$\max_i |\lambda_i| \leq 1 \quad \implies \quad \text{Das Verfahren ist stabil (notwendige und hinreichende Bedingung), d.h. es gilt (vgl. (22.9))}$$

$$\|\mathbf{y}^{j+1}\|_{(0,h)}^2 \leq \|\mathbf{y}^0\|_{(0,h)}^2 \quad \forall j.$$

$$\exists i : |\lambda_i| > 1 \quad \implies \quad \text{Das Verfahren ist instabil (d.h. für beschränkte Anfangswerte kann bei Gitterverfeinerung eine unbeschränkte Lösung existieren)}$$

2. Sind die Anfangswerte nicht periodisch (man kann dann nicht alle Lösungen der Differenzengleichung darstellen), so erhält man nur hinreichende Instabilitätsbedingungen, falls ein $|\lambda_i| > 1$ auftritt. Dann existiert eine explodierende Lösung. Dazu müssen die Nullstellen der charakteristischen Gleichung nicht einfach sein.

Damit erhält man für den praktischen Verlauf der Stabilitätsuntersuchung nach Neumann

1. Man setzt den Ansatz $y_k^j = \mu^k \lambda^j = e^{i\varphi k} \lambda^j$ in die Differenzengleichung ein.
Falls die Aufgabe periodische Anfangswerte besitzt, ist $\varphi = 2\pi\ell h$ zu setzen. Im nicht periodischen Fall ist φ frei (ggf. geeignet) wählbar: man untersucht dann ja "nur", ob etwas schief geht.
2. Bestimme $\lambda = \lambda(\varphi)$ so, daß der Ansatz eine Lösung der Differenzengleichung liefert.
Falls ein $|\lambda_i| > 0$ existiert \implies Instabilität
Falls alle $|\lambda_i| \leq 1$ und die Differenzengleichung in Zeitrichtung nur einfache Wurzeln hat \implies Stabilität.
Der Fall mehrfacher Wurzeln muß, ggf. in Abhängigkeit von φ , extra untersucht werden (vgl. Beispiel 4).
3. Falls in den Differenzengleichungen in Zeitrichtung ein τ auftritt, kann es auch zu Abschätzungen für hinreichend kleine τ kommen.

$$|\lambda| \leq 1 + c_0\tau, c_0 = \text{konst.} \quad \text{Neumannsches Stabilitätskriterium.}$$

Dann folgt

$$|\lambda^j| \leq (1 + c_0\tau)^j \leq e^{c_0\tau j} \leq e^{c_0T} := M \quad \text{für } \tau j \leq T$$

Beachte: $(1 + c_0\tau)$ ist der Anfang der Taylorreihe für $e^{c_0\tau j}$ und damit (vgl. (22.8))

$$|c_\ell^{(j)}| \leq |c_\ell^{(0)}| M, \quad \text{also} \quad \|\mathbf{y}^j\|_{(0,h)} \leq M \|\mathbf{y}^0\|_{(0,h)},$$

d.h. man erhält Stabilität auf jeder endlichen Zeitschicht.

Schwierigkeiten bei der Neumann'schen Stabilitätsanalyse

1. mehrfache Wurzeln,
2. bei impliziten Verfahren können rationale Bedingungen für λ auftreten (Formel-manipulationssysteme!, z.B. Maple oder Mathematika),
3. oft erhält man hinreichende Bedingungen für Instabilität statt hinreichender Bedingungen für Stabilität.

Beispiele zur Stabilitätsanalyse

Beispiel 1: Wärmeleitungsgleichung aus § 4.

Beispiel 2: Wellengleichung in § 21.

Beispiel 3: Wir betrachten ein implizites Verfahren für die Wellengleichung

$$\mathbf{y}_{\bar{t}\bar{t}} = a^2 \hat{\mathbf{y}}_{\bar{x}\bar{x}}.$$

Einsetzen des Ansatzes $y_k^j = e^{i\varphi k} \lambda^j$, $\varphi = 2\pi\ell h$ in die Wellengleichung liefert

$$\begin{aligned} (y_{\bar{x}\bar{x}})_k^{j+1} &= \frac{\lambda^{j+1}}{h^2} (e^{i\varphi(k+1)} - 2e^{i\varphi k} + e^{i\varphi(k-1)}) \\ &= \lambda^{j+1} e^{i\varphi k} \left(\frac{e^{i\varphi} - 2 + e^{-i\varphi}}{h^2} \right) \\ &= \frac{\lambda y_k^j}{h^2} (-4) \left(\frac{e^{i\frac{\varphi}{2}} - e^{-i\frac{\varphi}{2}}}{2i} \right)^2 \\ &= \frac{-4y_k^j \lambda}{h^2} \sin^2\left(\frac{\varphi}{2}\right) \\ (y_{\bar{t}\bar{t}})_k^j &= \frac{e^{i\varphi k}}{\tau^2} \frac{\lambda^j}{\lambda} (\lambda^2 - 2\lambda + 1) = y_k^j \frac{1}{\tau^2 \lambda} (1 - \lambda)^2. \end{aligned}$$

Damit folgt für die Differenzgleichung

$$\frac{(1 - \lambda)^2}{\tau^2 \lambda} = -\frac{4}{h^2} a^2 \lambda \sin^2\left(\frac{\varphi}{2}\right)$$

Mit $s = \sin \frac{\varphi}{2}$ und $\gamma = \frac{\tau|a|}{h}$ erhält man

$$\begin{aligned} (1 - \lambda)^2 &= -\gamma^2 \lambda^2 \cdot 4s^2 \quad \text{bzw.} \quad (1 + 4\gamma^2 s^2) \lambda^2 - 2\lambda + 1 = 0 \\ \lambda_{1,2} &= \frac{1}{2(1 + 4\gamma^2 s^2)} \left(2 \pm \sqrt{4 - 4(1 + 4\gamma^2 s^2)} \right) \\ &= \frac{1}{2(1 + 4\gamma^2 s^2)} \left(2 \pm \sqrt{-16\gamma^2 s^2} \right) \\ &= \frac{1 \pm \sqrt{-4\gamma^2 s^2}}{1 + 4\gamma^2 s^2} \quad \text{und daher (Betrag ausrechnen)} \\ |\lambda_{1,2}| &= \frac{1}{\sqrt{1 + 4\gamma^2 s^2}} < 1 \quad \text{unabhängig von } \gamma \end{aligned}$$

also: Stabilität für periodische Anfangswerte.

Beispiel 4: Stabilität bei mehrfachen Nullstellen

Wir betrachten nochmals die Wellengleichung mit periodischen Anfangswerten

$$u_{\bar{t}\bar{t}} = a^2 u_{\bar{x}\bar{x}}$$

Bekannt ist: $\gamma > 1 \implies$ Instabilität, $\gamma \leq 1$ notwendig für Stabilität.

Wir untersuchen den Fall $\gamma = 1$ für die Diskretisierung

$$\frac{y_k^{j+1} - 2y_k^j + y_k^{j-1}}{\tau^2} = a^2 \frac{y_{k+1}^j - 2y_k^j + y_{k-1}^j}{h^2}$$

nochmals genauer mit dem Ansatz

$$y_k^j = \lambda^j e^{i\varphi k}, \quad \varphi = 2\pi\ell h, \quad \ell = 1, \dots, n, \quad h = \frac{1}{n+1}, \quad \implies$$

$$\begin{aligned} \frac{(\lambda - 1)^2}{\tau^2 \lambda} &= a^2 \frac{e^{i\varphi} - 2 + e^{-i\varphi}}{h^2} = \frac{a^2}{h^2} \left(e^{i\frac{\varphi}{2}} - e^{-i\frac{\varphi}{2}} \right)^2 \\ &= -4 \sin^2\left(\frac{\varphi}{2}\right). \end{aligned}$$

Mit $\gamma = \frac{|a|\tau}{h}$, $s = \sin\left(\frac{\varphi h}{2}\right)$ folgt

$$\lambda^2 - 2\lambda + 1 = -4\gamma^2 s^2 \lambda$$

$$\lambda_{1,2} = 1 - 2\gamma^2 s^2 \pm 2\sqrt{\gamma^2 s^2 (\gamma^2 s^2 - 1)}.$$

Bei $\gamma = 1$ kann nur eine mehrfache Wurzel auftreten für $s^2 = 1$, also

$$\begin{aligned} \sin\left(\frac{2\pi\ell h}{2}\right) &= \sin(\pi\ell h) = \sin\left(\frac{\pi\ell}{n+1}\right) = \pm 1, \quad \text{d.h. für ganzzahliges } z \\ \frac{\pi\ell}{n+1} &= \frac{\pi}{2} + z\pi, \quad \iff \frac{\ell}{n+1} = \frac{1}{2} + z \iff \ell = \frac{n+1}{2}(1+2z) \end{aligned}$$

Die letzte Gleichung hat nur für ungerades n eine ganzzahlige Lösung für ℓ und n . Für geradzahliges n ist $s^2 \neq 1$, also $|\sin^2(\frac{\pi\ell}{n+1})| < 1$ und damit $\lambda_1 \neq \lambda_2$. Beide Nullstellen sind einfach und wegen $|\lambda_{1,2}| = 1$ (nachrechnen!) folgt somit:

Für n geradzahlig liegt Stabilität vor.

Für n nicht geradzahlig liefert die Lösungsformel $\lambda_1 = \lambda_2 = -1$. Nach Satz 22.1 erhält man der doppelten Nullstelle wegen die Lösungen $y_k^j = e^{i\varphi k}(-1)^j$ und $y_k^j = e^{i\varphi k} j(-1)^j$. Einsetzen in die Differenzgleichung bestätigt die Lösungseigenschaft. Damit ist $e^{i\varphi k} j(-1)^j$ eine Lösung mit polynomialen Wachstum.

Der Grenzfall $\gamma = 1$, n , gerade gehört nicht zum Stabilitätsbereich.

$\gamma = 1$ genügt also nicht für die Stabilität

Beispiel 5: Wir untersuchen eine Approximation 2. Ordnung in Orts- und Zeitrichtung für die Wärmeleitungsgleichung

$$\frac{y^{j+1} - y^{j-1}}{2\tau} = y_{\bar{x}x}$$

Mit dem Ansatz $y_k^j = e^{i\varphi k} \lambda^j$ folgt aus der Differentialgleichung

$$\begin{aligned} \frac{\lambda - \frac{1}{\lambda}}{2\tau} y_k^j &= \frac{1}{h^2} (e^{i\varphi} - 2 + e^{-i\varphi}) y_k^j \\ \frac{\lambda - \frac{1}{\lambda}}{2\tau} &= \frac{2 \cos \varphi - 2}{h^2} = -\frac{2 \sin^2 \frac{\varphi}{2}}{h^2} \cdot 2 = \frac{4}{h^2} \quad \text{mit } s = \sin^2 \frac{\varphi}{2} \\ \lambda^2 - 1 &= -8 \frac{\tau}{h^2} \lambda s^2 \end{aligned}$$

$\gamma := \frac{\tau}{h^2}$ ist eine Art Courant-Zahl für die Wärmeleitungsgleichung (vgl. (4.2))

$$\begin{aligned} \lambda^2 + 8\lambda\gamma s^2 - 1 &= 0 \\ \lambda_{1,2} &= \frac{1}{2} (-8\gamma s^2 \pm \sqrt{64\gamma^2 s^4 + 4}). \end{aligned}$$

Da $\gamma > 0$, folgt $|\lambda_2| > 0$, unabhängig von der Größe von γ

Instabilität, unabhängig vom Verhältnis τ zu h^2 . Verfahren unbrauchbar.

Daß die Mittelpunkregel, die zur Approximation von y_t verwandt wurde, nicht immer schlechte Ergebnisse liefern muß, zeigt das

Beispiel 6: Die Schrödingergleichung

$$-\frac{\hbar^2}{2m} \Delta u + Uu = i\hbar \frac{\partial u}{\partial t}$$

hierbei sind: \hbar = Planck'sches Wirkungsquantum

m = Masse

U = ein Potential.

Wir behandeln nur den 1D-Fall (ohne physikalische Konstanten) wie folgt

$$-i\alpha \frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} - \beta u, \quad \alpha, \beta > 0$$

mit der zeitlichen Approximation durch die Mittelpunktsregel (auf der Zeitschicht j)

$$-i\alpha \frac{y^{j+1} - y^{j-1}}{2\tau} = y_{\bar{x}\bar{x}} - \beta y.$$

Aus dem üblichen Ansatz (vgl. Vorseite, inclusive Abkürzungen) folgt

$$\begin{aligned}
 -i\alpha \frac{\lambda - \frac{1}{\lambda}}{2\tau} &= -\frac{4}{h^2} s^2 - \beta \\
 -i\alpha \lambda^2 - 8\gamma s^2 \lambda - 2\tau\beta \lambda &= 0, \quad \gamma = \frac{\tau}{h^2} \\
 \lambda^2 + i \underbrace{\left(\frac{8\gamma s^2}{\alpha} + \frac{2\tau\beta}{\alpha} \right)}_{=:d} \lambda - 1 &= 0 \\
 \lambda_{1,2} &= \frac{1}{2}(-di \pm \sqrt{-d^2 + 4}) \implies \\
 |\lambda_1 \lambda_2| &= 1
 \end{aligned}$$

und $|\lambda_i| = 1$ solange $d^2 \leq 4$ (konjugiert komplexe Wurzeln)

$$d^2 \leq 4 \iff 8\gamma s^2 + 2\tau\beta \leq 2\alpha$$

$s^2 = 1$ im schlimmsten Fall

$$8\frac{\tau}{h^2} + 2\beta \leq 2\alpha$$

Folge: Falls $\tau = O(h^2)$ erhält man Spielraum für eine stabile Rechnung trotz Verwendung der Mittelpunkregel.

Falls $d^2 > 4$ hat man Instabilität.

Diskretisierungsfehler für die 1D-Wellengleichung

$$\psi = u_{\bar{t}t} - a^2 u_{\bar{x}x}^\sigma - \varphi, \quad \text{wobei } u = \text{exakte Lösung}$$

Wir benutzen die bekannten Entwicklungen aus § 2 für $u \in C^{4,4}$.

$$\begin{aligned}
 u_{\bar{t}t} &= \frac{\partial^2 u}{\partial t^2} + \frac{h^2}{12} \frac{\partial^4 u}{\partial t^4} + O(\tau^4) \\
 u_{\bar{x}x} &= \frac{\partial^2 u}{\partial x^2} + \frac{h^2}{12} \frac{\partial^4 u}{\partial x^4} + O(h^4) \\
 u^{(\sigma)} &= \left(I + \sigma\tau^2 \frac{\partial^2}{\partial x^2} \right) u + O(\tau^4)
 \end{aligned}$$

Aus den letzten beiden Gleichungen folgt

$$\begin{aligned}
 u_{\bar{x}x}^{(\sigma)} &= \left(I + \sigma\tau^2 \frac{\partial^2}{\partial x^2} \right) \left(\frac{\partial^2 u}{\partial x^2} + \frac{h^2}{24} \frac{\partial^4 u}{\partial x^4} \right) u \\
 &= \frac{\partial^2 u}{\partial x^2} + \frac{h^2}{12} \frac{\partial^4 u}{\partial x^4} + \sigma\tau^2 \frac{\partial^4 u}{\partial x^2 \partial t^2} + O(\tau^4 + h^4)
 \end{aligned}$$

Damit erhält man (für $\varphi = f^j$, üblicherweise)

$$\psi = \left(\frac{\partial^2 u}{\partial t^2} - a^2 \frac{\partial^2 u}{\partial x^2} - f \right) + \frac{\tau^2}{12} \frac{\partial^4 u}{\partial t^4} - a^2 \left(\frac{h^2}{12} \frac{\partial^4 u}{\partial x^4} + \sigma \tau^2 \frac{\partial^4 u}{\partial x^2 \partial t^2} \right) + O(\tau^4 + h^4).$$

$$\psi = O(\tau^2 + h^2) \quad \text{falls } u \in C^{4,4}.$$

Eine Voraussetzung $u \in C^{6,6}$ ist jedoch nicht überflüssig, denn eine weitere Rechnung liefert

$$\psi = O(\tau^4 + h^4) \quad \text{falls } u \in C^{6,6} \text{ und}$$

$$\sigma = \frac{1 - \frac{1}{\gamma^2}}{12}, \quad \gamma = \frac{|a|\tau}{h}, \quad \varphi = f + \frac{1}{12} \left(\tau^2 \frac{\partial^2 f}{\partial t^2} + h^2 \frac{\partial^2 f}{\partial x^2} \right)$$

Bemerkung: Hier, wie im parabolischen und elliptischen Fall erhält man die Konvergenz aus Stabilität und Diskretisierungsordnung.

Beachte: Für $\gamma = 1$, also $\sigma = 0$ erhält man so ein explizites Verfahren 4. Ordnung (im Unterschied zum parabolischen Fall). Beachte jedoch Beispiel 4. Noch höhere Ordnungen sind möglich, falls noch höhere Differenzierbarkeitsordnung vorausgesetzt wird.

Folge: Der Bedarf nach σ ist im hyperbolischen Fall kleiner als im parabolischen Fall. Verfahren mit $\sigma > 0$ werden wichtig bei Eigenwertaufgaben, falls die Approximationen für die höheren Eigenwerte schlecht sind und wenn sich diese höheren Eigenwerte bei Schwingungsaufgaben bemerkbar machen.

Dann benutzt man implizite Verfahren, die automatisch eine Dämpfung des Einflusses der höheren Eigenwerte bewirken (ohne Beweis).

Bemerkung: Die Stabilitätsabschätzungen, die wir in § 4 und § 8 (Stabilität bzgl. der rechten Seite auf Grund der Stabilität bzgl. der Anfangswerte) hergeleitet haben, lassen sich auf 3-Schicht-Schemata (für hyperbolische und parabolische Differentialgleichungen) verallgemeinern (natürlich unter entsprechenden Voraussetzungen). Dadurch werden die Neumann'schen Stabilitäts-Analysen auf Grund anderer Voraussetzungen und auch wegen der möglichen Instabilitätsbeweise, die bei vielen Differentialgleichungen Diskretisierungen, die "eigentlich auf der Hand liegen", ausscheiden.

§ 23 Literatur

1. SAMARSKY, A.A.: *Theorie der Differenzenverfahren*, Akad. Verlagsgesellschaft Geest u. Portig, Leipzig 1984
2. MORTON, K.W. / MAYERS, D.F.: *Num. Solution of Partial Diff. Equ.*, Cambridge University Press 1994
3. ISERLES, A.:] *A first course in the Num. Analysis of Diff. Equ.*, Cambridge Texts in Appl. Math. 1996
4. GROSSMANN, CH. / ROOS, H.-G.: *Numerik partieller Differentialgleichungen*, Teubner Studienbücher 1993
5. RICHTMYER, R.D. / MORTON, K.W.: *Difference Methods for Initial-Value-Problems*, Interscience Publishers (Wiley and Sons) 1967
6. FORSYTHE, G.E. / WASOW, W.R.: *Finite Difference Methods for PDE*, Wiley and Sons 1960
7. HACKBUSCH, W.: *Theorie und Numerik elliptischer Dgln*, Teubner Studienbücher 1986
8. HACKBUSCH, W.: *Multi-Grid Methods and Application*, Springer 1986
9. TROTTEBERG, U. / OOSTERLEE, C. / SCHÜLLER, A.: *Multigrid*, Academic Press 2001
10. BRANDT, A. *A guide to multigrid development: Multigrid Methods I*, Lecture Notes in Mathematics 960 S 230-312, 1982
11. HACKBUSCH, W. / TROTTEBERG, U.: *Multigrid II*, Lecture Notes in Mathematics 1228, 1986

Index

- ε -Ungleichung, 27
- Courant-Levy-Friedrich-Zahl,, 177
- 2-Gitter-Verfahren, 127
- 9-Punkte-Prolongation, 142

- bedingt stabil, 29
- gewichtetes Skalarprodukt $(\cdot, \cdot)_{(0,h)}$, 12

- Abbaurrate q , 51
- Adjungierte Operatoren, 143
- alternierende Richtungen, 73
- Approximationseigenschaft, falls, 148
- Approximationsfehler, 31

- bedingt stabil, 29
- Bezeichnungen, 6

- CLF-Kriterium, 178
- CLF-Zahl, 177
- Crank-Nicolson Schema, 22

- D'Alembert'sche Lösungsformel, 174
- diskretes Maximumprinzip, 88
- Diskretisierungsfehler, 31, 47

- Einzelstschritt-Verfahren, 137
- energetische Norm, 17
- explizite Differenzenschema, 7

- Fehleranalyse, 122

- Gauß-Seidel-Iteration, 137
- Gitterfunktion, 4
- Gittergeschwindigkeit, 177
- Gitternummer, 128
- Glättung, 124
- Glättungseigenschaft, 148
- Glättungsnummer, 135, 136
- Glättungsoperator, 127
- Glättungsrate, 135

- implizites Euler Verfahren, 22

- instabil, 181
- irreduzibel, 112
- Iterationsmatrix, 110
- Iterationsmatrix des MGV, 159
- Iterationsmatrix des ZGV, 146, 147

- Konvergenzrate, 110
- Konvergenzsatz für das MGV, 161
- Korrekturgleichung, 120

- L-Matrizen, 85
- Lexikographische Ordnung:, 141
- lineare Konsistenzordnung, 4

- M-Matrix, 42
- Maximumprinzip für Matrizen, 88
- MGV($3, \mathbf{u}, \mathbf{b}$), 131

- Nachglättung, 127
- Nachiteration, 120
- nested iteration, 133
- Neumann'sche Stabilitätsanalyse, 180
- Neumannsches Stabilitätskriterium, 184
- Normvergleiche, 17
- numerisches Einflußgebiet, 177

- Oszillationsfreiheit, 66

- Parallelisierung, 139
- physikalisches Einflußgebiet, 177
- positiv definit, 12
- positiv semidefinit, 11
- positive Definitheit, 60
- Prolongation , 121

- Rayleigh-Quotienten., 10
- reduzibel, 112
- reell positiv definit, 56
- Restriktionsoperator , 121
- rotierte Lexikographische Ordnung:, 141

Schachbrettanordnung, 138
Shortley-Wellers-Approximation, 100, 103
Spektralnorm, 11
Splitting-Verfahren, 73
stabil, 29
Stabilität, 23
symmetrische Matrix, 60
symmetrisches Gauß-Seidel-Verfahren, 117

Tridiagonalalgorithmus, 38
Tridiagonalmatrix, 73
Tridiagonalverfahren, 39
truncation error, 31

unbedingt stabil, 28, 29

Verfahrensfehler, 35
Vierfarben-Ordnung, 140
Volles Mehrgitter-Verfahren, 133
Vorglättung, 127
vorwärtsgenommener Differenzenquotient, 4

Wellengeschwindigkeit, 177

zebra-line-ordering, 141
Zeitschicht, 4
Zentrale Differenzenquotienten, 5